

Introduction to Empirical Processes and Semiparametric Inference¹

Michael R. Kosorok

August 2006

¹©2006 SPRINGER SCIENCE+BUSINESS MEDIA, INC. All rights reserved. Permission is granted to print a copy of this preliminary version for non-commercial purposes but not to distribute it in printed or electronic form. The author would appreciate notification of typos and other suggestions for improvement.

Preface

The goal of this book is to introduce statisticians, and other researchers with a background in mathematical statistics, to empirical processes and semiparametric inference. These powerful research techniques are surprisingly useful for studying large sample properties of statistical estimates from realistically complex models as well as for developing new and improved approaches to statistical inference.

This book is more a textbook than a research monograph, although some new results are presented in later chapters. The level of the book is more introductory than the seminal work of van der Vaart and Wellner (1996). In fact, another purpose of this work is to help readers prepare for the mathematically advanced van der Vaart and Wellner text, as well as for the semiparametric inference work of Bickel, Klaassen, Ritov and Wellner (1997). These two books, along with Pollard (1990) and chapters 19 and 25 of van der Vaart (1998), formulate a very complete and successful elucidation of modern empirical process methods. The present book owes much by the way of inspiration, concept, and notation to these previous works. What is perhaps new is the introductory, gradual and unified way this book introduces the reader to the field.

The book consists of three parts. The first part is an overview which concisely covers the basic concepts in both empirical processes and semiparametric inference, while avoiding many technicalities. The second part is devoted to empirical processes, while the third part is devoted to semiparametric efficiency and inference. In each of the last two parts, the second chapter (after the introductory chapter) is devoted to the relevant mathematical concepts and techniques. For example, an overview of metric

spaces—which are necessary to the study of weak convergence—is included in the second chapter of the second part. Thus the book is largely self contained. In addition, a chapter devoted to case studies is included at the end of each of the three parts of the book. These case studies explore in detail practical examples which illustrate applications of theoretical concepts.

The impetus for this work came from a course the author gave in the Department of Statistics at the University of Wisconsin-Madison, during the Spring semester of 2001. Accordingly, the book is designed to be used as a text in a one or two semester sequence in empirical processes and semiparametric inference. In a one semester course, some of the material would need to be skipped. Students should have had at least half a year of graduate level probability as well as a year of graduate level mathematical statistics.

Contents

Preface	iii
I Overview	1
1 Introduction	3
2 An Overview of Empirical Processes	9
2.1 The Main Features	9
2.2 Empirical Process Techniques	13
2.2.1 Stochastic Convergence	13
2.2.2 Entropy for Glivenko-Cantelli and Donsker Theorems	16
2.2.3 Bootstrapping Empirical Processes	19
2.2.4 The Functional Delta Method	21
2.2.5 Z-Estimators	24
2.2.6 M-Estimators	28
2.3 Other Topics	30
2.4 Exercises	31
2.5 Notes	32
3 Overview of Semiparametric Inference	33
3.1 Semiparametric Models and Efficiency	33
3.2 Score Functions and Estimating Equations	37
3.3 Maximum Likelihood Estimation	42

3.4	Other Topics	45
3.5	Exercises	46
3.6	Notes	46
4	Case Studies I	47
4.1	Linear Regression	48
4.1.1	Mean Zero Residuals	48
4.1.2	Median Zero Residuals	50
4.2	Counting Process Regression	52
4.2.1	The General Case	53
4.2.2	The Cox Model	56
4.3	The Kaplan-Meier Estimator	58
4.4	Efficient Estimating Equations for Regression	60
4.4.1	Simple Linear Regression	64
4.4.2	A Poisson Mixture Regression Model	66
4.5	Partly Linear Logistic Regression	67
4.6	Exercises	69
4.7	Notes	70
II	Empirical Processes	71
5	Introduction to Empirical Processes	73
6	Preliminaries for Empirical Processes	77
6.1	Metric Spaces	77
6.2	Outer Expectation	84
6.3	Linear Operators and Functional Differentiation	89
6.4	Proofs	92
6.5	Exercises	95
6.6	Notes	98
7	Stochastic Convergence	99
7.1	Stochastic Processes in Metric Spaces	99
7.2	Weak Convergence	103
7.2.1	General Theory	103
7.2.2	Spaces of Bounded Functions	109
7.3	Other Modes of Convergence	111
7.4	Proofs	116
7.5	Exercises	121
7.6	Notes	122
8	Empirical Process Methods	123
8.1	Maximal Inequalities	124
8.1.1	Orlicz Norms and Maxima	124

8.1.2	Maximal Inequalities for Processes	127
8.2	The Symmetrization Inequality and Measurability	134
8.3	Glivenko-Cantelli Results	140
8.4	Donsker Results	143
8.5	Exercises	147
8.6	Notes	148
9	Entropy Calculations	149
9.1	Uniform Entropy	150
9.1.1	VC-Classes	150
9.1.2	BUEI-Classes	156
9.2	Bracketing Entropy	160
9.3	Glivenko-Cantelli Preservation	162
9.4	Donsker Preservation	165
9.5	Proofs	167
9.6	Exercises	170
9.7	Notes	171
10	Bootstrapping Empirical Processes	173
10.1	The Bootstrap for Donsker Classes	174
10.1.1	An Unconditional Multiplier Central Limit Theorem	175
10.1.2	Conditional Multiplier Central Limit Theorems . . .	177
10.1.3	Bootstrap Central Limit Theorems	181
10.1.4	Continuous Mapping Results	183
10.2	The Bootstrap for Glivenko-Cantelli Classes	187
10.3	A Simple Z-Estimator Master Theorem	190
10.4	Proofs	192
10.5	Exercises	198
10.6	Notes	199
11	Additional Empirical Process Results	201
11.1	Bounding Moments and Tail Probabilities	202
11.2	Sequences of Functions	204
11.3	Contiguous Alternatives	208
11.4	Sums of Independent but not Identically Distributed Stochastic Processes	211
11.4.1	Central Limit Theorems	211
11.4.2	Bootstrap Results	215
11.5	Function Classes Changing with n	216
11.6	Dependent Observations	220
11.7	Proofs	223
11.8	Exercises	226
11.9	Notes	226
12	The Functional Delta Method	227

12.1	Main Results and Proofs	227
12.2	Examples	229
12.2.1	Composition	229
12.2.2	Integration	230
12.2.3	Product Integration	234
12.2.4	Inversion	238
12.2.5	Other Mappings	241
12.3	Exercises	241
12.4	Notes	242
13	Z-Estimators	243
13.1	Consistency	244
13.2	Weak Convergence	245
13.2.1	The General Setting	245
13.2.2	Using Donsker Classes	246
13.2.3	A Master Theorem and the Bootstrap	247
13.3	Using the Delta Method	249
13.4	Exercises	253
13.5	Notes	253
14	M-Estimators	255
14.1	The Argmax Theorem	256
14.2	Consistency	258
14.3	Rate of Convergence	259
14.4	Regular Euclidean M-Estimators	261
14.5	Non-Regular Examples	263
14.5.1	A Change-Point Model	263
14.5.2	Monotone Density Estimation	268
14.6	Exercises	271
14.7	Notes	272
15	Case Studies II	273
III	Semiparametric Inference	275
16	Introduction to Semiparametric Inference	277
17	Preliminaries for Semiparametric Inference	279
18	Semiparametric Models and Efficiency	281
19	Score Functions and Estimating Equations	283
20	Maximum Likelihood Estimation and Inference	285

21 Semiparametric M-Estimation	287
22 Case Studies III	289
References	291
List of Symbols	297

Part I

Overview

1

Introduction

Both empirical processes and semiparametric inference techniques have become increasingly important tools for solving statistical estimation and inference problems. These tools are particularly important when the statistical model for the data at hand is *semiparametric*, in that it has one or more unknown component which is a function, measure or some other infinite dimensional quantity. Semiparametric models also typically have one or more finite-dimensional Euclidean parameters of particular interest. The term *nonparametric* is often reserved for semiparametric models with no Euclidean parameters. Empirical process methods are powerful techniques for evaluating the large sample properties of estimators based on semiparametric models, including consistency, distributional convergence, and validity of the bootstrap. Semiparametric inference tools complement empirical process methods by evaluating whether estimators make efficient use of the data.

Consider, for example, the semiparametric model

$$(1.1) \quad Y = \beta'Z + e,$$

where $\beta, Z \in \mathbb{R}^p$ are restricted to bounded sets, prime denotes transpose, (Y, Z) are the observed data, $E[e|Z] = 0$ and $E[e^2|Z] \leq K < \infty$ almost surely, and $E[ZZ']$ is positive definite. Given an independent and identically distributed (i.i.d.) sample of such data $(Y_i, Z_i), i = 1 \dots n$, we are interested in estimating β without having to further specify the joint distribution of (e, Z) . This is a very simple semiparametric model, and we definitely do

not need empirical process methods to verify that

$$\hat{\beta} = \left[\sum_{i=1}^n Z_i Z_i' \right]^{-1} \sum_{i=1}^n Z_i Y_i$$

is consistent for β and that $\sqrt{n}(\hat{\beta} - \beta)$ is asymptotically normal with bounded variance.

A deeper question is whether $\hat{\beta}$ achieves the lowest possible variance among all “reasonable” estimators. One important criteria for “reasonableness” is *regularity* which is satisfied by $\hat{\beta}$ and most standard \sqrt{n} consistent estimators and which we will define more precisely in chapter 3. A regular estimator is *efficient* if it achieves the lowest possible variance among regular estimators. Semiparametric inference tools are required to establish this kind of optimality. Unfortunately, $\hat{\beta}$ does not have the lowest possible variance among all *regular* estimators, unless we are willing to make some very strong assumptions. For instance, $\hat{\beta}$ has optimal variance if we are willing to assume, in addition to what we have already assumed, that e has a Gaussian distribution and is independent of Z . In this instance, the model is almost fully parametric (except that the distribution of Z remains unspecified). Returning to the more general model in the previous paragraph, there is a modification of $\hat{\beta}$ that does have the lowest possible variance among regular estimators, but computation of this modified estimator requires estimation of the function $z \mapsto E[e^2|Z = z]$. We will explore this particular example in greater detail in chapters 3 and 4.

There is an interesting semiparametric model part way between the fully parametric Gaussian residual model and the more general model which only assumes $E[e|Z] = 0$ almost surely. This alternative model assumes that the residual e and covariate Z are independent, with $E[e] = 0$ and $E[e^2] < \infty$, but no additional restrictions are placed on the residual distribution F . Unfortunately, $\hat{\beta}$ still does not have optimal variance. However, $\hat{\beta}$ is a very good estimator and may be good enough for most purposes. In this setting, it may be useful to estimate F to determine whether the residuals are Gaussian. One promising estimator is

$$(1.2) \quad \hat{F}(t) = n^{-1} \sum_{i=1}^n 1 \left\{ Y_i - \hat{\beta}' Z_i \leq t \right\},$$

where $1\{A\}$ is the indicator of A . Empirical process methods can be used to show that \hat{F} is uniformly consistent for F and that $\sqrt{n}(\hat{F} - F)$ converges “in distribution” in a uniform sense to a certain Gaussian quantity, provided f is uniformly bounded. Quotes are used here because the convergence in question involves random real functions rather than Euclidean random variables. This kind of convergence is called *weak convergence* and is a generalization of convergence in distribution which will be defined more precisely in chapter 2.

Now we will consider a more complex example. Let $N(t)$ be a counting process over the finite interval $[0, \tau]$ which is free to jump as long as $t \in [0, V]$, where $V \in (0, \tau]$ is a random time. A counting process, by definition, is nonnegative and piecewise constant with positive jumps of size 1. Typically, the process counts a certain kind of event, such as hospitalizations, for an individual (see chapter 1 of Fleming and Harrington, 1991, or chapter II of Andersen, Borgan, Gill and Keiding, 1993). Define also the “at-risk” process $Y(t) = 1\{V \geq t\}$. This process indicates whether an individual is at-risk at time $t-$ (just to the left of t) for a jump in N at time t . Suppose we also have baseline covariates $Z \in \mathbb{R}^p$, and, for all $t \in [0, \tau]$, we assume

$$(1.3) \quad \mathbb{E}\{N(t)|Z\} = \int_0^t \mathbb{E}\{Y(s)|Z\} e^{\beta'Z} d\Lambda(s),$$

for some $\beta \in \mathbb{R}^p$ and continuous nondecreasing function $\Lambda(t)$ with $\Lambda(0) = 0$ and $0 < \Lambda(\tau) < \infty$. The model (1.3) is a variant of the “multiplicative intensity model” (see definition 4.2.1 of Fleming and Harrington, 1991). Basically, we are assuming that the mean of the counting process is proportional to $e^{\beta'Z}$. We also need to assume that $\mathbb{E}\{Y(\tau)\} > 0$, $\mathbb{E}\{N^2(\tau)\} < \infty$, and that Z is restricted to a bounded set, but we do not otherwise restrict the distribution of N . Given an i.i.d. sample $(N_i, Y_i, Z_i), i = 1 \dots n$, we are interested in estimating β and Λ .

Under mild regularity conditions, the estimating equation

$$(1.4) \quad U_n(t, \beta) = n^{-1} \sum_{i=1}^n \int_0^t [Z_i - E_n(s, \beta)] dN_i(s),$$

where

$$E_n(s, \beta) = \frac{n^{-1} \sum_{i=1}^n Z_i Y_i(s) e^{\beta'Z_i}}{n^{-1} \sum_{i=1}^n Y_i(s) e^{\beta'Z_i}},$$

can be used for estimating β . The motivation for this estimating equation is that it arises as the score equation from the celebrated Cox partial likelihood (Cox, 1975) for either failure time data (where the counting process $N(t)$ simply indicates whether the failure time has occurred by time t) or the multiplicative intensity model under an independent increment assumption on N (See chapter 4 of Fleming and Harrington, 1991). Interestingly, this estimating equation can be shown to work under the more general model (1.3). Specifically, we can establish that (1.4) has an asymptotically unique zero $\hat{\beta}$ at $t = \tau$, that $\hat{\beta}$ is consistent for β , and that $\sqrt{n}(\hat{\beta} - \beta)$ is asymptotically mean zero Gaussian. Empirical process tools are needed to accomplish this. These same techniques can also establish that

$$\hat{\Lambda}(t) = \int_0^t \frac{n^{-1} \sum_{i=1}^n dN_i(s)}{n^{-1} \sum_{i=1}^n Y_i(s) e^{\hat{\beta}'Z_i}}$$

is uniformly consistent for $\Lambda(t)$ over $t \in [0, \tau]$ and that $\sqrt{n}(\hat{\Lambda} - \Lambda)$ converges weakly to a mean zero Gaussian quantity. These methods can also be used to construct valid confidence intervals and confidence bands for β and Λ .

At this point, efficiency of these estimators is difficult to determine because so little has been specified about the distribution of N . Consider, however, the special case of right-censored failure time data. In this setting, the observed data are $(V_i, d_i, Z_i), i = 1 \dots n$, where $V_i = T_i \wedge C_i$, $d_i = 1\{T_i \leq C_i\}$, T_i is a failure time of interest with integrated hazard function $e^{\beta' Z_i} \Lambda(t)$ given Z_i , and C_i is a right censoring time independent of T_i given Z_i with distribution not depending on β or Λ . Here, $N_i(t) = d_i 1\{V_i \leq t\}$ and $a \wedge b$ denotes the minimum of a and b . Semi-parametric inference techniques can now establish that both $\hat{\beta}$ and $\hat{\Lambda}$ are efficient. We will revisit this example in greater detail in chapters 3 and 4.

In both of the previous examples, the estimator of the infinite-dimensional parameter (F in the first example and Λ in the second) is \sqrt{n} consistent, but slower rates of convergence for the infinite dimensional part are also possible. As a third example, consider the partly linear logistic regression model described in Mammen and van de Geer (1997) and van der Vaart (1998, page 405). The observed data are n independent realizations of the random triplet (Y, Z, U) , where $Z \in \mathbb{R}^p$ and $U \in \mathbb{R}$ are covariates which are not linearly dependent, and Y is a dichotomous outcome with

$$(1.5) \quad \mathbb{E}\{Y|Z, U\} = \nu[\beta'Z + \eta(U)],$$

where $\beta \in \mathbb{R}^p$, Z is restricted to a bounded set, $U \in [0, 1]$, $\nu(t) = 1/(1+e^{-t})$, and η is an unknown smooth function. We assume, for some integer $k \geq 1$, that the first $k - 1$ derivatives of η exist and are absolutely continuous with $J^2(\eta) \equiv \int_0^1 [\eta^{(k)}(t)]^2 dt < \infty$, where superscript (k) denotes the k -th derivative. Given an i.i.d. sample $X_i = (Y_i, Z_i, U_i), i = 1 \dots n$, we are interested in estimating β and η .

The conditional density at $Y = y$ given the covariates $(Z, U) = (z, u)$ has the form

$$p_{\beta, \eta}(x) = \{\nu[\beta'z + \eta(u)]\}^y \{1 - \nu[\beta'z + \eta(u)]\}^{1-y}.$$

This cannot be used directly for defining a likelihood since for any $1 \leq n < \infty$ and fixed sample x_1, \dots, x_n , there exists a sequence of smooth functions $\{\hat{\eta}_m\}$ satisfying our criteria which converges to $\hat{\eta}$, where $\hat{\eta}(u_i) = \infty$ when $y_i = 1$ and $\hat{\eta}(u_i) = -\infty$ when $y_i = 0$. The issue is that requiring $J(\hat{\eta}) < \infty$ does not restrict $\hat{\eta}$ on any finite collection of points. There are a number of methods for addressing this problem, including requiring $J(\hat{\eta}) \leq M_n$ for each n , where $M_n \uparrow \infty$ at an appropriately slow rate, or using a series of increasingly complex spline approximations.

An important alternative is to use the penalized log-likelihood

$$(1.6) \quad \tilde{L}_n(\beta, \eta) = n^{-1} \sum_{i=1}^n \log p_{\beta, \eta}(X_i) - \hat{\lambda}_n^2 J^2(\eta),$$

where $\hat{\lambda}_n$ is a possibly data-dependent *smoothing parameter*. \tilde{L}_n is maximized over β and η to obtain the estimators $\hat{\beta}$ and $\hat{\eta}$. Large values of $\hat{\lambda}_n$ lead to very smooth but somewhat biased $\hat{\eta}$, while small values of $\hat{\lambda}_n$ lead to less smooth but less biased $\hat{\eta}$. The proper trade-off between smoothness and bias is usually best achieved by data-dependent schemes such as cross-validation. If $\hat{\lambda}_n$ is chosen to satisfy $\hat{\lambda}_n = o_p(n^{-1/4})$ and $\hat{\lambda}_n^{-1} = O_p(n^{k/(2k+1)})$, then both $\hat{\beta}$ and $\hat{\eta}$ are uniformly consistent and $\sqrt{n}(\hat{\beta} - \beta)$ converges to a mean zero Gaussian vector. Furthermore, $\hat{\beta}$ can be shown to be efficient even though $\hat{\eta}$ is not \sqrt{n} consistent. More about this example will be discussed in chapter 4.

These three examples illustrate the goals of empirical process and semiparametric inference research as well as hint at the power of these methods for solving statistical inference problems involving infinite-dimensional parameters. The goal of the first part of this book is to present the key ideas of empirical processes and semiparametric inference in a concise and heuristic way, without being distracted by technical issues, and to provide motivation to pursue the subject in greater depth as given in the remaining parts of the book. Even for those anxious to pursue the subject in depth, the broad view contained in this first part provides a valuable context for learning the details.

Chapter 2 presents an overview of empirical process methods and results, while chapter 3 presents an overview of semiparametric inference techniques. Several case studies illustrating these methods, including further details on the examples given above, are presented in chapter 4, which concludes the overview part. The empirical process part of the book (part II) will be introduced in chapter 5, while the semiparametric inference part (part III) will be introduced in chapter 16.

2

An Overview of Empirical Processes

This chapter presents an overview of the main ideas and techniques of empirical process research. The emphasis is on those concepts which directly impact statistical estimation and inference. The major distinction between empirical process theory and more standard asymptotics is that the random quantities studied have realizations as functions rather than real numbers or vectors. Proofs of results and certain details in definitions are postponed until part II of the book.

We begin by defining and sketching the main features and asymptotic concepts of empirical processes with a view towards statistical issues. An outline of the main empirical process techniques covered in this book is presented next. This chapter concludes with a discussion of several additional related topics which will not be pursued in later chapters.

2.1 The Main Features

A *stochastic process* is a collection of random variables $\{X_t, t \in T\}$ on the same probability space, indexed by an arbitrary index set T . An *empirical process* is a stochastic process based on a random sample. For example, consider a random sample X_1, \dots, X_n of i.i.d. real random variables with distribution F . The *empirical distribution function* is

$$(2.1) \quad \mathbb{F}_n(t) = n^{-1} \sum_{i=1}^n 1\{X_i \leq t\},$$

where the index t is allowed to vary over $T = \mathbb{R}$, the real line.

More generally, we can consider a random sample X_1, \dots, X_n of independent draws from a probability measure P on an arbitrary sample space \mathcal{X} . We define the *empirical measure* to be $\mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$, where δ_x is the measure which assigns mass 1 at x and zero elsewhere. For a measurable function $f : \mathcal{X} \mapsto \mathbb{R}$, we denote $\mathbb{P}_n f = n^{-1} \sum_{i=1}^n f(X_i)$. For any class \mathcal{F} of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$, an empirical process $\{\mathbb{P}_n f, f \in \mathcal{F}\}$ can be defined. This simple approach can generate a surprising variety of empirical processes, many of which we will consider in later sections in this chapter as well as in part II.

Setting $\mathcal{X} = \mathbb{R}$, we can now re-express \mathbb{F}_n as the empirical process $\{\mathbb{P}_n f, f \in \mathcal{F}\}$, where $\mathcal{F} = \{1\{x \leq t\}, t \in \mathbb{R}\}$. Thus one can view the stochastic process \mathbb{F}_n as indexed by either $t \in \mathbb{R}$ or $f \in \mathcal{F}$. We will use either indexing approach, depending on which is most convenient for the task at hand. However, because of its generality, indexing empirical processes by classes of functions will be the primary approach taken throughout this book.

By the law of large numbers, we know that

$$(2.2) \quad \mathbb{F}_n(t) \xrightarrow{\text{as}} F(t)$$

for each $t \in \mathbb{R}$, where $\xrightarrow{\text{as}}$ denotes almost sure convergence. A primary goal of empirical process research is to study empirical processes as random functions over the associated index set. Each realization of one of these random functions is a *sample path*. To this end, Glivenko (1933) and Cantelli (1933) demonstrated that (2.2) could be strengthened to

$$(2.3) \quad \sup_{t \in \mathbb{R}} |\mathbb{F}_n(t) - F(t)| \xrightarrow{\text{as}} 0.$$

Another way of saying this is that the sample paths of F_n get uniformly closer to F as $n \rightarrow \infty$. Returning to general empirical processes, a class \mathcal{F} of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$, is said to be a *P-Glivenko-Cantelli* class if

$$(2.4) \quad \sup_{f \in \mathcal{F}} |\mathbb{P}_n f - P f| \xrightarrow{\text{as}^*} 0,$$

where $P f = \int_{\mathcal{X}} f(x) P(dx)$ and $\xrightarrow{\text{as}^*}$ is a mode of convergence slightly stronger than $\xrightarrow{\text{as}}$ but which will not be precisely defined until later in this chapter (both modes of convergence are equivalent in the setting of (2.3)). Sometimes the P in P -Glivenko-Cantelli can be dropped if the context is clear.

Returning to \mathbb{F}_n , we know by the central limit theorem that for each $t \in \mathbb{R}$

$$G_n(t) \equiv \sqrt{n} [\mathbb{F}_n(t) - F(t)] \rightsquigarrow G(t),$$

where \rightsquigarrow denotes convergence in distribution and $G(t)$ is a mean zero normal random variable with variance $F(t)[1 - F(t)]$. In fact, we know that G_n , simultaneously for all t in a finite set $T_k = \{t_1, \dots, t_k\} \in \mathbb{R}$, will converge in distribution to a mean zero multivariate normal vector $G = \{G(t_1), \dots, G(t_k)\}'$, where

$$(2.5) \quad \text{cov}[G(s), G(t)] = \text{E}[G(s)G(t)] = F(s \wedge t) - F(s)F(t)$$

for all $s, t \in T_k$, and where $a \wedge b$ is the minimum of a and b .

Much more can be said. Donsker (1952) showed that the sample paths of G_n , as functions on \mathbb{R} , converge in distribution to a certain stochastic process G . *Weak convergence* is the generalization of convergence in distribution from vectors of random variables to sample paths of stochastic processes. Donsker's result can be stated succinctly as $G_n \rightsquigarrow G$ in $\ell^\infty(\mathbb{R})$, where, for any index set T , $\ell^\infty(T)$ is the collection of all bounded functions $f : T \mapsto \mathbb{R}$. $\ell^\infty(T)$ is used in settings like this to remind us that we are thinking of distributional convergence in terms of the sample paths.

The limiting process G is a mean zero *Gaussian process* with $\text{E}[G(s)G(t)] = (2.5)$ for every $s, t \in \mathbb{R}$. A Gaussian process is a stochastic process $\{Z_t, t \in T\}$, where for every finite $T_k \subset T$, $\{Z_t, t \in T_k\}$ is multivariate normal, and where all sample paths are continuous in a certain sense which will be made more explicit later in this chapter. The process G can be written $G(t) = \mathbb{B}(F(t))$, where \mathbb{B} is a standard Brownian bridge on the unit interval. The process \mathbb{B} has covariance $s \wedge t - st$ and is equivalent to the process $\mathbb{W}(t) - t\mathbb{W}(1)$, for $t \in [0, 1]$, where \mathbb{W} is a *standard Brownian motion* process. The standard Brownian motion is a Gaussian process on $[0, \infty)$ with continuous sample paths, with $\mathbb{W}(0) = 0$, and with covariance $s \wedge t$. Both \mathbb{B} and \mathbb{W} are important examples of Gaussian processes.

Returning again to general empirical processes, define the random measure $\mathbb{G}_n = \sqrt{n}(\mathbb{P}_n - P)$, and, for any class \mathcal{F} of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$, let \mathbb{G} be a mean zero Gaussian process indexed by \mathcal{F} , with covariance $\text{E}[f(X)g(X)] - \text{E}f(X)\text{E}g(X)$ for all $f, g \in \mathcal{F}$, and having appropriately continuous sample paths. Both \mathbb{G}_n and \mathbb{G} can be thought of as being indexed by \mathcal{F} . We say that \mathcal{F} is P -Donsker if $\mathbb{G}_n \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$. The P and/or the $\ell^\infty(\mathcal{F})$ may be dropped if the context is clear. Donsker's (1952) theorem tells us that $\mathcal{F} = \{1\{x \leq t\}, t \in \mathbb{R}\}$ is Donsker for all probability measures which are based on some real distribution function F . With $f(x) = 1\{x \leq t\}$ and $g(x) = 1\{x \leq s\}$,

$$\text{E}[f(X)g(X)] - \text{E}f(X)\text{E}g(X) = F(s \wedge t) - F(s)F(t).$$

For this reason, \mathbb{G} is also referred to as a Brownian bridge.

Suppose we are interested in forming confidence bands for F over some subset $T \subset \mathbb{R}$. Because $\mathcal{F} = \{1\{x \leq t\}, t \in \mathbb{R}\}$ is Glivenko-Cantelli, we can uniformly consistently estimate the covariance $\sigma(s, t) = F(s \wedge t) - F(s)F(t)$ of G with $\hat{\sigma}(s, t) = \mathbb{F}_n(s \wedge t) - \mathbb{F}_n(s)\mathbb{F}_n(t)$. While such a covariance could be

used to form confidence bands when H is finite, it is of little use when H is infinite, such as when H is a subinterval of \mathbb{R} . In this case, it is preferable to make use of the Donsker result for G_n . Let $U_n = \sup_{t \in T} |G_n(t)|$. The *continuous mapping theorem* tells us that whenever a process $\{Z_n(t), t \in T\}$ converges weakly to a tight limiting process $\{Z(t), t \in T\}$ in $\ell^\infty(T)$, then $h(Z_n) \rightsquigarrow h(Z)$ in $h(\ell^\infty(T))$ for any continuous map h . In our setting $U_n = h(G_n)$, where $h(g) = \sup_{t \in T} |g(t)|$, for any $g \in \ell^\infty(\mathbb{R})$, is a continuous real function. Thus the continuous mapping theorem tells us that $U_n \rightsquigarrow U = \sup_{t \in \mathbb{R}} |G(t)|$. When F is continuous and $T = \mathbb{R}$, $U = \sup_{t \in [0,1]} |\mathbb{B}(t)|$ has a known distribution from which it is easy to compute quantiles. If we let u_p be the p -th quantile of U , then an asymptotically valid symmetric $1 - \alpha$ level confidence band for F is $\mathbb{F}_n \pm u_{1-\alpha}/\sqrt{n}$.

An alternative is to construct confidence bands based on a large number of bootstraps of \mathbb{F}_n . The bootstrap for \mathbb{F}_n can be written as $\hat{\mathbb{F}}_n(t) = n^{-1} \sum_{i=1}^n W_{ni} \mathbf{1}\{X_i \leq t\}$, where (W_{n1}, \dots, W_{nn}) is a multinomial random n -vector, with probabilities $1/n, \dots, 1/n$ and number of trials n , and which is independent of the data X_1, \dots, X_n . The conditional distribution of $\hat{G}_n = \sqrt{n}(\hat{\mathbb{F}}_n - \mathbb{F}_n)$ given X_1, \dots, X_n can be shown to converge weakly to the distribution of G in $\ell^\infty(\mathbb{R})$. Thus the bootstrap is an asymptotically valid way to obtain confidence bands for F .

Returning to the general empirical process set-up, let \mathcal{F} be a Donsker class and suppose we wish to construct confidence bands for $Ef(X)$ which are simultaneously valid for all $f \in \mathcal{H} \subset \mathcal{F}$. Provided certain second moment conditions hold on \mathcal{F} , the estimator $\hat{\sigma}(f, g) = \mathbb{P}_n[f(X)g(X)] - \mathbb{P}_n f(X)\mathbb{P}_n g(X)$ is consistent for $\sigma(f, g) = E[f(X)g(X)] - Ef(X)Eg(X)$ uniformly over all $f, g \in \mathcal{F}$. As with the empirical distribution function estimator, this covariance is enough to form confidence bands provided \mathcal{H} is finite. Fortunately, the bootstrap is always asymptotically valid when \mathcal{F} is Donsker and can therefore be used for infinite \mathcal{H} . More precisely, if $\hat{G}_n = \sqrt{n}(\hat{\mathbb{P}}_n - \mathbb{P}_n)$, where $\hat{\mathbb{P}}_n f = n^{-1} \sum_{i=1}^n W_{ni} f(X_i)$ and (W_{n1}, \dots, W_{nn}) is defined as before, then the conditional distribution of \hat{G}_n given the data converges weakly to \mathbb{G} in $\ell^\infty(\mathcal{F})$. Since this is true for all of \mathcal{F} , it is certainly true for any $\mathcal{H} \subset \mathcal{F}$. The bootstrap result for \mathbb{F}_n is clearly a special case of this more general result.

Many important statistics based on i.i.d. data cannot be written as empirical processes, but they can frequently be written in the form $\phi(\mathbb{P}_n)$, where \mathbb{P}_n is indexed by some \mathcal{F} and ϕ is a smooth map from $\ell^\infty(\mathcal{F})$ to some set B (possibly infinite-dimensional). Consider, for example, the quantile process $\xi_n(p) = \mathbb{F}_n^{-1}(p)$ for $p \in [a, b]$, where $H^{-1}(p) = \inf\{t : H(t) \geq p\}$ for a distribution function H and $0 < a < b < 1$. Here, $\xi_n = \phi(\mathbb{F}_n)$, where ϕ maps a distribution function H to H^{-1} . When the underlying distribution F is continuous over $N = [H^{-1}(a) - \epsilon, H^{-1}(b) + \epsilon] \subset [0, 1]$, for some $\epsilon > 0$, with continuous density f such that $0 < \inf_{t \in N} f(t) \leq \sup_{t \in N} f(t) < \infty$, then $\sqrt{n}(\xi_n(p) - \xi_p)$, where $\xi_p = F^{-1}(p)$, is uniformly

asymptotically equivalent to $-G_n(F^{-1}(p))/f(F^{-1}(p))$ and hence converges weakly to $G(F^{-1}(p))/f(F^{-1}(p))$ in $\ell^\infty([a, b])$. (Because the process G is symmetric around zero, both $-G$ and G have the same distribution.) The above weak convergence result is a special case of the *functional delta-method* principle which states that $\sqrt{n}[\phi(\mathbb{P}_n) - \phi(P)]$ converges weakly in B to $\phi'(\mathbb{G})$, whenever \mathcal{F} is Donsker and ϕ has a “Hadamard derivative” ϕ' which will be defined more precisely later in this chapter.

Many additional statistics can be written as zeros or maximizers of certain data-dependent processes. The former are known as *Z-estimators* and the latter as *M-estimators*. Consider the linear regression example given in chapter 1. Since $\hat{\beta}$ is the zero of $U_n(\beta) = \mathbb{P}_n[X(Y - X'\beta)]$, $\hat{\beta}$ is a Z-estimator. In contrast, the penalized likelihood estimators $(\hat{\beta}, \hat{\eta})$ in the partly linear logistic regression example of the same chapter are M-estimators since they are maximizers of $\tilde{L}(\beta, \eta)$ given in (1.6). As is the case with U_n and \tilde{L}_n , the data-dependent objective functions used in Z- and M- estimation are often empirical processes, and thus empirical process methods are frequently required when studying the large sample properties of the associated statistics.

The key attribute of empirical processes is that they are random functions—or stochastic processes—based on a random data sample. The main asymptotic issue is studying the limiting behavior of these processes in terms of their sample paths. Primary achievements in this direction are Glivenko–Cantelli results which extend the law of large numbers, Donsker results which extend the central limit theorem, the validity of the bootstrap for Donsker classes, and the functional delta method.

2.2 Empirical Process Techniques

In this section, we expand on several important techniques used in empirical processes. We first define and discuss several important kinds of stochastic convergence, including convergence in probability as well as almost sure and weak convergence. We then introduce the concept of entropy and introduce several Glivenko–Cantelli and Donsker theorems based on entropy. The empirical bootstrap and functional delta method are described next. A brief outline of Z- and M- estimator methods are then presented. This section is essentially a review in miniature of the main points covered in Part II of this book, with a minimum of technicalities.

2.2.1 Stochastic Convergence

When discussing convergence of stochastic processes, there is always a *metric space* (\mathbb{D}, d) implicitly involved, where \mathbb{D} is the space of possible values for the processes and d is a *metric* (distance measure), satisfying $d(x, y) \geq 0$,

$d(x, y) = d(y, x)$, $d(x, z) \leq d(x, y) + d(y, z)$, and $d(x, y) = 0$ if and only if $x = y$, for all $x, y, z \in \mathbb{D}$. Frequently, $\mathbb{D} = \ell^\infty(T)$, where T is the index set for the processes involved, and d is the uniform distance on \mathbb{D} , i.e., $d(x, y) = \sup_{t \in T} |x(t) - y(t)|$ for any $x, y \in \mathbb{D}$. We are primarily interested in the convergence properties of the sample paths of stochastic processes. Weak convergence, or convergence in distribution, of a stochastic process X_n happens when the sample paths of X_n begin to behave in distribution, as $n \rightarrow \infty$, more and more like a specific random process X . When X_n and X are *Borel measurable*, weak convergence is equivalent to saying that $Ef(X_n) \rightarrow Ef(X)$ for every bounded, continuous function $f : \mathbb{D} \mapsto \mathbb{R}$, where the notation $f : A \mapsto B$ means that f is a mapping from A to B , and where continuity is in terms of d . Hereafter, we will let $C_b(\mathbb{D})$ denote the space of bounded, continuous maps $f : \mathbb{D} \mapsto \mathbb{R}$. We will define Borel measurability in detail later in part II, but, for now, it is enough to say that lack of this property means that there are certain important subsets $A \subset \mathbb{D}$ where the probability that $X_n \in A$ is not defined.

In many statistical applications, X_n may not be Borel measurable. To resolve this problem, we need to introduce the notion of *outer expectation* for arbitrary maps $T : \Omega \mapsto \bar{\mathbb{R}} \equiv [-\infty, \infty]$, where Ω is the sample space. T is not necessarily a random variable because it is not necessarily Borel measurable. The outer expectation of T , denoted E^*T , is the infimum over all EU , where $U : \Omega \mapsto \mathbb{R}$ is measurable, $U \geq T$, and EU exists. For EU to exist, it must not be indeterminate, although it can be $\pm\infty$, provided the sign is clear. We analogously define inner expectation: $E_*T = -E^*[-T]$. There also exists a measurable function $T^* : \Omega \mapsto \mathbb{R}$, called the *minimal measurable majorant*, satisfying $T^*(\omega) \geq T(\omega)$ for all $\omega \in \Omega$ and which is almost surely the smallest measurable function $\geq T$. Furthermore, when $E^*T < \infty$, $E^*T = ET^*$. The *maximal measurable minorant* is $T_* = -(-T)^*$. We also define outer probability for possibly nonmeasurable sets: $P^*(A)$ as the infimum over all $P(B)$ with $A \subset B \subset \Omega$ and B a Borel measurable set. Inner probability is defined as $P_*(A) = 1 - P^*(\Omega - A)$. This use of outer measure permits defining weak convergence, for possibly nonmeasurable X_n , as $E^*f(X_n) \mapsto Ef(X)$ for all $f \in C_b(\mathbb{D})$. We denote this convergence by $X_n \rightsquigarrow X$. Notice that we require the limiting process X to be measurable. This definition of weak convergence also carries with it an implicit measurability requirement on X_n : $X_n \rightsquigarrow X$ implies that X_n is *asymptotically measurable*, in that $E^*f(X_n) - E_*f(X_n) \rightarrow 0$, for every $f \in C_b(\mathbb{D})$.

We now consider convergence in probability and almost surely. We say X_n converges to X in probability if $P\{d(X_n, X)^* > \epsilon\} \rightarrow 0$ for every $\epsilon > 0$, and we denote this $X_n \xrightarrow{P} X$. We say that X_n converges outer almost surely to X if there exists a sequence Δ_n of measurable random variables with $d(X_n, X) \leq \Delta_n$ for all n and with $P\{\limsup_{n \rightarrow \infty} \Delta_n = 0\} = 1$. We denote this kind of convergence $X_n \xrightarrow{\text{as}^*} X$. While these modes of convergence are slightly different than the standard ones, they are identical when all the

quantities involved are measurable. The properties of the standard modes are also generally preserved in these new modes. The major difference is that these new modes can accommodate many situations in statistics and in other fields which could not be as easily accommodated with the standard ones. As much as possible, we will keep measurability issues suppressed throughout this book, except where it is necessary for clarity. From this point on, the metric d of choice will be the uniform metric unless noted otherwise.

For almost all of the weak convergence applications in this book, the limiting quantity X will be *tight*, in the sense that the sample paths of X will have a certain minimum amount of smoothness. To be more precise, for an index set T , let ρ be a *semimetric* on T , in that ρ has all the properties of a metric except that $\rho(s, t) = 0$ does not necessarily imply $s = t$. We say that T is totally bounded by ρ if for every $\epsilon > 0$, there exists a finite collection $T_\epsilon = \{t_1, \dots, t_k\} \subset T$ such that for all $t \in T$, we have $\rho(t, s) \leq \epsilon$ for some $s \in T_\epsilon$. Now define $UC(T, \rho)$ to be the subset of $\ell^\infty(T)$ where each $x \in UC(T, \rho)$ satisfies

$$\lim_{\delta \downarrow 0} \sup_{s, t \in T \text{ with } \rho(s, t) \leq \delta} |x(t) - x(s)| = 0.$$

The “ UC ” refers to uniform continuity. The stochastic process X is tight if $X \in UC(T, \rho)$ almost surely for some ρ for which T is totally bounded. If X is a Gaussian process, then ρ can be chosen as $\rho(s, t) = (\text{var}[X(s) - X(t)])^{1/2}$. Tight Gaussian processes will be the most important limiting processes considered in this book.

Two conditions need to be met in order for X_n to converge weakly in $\ell^\infty(T)$ to a tight X . This is summarized in the following theorem which we present now but prove later in chapter 7:

THEOREM 2.1 *X_n converges weakly to a tight X in $\ell^\infty(T)$ if and only if:*

- (i) *For all finite $\{t_1, \dots, t_k\} \subset T$, the multivariate distribution of $\{X_n(t_1), \dots, X_n(t_k)\}$ converges to that of $\{X(t_1), \dots, X(t_k)\}$.*
- (ii) *There exists a semimetric ρ for which T is totally bounded and*

$$(2.6) \quad \lim_{\delta \downarrow 0} \limsup_{n \rightarrow \infty} \mathbf{P}^* \left\{ \sup_{s, t \in T \text{ with } \rho(s, t) < \delta} |X_n(s) - X_n(t)| > \epsilon \right\} = 0,$$

for all $\epsilon > 0$.

Condition (i) is convergence of all finite dimensional distributions and condition (ii) implies *asymptotic tightness*. In many applications, condition (i) is not hard to verify while condition (ii) is much more difficult.

In the empirical process setting based on i.i.d. data, we are interested in establishing that $\mathbb{G}_n \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$, where \mathcal{F} is some class of measurable

functions $f : \mathcal{X} \mapsto \mathbb{R}$, and where \mathcal{X} is the sample space. When $Ef^2(X) < \infty$ for all $f \in \mathcal{F}$, condition (i) above is automatically satisfied by the standard central limit theorem, whereas establishing condition (ii) is much more work and is the primary motivator behind the development of much of modern empirical process theory. Whenever \mathcal{F} is Donsker, the limiting process \mathbb{G} is always a tight Gaussian process, and \mathcal{F} is totally bounded by the semimetric $\rho(f, g) = \{\text{var}[f(X) - g(X)]\}^{1/2}$. Thus conditions (i) and (ii) of theorem 2.1 are both satisfied with $T = \mathcal{F}$, $X_n(f) = \mathbb{G}_n f$, and $X(f) = \mathbb{G}f$, for all $f \in \mathcal{F}$.

Another important result is the *continuous mapping theorem*. This theorem states that if $g : \mathbb{D} \mapsto \mathbb{E}$ is continuous at every point of a set $\mathbb{D}_0 \subset \mathbb{D}$, and if $X_n \rightsquigarrow X$, where X takes all its values in \mathbb{D}_0 , then $g(X_n) \rightsquigarrow g(X)$. For example, if \mathcal{F} is a Donsker class, then $\sup_{f \in \mathcal{F}} |\mathbb{G}_n f|$ has the same limiting distribution as $\sup_{f \in \mathcal{F}} |\mathbb{G}f|$, since the supremum map is uniformly continuous, i.e., $|\sup_{f \in \mathcal{F}} |x(f)| - \sup_{f \in \mathcal{F}} |y(f)|| \leq \sup_{f \in \mathcal{F}} |x(f) - y(f)|$ for all $x, y \in \ell^\infty(\mathcal{F})$. This fact can be used to construct confidence bands for Pf . The continuous mapping theorem has many other practical uses which we will utilize at various points throughout this book.

2.2.2 Entropy for Glivenko-Cantelli and Donsker Theorems

The major challenge in obtaining Glivenko-Cantelli or Donsker theorems for classes of functions \mathcal{F} is to somehow show that going from pointwise convergence to uniform convergence is feasible. Clearly the complexity, or *entropy*, of \mathcal{F} plays a major role. The easiest entropy to introduce is *entropy with bracketing*. For $1 \leq r < \infty$, Let $L_r(P)$ denote the collection of functions $g : \mathcal{X} \mapsto \mathbb{R}$ such that $\|g\|_{r,P} \equiv [\int_{\mathcal{X}} |g(x)|^r dP(x)]^{1/r} < \infty$. An ϵ -bracket in $L_r(P)$ is a pair of functions $l, u \in L_r(P)$ with $P\{l(X) \leq u(X)\} = 1$ and with $\|l - u\|_{r,P} \leq \epsilon$. A function $f \in \mathcal{F}$ lies in the bracket l, u if $P\{l(X) \leq f(X) \leq u(X)\} = 1$. The *bracketing number* $N_{[]}(\epsilon, \mathcal{F}, L_r(P))$ is the minimum number of ϵ -brackets in $L_r(P)$ needed to ensure that every $f \in \mathcal{F}$ lies in at least one bracket. The logarithm of the bracketing number is the entropy with bracketing. The following is one of the simplest Glivenko-Cantelli theorems (the proof is deferred until part II):

THEOREM 2.2 *Let \mathcal{F} be a class of measurable functions and suppose that $N_{[]}(\epsilon, \mathcal{F}, L_1(P)) < \infty$ for every $\epsilon > 0$. Then \mathcal{F} is P -Glivenko-Cantelli.*

Consider, for example, the empirical distribution function \mathbb{F}_n based on an i.i.d. sample X_1, \dots, X_n of real random variables with distribution F (which defines the probability measure P on $\mathcal{X} = \mathbb{R}$). In this setting, \mathbb{F}_n is the empirical process \mathbb{G}_n with class $\mathcal{F} = \{1\{x \leq t\}, t \in \mathbb{R}\}$. For any $\epsilon > 0$, a finite collection of real numbers $-\infty = t_1 < t_2 < \dots < t_k = \infty$ can be found so that $F(t_j-) - F(t_{j-1}) \leq \epsilon$ for all $1 < j \leq k$, $F(t_1) = 0$ and $F(t_k-) = 1$, where $H(t-) = \lim_{s \uparrow t} H(s)$ when such a limit exists.

This can always be done in such a way that $k \leq 2 + 1/\epsilon$. Consider the collection of brackets $\{(l_j, u_j), 1 < j \leq k\}$, with $l_j(x) = 1\{x \leq t_{j-1}\}$ and $u_j(x) = 1\{x < t_j\}$ (notice that u_j is not in \mathcal{F}). Now each $f \in \mathcal{F}$ is in at least one bracket and $\|u_j - l_j\|_{P,1} = F(t_j) - F(t_{j-1}) \leq \epsilon$ for all $1 < j \leq k$. Thus $N_{[]}(\epsilon, \mathcal{F}, L_1(P)) < \infty$ for every $\epsilon > 0$, and the conditions of theorem 2.2 are met.

Donsker theorems based on entropy with bracketing require more stringent conditions on the number of brackets needed to cover \mathcal{F} . The *bracketing integral*,

$$J_{[]}(\delta, \mathcal{F}, L_r(P)) \equiv \int_0^\delta \sqrt{\log N_{[]}(\epsilon, \mathcal{F}, L_r(P))} d\epsilon,$$

needs to be bounded for $r = 2$ and $\delta = \infty$ to establish that \mathcal{F} is Donsker. Hence the bracketing entropy is permitted to go to ∞ as $\epsilon \downarrow 0$, but not too quickly. For most of the classes \mathcal{F} of interest, the entropy does go to ∞ as $\epsilon \downarrow 0$. However, a surprisingly large number of these classes satisfy the conditions of theorem 2.3 below, our first Donsker theorem (which we prove in chapter 8):

THEOREM 2.3 *Let \mathcal{F} be a class of measurable functions with $J_{[]}(\infty, \mathcal{F}, L_2(P)) < \infty$. Then \mathcal{F} is P -Donsker.*

Returning again to the empirical distribution function example, we have for the ϵ -brackets used previously that $\|u_j - l_j\|_{P,2} = (\|u_j - l_j\|_{P,1})^{1/2} \leq \epsilon^{1/2}$. Hence the minimum number of L_2 ϵ -brackets needed to cover \mathcal{F} is bounded by $1 + 1/\epsilon^2$, since an L_1 ϵ^2 -bracket is an L_2 ϵ -bracket. For $\epsilon > 1$, the number of brackets needed is just 1. $J_{[]}(\infty, \mathcal{F}, L_2(P))$ will therefore be finite if $\int_0^1 \sqrt{\log(1 + 1/\epsilon^2)} d\epsilon < \infty$. Using the fact that $\log(1+a) \leq 1 + \log(a)$ for $a \geq 1$ and the variable substitution $u = 1 + \log(1/\epsilon^2)$, we obtain that this integral is bounded by $\int_0^\infty u^{1/2} e^{-u/2} du = \sqrt{2\pi}$. Thus the conditions of theorem 2.3 are easily satisfied. We now give two other examples of classes with bounded $L_r(P)$ bracketing integral. Parametric classes of the form $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ work, provided Θ is a bounded subset of \mathbb{R}^p and there exists an $m \in L_r(P)$ such that $|f_{\theta_1}(x) - f_{\theta_2}(x)| \leq m(x)\|\theta_1 - \theta_2\|$ for all $\theta_1, \theta_2 \in \Theta$. Here, $\|\cdot\|$ is the standard Euclidean norm on \mathbb{R}^p . The class \mathcal{F} of all monotone functions $f : \mathbb{R} \mapsto [0, 1]$ also works for all $1 \leq r < \infty$ and all probability measures P .

Entropy calculations for other classes that arise in statistical applications can be difficult. However, there are a number of techniques for doing this which are not difficult to apply in practice and which we will explore briefly later on in this section. Unfortunately, there are also many classes \mathcal{F} for which entropy with bracketing does not work at all. An alternative which can be useful in such settings is entropy based on *covering numbers*. The covering number $N(\epsilon, \mathcal{F}, L_r(Q))$ is the minimum number of $L_r(Q)$ ϵ -balls needed to cover \mathcal{F} , where an $L_r(Q)$ ϵ -ball around a function $g \in L_r(Q)$ is

the set $\{h \in L_r(Q) : \|h - g\|_{Q,r} < \epsilon\}$. For a collection of balls to cover \mathcal{F} , all elements of \mathcal{F} must be included in at least one of the balls, but it is not necessary that the centers of the balls be contained in \mathcal{F} . The *entropy* is the logarithm of the covering number. The bracketing entropy conditions in theorems 2.2 and 2.3 can be replaced by conditions based on the *uniform covering numbers*

$$(2.7) \quad \sup_Q N(\epsilon \|F\|_{Q,r}, \mathcal{F}, L_r(Q)),$$

where $F : \mathcal{X} \mapsto \mathbb{R}$ is an *envelope* for \mathcal{F} , meaning that $|f(x)| \leq F(x)$ for all $x \in \mathcal{X}$ and all $f \in \mathcal{F}$, and where the supremum is taken over all finitely discrete probability measures Q with $\|F\|_{Q,r} > 0$. A finitely discrete probability measure on \mathcal{X} puts mass only at a finite number of points in \mathcal{X} . Notice that the uniform covering number does not depend on the probability measure P for the observed data. The *uniform entropy integral* is

$$J(\delta, \mathcal{F}, L_r) = \int_0^\delta \sqrt{\log \sup_Q N(\epsilon \|F\|_{Q,r}, \mathcal{F}, L_r(Q))} d\epsilon,$$

where the supremum is taken over the same set used in (2.7).

The following two theorems (given without proof) are Glivenko-Cantelli and Donsker results for uniform entropy:

THEOREM 2.4 *Let \mathcal{F} be an appropriately measurable class of measurable functions with $\sup_Q N(\epsilon \|F\|_{1,Q}, \mathcal{F}, L_1(Q)) < \infty$ for every $\epsilon > 0$, where the supremum is taken over the same set used in (2.7). If $P^*F < \infty$, then \mathcal{F} is P -Glivenko-Cantelli.*

THEOREM 2.5 *Let \mathcal{F} be an appropriately measurable class of measurable functions with $J(1, \mathcal{F}, L_2) < \infty$. If $P^*F^2 < \infty$, then \mathcal{F} is P -Donsker.*

Discussion of the ‘‘appropriately measurable’’ condition will be postponed until part II, but suffice it to say that it is satisfied for many function classes of interest in statistical applications.

An important collection of function classes \mathcal{F} , which satisfies $J(1, \mathcal{F}, L_r) < \infty$ for any $1 \leq r < \infty$, are the *Vapnik-Červonenkis* classes, or VC classes. Many classes of interest in statistics are VC, including the class of indicator functions explored earlier in the empirical distribution function example and also vector space classes. A vector space class \mathcal{F} has the form $\{\sum_{i=1}^k \lambda_i f_i(x), (\lambda_1, \dots, \lambda_k) \in \mathbb{R}^k\}$ for fixed functions f_1, \dots, f_k . We will postpone further definition and discussion of VC classes until part II.

The important thing to know at this point is that one does not need to calculate entropy for each new problem. There are a number of easy methods which can be used to determine whether a given class is Glivenko-Cantelli or Donsker based on whether the class is built up of other, well-known classes. For example, subsets of Donsker classes are Donsker since

condition (ii) of theorem 2.1 is clearly satisfied for any subset of T if it is satisfied for T . One can also use theorem 2.1 to show that finite unions of Donsker classes are Donsker. When \mathcal{F} and \mathcal{G} are Donsker, the following are also Donsker: $\{f \wedge g : f \in \mathcal{F}, g \in \mathcal{G}\}$, $\{f \vee g : f \in \mathcal{F}, g \in \mathcal{G}\}$, where \vee denotes maximum, and $\{f + g : f \in \mathcal{F}, g \in \mathcal{G}\}$. If \mathcal{F} and \mathcal{G} are bounded Donsker classes, then $\{fg : f \in \mathcal{F}, g \in \mathcal{G}\}$ is Donsker. Also, Lipschitz continuous functions of Donsker classes are Donsker. Furthermore, if \mathcal{F} is Donsker, then it is also Glivenko-Cantelli. These, and many other tools for verifying that a given class is Glivenko-Cantelli or Donsker, will be discussed in greater detail in chapter 9.

2.2.3 Bootstrapping Empirical Processes

An important aspect of inference for empirical processes is to be able to obtain covariance and confidence band estimates. The limiting covariance for a P -Donsker class \mathcal{F} is $\sigma : \mathcal{F} \times \mathcal{F} \mapsto \mathbb{R}$, where $\sigma(f, g) = Pf g - PfPg$. The covariance estimate $\hat{\sigma} : \mathcal{F} \times \mathcal{F} \mapsto \mathbb{R}$, where $\hat{\sigma}(f, g) = \mathbb{P}_n f g - \mathbb{P}_n f \mathbb{P}_n g$, is uniformly consistent for σ outer almost surely if and only if $P^* \left[\sup_{f \in \mathcal{F}} (f(X) - Pf)^2 \right] < \infty$. This will be proved later in part II. However, this is only of limited use since critical values for confidence bands cannot in general be determined from the covariance when \mathcal{F} is not finite. The bootstrap is an effective alternative.

As mentioned earlier, some care must be taken to ensure that the concept of weak convergence makes sense when the statistics of interest may not be measurable. This issue becomes more delicate with bootstrap results which involve convergence of conditional laws given the observed data. In this setting, there are two sources of randomness, the observed data and the resampling done by the bootstrap. For this reason, convergence of conditional laws is assessed in a slightly different manner than regular weak convergence. An important result is that $X_n \rightsquigarrow X$ in the metric space (\mathbb{D}, d) if and only if

$$(2.8) \quad \sup_{f \in BL_1} |E^* f(X_n) - Ef(X)| \rightarrow 0,$$

where BL_1 is the space of functions $f : \mathbb{D} \mapsto \mathbb{R}$ with Lipschitz norm bounded by 1, i.e., $\|f\|_\infty \leq 1$ and $|f(x) - f(y)| \leq d(x, y)$ for all $x, y \in \mathbb{D}$, and where $\|\cdot\|_\infty$ is the uniform norm.

We can now use this alternative definition of weak convergence to define convergence of the conditional limit laws of bootstraps. Let \hat{X}_n be a sequence of bootstrapped processes in \mathbb{D} with random weights which we will denote M . For some tight process X in \mathbb{D} , we use the notation $\hat{X}_n \overset{P}{\rightsquigarrow}_M X$ to mean that $\sup_{h \in BL_1} \left| E_M h(\hat{X}_n) - Eh(X) \right| \xrightarrow{P} 0$ and $E_M h(\hat{X}_n)^* - E_M h(\hat{X}_n)_* \xrightarrow{P} 0$, for all $h \in BL_1$, where the subscript M in the expectations indicates

conditional expectation over the weights M given the remaining data, and where $h(\hat{X}_n)^*$ and $h(\hat{X}_n)_*$ denote measurable majorants and minorants with respect to the joint data (including the weights M). We use the notation $\hat{X}_n \xrightarrow[M]{\text{as}^*} X$ to mean the same thing except with all $\overset{P}{\rightarrow}$'s replaced by $\overset{\text{as}^*}{\rightarrow}$'s. Note that the $h(\hat{X}_n)$ inside of the supremum does not have an asterisk: this is because Lipschitz continuous function of the bootstrapped processes we will study in this book will always be measurable functions of the random weights when conditioning on the data.

As mentioned previously, the bootstrap empirical measure can be defined as $\hat{\mathbb{P}}_n f = n^{-1} \sum_{i=1}^n W_{ni} f(X_i)$, where $\vec{W}_n = (W_{n1}, \dots, W_{nn})$ is a multinomial vector with probabilities $(1/n, \dots, 1/n)$ and number of trials n , and where \vec{W}_n is independent of the data sequence $\vec{X} = (X_1, X_2, \dots)$. We can now define a useful and simple alternative to this standard non-parametric bootstrap. Let $\vec{\xi} = (\xi_1, \xi_2, \dots)$ be an infinite sequence of non-negative i.i.d. random variables, also independent of \vec{X} , which have mean $0 < \mu < \infty$ and variance $0 < \tau^2 < \infty$, and which satisfy $\|\xi\|_{2,1} < \infty$, where $\|\xi\|_{2,1} = \int_0^\infty \sqrt{P(|\xi| > x)} dx$. This last condition is slightly stronger than bounded second moment but is implied whenever the $2 + \epsilon$ moment exists for any $\epsilon > 0$. We can now define a *multiplier bootstrap* empirical measure $\tilde{\mathbb{P}}_n f = n^{-1} \sum_{i=1}^n (\xi_i / \bar{\xi}_n) f(X_i)$, where $\bar{\xi}_n = n^{-1} \sum_{i=1}^n \xi_i$ and $\tilde{\mathbb{P}}_n$ is defined to be zero if $\bar{\xi}_n = 0$. Note that the weights add up to n for both bootstraps. When ξ_1 has a standard exponential distribution, for example, the moment conditions are clearly satisfied, and the resulting multiplier bootstrap has Dirichlet weights.

Under these conditions, we have the following two theorems (which we prove in part II), for convergence of the bootstrap, both in probability and outer almost surely. Let $\hat{\mathbb{G}}_n = \sqrt{n}(\hat{\mathbb{P}}_n - \mathbb{P}_n)$, $\tilde{\mathbb{G}}_n = \sqrt{n}(\mu/\tau)(\tilde{\mathbb{P}}_n - \mathbb{P}_n)$, and \mathbb{G} be the standard Brownian bridge in $\ell^\infty(\mathcal{F})$.

THEOREM 2.6 *The following are equivalent:*

- (i) \mathcal{F} is P -Donsker.
- (ii) $\hat{\mathbb{G}}_n \xrightarrow[W]{P} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ and the sequence $\hat{\mathbb{G}}_n$ is asymptotically measurable.
- (iii) $\tilde{\mathbb{G}}_n \xrightarrow[\xi]{P} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ and the sequence $\tilde{\mathbb{G}}_n$ is asymptotically measurable.

THEOREM 2.7 *The following are equivalent:*

- (i) \mathcal{F} is P -Donsker and $P^* [\sup_{f \in \mathcal{F}} (f(X) - Pf)^2] < \infty$.
- (ii) $\hat{\mathbb{G}}_n \xrightarrow[W]{\text{as}^*} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$.
- (iii) $\tilde{\mathbb{G}}_n \xrightarrow[\xi]{\text{as}^*} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$.

According to theorem 2.7, the almost sure consistency of the bootstrap requires the same moment condition required for almost sure uniform consistency of the covariance estimator $\hat{\sigma}$. In contrast, the consistency in probability of the bootstrap given in theorem 2.6 only requires that \mathcal{F} is Donsker. Thus consistency in probability of the bootstrap empirical process is an automatic consequence of weak convergence in the first place. Fortunately, consistency in probability is adequate for most statistical applications, since this implies that confidence bands constructed from the bootstrap are asymptotically valid. This follows because, as we will also establish in part II, whenever the conditional law of a bootstrapped quantity (say \hat{X}_n) in a normed space (with norm $\|\cdot\|$) converges to a limiting law (say of X), either in probability or outer almost surely, then the conditional law of $\|\hat{X}_n\|$ converges to that of $\|X\|$ under mild regularity conditions. We will also establish a slightly more general in-probability continuous mapping theorem for the bootstrap when the continuous map g is real valued.

Suppose we wish to construct a $1 - \alpha$ level confidence band for $\{Pf, f \in \mathcal{F}\}$, where \mathcal{F} is P -Donsker. We can obtain a large number, say N , bootstrap realizations of $\sup_{f \in \mathcal{F}} |\hat{G}_n f|$ to estimate the $1 - \alpha$ quantile of $\sup_{f \in \mathcal{F}} |Gf|$. If we call this estimate $\hat{c}_{1-\alpha}$, then theorem 2.6 tells us that $\{\mathbb{P}_n f \pm \hat{c}_{1-\alpha}, f \in \mathcal{F}\}$ has coverage $1 - \alpha$ for large enough n and N . For a more specific example, consider estimating $F(t_1, t_2) = P\{Y_1 \leq t_1, Y_2 \leq t_2\}$, where $X = (Y_1, Y_2)$ has an arbitrary bivariate distribution. We can estimate $F(t_1, t_2)$ with $\hat{F}_n(t_1, t_2) = n^{-1} \sum_{i=1}^n 1\{Y_{1i} \leq t_1, Y_{2i} \leq t_2\}$. This is the same as estimating $\{Pf, f \in \mathcal{F}\}$, where $\mathcal{F} = \{f(x) = 1\{y_1 \leq t_1, y_2 \leq t_2\} : t_1, t_2 \in \mathbb{R}\}$. This is a bounded Donsker class since $\mathcal{F} = \{f_1 f_2 : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}$, where $\mathcal{F}_j = \{1\{y_j \leq t\}, t \in \mathbb{R}\}$ is a bounded Donsker class for $j = 1, 2$. We thus obtain consistency in probability of the bootstrap. We also obtain outer almost sure consistency of the bootstrap by theorem 2.7, since \mathcal{F} is bounded by 1.

2.2.4 The Functional Delta Method

Suppose X_n is a sequence of random variables with $\sqrt{n}(X_n - \theta) \rightsquigarrow X$ for some $\theta \in \mathbb{R}^p$, and the function $\phi : \mathbb{R}^p \mapsto \mathbb{R}^q$ has a derivative $\phi'(\theta)$ at θ . The standard delta method now tells us that $\sqrt{n}(\phi(X_n) - \phi(\theta)) \rightsquigarrow \phi'(\theta)X$. However, many important statistics based on i.i.d. data involve maps from empirical processes to spaces of functions, and hence cannot be handled by the standard delta method. A simple example is the map ϕ_ξ which takes cumulative distribution functions H and computes $\{\xi_p, p \in [a, b]\}$, where $\xi_p = H^{-1}(p) = \inf\{t : H(t) \geq p\}$ and $[a, b] \subset (0, 1)$. The sample p -th quantile is then $\hat{\xi}_n(p) = \phi_\xi(\mathbb{F}_n)(p)$. Although the standard delta method cannot be used here, the functional delta method can be.

Before giving the main functional delta method results, we need to define derivatives for functions between normed spaces \mathbb{D} and \mathbb{E} . A normed space

is a metric space (\mathbb{D}, d) , where $d(x, y) = \|x - y\|$, for any $x, y \in \mathbb{D}$, and where $\|\cdot\|$ is a norm. A norm satisfies $\|x + y\| \leq \|x\| + \|y\|$, $\|\alpha x\| = |\alpha| \|x\|$, $\|x\| \geq 0$, and $\|x\| = 0$ if and only if $x = 0$, for all $x, y \in \mathbb{D}$ and all complex numbers α . A map $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ is *Gateaux-differentiable* at $\theta \in \mathbb{D}$, if for every fixed $h \in \mathbb{D}$ with $\theta + th \in \mathbb{D}_\phi$ for all $t > 0$ small enough, there exists an element $\phi'_\theta(h) \in \mathbb{E}$ such that

$$\frac{\phi(\theta + th) - \phi(\theta)}{t} \rightarrow \phi'_\theta(h)$$

as $t \downarrow 0$. For the functional delta method, however, we need ϕ to have the stronger property of being *Hadamard-differentiable*. A map $\phi : \mathbb{D}_\phi \mapsto \mathbb{E}$ is Hadamard-differentiable at $\theta \in \mathbb{D}$, tangentially to a set $\mathbb{D}_0 \subset \mathbb{D}$, if there exists a continuous linear map $\phi'_\theta : \mathbb{D} \mapsto \mathbb{E}$ such that

$$\frac{\phi(\theta + t_n h_n) - \phi(\theta)}{t_n} \rightarrow \phi'_\theta(h),$$

as $n \rightarrow \infty$, for all converging sequences $t_n \rightarrow 0$ and $h_n \rightarrow h \in \mathbb{D}_0$, with $h_n \in \mathbb{D}$ and $\theta + t_n h_n \in \mathbb{D}_\phi$ for all $n \geq 1$ sufficiently large.

For example, let $\mathbb{D} = D[0, 1]$, where DA , for any interval $A \subset \mathbb{R}$, is the space of cadlag (right-continuous with left-hand limits) real functions on A with the uniform norm. Let $\mathbb{D}_\phi = \{f \in D[0, 1] : |f| > 0\}$. Consider the function $\phi : \mathbb{D}_\phi \mapsto \mathbb{E} = D[0, 1]$ defined by $\phi(g) = 1/g$. Notice that for any $\theta \in \mathbb{D}_\phi$, we have, for any converging sequences $t_n \downarrow 0$ and $h_n \rightarrow h \in \mathbb{D}$, with $h_n \in \mathbb{D}$ and $\theta + t_n h_n \in \mathbb{D}_\phi$ for all $n \geq 1$,

$$\frac{\phi(\theta + t_n h_n) - \phi(\theta)}{t_n} = \frac{1}{t_n(\theta + t_n h_n)} - \frac{1}{t_n \theta} = -\frac{h_n}{\theta(\theta + t_n h_n)} \rightarrow -\frac{h}{\theta^2},$$

where we have suppressed the argument in g for clarity. Thus ϕ is Hadamard-differentiable, tangentially to \mathbb{D} , with $\phi'_\theta(h) = -h/\theta^2$.

Sometimes Hadamard differentiability is also called *compact differentiability*. Another important property of this kind of derivative is that it satisfies a chain rule, in that compositions of Hadamard-differentiable functions are also Hadamard-differentiable. Details on this and several other aspects of functional differentiation will be postponed until part II. We have the following important result (the proof of which will be given in part II):

THEOREM 2.8 *For normed spaces \mathbb{D} and \mathbb{E} , let $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ be Hadamard-differentiable at θ tangentially to $\mathbb{D}_0 \subset \mathbb{D}$. Assume that $r_n(X_n - \theta) \rightsquigarrow X$ for some sequence of constants $r_n \rightarrow \infty$, where X_n takes its values in \mathbb{D}_ϕ , and X is a tight process taking its values in \mathbb{D}_0 . Then $r_n(\phi(X_n) - \phi(\theta)) \rightsquigarrow \phi'_\theta(X)$.*

Consider again the quantile map ϕ_ξ , and let the distribution function F be absolutely continuous over $N = [u, v] = [F^{-1}(a) - \epsilon, F^{-1}(b) + \epsilon]$, for some $\epsilon > 0$, with continuous density f such that $0 < \inf_{t \in N} f(t) \leq$

$\sup_{t \in N} f(t) < \infty$. Also let $\mathbb{D}_1 \subset D[u, v]$ be the space of all distribution functions restricted to $[u, v]$. We will now argue that ϕ_ξ is Hadamard-differentiable at F tangentially to $C[u, v]$, where for any interval $A \subset \mathbb{R}$, CA is the space of continuous real functions on A . Let $t_n \rightarrow 0$ and $\{h_n\} \in D[u, v]$ converge uniformly to $h \in C[u, v]$ such that $F + t_n h_n \in \mathbb{D}_1$ for all $n \geq 1$, and denote $\xi_p = F^{-1}(p)$, $\xi_{pn} = (F + t_n h_n)^{-1}(p)$, $\xi_{pn}^N = (\xi_{pn} \vee u) \wedge v$, and $\epsilon_{pn} = t_n^2 \wedge (\xi_{pn}^N - u)$. The reason for the modification ξ_{pn}^N is to ensure that the quantile estimate is contained in $[u, v]$ and hence also $\epsilon_{pn} \geq 0$. Thus there exists an $n_0 < \infty$, such that for all $n \geq n_0$, $(F + t_n h_n)(u) < a$, $(F + t_n h_n)(v) > b$, $\epsilon_{pn} > 0$ and $\xi_{pn}^N = \xi_{pn}$ for all $p \in [a, b]$, and therefore

$$(2.9) \quad (F + t_n h_n)(\xi_{pn}^N - \epsilon_{pn}) \leq F(\xi_p) \leq (F + t_n h_n)(\xi_{pn}^N)$$

for all $p \in [a, b]$, since $(F + t_n h_n)^{-1}(p)$ is the smallest x satisfying $(F + t_n h_n)(x) \geq p$ and $F(\xi_p) = p$.

Since $F(\xi_{pn}^N - \epsilon_{pn}) = F(\xi_{pn}^N) + O(\epsilon_{pn})$, $h_n(\xi_{pn}^N) - h(\xi_{pn}^N) = o(1)$, and $h_n(\xi_{pn}^N - \epsilon_{pn}) - h(\xi_{pn}^N - \epsilon_{pn}) = o(1)$, where O and o are uniform over $p \in [a, b]$ (here and for the remainder of our argument), we have that (2.9) implies

$$(2.10) \quad \begin{aligned} F(\xi_{pn}^N) + t_n h(\xi_{pn}^N - \epsilon_{pn}) + o(t_n) &\leq F(\xi_p) \\ &\leq F(\xi_{pn}^N) + t_n h(\xi_{pn}^N) + o(t_n). \end{aligned}$$

But this implies that $F(\xi_{pn}^N) + O(t_n) \leq F(\xi_p) \leq F(\xi_{pn}^N) + O(t_n)$, which implies that $|\xi_{pn} - \xi_p| = O(t_n)$. This, together with (2.10) and the fact that h is continuous, implies that $F(\xi_{pn}) - F(\xi_p) = -t_n h(\xi_p) + o(t_n)$. This now yields

$$\frac{\xi_{pn} - \xi_p}{t_n} = -\frac{h(\xi_p)}{f(\xi_p)} + o(1),$$

and the desired Hadamard-differentiability of ϕ_ξ follows, with derivative $\phi'_F(h) = \{-h(F^{-1}(p))/f(F^{-1}(p)), p \in [a, b]\}$.

The functional delta method also applies to the bootstrap. Consider the sequence of random elements $\mathbb{X}_n(X_n)$ in a normed space \mathbb{D} , and assume that $r_n(\mathbb{X}_n - \mu) \rightsquigarrow \mathbb{X}$, where \mathbb{X} is tight in \mathbb{D} , for some sequence of constants $0 < r_n \rightsquigarrow \infty$. Here, \mathbb{X}_n is a generic empirical process based on the data sequence $\{X_n, n \geq 1\}$, and is not restricted to i.i.d. data. Now assume we have a bootstrap of \mathbb{X}_n , $\hat{\mathbb{X}}_n(X_n, W_n)$, where $W = \{W_n\}$ is a sequence of random bootstrap weights which are independent of X_n . Also assume $\hat{\mathbb{X}}_n \xrightarrow[W]{P} \mathbb{X}$. We have the following bootstrap result:

THEOREM 2.9 *For normed spaces \mathbb{D} and \mathbb{E} , let $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ be Hadamard-differentiable at μ tangentially to $\mathbb{D}_0 \subset \mathbb{D}$, with derivative ϕ'_μ . Let \mathbb{X}_n and $\hat{\mathbb{X}}_n$ have values in \mathbb{D}_ϕ , with $r_n(\mathbb{X}_n - \mu) \rightsquigarrow \mathbb{X}$, where \mathbb{X} is tight and takes its values in \mathbb{D}_0 , the maps $W_n \mapsto \hat{\mathbb{X}}_n$ are appropriately measurable, and where $r_n c(\hat{\mathbb{X}}_n - \mathbb{X}_n) \xrightarrow[W]{P} \mathbb{X}$. Then $r_n c(\phi(\hat{\mathbb{X}}_n) - \phi(\mathbb{X}_n)) \xrightarrow[W]{P} \phi'_\mu(\mathbb{X})$.*

We will postpone until part II a more precise discussion of what “appropriately measurable” means in this context.

When \mathbb{X}_n in the previous theorem is the empirical process \mathbb{P}_n indexed by a Donsker class \mathcal{F} and $r_n = \sqrt{n}$, the results of theorem 2.6 apply with $\mu = P$ for either the nonparametric or multiplier bootstrap weights. Moreover, the above measurability condition also holds (this will be verified in chapter 12). Thus the bootstrap is automatically valid for Hadamard-differentiable functions applied to empirical processes indexed by Donsker classes. As a simple example, bootstraps of the quantile process $\{\hat{\xi}_n(p), p \in [a, b] \subset (0, 1)\}$ are valid, provided the conditions given in the example following theorem 2.8 for the density f over the interval N are satisfied. This can be used, for example, to create asymptotically valid confidence bands for $\{F^{-1}(p), p \in [a, b]\}$. There are also results for outer almost sure conditional convergence of the conditional laws of the bootstrapped process $r_n(\phi(\hat{\mathbb{X}}_n) - \phi(\mathbb{X}_n))$, but this requires stronger conditions on the differentiability of ϕ , and we will not pursue this further in this book.

2.2.5 Z-Estimators

A Z-estimator $\hat{\theta}_n$ is the approximate zero of a data-dependent function. To be more precise, let the parameter space be Θ and let $\Psi_n : \Theta \mapsto \mathbb{L}$ be a data-dependent function between two normed spaces, with norms $\|\cdot\|$ and $\|\cdot\|_{\mathbb{L}}$, respectively. If $\|\Psi_n(\hat{\theta}_n)\|_{\mathbb{L}} \xrightarrow{P} 0$, then $\hat{\theta}_n$ is a Z-estimator. The main statistical issues for such estimators are consistency, asymptotic normality and validity of the bootstrap. Usually, Ψ_n is an estimator of a fixed function $\Psi : \Theta \mapsto \mathbb{L}$ with $\Psi(\theta_0) = 0$ for some parameter of interest $\theta_0 \in \Theta$. We save the proof of the following theorem as an exercise:

THEOREM 2.10 *Let $\Psi(\theta_0) = 0$ for some $\theta_0 \in \Theta$, and assume $\|\Psi(\theta_n)\|_{\mathbb{L}} \rightarrow 0$ implies $\|\theta_n - \theta_0\| \rightarrow 0$ for any sequence $\{\theta_n\} \in \Theta$ (this is an “identifiability” condition). Then*

- (i) *If $\|\Psi_n(\hat{\theta}_n)\|_{\mathbb{L}} \xrightarrow{P} 0$ for some sequence of estimators $\hat{\theta}_n \in \Theta$ and $\sup_{\theta \in \Theta} \|\Psi_n(\theta) - \Psi(\theta)\|_{\mathbb{L}} \xrightarrow{P} 0$, then $\|\hat{\theta}_n - \theta_0\| \xrightarrow{P} 0$.*
- (ii) *If $\|\Psi_n(\hat{\theta}_n)\|_{\mathbb{L}} \xrightarrow{\text{as*}} 0$ for some sequence of estimators $\hat{\theta}_n \in \Theta$ and $\sup_{\theta \in \Theta} \|\Psi_n(\theta) - \Psi(\theta)\|_{\mathbb{L}} \xrightarrow{\text{as*}} 0$, then $\|\hat{\theta}_n - \theta_0\| \xrightarrow{\text{as*}} 0$.*

Consider, for example, estimating the survival function for right-censored failure time data. In this setting, we observe $X = (U, \delta)$, where $U = T \wedge C$, $\delta = 1\{T \leq C\}$, T is a failure time of interest with distribution function F_0 and survival function $S_0 = 1 - F_0$ with $S_0(0) = 1$, and C is a censoring time with distribution and survival functions G and $L = 1 - G$, respectively, with $L(0) = 1$. For a sample of n observations $\{X_i, i = 1 \dots n\}$, let $\{\tilde{T}_j, j = 1 \dots m_n\}$ be the unique observed failure times. The *Kaplan-Meier estimator*

\hat{S}_n of S_0 is then given by

$$\hat{S}_n(t) = \prod_{j: \tilde{T}_j \leq t} \left(1 - \frac{\sum_{i=1}^n \delta_i 1\{U_i = \tilde{T}_j\}}{\sum_{i=1}^n 1\{U_i \geq \tilde{T}_j\}} \right).$$

Consistency and other properties of this estimator can be demonstrated via standard continuous-time martingale arguments (Fleming and Harrington, 1991; Andersen, Borgun, Keiding and Gill, 1993); however, it is instructive to use empirical process arguments for Z-estimators.

Let $\tau < \infty$ satisfy $L(\tau-)S_0(\tau-) > 0$, and let Θ be the space of all survival functions S with $S(0) = 1$ and restricted to $[0, \tau]$. We will use the uniform norm $\|\cdot\|_\infty$ on Θ . After some algebra, the Kaplan-Meier estimator can be shown to be the solution of $\Psi_n(\hat{S}_n) = 0$, where $\Psi_n : \Theta \mapsto \Theta$ has the form $\Psi_n(S)(t) = \mathbb{P}_n \psi_{S,t}$, where

$$\psi_{S,t}(X) = 1\{U > t\} + (1 - \delta)1\{U \leq t\}1\{S(U) > 0\} \frac{S(t)}{S(U)} - S(t).$$

This is Efron's (1967) "self-consistency" expression for the Kaplan-Meier. For the fixed function Ψ , we use $\Psi(S)(t) = P\psi_{S,t}$. Somewhat surprisingly, the class of function $\mathcal{F} = \{\psi_{S,t} : S \in \Theta, t \in [0, \tau]\}$ is P -Donsker. To see this, first note that the class \mathcal{M} of monotone functions $f : [0, \tau] \mapsto [0, 1]$ of the real random variable U has bounded entropy (with bracketing) integral, which fact we establish later in part II. Now the class of functions $\mathcal{M}_1 = \{\tilde{\psi}_{S,t} : S \in \Theta, t \in [0, \tau]\}$, where

$$\tilde{\psi}_{S,t}(U) = 1\{U > t\} + 1\{U \leq t\}1\{S(U) > 0\} \frac{S(t)}{S(U)},$$

is a subset of \mathcal{M} , since $\tilde{\psi}_{S,t}(U)$ is monotone in U on $[0, \tau]$ and takes values only in $[0, 1]$ for all $S \in \Theta$ and $t \in [0, \tau]$. Note that $(1\{U \leq t\} : t \in [0, \tau])$ is also Donsker (as argued previously), and so is $\{\delta\}$ (trivially) and $\{S(t) : S \in \Theta, t \in [0, \tau]\}$, since any class of fixed functions is always Donsker. Since all of these Donsker classes are bounded, we now have that \mathcal{F} is Donsker since sums and products of bounded Donsker classes are also Donsker. Since Donsker classes are also Glivenko-Cantelli, we have that $\sup_{S \in \Theta} \|\Psi_n(S) - \Psi(S)\|_\infty \xrightarrow{\text{as}^*} 0$. If we can establish the identifiability condition for Ψ , the outer almost sure version of theorem 2.10 gives us that $\|\hat{S}_n - S_0\|_\infty \xrightarrow{\text{as}^*} 0$.

After taking expectations, the function Ψ can be shown to have the form

$$(2.11) \quad \Psi(S)(t) = P\psi_{S,t} = S_0(t)L(t) + \int_0^t \frac{S_0(u)}{S(u)} dG(u)S(t) - S(t).$$

Thus, if we make the substitution $\epsilon_n(t) = S_0(t)/S_n(t) - 1$, $\Psi(S_n)(t) \rightarrow 0$ uniformly over $t \in [0, \tau]$ implies that $u_n(t) = \epsilon_n(t)L(t) + \int_0^t \epsilon_n(u)dG(u) \rightarrow 0$

uniformly over the same interval. By solving this integral equation, we obtain $\epsilon_n(t) = u_n(t)/L(t-) - \int_0^{t-} [L(s)L(s-)]^{-1} u_n(s) dG(s)$, which implies $\epsilon_n(t) \rightarrow 0$ uniformly, since $L(t-) \geq L(\tau-) > 0$. Thus $\|S_n - S_0\|_\infty \rightarrow 0$, implying the desired identifiability.

We now consider weak convergence of Z-estimators. Let Ψ_n, Ψ, Θ and \mathbb{L} be as at the beginning of this section. We have the following master theorem for Z-estimators, the proof of which will be given in part II:

THEOREM 2.11 *Assume that $\Psi(\theta_0) = 0$ for some θ_0 in the interior of Θ , $\sqrt{n}\Psi_n(\hat{\theta}_n) \xrightarrow{P} 0$, and $\|\hat{\theta}_n - \theta_0\| \xrightarrow{P} 0$ for the random sequence $\{\hat{\theta}_n\} \in \Theta$. Assume also that $\sqrt{n}(\Psi_n - \Psi)(\theta_0) \rightsquigarrow Z$, for some tight random Z , and that*

$$(2.12) \quad \frac{\left\| \sqrt{n}(\Psi_n(\hat{\theta}_n) - \Psi(\hat{\theta}_n)) - \sqrt{n}(\Psi_n(\theta_0) - \Psi(\theta_0)) \right\|_{\mathbb{L}}}{1 + \sqrt{n}\|\hat{\theta}_n - \theta_0\|} \xrightarrow{P} 0.$$

If $\theta \mapsto \Psi(\theta)$ is Fréchet-differentiable at θ_0 (defined below) with continuously-invertible (also defined below) derivative $\dot{\Psi}_{\theta_0}$, then

$$(2.13) \quad \left\| \sqrt{n}\dot{\Psi}_{\theta_0}(\hat{\theta}_n - \theta_0) + \sqrt{n}(\Psi_n - \Psi)(\theta_0) \right\|_{\mathbb{L}} \xrightarrow{P} 0$$

and thus $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow -\dot{\Psi}_{\theta_0}^{-1}(Z)$.

Fréchet-differentiability of a map $\phi : \Theta \subset \mathbb{D} \mapsto \mathbb{L}$ at $\theta \in \Theta$ is stronger than Hadamard-differentiability, in that it means there exists a continuous, linear map $\phi'_\theta : \mathbb{D} \mapsto \mathbb{L}$ with

$$(2.14) \quad \frac{\|\phi(\theta + h_n) - \phi(\theta) - \phi'_\theta(h_n)\|_{\mathbb{L}}}{\|h_n\|} \rightarrow 0$$

for all sequences $\{h_n\} \subset \mathbb{D}$ with $\|h_n\| \rightarrow 0$ and $\theta + h_n \in \Theta$ for all $n \geq 1$. *Continuous invertibility* of an operator $A : \Theta \mapsto \mathbb{L}$ essentially means A is invertible with the property that for a constant $c > 0$ and all $\theta_1, \theta_2 \in \Theta$,

$$(2.15) \quad \|A(\theta_1) - A(\theta_2)\|_{\mathbb{L}} \geq c\|\theta_1 - \theta_2\|.$$

An operator is a map between spaces of function, such as the maps Ψ and Ψ_n . We will postpone further discussion of operators and continuous invertibility until part II.

Returning to our Kaplan-Meier example, with $\Psi_n(S)(t) = \mathbb{P}_n \psi_{S,t}$ and $\Psi(S)(t) = P \psi_{S,t}$ as before, note that since $\mathcal{F} = \{\psi_{S,t}, S \in \Theta, t \in [0, \tau]\}$ is Donsker, we easily have that $\sqrt{n}(\Psi_n - \Psi)(\theta_0) \rightsquigarrow Z$, for $\theta_0 = S_0$ and some tight random Z . We also have that for any $\{S_n\} \in \Theta$ converging uniformly to S_0 ,

$$\begin{aligned} \sup_{t \in [0, \tau]} P(\psi_{S_n, t} - \psi_{S_0, t})^2 &\leq 2 \sup_{t \in [0, \tau]} \int_0^t \left[\frac{S_n(t)}{S_n(u)} - \frac{S_0(t)}{S_0(u)} \right]^2 S_0(u) dG(u) \\ &\quad + 2 \sup_{t \in [0, \tau]} (S_n(t) - S_0(t))^2 \\ &\rightarrow 0. \end{aligned}$$

This can be shown to imply (2.12). After some analysis, Ψ can be shown to be Fréchet-differentiable at S_0 , with derivative

$$(2.16) \quad \dot{\Psi}_{\theta_0}(h)(t) = - \int_0^t \frac{S_0(t)h(u)}{S_0(u)} dG(u) - L(t)h(t),$$

for all $h \in D[0, \tau]$, having continuous inverse

$$(2.17) \quad \begin{aligned} \dot{\Psi}_{\theta_0}^{-1}(a)(t) &= -S_0(t) \\ &\times \left\{ a(0) + \int_0^t \frac{1}{L(u-)S_0(u-)} \left[da(u) + \frac{a(u)dF_0(u)}{S_0(u)} \right] \right\}, \end{aligned}$$

for all $a \in D[0, \tau]$. Thus all of the conditions of theorem 2.11 are satisfied, and we obtain the desired weak convergence of $\sqrt{n}(\hat{S}_n - S_0)$ to a tight, mean zero Gaussian process. The covariance of this process is

$$V(s, t) = S_0(s)S_0(t) \int_0^{s \wedge t} \frac{dF_0(u)}{L(u-)S_0(u)S_0(u-)},$$

which can be derived after lengthy but straightforward calculations (which we omit).

Returning to general Z-estimators, there are a number of methods for showing that the conditional law of a bootstrapped Z-estimator, given the observed data, converges to the limiting law of the original Z-estimator. One important approach which is applicable to non-i.i.d. data involves establishing Hadamard-differentiability of the map ϕ which extracts a zero from the function Ψ . We will explore this approach in part II. We close this section with a simple bootstrap result for the setting where $\Psi_n(\theta)(h) = \mathbb{P}_n \psi_{\theta, h}$ and $\Psi(\theta)(h) = P \psi_{\theta, h}$, for random and fixed real maps indexed by $\theta \in \Theta$ and $h \in \mathcal{H}$. Assume that $\Psi(\theta_0)(h) = 0$ for some $\theta_0 \in \Theta$ and all $h \in \mathcal{H}$, that $\sup_{h \in \mathcal{H}} |\Psi(\theta_n)(h)| \rightarrow 0$ implies $\|\theta_n - \theta_0\| \rightarrow 0$ for any sequence $\{\theta_n\} \in \Theta$, and that Ψ is Fréchet-differentiable with continuously invertible derivative $\dot{\Psi}_{\theta_0}$. Also assume that $\mathcal{F} = \{\psi_{\theta, h} : \theta \in \Theta, h \in \mathcal{H}\}$ is P -G-C with $\sup_{\theta \in \Theta, h \in \mathcal{G}} P|\psi_{\theta, h}| < \infty$. Furthermore, assume that $\mathcal{G} = \{\psi_{\theta, h} : \theta \in \Theta, \|\theta - \theta_0\| \leq \delta, h \in \mathcal{H}\}$, where $\delta > 0$, is P -Donsker and that $\sup_{h \in \mathcal{H}} P(\psi_{\theta_n, h} - \psi_{\theta_0, h})^2 \rightarrow 0$ for any sequence $\{\theta_n\} \in \Theta$ with $\|\theta_n - \theta_0\| \rightarrow 0$. Then, using arguments similar to those used in the Kaplan-Meier example and with the help of theorems 2.10 and 2.11, we have that if $\hat{\theta}_n$ satisfies $\sup_{h \in \mathcal{H}} |\sqrt{n} \Psi_n(\hat{\theta}_n)| \xrightarrow{P} 0$, then $\|\hat{\theta}_n - \theta_0\| \xrightarrow{P} 0$ and $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow -\dot{\Psi}_{\theta_0}^{-1}(Z)$, where Z is the tight limiting distribution of $\sqrt{n}(\Psi_n(\theta_0) - \Psi(\theta_0))$.

Let $\Psi_n^\circ(\theta)(h) = \mathbb{P}_n^\circ \psi_{\theta, h}$, where \mathbb{P}_n° is either the nonparametric bootstrap $\hat{\mathbb{P}}_n$ or the multiplier bootstrap $\tilde{\mathbb{P}}_n$ defined in section 2.2.3, and define the bootstrap estimator $\hat{\theta}_n^\circ \in \Theta$ to be a minimizer of $\sup_{h \in \mathcal{H}} |\Psi_n^\circ(\theta)(h)|$ over $\theta \in \Theta$. We will prove in part II that these conditions are more than enough

to ensure that $\sqrt{n}(\hat{\theta}_n^\circ - \hat{\theta}_n) \overset{P}{\rightsquigarrow} -\dot{\Psi}_{\theta_0}^{-1}(Z)$, where W refers to either the nonparametric or multiplier bootstrap weights. Thus the bootstrap is valid. These conditions for the bootstrap are satisfied in the Kaplan-Meier example, for either the nonparametric or multiplier weights, thus enabling the construction of confidence bands for $S_0(t)$ over $t \in [0, \tau]$.

2.2.6 *M-Estimators*

An M-estimator $\hat{\theta}_n$ is the approximate maximum of a data-dependent function. To be more precise, let the parameter set be a metric space (Θ, d) and let $M_n : \Theta \mapsto \mathbb{R}$ be a data-dependent real function. If $M_n(\hat{\theta}_n) = \sup_{\theta \in \Theta} M_n(\theta) - o_p(1)$, then $\hat{\theta}_n$ is an M-estimator. Maximum likelihood and least-squares (after changing the sign of the objective function) estimators are some of the most important examples, but there are many other examples as well. As with Z-estimators, the main statistical issues for M-estimators are consistency, weak convergence and validity of the bootstrap. Unlike Z-estimators, the rate of convergence for M-estimators is not necessarily \sqrt{n} , even for i.i.d. data, and finding the right rate can be quite challenging.

For establishing consistency, M_n is often an estimator of a fixed function $M : \Theta \mapsto \mathbb{R}$. We now present the following consistency theorem (the proof of which is deferred to part II):

THEOREM 2.12 *Assume for some $\theta_0 \in \Theta$ that $\liminf_{n \rightarrow \infty} M(\theta_n) \geq M(\theta_0)$ implies $d(\theta_n, \theta_0) \rightarrow 0$ for any sequence $\{\theta_n\} \in \Theta$ (this is another identifiability condition). Then, for a sequence of estimators $\hat{\theta}_n \in \Theta$,*

(i) *If $M_n(\hat{\theta}_n) = \sup_{\theta \in \Theta} M_n(\theta) - o_p(1)$ and $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0$, then $d(\hat{\theta}_n, \theta_0) \xrightarrow{P} 0$.*

(ii) *If $M_n(\hat{\theta}_n) = \sup_{\theta \in \Theta} M_n(\theta) - o_{as^*}(1)$ and $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{as^*} 0$, then $d(\hat{\theta}_n, \theta_0) \xrightarrow{as^*} 0$.*

Suppose, for now, we know that the rate of convergence for the M-estimator $\hat{\theta}_n$ is r_n , or, in other words, we know that $Z_n = r_n(\hat{\theta}_n - \theta_0) = O_p(1)$. Z_n can now be re-expressed as the approximate maximum of the criterion function $h \mapsto H_n(h) = M_n(\theta_0 + h/r_n)$ for h ranging over some metric space \mathbb{H} . If the argmax of H_n over bounded subsets of \mathbb{H} can now be shown to converge weakly to the argmax of a tight limiting process H over the same bounded subsets, then Z_n converges weakly to $\operatorname{argmax}_{h \in \mathbb{H}} H(h)$.

We will postpone the technical challenges associated with determining these rates of convergence until part II, and restrict ourselves to an interesting special case involving Euclidean parameters, where the rate is known to be \sqrt{n} . The proof of the following theorem is also deferred to part II:

THEOREM 2.13 *Let X_1, \dots, X_n be i.i.d. with sample space \mathcal{X} and law P , and let $m_\theta : \mathcal{X} \mapsto \mathbb{R}$ be measurable functions indexed by θ ranging over an open subset of Euclidean space $\Theta \subset \mathbb{R}^p$. Let θ_0 be a bounded point of maximum of Pm_θ in the interior of Θ , and assume for some neighborhood $\Theta_0 \subset \Theta$ including θ_0 , that there exists measurable functions $\dot{m} : \mathcal{X} \mapsto \mathbb{R}$ and $\dot{m}_{\theta_0} : \mathcal{X} \mapsto \mathbb{R}^p$ satisfying*

$$(2.18) \quad |m_{\theta_1}(x) - m_{\theta_2}(x)| \leq \dot{m}(x)\|\theta_1 - \theta_2\|,$$

$$(2.19) \quad P[m_\theta - m_{\theta_0} - (\theta - \theta_0)'\dot{m}_{\theta_0}]^2 = o(\|\theta - \theta_0\|^2),$$

$P\dot{m}^2 < \infty$, and $P\|\dot{m}_{\theta_0}\|^2 < \infty$, for all $\theta_1, \theta_2, \theta \in \Theta_0$. Assume also that $M(\theta) = Pm_\theta$ admits a second order Taylor expansion with nonsingular second derivative matrix V . Denote $M_n(\theta) = \mathbb{P}_n m_\theta$, and assume the approximate maximizer $\hat{\theta}_n$ satisfies $M_n(\hat{\theta}_n) = \operatorname{argmax}_{\theta \in \Theta} M_n(\theta) - o_p(n^{-1})$ and $\|\hat{\theta}_n - \theta_0\| \xrightarrow{P} 0$. Then $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow -V^{-1}Z$, where Z is the limiting Gaussian distribution of $\mathbb{G}_n \dot{m}_{\theta_0}$.

Consider, for example, least absolute deviation regression. In this setting, we have i.i.d. random vectors U_1, \dots, U_n in \mathbb{R}^p and random errors e_1, \dots, e_n , but we observe only the data $X_i = (Y_i, U_i)$, where $Y_i = \theta_0' U_i + e_i$, $i = 1 \dots n$. The least-absolute-deviation estimator $\hat{\theta}_n$ minimizes the function $\theta \mapsto \mathbb{P}_n \tilde{m}_\theta$, where $\tilde{m}_\theta(X) = |Y - \theta' U|$. Since a minimizer of a criterion function M_n is also a maximizer of $-M_n$, M-estimation methods can be used in this context with only a change in sign. Although boundedness of the parameter space Θ is not necessary for this regression setting, we restrict—for ease of discourse— Θ to be a bounded, open subset of \mathbb{R}^p containing θ_0 . We also assume that the distribution of the errors e_i has median zero and positive density at zero, which we denote $f(0)$, and that $P[UU']$ is positive definite.

Note that since we are not assuming $E|e_i| < \infty$, it is possible that $P\tilde{m}_\theta = \infty$ for all $\theta \in \Theta$. Since minimizing $\mathbb{P}_n \tilde{m}_\theta$ is the same as minimizing $\mathbb{P}_n m_\theta$, where $m_\theta = \tilde{m}_\theta - \tilde{m}_{\theta_0}$, we will use $M_n(\theta) = \mathbb{P}_n m_\theta$ as our criterion function hereafter (without modifying the estimator $\hat{\theta}_n$). By the definition of Y , $m_\theta(X) = |e - (\theta - \theta_0)' U| - |e|$, and we now have that $Pm_\theta \leq \|\theta - \theta_0\| (E\|U\|^2)^{1/2} < \infty$ for all $\theta \in \Theta$. Since

$$(2.20) \quad |m_{\theta_1}(x) - m_{\theta_2}(x)| \leq \|\theta_1 - \theta_2\| \times \|u\|,$$

it is not hard to show that the class of function $\{m_\theta : \theta \in \Theta\}$ is P -Glivenko-Cantelli. It can also be shown that $Pm_\theta \geq 0$ with equality only when $\theta = \theta_0$. Hence theorem 2.12, part (ii), yields that $\hat{\theta}_n \xrightarrow{\text{as}^*} \theta_0$.

Now we consider $M(\theta) = Pm_\theta$. By conditioning on U , one can show after some analysis that $M(\theta)$ is two times continuously differentiable, with second derivative $V = 2f(0)P[UU']$ at θ_0 . Note that (2.20) satisfies condition (2.18); and with $\dot{m}_\theta(X) = -U \operatorname{sign}(e)$, we also have that

$$|m_\theta(X) - m_{\theta_0}(X) - (\theta - \theta_0)'\dot{m}_\theta(X)| \leq 1 \{|e| \leq |(\theta - \theta_0)' U|\} [(\theta - \theta_0)' U]^2$$

satisfies condition (2.19). Thus all the conditions of theorem 2.13 are satisfied. Hence $\sqrt{n}(\hat{\theta}_n - \theta_0)$ is asymptotically mean zero normal, with variance $V^{-1}P[\dot{m}_{\theta_0}\dot{m}'_{\theta_0}]V^{-1} = (P[UU'])^{-1} / (4f^2(0))$. This variance is not difficult to estimate from the data, but we postpone presenting the details.

Another technique for obtaining weak convergence of M-estimators which are \sqrt{n} consistent, is to first establish consistency and then take an appropriate derivative of the criterion function $M_n(\theta)$, $\Psi_n(\theta)(h)$, for h ranging over some index set H , and apply Z-estimator techniques to Ψ_n . This works because the derivative of a smooth criterion function at an approximate maximizer is approximately zero. This approach facilitates establishing the validity of the bootstrap since such validity is often easier to obtain for Z-estimators than for M-estimators. This approach is also applicable to certain nonparametric maximum likelihood estimators which we will consider in part III.

2.3 Other Topics

In addition to the empirical process topics outlined in the previous sections, we will cover a few other related topics in part II, including results for sums of independent but not identically distributed stochastic processes and, briefly, for dependent but stationary processes. However, there are a number of interesting empirical process topics we will not pursue in later chapters, including general results for convergence of nets. In the remainder of this section, we briefly outline a few additional topics not covered later which involve sequences of empirical processes based on i.i.d. data. For simplicity, we will primarily restrict ourselves to the empirical process $G_n = \sqrt{n}(\mathbb{F}_n - F)$, although many of these results have extensions which apply to more general empirical processes.

The law of the iterated logarithm for G_n states that

$$(2.21) \quad \limsup_{n \rightarrow \infty} \frac{\|G_n\|_{\infty}}{\sqrt{2 \log \log n}} \leq \frac{1}{2}, \quad \text{a.s.},$$

with equality if $1/2$ is in the range of F , where $\|\cdot\|_{\infty}$ is the uniform norm. This can be generalized to empirical processes on P -Donsker classes \mathcal{F} which have a measurable envelope with bounded second moment (Dudley and Philipp, 1983):

$$\limsup_{n \rightarrow \infty} \frac{[\sup_{f \in \mathcal{F}} |G_n(f)|]^*}{\sqrt{(2 \log \log n) \sup_{f \in \mathcal{F}} |P(f - Pf)^2|}} \leq 1, \quad \text{a.s.}$$

Result (2.21) can be further strengthened to Strassen's (1964) theorem, which states that on a set with probability 1, the set of all limiting paths of $\sqrt{1/(2 \log \log n)}G_n$ is exactly the set of all functions of the form $h(F)$,

where $h(0) = h(1) = 0$ and h is absolutely continuous with derivative h' satisfying $\int_0^1 [h'(s)]^2 ds \leq 1$. While the previous results give upper bounds on $\|G_n\|_\infty$, it is also known that

$$\liminf_{n \rightarrow \infty} \sqrt{2 \log \log n} \|G_n\|_\infty = \frac{\pi}{2}, \text{ a.s.},$$

implying that the smallest uniform distance between \mathbb{F}_n and F is at least $O(1/\sqrt{n \log \log n})$.

A topic of interest regarding Donsker theorems is the closeness of the empirical process sample paths to the limiting Brownian bridge sample paths. The strongest result on this question for the empirical process G_n is the KMT construction, named after Komlós, Major and Tusnády (1975, 1976). The KMT construction states that there exists fixed positive constants a , b , and c , and a sequence of standard Brownian bridges $\{\mathbb{B}_n\}$, such that

$$\mathbb{P} \left(\|G_n - \mathbb{B}_n(F)\|_\infty > \frac{a \log n + x}{\sqrt{n}} \right) \leq b e^{-cx},$$

for all $x > 0$ and $n \geq 1$. This powerful result can be shown to imply both

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{\sqrt{n}}{\log n} \|G_n - \mathbb{B}_n(F)\|_\infty &< \infty, \text{ a.s., and} \\ \limsup_{n \rightarrow \infty} \mathbb{E} \left[\frac{\sqrt{n}}{\log n} \|G_n - \mathbb{B}_n(F)\|_\infty \right]^m &< \infty, \end{aligned}$$

for all $0 < m < \infty$. These results are called *strong approximations* and have applications in statistics, such as in the construction of confidence bands for kernel density estimators (see, for example, Bickel and Rosenblatt, 1973).

2.4 Exercises

2.4.1. Let X, Y be a pair of real random numbers with joint distribution P . Compute upper bounds for $N_{[]}(\epsilon, \mathcal{F}, L_r(P))$, for $r = 1, 2$, where $\mathcal{F} = \{1\{X \leq s, Y \leq t\} : s, t \in \mathbb{R}\}$.

2.4.2. Prove theorem 2.10.

2.4.3. Consider the Z-estimation framework for the Kaplan-Meier estimator discussed in section 2.2.5. Let $\Psi(S)(t)$ be as defined in (2.11). Show that Ψ is Fréchet-differentiable at S_0 , with derivative $\dot{\Psi}_{\theta_0}(h)(t)$ given by (2.16), for all $h \in D[0, \tau]$.

2.4.4. Continuing with the set-up of the previous problem, show that $\dot{\Psi}_{\theta_0}$ is continuously invertible, with inverse $\dot{\Psi}_{\theta_0}^{-1}$ given in (2.17). The following approach may be easiest: First show that for any $a \in D[0, \tau]$,

$h(t) = \dot{\Psi}_{\theta_0}^{-1}(a)(t)$ satisfies $\dot{\Psi}_{\theta_0}(h)(t) = a(t)$. The following identity may be helpful:

$$d \left[\frac{a(t)}{S_0(t)} \right] = \frac{da(t)}{S_0(t-)} + \frac{a(t)dF_0(t)}{S_0(t-)S_0(t)}.$$

Now show that there exists an $M < \infty$ such that $\|\dot{\Psi}_{\theta_0}^{-1}(a)\| \leq M\|a\|$, where $\|\cdot\|$ is the uniform norm. This then implies that there exists a $c > 0$ such that $\|\dot{\Psi}_{\theta_0}(h)\| \geq c\|h\|$.

2.5 Notes

Theorem 2.1 is a composite of theorems 1.5.4 and 1.5.7 of van der Vaart and Wellner (1996) (hereafter abbreviated VW). Theorems 2.2, 2.3, 2.4 and 2.5 correspond to theorems 19.4, 19.5, 19.13 and 19.14, respectively, of van der Vaart (1998). The if and only if implications of (2.8) are described in VW, page 73. The implications (i) \Leftrightarrow (ii) in theorems 2.6 and 2.7 are given in theorems 3.6.1 and 3.6.2, respectively, of VW. Theorems 2.8 and 2.11 correspond to theorems 3.9.4 and 3.3.1, of VW, while theorem 2.13 comes from example 3.2.22 of VW.

3

Overview of Semiparametric Inference

This chapter presents an overview of the main ideas and techniques of semiparametric inference, with particular emphasis on semiparametric efficiency. The major distinction between this kind of efficiency and the standard notion of efficiency for parametric maximum likelihood estimators—as expressed in the Cramér-Rao lower bound—is the presence of an infinite-dimensional nuisance parameter in semiparametric models. Proofs and other technical details will generally be postponed until part III.

In the first section, we define and sketch the main features of semiparametric models and semiparametric efficiency. The second section discusses efficient score functions and estimating equations and their connection to efficient estimation. The third section discusses nonparametric maximum likelihood estimation, the main tool for constructing efficient estimators. The fourth and final section briefly discusses several additional related topics, including variance estimation and confidence band construction for efficient estimators.

3.1 Semiparametric Models and Efficiency

A *statistical model* is a collection of probability measures $\{P \in \mathcal{P}\}$ on a sample space \mathcal{X} . Such models can be expressed in the form $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, where Θ is some parameter space. *Semiparametric models* are statistical models where Θ has one or more infinite-dimensional component. For example, the parameter space for the linear regression model (1.1), where

$Y = \beta'Z + e$, consists of two components, a subset of p -dimensional Euclidean space (for the regression parameter β) and the infinite-dimensional space of all joint distribution functions of (e, Z) with $E[e|Z] = 0$ and $E[e^2|Z] \leq K < \infty$ almost surely.

The goal of *semiparametric inference* is to construct optimal estimators and test statistics for evaluating semiparametric model parameters. The parameter component of interest can be succinctly expressed as a function of the form $\psi : \mathcal{P} \mapsto \mathbb{D}$, where ψ extracts the component of interest and takes values in \mathbb{D} . For now, we assume \mathbb{D} is finite dimensional. As an illustration, if we are interested in the unconditional variance of the residual errors in the regression model (1.1), ψ would be the map $\psi(P) = \int_{\mathbb{R}} e^2 dF(e)$, where $F(t) = P[e \leq t]$ is the unconditional residual distribution component of P . Throughout this book, the statistics of interest will be based on an i.i.d. sample, X_1, \dots, X_n , of realizations from some $P \in \mathcal{P}$.

An estimator T_n of the parameter $\psi(P)$, based on such a sample, is *efficient* if the limiting variance V of $\sqrt{n}(T_n - \psi(P))$ is the smallest possible among all *regular* estimators of $\psi(P)$. The inverse of V is the *information* for T_n . Regularity will be defined more explicitly later in this section, but suffice it to say for now that the limiting distribution of $\sqrt{n}(T_n - \psi(P))$ (as $n \rightarrow \infty$), for a regular estimator T_n , is continuous in P . Note that as we are changing P , the distribution of T_n changes as well as the parameter $\psi(P)$. Not all estimators are regular, but most commonly used estimators in statistics are. Optimality of test statistics is closely related to efficient estimation, in that the most powerful test statistics for a hypothesis about a parameter are usually based on efficient estimators for that parameter.

The optimal efficiency for estimators of a parameter $\psi(P)$ depends in part on the complexity of the model \mathcal{P} . Estimation under the model \mathcal{P} is more taxing than estimation under any parametric submodel $\mathcal{P}_0 = \{P_\theta : \theta \in \Theta_0\} \subset \mathcal{P}$, where Θ_0 is finite dimensional. Thus the information for estimation under model \mathcal{P} is worse than the information under any parametric submodel \mathcal{P}_0 . If the information for the regular estimator T_n is equal to the minimum of the information over all efficient estimators for all parametric submodels \mathcal{P}_0 , then T_n is *semiparametric efficient*. For semiparametric models, this minimizer is the best possible, since the only models with more information are parametric models. A parametric model which achieves this minimum, if such a model exists, is called a *least favorable* or *hardest* submodel. Note that efficient estimators for parametric models are trivially semiparametric efficient since such models are their own parametric submodels.

Fortunately, finding the minimum information over parametric submodels usually only requires consideration of one-dimensional parametric submodels $\{P_t : t \in N_\epsilon\}$ surrounding representative distributions $P \in \mathcal{P}$, where $N_\epsilon = [0, \epsilon)$ for some $\epsilon > 0$, $P_0 = P$, and $P_t \in \mathcal{P}$ for all $t \in N_\epsilon$. If \mathcal{P} has a dominating measure μ , then each $P \in \mathcal{P}$ can be expressed as a density p . In this case, we require the submodels around a representative density p

to be smooth enough so that the real function $g(x) = \partial \log p_t(x)/(\partial t)|_{t=0}$ exists with $\int_{\mathcal{X}} g^2(x)p(x)\mu(dx) < \infty$. This idea can be restated in sufficient generality to allow for models which may not be dominated. In this more general case, we require

$$(3.1) \quad \int \left[\frac{(dP_t(x))^{1/2} - (dP(x))^{1/2}}{t} - \frac{1}{2}g(x)(dP(x))^{1/2} \right]^2 \rightarrow 0,$$

as $t \downarrow 0$. In this setting, we say that the submodel $\{P_t : t \in N_\epsilon\}$ is *differentiable in quadratic mean* at $t = 0$, with *score function* $g : \mathcal{X} \mapsto \mathbb{R}$.

In evaluating efficiency, it is necessary to consider many such one-dimensional submodels surrounding the representative P , each with a different score function. Such a collection of score functions is called a *tangent set* of the model \mathcal{P} at P , and is denoted $\dot{\mathcal{P}}_P$. Because $Pg = 0$ and $Pg^2 < \infty$ for any $g \in \dot{\mathcal{P}}_P$, such tangent sets are subsets of $L_2^0(P)$, the space of all functions $h : \mathcal{X} \mapsto \mathbb{R}$ with $Ph = 0$ and $Ph^2 < \infty$. Note that, as we have done here, we will sometimes omit function arguments for simplicity, provided the context is clear. When the tangent set is closed under linear combinations, it is called a *tangent space*. Usually, one can take the closed linear span (the closure under linear combinations) of a tangent set to make a tangent space.

Consider $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, where Θ is an open subset of \mathbb{R}^k . Assume that \mathcal{P} is dominated by μ , and that the classical score function $\dot{\ell}_\theta(x) = \partial \log p_\theta(x)/(\partial \theta)$ exists with $P_\theta[\dot{\ell}_\theta \dot{\ell}'_\theta]$ bounded. Now for each $h \in \mathbb{R}^k$, let $\epsilon > 0$ be small enough so that $\{P_t : t \in N_\epsilon\} \subset \mathcal{P}$, where for each $t \in N_\epsilon$, $P_t = P_{\theta+th}$. One can show that each of these one-dimensional submodels satisfy (3.1), with $g = h' \dot{\ell}_\theta$, resulting in the tangent space $\dot{\mathcal{P}}_{P_\theta} = \{h' \dot{\ell}_\theta : h \in \mathbb{R}^k\}$. Thus there is a simple connection between the classical score function and the more general idea of tangent sets and tangent spaces.

Continuing with the parametric setting, the estimator $\hat{\theta}$ is efficient for estimating θ if it is regular with information achieving the Cramér-Rao lower bound $P[\dot{\ell}_\theta \dot{\ell}'_\theta]$. Thus the tangent set for the model contains information about the optimal efficiency. This is also true for semiparametric models in general, although the relationship between tangent sets and the optimal information is more complex.

Consider estimation of the parameter $\psi(P) \in \mathbb{R}^k$ for the semiparametric model \mathcal{P} . For any estimator T_n of $\psi(P)$, if $\sqrt{n}(T_n - \psi(P)) = \sqrt{n}\mathbb{P}_n \dot{\psi}_P + o_p(1)$, where $o_p(1)$ denotes a quantity going to zero in probability, then $\dot{\psi}_P$ is an *influence function* for $\psi(P)$ and T_n is *asymptotically linear*. For a given tangent set $\dot{\mathcal{P}}_P$, assume for each submodel $\{P_t : t \in N_\epsilon\}$ satisfying (3.1) with some $g \in \dot{\mathcal{P}}_P$ and some $\epsilon > 0$, that $d\psi(P_t)/(dt)|_{t=0} = \psi_P(g)$ for some linear map $\psi_P : L_2^0(P) \mapsto \mathbb{R}^k$. In this setting, we say that ψ is differentiable at P relative to $\dot{\mathcal{P}}_P$. When $\dot{\mathcal{P}}_P$ is a linear space, there exists a measurable function $\tilde{\psi}_P : \mathcal{X} \mapsto \mathbb{R}^k$ such that $\dot{\psi}_P(g) = P[\tilde{\psi}_P(X)g(X)]$, for each $g \in \dot{\mathcal{P}}_P$. The function $\tilde{\psi}_P \in \dot{\mathcal{P}}_P \subset L_2^0(P)$ is unique and is called the

efficient influence function for the parameter ψ in the model P (relative to the tangent space $\dot{\mathcal{P}}_P$). Here, and throughout the book, we abuse notation slightly by declaring that a random vector is in a given linear space if and only if each component of the vector is. Note that $\tilde{\psi}_P \in \dot{\mathcal{P}}_P$. Frequently, the efficient influence function can be found by taking a candidate influence function $\check{\psi}_P \in L_2^0(P)$ and projecting it onto $\dot{\mathcal{P}}_P$ to obtain $\tilde{\psi}_P$. If $\sqrt{n}(T_n - \psi(P))$ is asymptotically equivalent to $\sqrt{n}\mathbb{P}_n\tilde{\psi}_P$, then T_n can be shown to be semiparametric efficient (which we refer to hereafter simply as “efficient”).

Consider again the parametric example, with $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, where $\Theta \subset \mathbb{R}^k$. Suppose the parameter of interest is $\psi(P_\theta) = f(\theta)$, where $f : \mathbb{R}^k \mapsto \mathbb{R}^d$ has derivative \dot{f}_θ at θ , and that the Fisher information matrix $I_\theta = P[\dot{\ell}_\theta \dot{\ell}'_\theta]$ is invertible. It is not hard to show that $\dot{\psi}_P(g) = \dot{f}_\theta I_\theta^{-1} P[\dot{\ell}_\theta(X)g(X)]$, and thus $\tilde{\psi}_P(x) = \dot{f}_\theta I_\theta^{-1} \dot{\ell}_\theta(x)$. Any estimator T_n for which $\sqrt{n}(T_n - f(\theta))$ is asymptotically equivalent to $\sqrt{n}\mathbb{P}_n[\dot{f}_\theta I_\theta^{-1} \dot{\ell}_\theta]$, has asymptotic variance equal to the Cramér-Rao lower bound $\dot{f}_\theta I_\theta^{-1} \dot{f}'_\theta$.

Returning to the semiparametric setting, any one-dimensional submodel $\{P_t : t \in N_\epsilon\}$, satisfying (3.1) for the score function $g \in \dot{\mathcal{P}}_P$ and some $\epsilon > 0$, is a parametric model with parameter t . The Fisher information for t , evaluated at $t = 0$, is Pg^2 . Thus the Cramér-Rao lower bound for estimating a univariate parameter $\psi(P)$ based on this model is $(P[\tilde{\psi}_P g])^2 / Pg^2$, since $d\psi(P_t)/(dt)|_{t=0} = P[\tilde{\psi}_P g]$. Provided that $\tilde{\psi}_P \in \dot{\mathcal{P}}_P$ and all necessary derivatives exist, the maximum Cramér-Rao lower bound over all such submodels in $\dot{\mathcal{P}}_P$ is thus $P\tilde{\psi}_P^2$. For more general Euclidean parameters $\psi(P)$, this lower bound on the asymptotic variance is $P[\tilde{\psi}_P \tilde{\psi}'_P]$.

Hence, for $P[\tilde{\psi}_P \tilde{\psi}'_P]$ to be the upper bound for all parametric submodels, the tangent set must be sufficiently large. Obviously, the tangent set must also be restricted to score functions which reflect valid submodels. In addition, the larger the tangent set, the fewer the number of regular estimators. To see this, we will now provide a more precise definition of regularity. Let $P_{t,g}$ denote a submodel $\{P_t : t \in N_\epsilon\}$ satisfying (3.1) for the score g and some $\epsilon > 0$. T_n is regular for $\psi(P)$ if the limiting distribution of $\sqrt{n}(T_n - \psi(P_{1/\sqrt{n},g}))$, over the sequence of distributions $P_{1/\sqrt{n},g}$, exists and is constant over all $g \in \dot{\mathcal{P}}_P$. Thus the tangent set must be chosen to be large enough but not too large. Fortunately, most estimators in common use for semiparametric inference are regular for large tangent sets, and thus there is usually quite a lot to be gained by making the effort to obtain an efficient estimator. Once the tangent set has been identified, the corresponding efficient estimator T_n of $\psi(P)$ will always be asymptotically linear with influence function equal to the efficient influence function $\tilde{\psi}_P$.

Consider, for example, the unrestricted model \mathcal{P} of all distributions on \mathcal{X} . Suppose we are interested in estimating $\psi(P) = Pf$ for some $f \in L_2(P)$,

the space of all measurable functions h with $Ph^2 < \infty$. For bounded $g \in L_2^0(P)$, the one-dimensional submodel $\{P_t : dP_t = (1 + tg)dP, t \in N_\epsilon\} \subset \mathcal{P}$ for ϵ small enough. Furthermore, (3.1) is satisfied with $\partial\psi(P_t)/(\partial t)|_{t=0} = P[fg]$. It is not hard to show, in fact, that one-dimensional submodels satisfying (3.1) with $\partial\psi(P_t)/(\partial t)|_{t=0} = P[fg]$ exist for all $g \in L_2^0(P)$. Thus $\dot{\mathcal{P}}_P = L_2^0(P)$ is a tangent set for the unrestricted model and $\psi_P(x) = f(x) - Pf$ is the corresponding efficient influence function. Since $\sqrt{n}(\mathbb{P}_n f - \psi(P)) = \sqrt{n}\mathbb{P}_n \tilde{\psi}_P$, $\mathbb{P}_n f$ is efficient for estimating Pf . Note that, in this unrestricted model, $\dot{\mathcal{P}}_P$ is the maximal possible tangent set, since all tangent sets must be subsets of $L_2^0(P)$. In general, the size of a tangent set reflects the amount of restrictions placed on a model, in that larger tangent sets reflect fewer restrictions.

Efficiency can also be established for infinite dimensional parameters $\psi(P)$ when \sqrt{n} consistent regular estimators for $\psi(P)$ exist. There are a number of ways of expressing efficiency in this context, but we will only mention the convolution approach here. The convolution theorem states that for any regular estimator T_n of $\psi(P)$, $\sqrt{n}(T_n - \psi(P))$ has a weak limiting distribution which is the convolution of a Gaussian process Z and an independent process M , where Z has the same limiting distribution as $\sqrt{n}\mathbb{P}_n \tilde{\psi}_P$. In other words, an inefficient estimator always has an asymptotically non-negligible independent noise process M added to the efficient estimator distribution. A regular estimator T_n for which M is zero is efficient. Occasionally, we will use the term *uniformly efficient* when it is helpful to emphasize the fact that $\psi(P)$ is infinite dimensional. If $\psi(P)$ is indexed by $\{t \in T\}$, and if $T_n(t)$ is efficient for $\psi(P)(t)$ for each $t \in T$, then it can be shown that weak convergence of $\sqrt{n}(T_n - \psi(P))$ to a tight, mean zero Gaussian process implies uniform efficiency of T_n . Another important fact is that if T_n is an efficient estimator for $\psi(P)$ and ϕ is a suitable Hadamard-differentiable function, then $\phi(T_n)$ is an efficient estimator for $\phi(\psi(P))$. We will make these results more explicit in part III.

3.2 Score Functions and Estimating Equations

A parameter $\psi(P)$ of particular interest is the parametric component θ of a semiparametric model $\{P_{\theta,\eta} : \theta \in \Theta, \eta \in H\}$, where Θ is an open subset of \mathbb{R}^k and H is an arbitrary subset that may be infinite dimensional. Tangent sets can be used to develop an efficient estimator for $\psi(P_{\theta,\eta}) = \theta$ through the formation of an *efficient score function*. In this setting, we consider submodels of the form $\{P_{\theta+ta,\eta_t}, t \in N_\epsilon\}$ which are differentiable in quadratic mean with score function $\partial \log dP_{\theta+ta,\eta_t}/(\partial t)|_{t=0} = a' \dot{\ell}_{\theta,\eta} + g$, where $a \in \mathbb{R}^k$, $\dot{\ell}_{\theta,\eta} : \mathcal{X} \mapsto \mathbb{R}^k$ is the ordinary score for θ when η is fixed, and where $g : \mathcal{X} \mapsto \mathbb{R}$ is an element of a tangent set $\dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$ for the submodel $\mathcal{P}_\theta = \{P_{\theta,\eta} : \eta \in H\}$ (holding θ fixed). This tangent set is the *tangent set for*

η and should be rich enough to reflect all parametric submodels of \mathcal{P}_θ . The tangent set for the full model is $\dot{\mathcal{P}}_{P_{\theta,\eta}} = \left\{ a' \dot{\ell}_{\theta,\eta} + g : a \in \mathbb{R}^k, g \in \dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)} \right\}$.

While $\psi(P_{\theta+ta,\eta_t}) = \theta + ta$ is clearly differentiable with respect to t , we also need, as in the previous section, that there exists a function $\tilde{\psi}_{\theta,\eta} : \mathcal{X} \mapsto \mathbb{R}^k$ such that

$$(3.2) \quad \left. \frac{\partial \psi(P_{\theta+ta,\eta_t})}{\partial t} \right|_{t=0} = a = P \left[\tilde{\psi}_{\theta,\eta} \left(\dot{\ell}'_{\theta,\eta} a + g \right) \right],$$

for all $a \in \mathbb{R}^k$ and all $g \in \dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$. After setting $a = 0$, we see that such a function must be uncorrelated with all of the elements of $\dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$.

Define $\Pi_{\theta,\eta}$ to be the orthogonal projection onto the closed linear span of $\dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$ in $L_2^0(P_{\theta,\eta})$. We will describe how to obtain such projections in detail in part III, but suffice it to say that for any $h \in L_2^0(P_{\theta,\eta})$, $h = h - \Pi_{\theta,\eta}h + \Pi_{\theta,\eta}h$, where $\Pi_{\theta,\eta}h \in \dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$ but $P[(h - \Pi_{\theta,\eta}h)g] = 0$ for all $g \in \dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$. The *efficient score function* for θ is $\tilde{\ell}_{\theta,\eta} = \dot{\ell}_{\theta,\eta} - \Pi_{\theta,\eta}\dot{\ell}_{\theta,\eta}$, while the *efficient information matrix* for θ is $\tilde{I}_{\theta,\eta} = P \left[\tilde{\ell}_{\theta,\eta} \tilde{\ell}'_{\theta,\eta} \right]$.

Provided that $\tilde{I}_{\theta,\eta}$ is nonsingular, the function $\tilde{\psi}_{\theta,\eta} = \tilde{I}_{\theta,\eta}^{-1} \tilde{\ell}_{\theta,\eta}$ satisfies (3.2) for all $a \in \mathbb{R}^k$ and all $g \in \dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$. Thus the functional (parameter) $\psi(P_{\theta,\eta}) = \theta$ is differentiable at $P_{\theta,\eta}$ relative to the tangent set $\dot{\mathcal{P}}_{P_{\theta,\eta}}$, with efficient influence function $\tilde{\psi}_{\theta,\eta}$. Hence the search for an efficient estimator of θ is over if one can find an estimator T_n satisfying $\sqrt{n}(T_n - \theta) = \sqrt{n}\mathbb{P}_n \tilde{\psi}_{\theta,\eta} + o_P(1)$. Note that $\tilde{I}_{\theta,\eta} = I_{\theta,\eta} - P \left[\Pi_{\theta,\eta} \dot{\ell}_{\theta,\eta} \left(\Pi_{\theta,\eta} \dot{\ell}_{\theta,\eta} \right)' \right]$, where $I_{\theta,\eta} = P \left[\dot{\ell}_{\theta,\eta} \dot{\ell}'_{\theta,\eta} \right]$. An intuitive justification for the form of the efficient score is that some information for estimating θ is lost due to a lack of knowledge about η . The amount subtracted off of the efficient score, $\Pi_{\theta,\eta} \dot{\ell}_{\theta,\eta}$, is the minimum possible amount for regular estimators when η is unknown.

Consider again the semiparametric regression model (1.1), where $Y = \beta'Z + e$, $\mathbb{E}[e|Z] = 0$ and $\mathbb{E}[e^2|Z] \leq K < \infty$ almost surely, and where we observe (Y, Z) , with the joint density η of (e, Z) satisfying $\int_{\mathbb{R}} e\eta(e, Z)de = 0$ almost surely. Assume η has partial derivative with respect to the first argument, $\dot{\eta}_1$, satisfying $\dot{\eta}_1/\eta \in L_2(P_{\beta,\eta})$, and hence $\dot{\eta}_1/\eta \in L_2^0(P_{\beta,\eta})$, where $P_{\beta,\eta}$ is the joint distribution of (Y, Z) . The Euclidean parameter of interest in this semiparametric model is $\theta = \beta$. The score for β , assuming η is known, is $\dot{\ell}_{\beta,\eta} = -Z(\dot{\eta}_1/\eta)(Y - \beta'Z, Z)$, where we use the shorthand $(f/g)(u, v) = f(u, v)/g(u, v)$ for ratios of functions.

One can show that the tangent set $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$ for η is the subset of $L_2^0(P_{\beta,\eta})$ which consists of all functions $g(e, Z) \in L_2^0(P_{\beta,\eta})$ which satisfy

$$\mathbb{E}[eg(e, Z)|Z] = \frac{\int_{\mathbb{R}} eg(e, Z)\eta(e, Z)de}{\int_{\mathbb{R}} \eta(e, Z)de} = 0,$$

almost surely. One can also show that this set is the orthocomplement in $L_2^0(P_{\beta,\eta})$ of all functions of the form $ef(Z)$, where f satisfies $P_{\beta,\eta}f^2(Z) < \infty$. This means that $\tilde{\ell}_{\beta,\eta} = (I - \Pi_{\beta,\eta})\dot{\ell}_{\beta,\eta}$ is the projection in $L_2^0(P_{\beta,\eta})$ of $-Z(\dot{\eta}_1/\eta)(e, Z)$ onto $\{ef(Z) : P_{\beta,\eta}f^2(Z) < \infty\}$, where I is the identity. Thus

$$\tilde{\ell}_{\beta,\eta}(Y, Z) = \frac{-Ze \int_{\mathbb{R}} \dot{\eta}_1(e, Z)ede}{P_{\beta,\eta}[e^2|Z]} = -\frac{Ze(-1)}{P_{\beta,\eta}[e^2|Z]} = \frac{Z(Y - \beta'Z)}{P_{\beta,\eta}[e^2|Z]},$$

where the second-to-last step follows from the identity $\int_{\mathbb{R}} \dot{\eta}_1(e, Z)ede = \partial \int_{\mathbb{R}} \eta(te, Z)de / (\partial t)|_{t=1}$, and the last step follows since $e = Y - \beta'Z$. When the function $z \mapsto P_{\beta,\eta}[e^2|Z = z]$ is non-constant in z , $\tilde{\ell}_{\beta,\eta}(Y, Z)$ is not proportional to $Z(Y - \beta'Z)$, and the estimator $\hat{\beta}$ defined in chapter 1 will not be efficient. We will discuss efficient estimation for this model in greater detail in chapter 4.

Two very useful tools for computing efficient scores are score and information operators. Although we will provide more precise definitions in parts II and III, operators are maps between spaces of functions. Returning to the generic semiparametric model $\{P_{\theta,\eta} : \theta \in \Theta, \eta \in H\}$, sometimes it is easier to represent an element g in the tangent set for η , $\dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)}$, as $B_{\theta,\eta}b$, where b is an element of another set \mathbb{H}_η and $B_{\theta,\eta}$ is an operator satisfying $\dot{\mathcal{P}}_{P_{\theta,\eta}}^{(\eta)} = \{B_{\theta,\eta}b : b \in \mathbb{H}_\eta\}$. Such an operator is a score operator. The adjoint of the score operator $B_{\theta,\eta} : \mathbb{H}_\eta \mapsto L_2^0(P_{\theta,\eta})$ is another operator $B_{\theta,\eta}^* : L_2^0(P_{\theta,\eta}) \mapsto \overline{\text{lin}} \mathbb{H}_\eta$ which is similar in spirit to a transpose for matrices. Here we use $\overline{\text{lin}} \mathbb{A}$ to denote *closed linear span* (the linear space consisting of all linear combinations) of \mathbb{A} . Additional details on adjoints and methods for computing them will be described in part III. One can now define the information operator $B_{\theta,\eta}^* B_{\theta,\eta} : H_\eta \mapsto \overline{\text{lin}} H_\eta$. If $B_{\theta,\eta}^* B_{\theta,\eta}$ has an inverse, then it can be shown that the efficient score for θ has the form $\tilde{\ell}_{\theta,\eta} = \left(I - B_{\theta,\eta} \left[B_{\theta,\eta}^* B_{\theta,\eta} \right]^{-1} B_{\theta,\eta}^* \right) \dot{\ell}_{\theta,\eta}$.

To illustrate these methods, consider the Cox model for right-censored data introduced in chapter 1. In this setting, we observe a sample of n realizations of $X = (V, d, Z)$, where $V = T \wedge C$, $d = 1\{V = T\}$, $Z \in \mathbb{R}^k$ is a covariate vector, T is a failure time, and C is a censoring time. We assume that T and C are independent given Z , that T given Z has integrated hazard function $e^{\beta'Z} \Lambda(t)$ for β in an open subset $B \subset \mathbb{R}^k$ and Λ is continuous and monotone increasing with $\Lambda(0) = 0$, and that the censoring distribution does not depend on β or Λ (ie., censoring is uninformative). Define the counting and at-risk processes $N(t) = 1\{V \leq t\}d$ and $Y(t) = 1\{V \geq t\}$, and let $M(t) = N(t) - \int_0^t Y(s)e^{\beta'Z} d\Lambda(s)$. For some $0 < \tau < \infty$ with $P\{C \geq \tau\} > 0$, let H be the set of all Λ 's satisfying our criteria with $\Lambda(\tau) < \infty$. Now the set of models \mathcal{P} is indexed by $\beta \in B$ and $\Lambda \in H$. We let $P_{\beta,\Lambda}$ be the distribution of (V, d, Z) corresponding to the given parameters.

The likelihood for a single observation is thus proportional to $p_{\beta,\Lambda}(X) = \left[e^{\beta'Z} \lambda(V) \right]^d \exp \left[-e^{\beta'Z} \Lambda(V) \right]$, where λ is the derivative of Λ . Now let $L_2(\Lambda)$ be the set of measurable functions $b : [0, \tau] \mapsto \mathbb{R}$ with $\int_0^\tau b^2(s) d\Lambda(s) < \infty$. If $b \in L_2(\Lambda)$ is bounded, then $\Lambda_t(s) = \int_0^s e^{tb(u)} d\Lambda(u) \in H$ for all t . The score function $\partial \log p_{\beta+ta, \Lambda_t}(X) / (\partial t)|_{t=0}$ is thus $\int_0^\tau [a'Z + b(s)] dM(s)$, for any $a \in \mathbb{R}^k$. The score function for β is therefore $\dot{\ell}_{\beta,\Lambda}(X) = ZM(\tau)$, while the score function for Λ is $\int_0^\tau b(s) dM(s)$. In fact, one can show that there exists one-dimensional submodels Λ_t such that $\log p_{\beta+ta, \Lambda_t}$ is differentiable with score $a' \dot{\ell}_{\beta,\Lambda}(X) + \int_0^\tau b(s) dM(s)$, for any $b \in L_2(\Lambda)$ and $a \in \mathbb{R}^k$.

The operator $B_{\beta,\Lambda} : L_2(\Lambda) \mapsto L_2^0(P_{\beta,\Lambda})$, given by $B_{\beta,\Lambda}(b) = \int_0^\tau b(s) dM(s)$, is the score operator which generates the tangent set for Λ , $\dot{\mathcal{P}}_{P_{\beta,\Lambda}}^{(\Lambda)} \equiv \{B_{\beta,\Lambda}b : b \in L_2(\Lambda)\}$. It can be shown that this tangent space spans all square-integrable score functions for Λ generated by parametric submodels. The adjoint operator can be shown to be $B_{\beta,\Lambda}^* : L_2(P_{\beta,\Lambda}) \mapsto L_2(\Lambda)$, where $B_{\beta,\Lambda}^*(g)(t) = P_{\beta,\Lambda}[g(X)dM(t)]/d\Lambda(t)$. The information operator $B_{\beta,\Lambda}^* B_{\beta,\Lambda} : L_2(\Lambda) \mapsto L_2(\Lambda)$ is thus

$$B_{\beta,\Lambda}^* B_{\beta,\Lambda}(b)(t) = \frac{P_{\beta,\Lambda} \left[\int_0^\tau b(s) dM(s) dM(u) \right]}{d\Lambda(u)} = P_{\beta,\Lambda} \left[Y(t) e^{\beta'Z} \right] b(t),$$

using martingale methods.

Since $B_{\beta,\Lambda}^* \left(\dot{\ell}_{\beta,\Lambda} \right) (t) = P_{\beta,\Lambda} \left[ZY(t) e^{\beta'Z} \right]$, we have that the efficient score for β is

$$\begin{aligned} (3.3) \quad \tilde{\ell}_{\beta,\Lambda} &= \left(I - B_{\beta,\Lambda} \left[B_{\beta,\Lambda}^* B_{\beta,\Lambda} \right]^{-1} B_{\beta,\Lambda}^* \right) \dot{\ell}_{\beta,\Lambda} \\ &= \int_0^\tau \left\{ Z - \frac{P_{\beta,\Lambda} \left[ZY(t) e^{\beta'Z} \right]}{P_{\beta,\Lambda} \left[Y(t) e^{\beta'Z} \right]} \right\} dM(t). \end{aligned}$$

When $\tilde{I}_{\beta,\Lambda} \equiv P_{\beta,\Lambda} \left[\tilde{\ell}_{\beta,\Lambda} \tilde{\ell}_{\beta,\Lambda}' \right]$ is positive definite, the resulting efficient influence function is $\tilde{\psi}_{\beta,\Lambda} \equiv \tilde{I}_{\beta,\Lambda}^{-1} \tilde{\ell}_{\beta,\Lambda}$. Since the estimator $\hat{\beta}_n$ obtained from maximizing the *partial likelihood*

$$(3.4) \quad \tilde{L}_n(\beta) = \prod_{i=1}^n \left(\frac{e^{\beta'Z_i}}{\sum_{j=1}^n \mathbf{1}\{V_j \geq V_i\} e^{\beta'Z_j}} \right)^{d_i}$$

can be shown to satisfy $\sqrt{n}(\hat{\beta}_n - \beta) = \sqrt{n} \mathbb{P}_n \tilde{\psi}_{\beta,\Lambda} + o_p(1)$, this estimator is efficient.

Returning to our discussion of score and information operators, these operators are also useful for generating scores for the entire model, not just for the nuisance component. With semiparametric models having score functions of the form $a' \dot{\ell}_{\theta,\eta} + B_{\theta,\eta} b$, for $a \in \mathbb{R}^k$ and $b \in \mathbb{H}_\eta$, we can define

a new operator $A_{\beta,\eta} : \{(a, b) : a \in \mathbb{R}^k, b \in \text{lin } \mathbb{H}_\eta\} \mapsto L_2^0(P_{\theta,\eta})$ where $A_{\beta,\eta}(a, b) = a'\dot{\ell}_{\theta,\eta} + B_{\theta,\eta}b$. More generally, we can define the score operator $A_\eta : \text{lin } \mathbb{H}_\eta \mapsto L_2(P_\eta)$ for the model $\{P_\eta : \eta \in H\}$, where H indexes the entire model and may include both parametric and nonparametric components, and where $\text{lin } \mathbb{H}_\eta$ indexes directions in H . Let the parameter of interest be $\psi(P_\eta) = \chi(\eta) \in \mathbb{R}^k$. We assume there exists a linear operator $\dot{\chi} : \text{lin } \mathbb{H}_\eta \mapsto \mathbb{R}^k$ such that, for every $b \in \text{lin } \mathbb{H}_\eta$, there exists a one-dimensional submodel $\{P_{\eta_t} : \eta_t \in H, t \in N_\epsilon\}$ satisfying

$$\int \left[\frac{(dP_{\eta_t})^{1/2} - (dP_\eta)^{1/2}}{t} - \frac{1}{2}A_\eta b(dP_\eta)^{1/2} \right]^2 \rightarrow 0,$$

as $t \downarrow 0$, and $\partial\chi(\eta_t)/(\partial t)|_{t=0} = \dot{\chi}(b)$.

We require \mathbb{H}_η to have a suitable *inner product* $\langle \cdot, \cdot \rangle_\eta$, where an inner product is an operation on elements in \mathbb{H}_η with the property that $\langle a, b \rangle_\eta = \langle b, a \rangle_\eta$, $\langle a + b, c \rangle_\eta = \langle a, c \rangle_\eta + \langle b, c \rangle_\eta$, $\langle a, a \rangle_\eta \geq 0$, and $\langle a, a \rangle_\eta = 0$ if and only if $a = 0$, for all $a, b, c \in \mathbb{H}_\eta$. The efficient influence function is the solution $\tilde{\psi}_{P_\eta} \in \overline{R}(A_\eta) \subset L_2^0(P_\eta)$ of

$$(3.5) \quad A_\eta^* \tilde{\psi}_{P_\eta} = \tilde{\chi}_\eta,$$

where R denotes range, \overline{B} denotes closure of the set B , A_η^* is the adjoint of A_η , and $\tilde{\chi}_\eta \in \mathbb{H}_\eta$ satisfies $\langle \tilde{\chi}_\eta, b \rangle_\eta = \dot{\chi}_\eta(b)$ for all $b \in \mathbb{H}_\eta$. Methods for obtaining such a $\tilde{\chi}_\eta$ will be given in part III. When $A_\eta^*A_\eta$ is invertible, then the solution to (3.5) can be written $\tilde{\psi}_{P_\eta} = A_\eta (A_\eta^*A_\eta)^{-1} \tilde{\chi}_\eta$. In chapter 4, We will illustrate this approach to derive efficient estimators for all parameters of the Cox model.

Returning to the semiparametric model setting, where $\mathcal{P} = \{P_{\theta,\eta} : \theta \in \Theta, \eta \in H\}$, Θ is an open subset of \mathbb{R}^k , and H is a set, the efficient score can be used to derive *estimating equations* for computing efficient estimators of θ . An estimating equation is a data dependent function $\Psi_n : \Theta \mapsto \mathbb{R}^k$ for which an approximate zero yields a Z-estimator for θ . When $\Psi_n(\tilde{\theta})$ has the form $\mathbb{P}_n \hat{\ell}_{\tilde{\theta},n}$, where $\hat{\ell}_{\tilde{\theta},n}(X|X_1, \dots, X_n)$ is a function for the generic observation X which depends on the sample data X_1, \dots, X_n , we have the following estimating equation result (the proof of which will be given in part III):

THEOREM 3.1 *Suppose that the model $\{P_{\theta,\eta} : \theta \in \Theta\}$, where $\Theta \subset \mathbb{R}^k$, is differentiable in quadratic mean with respect to θ at (θ, η) and let the efficient information matrix $\tilde{I}_{\theta,\eta}$ be nonsingular. Let $\hat{\theta}_n$ satisfy $\sqrt{n}\mathbb{P}_n \hat{\ell}_{\hat{\theta}_n,n} = o_p(1)$ and be consistent for θ . Also assume that the following conditions*

hold:

$$(3.6) \quad \sqrt{n}(\mathbb{P}_n - P_{\theta,\eta}) \left(\hat{\ell}_{\hat{\theta}_n,n} - \tilde{\ell}_{\theta,\eta} \right) \xrightarrow{P} 0,$$

$$(3.7) \quad P_{\hat{\theta}_n,\eta} \hat{\ell}_{\hat{\theta}_n,n} = o_p(n^{-1/2} + \|\hat{\theta}_n - \theta\|),$$

$$(3.8) \quad P_{\theta,\eta} \left\| \hat{\ell}_{\hat{\theta}_n,n} - \tilde{\ell}_{\theta,\eta} \right\|^2 \xrightarrow{P} 0, \quad P_{\hat{\theta}_n,\eta} \left\| \hat{\ell}_{\hat{\theta}_n,n} \right\|^2 = O_p(1).$$

Then $\hat{\theta}_n$ is asymptotically efficient at (θ, η) .

Returning to the Cox model example, the profile likelihood score is the partial likelihood score $\Psi_n(\tilde{\beta}) = \mathbb{P}_n \hat{\ell}_{\tilde{\beta},n}$, where

$$(3.9) \quad \hat{\ell}_{\tilde{\beta},n}(X = (V, d, Z) | X_1, \dots, X_n) = \int_0^\tau \left\{ Z - \frac{\mathbb{P}_n [ZY(t)e^{\tilde{\beta}'Z}]}{\mathbb{P}_n [Y(t)e^{\tilde{\beta}'Z}]} \right\} dM_{\tilde{\beta}}(t),$$

and $M_{\tilde{\beta}}(t) = N(t) - \int_0^t Y(u)e^{\tilde{\beta}'Z} d\Lambda(u)$. We will show in chapter 4 that all the conditions of theorem 3.1 are satisfied for the root of $\Psi_n(\tilde{\beta}) = 0$, $\hat{\beta}_n$, and thus the partial likelihood yields efficient estimation of β .

Returning to the general semiparametric setting, even if an estimating equation Ψ_n is not close enough to $\mathbb{P}_n \tilde{\ell}_{\theta,\eta}$ to result in an efficient estimator, frequently the estimator will still result in a \sqrt{n} -consistent estimator which is precise enough to be useful. In some cases, the computational effort needed to obtain an efficient estimator may be too great a cost, and one must settle for an inefficient estimating equation that works. Even in these settings, some modifications in the estimating equation can often be made which improve efficiency while maintaining computability. This issue will be explored in greater detail in part III.

3.3 Maximum Likelihood Estimation

The most common approach to efficient estimation is based on modifications of maximum likelihood estimation which lead to efficient estimates. These modifications, which we will call “likelihoods,” are generally not really likelihoods (products of densities) because of complications resulting from the presence of an infinite dimensional nuisance parameter. Consider estimating an unknown real density $f(x)$ from an i.i.d. sample X_1, \dots, X_n . The likelihood is $\prod_{i=1}^n f(X_i)$, and the maximizer over all densities has arbitrarily high peaks at the observations, with zero at the other values, and is therefore not a density. This problem can be fixed by using an empirical likelihood $\prod_{i=1}^n p_i$, where p_1, \dots, p_n are the masses assigned to the observations indexed by $i = 1, \dots, n$ and are constrained to satisfy $\sum_{i=1}^n p_i = 1$.

This leads to the empirical distribution function estimator, which is known to be fully efficient.

Consider again the Cox model for right-censored data explored in the previous section. The density for a single observation $X = (V, d, Z)$ is proportional to $\left[e^{\beta'Z}\lambda(V)\right]^d \exp\left[-e^{\beta'Z}\Lambda(V)\right]$. Maximizing the likelihood based on this density will result in the same problem raised in the previous paragraph. A likelihood that works is the following, which assigns mass only at observed failure times:

$$(3.10) \quad L_n(\beta, \Lambda) = \prod_{i=1}^n \left[e^{\beta'Z_i} \Delta\Lambda(V_i) \right]^{d_i} \exp \left[-e^{\beta'Z_i} \Lambda(V_i) \right],$$

where $\Delta\Lambda(t)$ is the jump size of Λ at t . For each value of β , one can maximize or *profile* $L_n(\beta, \Lambda)$ over the “nuisance” parameter Λ to obtain the profile likelihood $pL_n(\beta)$, which for the Cox model is $\exp\left[-\sum_{i=1}^n d_i\right]$ times the partial likelihood (3.4). Let $\hat{\beta}$ be the maximizer of $pL_n(\beta)$. Then the maximizer $\hat{\Lambda}$ of $L_n(\hat{\beta}, \Lambda)$ is the “Breslow estimator”

$$\hat{\Lambda}(t) = \int_0^t \frac{\mathbb{P}_n dN(s)}{\mathbb{P}_n \left[Y(s) e^{\hat{\beta}'Z} \right]}.$$

We will see in chapter 4 that $\hat{\beta}$ and $\hat{\Lambda}$ are both efficient.

Another useful class of likelihood variants are *penalized likelihoods*. Penalized likelihoods add a penalty term in order to maintain an appropriate level of smoothness for one or more of the nuisance parameters. This method is used in the partly linear logistic regression model described in chapter 1. Other methods of generating likelihood variants that work are possible. The basic idea is that using the likelihood principle to guide estimation of semiparametric models often leads to efficient estimators for the model components which are \sqrt{n} consistent. Because of the richness of this approach to estimation, one needs to verify for each new situation that a likelihood-inspired estimator is consistent, efficient and well-behaved for moderate sample sizes. Verifying efficiency usually entails demonstrating that the estimator satisfies the efficient score equation described in the previous section.

Unfortunately, there is no guarantee that the efficient score is a derivative of the log likelihood along some submodel. A way around this problem is to use *approximately least-favorable submodels*. This is done by finding a function $\eta_t(\theta, \eta)$ such that $\eta_0(\theta, \eta) = \eta$, for all $\theta \in \Theta$ and $\eta \in H$, where $\eta_t(\theta, \eta) \in H$ for all t small enough, and such that $\tilde{\kappa}_{\theta_0, \eta_0} = \tilde{\ell}_{\theta_0, \eta_0}$, where $\tilde{\kappa}_{\theta, \eta}(x) = \partial l_{\theta+t, \eta_t(\theta, \eta)}(x) / (\partial t)|_{t=0}$, $l_{\theta, \eta}(x)$ is the log-likelihood for the observed value x at the parameters (θ, η) , and where (θ_0, η_0) are the true parameter values. Note that we require $\tilde{\kappa}_{\theta, \eta} = \tilde{\ell}_{\theta, \eta}$ only when $(\theta, \eta) = (\theta_0, \eta_0)$. If $(\hat{\theta}_n, \hat{\eta}_n)$ is the maximum likelihood estimator, ie., the maximizer of $\mathbb{P}_n l_{\theta, \eta}$,

then the function $t \mapsto \mathbb{P}_n l_{\hat{\theta}_n + t, \eta_t}(\hat{\theta}_n, \hat{\eta}_n)$ is maximal at $t = 0$, and thus $(\hat{\theta}_n, \hat{\eta}_n)$ is a zero of $\mathbb{P}_n \tilde{\kappa}_{\hat{\theta}, \hat{\eta}}$. Now if $\hat{\theta}_n$ and $\hat{\ell}_{\hat{\theta}, \hat{\eta}_n} = \tilde{\kappa}_{\hat{\theta}, \hat{\eta}_n}$ satisfy the conditions of theorem 3.1 at $(\theta, \eta) = (\theta_0, \eta_0)$, then the maximum likelihood estimator $\hat{\theta}_n$ is asymptotically efficient at (θ_0, η_0) . We will explore in part III how this flexibility is helpful for certain models. Note that one should usually check first whether $\eta_t(\theta, \eta) = \eta$ works before trying more complicated functional forms.

Consider now the special case that both θ and η are \sqrt{n} consistent in the model $\{P_{\theta, \eta} : \theta \in \Theta, \eta \in H\}$, where $\Theta \subset \mathbb{R}^k$. Let $l_{\theta, \eta}$ be the log-likelihood for a single observation, and let $\hat{\theta}_n$ and $\hat{\eta}_n$ be the corresponding maximum likelihood estimators. Since θ is finite-dimensional, the log-likelihood can be varied with respect to θ in the usual way so that the maximum likelihood estimators satisfy $\mathbb{P}_n \dot{\ell}_{\hat{\theta}_n, \hat{\eta}_n} = 0$.

In contrast, varying η in the log-likelihood is more complex. We can typically use a subset of the submodels $t \mapsto \eta_t$ used for defining the tangent set and information for η in the model. This is easiest when the score for η is expressed as a score operator $B_{\theta, \eta}$ working on a set of indices $h \in \mathcal{H}$, as described in section 3.2. The likelihood equation for η is then usually of the form $\mathbb{P}_n B_{\hat{\theta}_n, \hat{\eta}_n} h - P_{\hat{\theta}_n, \hat{\eta}_n} B_{\hat{\theta}_n, \hat{\eta}_n} h = 0$ for all $h \in \mathcal{H}$. Note that we have forced the scores to be mean zero by subtracting off the mean rather than simply assuming $P_{\theta, \eta} B_{\theta, \eta} h = 0$. The approach is valid if there exists some path $t \mapsto \eta_t(\theta, \eta)$, with $\eta_0(\theta, \eta) = \eta$, such that

$$(3.11) \quad B_{\theta, \eta} h(x) - P_{\theta, \eta} B_{\theta, \eta} h = \partial l_{\theta + t, \eta_t(\theta, \eta)}(x) / (\partial t)|_{t=0}$$

for all x in the sample space \mathcal{X} . We assume (3.11) is valid for all $h \in \mathcal{H}$, where \mathcal{H} is chosen so that $B_{\theta, \eta} h(x) - P_{\theta, \eta} B_{\theta, \eta} h$ is uniformly bounded on \mathcal{X} for all $x \in \mathcal{X}$, $(\theta, \eta) \in \mathbb{R}^k \times H$.

We can now express $(\hat{\theta}_n, \hat{\eta}_n)$ as a Z-estimator with estimating function $\Psi_n : \mathbb{R}^k \times \mathcal{H} \mapsto \mathbb{R}^k \times \ell^\infty(\mathcal{H})$, where $\Psi_n = (\Psi_{n1}, \Psi_{n2})$, with $\Psi_{n1}(\theta, \eta) = \mathbb{P}_n \dot{\ell}_{\theta, \eta}$ and $\Psi_{n2}(\theta, \eta) = \mathbb{P}_n B_{\theta, \eta} h - P_{\theta, \eta} B_{\theta, \eta} h$, for all $h \in \mathcal{H}$. The expectation of these maps under the true parameter (θ_0, η_0) is the deterministic map $\Psi = (\Psi_1, \Psi_2)$, where $\Psi_1(\theta, \eta) = P_{\theta_0, \eta_0} \dot{\ell}_{\theta, \eta}$ and $\Psi_2(\theta, \eta) = P_{\theta_0, \eta_0} B_{\theta, \eta} h - P_{\theta, \eta} B_{\theta, \eta} h$, for all $h \in \mathcal{H}$. We have constructed these estimating equations so that the maximum likelihood estimators and true parameters satisfy $\Psi_n(\hat{\theta}_n, \hat{\eta}_n) = 0 = \Psi(\theta_0, \eta_0)$. Provided H is a subset of a normed space, we can use the Z-estimator master theorem (theorem 2.11) to obtain weak convergence of $\sqrt{n}(\hat{\theta}_n - \theta_0, \hat{\eta}_n - \eta_0)$:

COROLLARY 3.2 *Suppose that $\dot{\ell}_{\theta, \eta}$ and $B_{\theta, \eta} h$, with h ranging over \mathcal{H} and with (θ, η) ranging over a neighborhood of (θ_0, η_0) , are contained in a P_{θ_0, η_0} -Donsker class, and that both $P_{\theta_0, \eta_0} \left\| \dot{\ell}_{\theta, \eta} - \dot{\ell}_{\theta_0, \eta_0} \right\|^2 \xrightarrow{P} 0$ and $\sup_{h \in \mathcal{H}} P_{\theta_0, \eta_0} |B_{\theta, \eta} h - B_{\theta_0, \eta_0} h|^2 \xrightarrow{P} 0$, as $(\theta, \eta) \rightarrow (\theta_0, \eta_0)$. Also assume that Ψ is Fréchet-differentiable at (θ_0, η_0) with derivative $\dot{\Psi}_0 : \mathbb{R}^k \times \text{lin } H \mapsto \mathbb{R}^k \times \ell^\infty(\mathcal{H})$*

that is continuously-invertible and onto its range, with inverse $\dot{\Psi}_0^{-1} : \mathbb{R}^k \times \ell^\infty(\mathcal{H}) \mapsto \mathbb{R}^k \times \text{lin } H$. Then, provided $(\hat{\theta}_n, \hat{\eta}_n)$ is consistent for (θ_0, η_0) and $\Psi_n(\hat{\theta}_n, \hat{\eta}_n) = o_p(n^{-1/2})$ (uniformly over $\mathbb{R}^k \times \ell^\infty(\mathcal{H})$), $(\hat{\theta}_n, \hat{\eta}_n)$ is efficient at (θ_0, η_0) and $\sqrt{n}(\hat{\theta}_n - \theta_0, \hat{\eta}_n - \eta_0) \rightsquigarrow -\dot{\Psi}_0^{-1}Z$, where Z is the Gaussian limiting distribution of $\sqrt{n}\Psi_n(\theta_0, \eta_0)$.

The proof of this will be given later in part III. A function $f : U \mapsto V$ is onto if, for every $v \in V$, there exists a $u \in U$ with $v = f(u)$. Note that while \mathcal{H} can be a subset of the tangent set used for information calculations, it must be rich enough to ensure that the inverse of $\dot{\Psi}_0$ exists.

The efficiency in corollary 3.2 can be shown to follow from the score operator calculations given in section 3.2, but we will postpone further details until part III. As was done at the end of section 3.2, the above discussion can be completely re-expressed in terms of a single parameter model $\{P_\eta : \eta \in H\}$ with a single score operator $A_\eta : \mathcal{H}_\eta \mapsto L_2(P_\eta)$, where H is a richer parameter set, including, for example, both Θ and H as defined in the previous paragraphs, and where \mathcal{H}_η is similarly enriched to include the tangent sets for all subcomponents of the model.

3.4 Other Topics

Other topics of importance include frequentist and Bayesian methods for constructing confidence sets. We will focus primarily on frequentist approaches in this book and only briefly discuss Bayesian methods. While the bootstrap is generally valid in the setting of corollary 3.2, it is unclear that this remains true when the nuisance parameter converges at a rate slower than \sqrt{n} , even if interest is limited to the parametric component. Even when the bootstrap is valid, it may be excessively cumbersome to re-estimate the entire model for many bootstrapped data sets. We will explore this issue in more detail in part III. We now present one approach for hypothesis testing and variance estimation for the parametric component θ of the semiparametric model $\{P_{\theta, \eta} : \theta \in \Theta, \eta \in H\}$. This approach is often valid even when the nuisance parameter is not \sqrt{n} consistent.

Murphy and van der Vaart demonstrated that under reasonable regularity conditions, the log-profile likelihood, $pl_n(\theta)$, (profiling over the nuisance parameter) admits the following expansion about the maximum likelihood estimator for the parametric component $\hat{\theta}_n$:

$$(3.12) \quad \log pl_n(\tilde{\theta}_n) = \log pl_n(\hat{\theta}_n) - \frac{1}{2}n(\tilde{\theta}_n - \hat{\theta}_n)' \tilde{I}_{\theta_0, \eta_0}(\tilde{\theta}_n - \hat{\theta}_n) + o_p(\sqrt{n}\|\tilde{\theta}_n - \theta_0\| + 1)^2,$$

for any estimator $\tilde{\theta}_n \xrightarrow{P} \theta_0$ (Murphy and van der Vaart, 2000). This can be shown to lead naturally to chi-square tests of full versus reduced models.

Furthermore, this result demonstrates that the curvature of the log-partial likelihood can serve as a consistent estimator of the efficient information for θ at θ_0 , $\tilde{I}_{\theta_0, \eta_0}$, and thereby permit estimation of the limiting variance of $\sqrt{n}(\hat{\theta}_n - \theta_0)$. We will discuss this method, along with several other methods of inference, in greater detail toward the end of part III.

3.5 Exercises

3.5.1. Consider the paragraphs leading up to the efficient score for β in the Cox model, given in expression (3.3). Show that the tangent set for the nuisance parameter Λ , $\dot{\mathcal{P}}_{P_{\beta, \Lambda}}^{(\Lambda)}$, spans all square-integrable score functions for Λ generated by parametric submodels (where the parameters for Λ are independent of β).

3.5.2. Let $L_n(\beta, \Lambda)$ be the Cox model likelihood given in (3.10). Show that the corresponding profile likelihood for β , $pL_n(\beta)$, obtained by maximizing over Λ , equals $\exp[-\sum_{i=1}^n d_i]$ times the partial likelihood (3.4).

3.6 Notes

The linear regression example was partially inspired by example 25.28 of van der Vaart (1998), and theorem 3.1 is a generalization of his theorem 25.57 based on his condition (25.28).

4

Case Studies I

We now expand upon several examples introduced in chapters 1–3 to more fully illustrate the methods and theory we have outlined thus far. Certain technical aspects which involve concepts introduced later in the book will be glossed over to avoid getting bogged down with details. The main objective of this chapter is to initiate an appreciation for what empirical processes and efficiency calculations can accomplish.

The first example is linear regression with either mean zero or median zero residuals. In addition to efficiency calculations for model parameters, empirical processes are needed for inference on the distribution of the residuals. The second example is counting process regression for both general counting processes and the Cox model for right-censored failure time data. Empirical processes will be needed for parameter inference, and efficiency will be established under the Cox proportional hazards model for maximum likelihood estimation of both the regression parameters and the baseline hazard. The third example is the Kaplan-Meier estimator of the survival function for right-censored failure time data. Since weak convergence of the Kaplan-Meier has already been established in chapter 2 using empirical processes, the focus in this chapter will be on efficiency calculations. The fourth example considers estimating equations for general regression models when the residual variance may be a function of the covariates. Estimation of the variance function is needed for efficient estimation. We also consider optimality of a certain class of estimating equations. The general results are illustrated with both simple linear regression and a Poisson mixture regression model. In the latter case, the mixture distribution is not \sqrt{n} consistent in the uniform norm, the proof of which fact we omit. The

fifth, and final, example is partly linear logistic regression. The emphasis will be on efficient estimation of the parametric regression parameter. For the last two examples, both empirical processes and efficiency calculations will be needed.

4.1 Linear Regression

The semiparametric linear regression model is $Y = \beta'Z + e$, where we observe $X = (Y, Z)$ and assume $E\|Z\|^2 < \infty$, $E[e|Z] = 0$ and $E[e^2|Z] \leq K < \infty$ almost surely, Z includes the constant 1, and $E[ZZ']$ is full rank. The model for X is $\{P_{\beta,\eta} : \beta \in \mathbb{R}^k, \eta \in H\}$, where η is the joint density of the residual e and covariate Z with partial derivative with respect to the first argument $\dot{\eta}_1$ which we assume satisfies $\dot{\eta}/\eta \in L_2(P_{\beta,\eta})$ and hence $\dot{\eta}/\eta \in L_2^0(P_{\beta,\eta})$. We consider this model under two assumptions on η : the first is that the residuals have conditional mean zero, i.e., $\int_{\mathbb{R}} u\eta(u, Z)du = 0$ almost surely; and the second is that the residuals have median zero, i.e., $\int_{\mathbb{R}} \text{sign}(u)\eta(u, Z)du = 0$ almost surely.

4.1.1 Mean Zero Residuals

We have already demonstrated in section 3.2 that the usual least squares estimator $\hat{\beta} = [\mathbb{P}_n ZZ']^{-1} \mathbb{P}_n ZY$ is \sqrt{n} consistent but not always efficient for β when the only assumption we are willing to make is that the residuals have mean zero conditional on the covariates. The basic argument for this was taken from the form of the efficient score $\tilde{\ell}_{\beta,\eta}(Y, Z) = Z(Y - \beta'Z)/P_{\beta,\eta}[e^2|Z]$ which yields a distinctly different estimator than $\hat{\beta}$ when $z \mapsto P_{\beta,\eta}[e^2|Z = z]$ is non-constant in z . In section 4.4 of this chapter, we will describe a data-driven procedure for efficient estimation in this context based on approximating the efficient score.

For now, however, we turn our attention to efficient estimation when we also assume that the covariates are independent of the residual. Accordingly, we denote η to be the density of the residual e and $\dot{\eta}$ to be the derivative of η . We discussed this model in chapter 1 and pointed at that $\hat{\beta}$ is still not efficient in this setting. We also claimed that an empirical estimator \hat{F} of the residual distribution, based on the residuals $Y_1 - \hat{\beta}'Z_1, \dots, Y_n - \hat{\beta}'Z_n$, had the property that $\sqrt{n}(\hat{F} - F)$ converged in a uniform sense to a certain Gaussian quantity. We now derive the efficient score for estimating β for this model and sketch a proof of the claimed convergence of $\sqrt{n}(\hat{F} - F)$. We assume that η is continuously differentiable with $P_{\beta,\eta}(\dot{\eta}/\eta)^2 < \infty$.

Using techniques described in section 3.2, it is not hard to verify that the tangent set for the full model is the linear span of $-(\dot{\eta}/\eta)(e)a'Z + b(e)$, as a spans \mathbb{R}^k and b spans the tangent set $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$ for η , consisting of all functions in $L_2^0(e)$ which are orthogonal to e , where we use $L_2^0(U)$ to denote all mean

zero real functions f of the random variable U for which $Pf^2(U) < \infty$. The structure of the tangent set follows from the constraint that $\int_{\mathbb{R}} e\eta(e)de = 0$.

The projection of the usual score for β , $\dot{\ell}_{\beta,\eta} \equiv -(\dot{\eta}/\eta)(e)Z$, onto $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$ is a function $h \in \dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$ such that $-(\dot{\eta}/\eta)(e)Z - h(e)$ is uncorrelated with $b(e)$ for all $b \in \dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$. We will now verify that $h(e) = -(\dot{\eta}/\eta)(e)\mu - e\mu/\sigma^2$, where $\mu \equiv E[Z]$ and $\sigma^2 = E[e^2]$. It is easy to see that h is square integrable. Moreover, since $\int_{\mathbb{R}} e\dot{\eta}(e)de = -1$, h has mean zero. Thus $h \in \dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$. It also follows that $-(\dot{\eta}/\eta)(e)Z - h(e)$ is orthogonal to $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$, after noting that $-(\dot{\eta}/\eta)(e)Z - h(e) = -(\dot{\eta}/\eta)(e)(Z - \mu) + e\mu/\sigma^2 \equiv \tilde{\ell}_{\beta,\eta}$ is orthogonal to all square-integrable mean zero functions $b(e)$ which satisfy $P_{\beta,\eta}b(e)e = 0$.

Thus the efficient information is

$$(4.1) \quad \tilde{I}_{\beta,\eta} \equiv P_{\beta,\eta} \left[\tilde{\ell}_{\beta,\eta} \tilde{\ell}'_{\beta,\eta} \right] = P_{\beta,\eta} \left[\left\{ \frac{\dot{\eta}}{\eta}(e) \right\}^2 (Z - \mu)(Z - \mu)' \right] + \frac{\mu\mu'}{\sigma^2}.$$

Since

$$1 = \left(\int_{\mathbb{R}} e\dot{\eta}(e)de \right)^2 = \left(\int_{\mathbb{R}} e \frac{\dot{\eta}}{\eta}(e)\eta(e)de \right)^2 \leq \sigma^2 P_{\beta,\eta} \left\{ \frac{\dot{\eta}}{\eta}(e) \right\}^2,$$

we have that

$$(4.1) \geq P_{\beta,\eta} \left[\frac{(Z - \mu)(Z - \mu)'}{\sigma^2} \right] + \frac{\mu\mu'}{\sigma^2} = \frac{P[ZZ']}{\sigma^2},$$

where, for two $k \times k$ symmetric matrices A and B , $A \geq B$ means that $A - B$ is positive semidefinite. Thus the efficient estimator can have strictly lower variance than the least-squares estimator $\hat{\beta}$.

Developing a procedure for calculating an asymptotically efficient estimator appears to require nonparametric estimation of $\dot{\eta}/\eta$. We will show in the empirical process case studies of chapter 15 how to accomplish this via the residual distribution estimator \hat{F} defined above. We now turn our attention to verifying that $\sqrt{n}(\hat{F} - F)$ converges weakly to a tight, mean zero Gaussian process. This result can also be used to check normality of the residuals. Recall that if the residuals are Gaussian, the least-squares estimator is fully efficient.

Now, using $P = P_{\beta,\eta}$,

$$\begin{aligned} \sqrt{n}(\hat{F}(v) - F(v)) &= \sqrt{n} \left[\mathbb{P}_n 1\{Y - \hat{\beta}'Z \leq v\} - P 1\{Y - \beta'Z \leq v\} \right] \\ &= \sqrt{n}(\mathbb{P}_n - P) 1\{Y - \hat{\beta}'Z \leq v\} \\ &\quad + \sqrt{n}P \left[1\{Y - \hat{\beta}'Z \leq v\} - 1\{Y - \beta'Z \leq v\} \right] \\ &= U_n(v) + V_n(v). \end{aligned}$$

We will show in part II that $\{1\{Y - b'Z \leq v\} : b \in \mathbb{R}^k, v \in \mathbb{R}\}$ is a VC (and hence Donsker) class of functions. Thus, since

$$\sup_{v \in \mathbb{R}} P \left[1\{Y - \hat{\beta}'Z \leq v\} - 1\{Y - \beta'Z \leq v\} \right]^2 \xrightarrow{P} 0,$$

we have that $U_n(v) = \sqrt{n}(\mathbb{P}_n - P)1\{Y - \beta'Z \leq v\} + \epsilon_n(v)$, where $\sup_{v \in \mathbb{R}} |\epsilon_n(v)| \xrightarrow{P} 0$. It is not difficult to show that

$$(4.2) \quad V_n(v) = P \left[\int_v^{v+(\hat{\beta}-\beta)'Z} \eta(u) du \right].$$

We leave it as an exercise to show that for any $u, v \in \mathbb{R}$,

$$(4.3) \quad |\eta(u) - \eta(v)| \leq |F(u) - F(v)|^{1/2} \left(P \left\{ \frac{\dot{\eta}(e)}{\eta} \right\}^2 \right)^{1/2}.$$

Thus η is both bounded and equicontinuous, and thus by (4.2), $V_n(v) = \sqrt{n}(\hat{\beta} - \beta)' \mu \eta(v) + \epsilon'_n(v)$, where $\sup_{v \in \mathbb{R}} |\epsilon'_n(v)| \xrightarrow{P} 0$. Hence \hat{F} is asymptotically linear with influence function

$$\check{\psi}(v) = 1\{e \leq v\} - F(v) + eZ' \{P[ZZ']\}^{-1} \mu \eta(v),$$

and thus, since this influence function is a Donsker class (as the sum of two obviously Donsker classes), $\sqrt{n}(\hat{F} - F)$ converges weakly to a tight, mean zero Gaussian process.

4.1.2 Median Zero Residuals

We have already established in section 2.2.6 that when the residuals have median zero and are independent of the covariates, then the least-absolute-deviation estimator $\hat{\beta} \equiv \operatorname{argmin}_{b \in \mathbb{R}^k} \mathbb{P}_n |Y - b'Z|$ is consistent for β and $\sqrt{n}(\hat{\beta} - \beta)$ is asymptotically linear with influence function

$$\check{\psi} = \{2\eta(0)P[ZZ']\}^{-1} Z \operatorname{sign}(e).$$

Thus $\sqrt{n}(\hat{\beta} - \beta)$ is asymptotically mean zero Gaussian with covariance equal to $\{4\eta^2(0)P[ZZ']\}^{-1}$. In this section, we will study efficiency for this model and show that $\hat{\beta}$ is not in general fully efficient. Before doing this, however, we will briefly study efficiency in the more general model where we only assume $E[\operatorname{sign}(e)|Z] = 0$ almost surely.

Under this more general model, the joint density of (e, Z) , η , must satisfy $\int_{\mathbb{R}} \operatorname{sign}(e) \eta(e, Z) de = 0$ almost surely. As we did when we studied the conditionally mean zero residual case in section 3.2, assume η has partial

derivative with respect to the first argument, $\dot{\eta}_1$, satisfying $\dot{\eta}_1/\eta \in L_2(P_{\beta,\eta})$. Clearly, $(\dot{\eta}_1/\eta)(e, Z)$ also has mean zero. The score for β , assuming η is known, is $\dot{\ell}_{\beta,\eta} = -Z(\dot{\eta}_1/\eta)(Y - \beta'Z, Z)$.

Similar to what was done in section 3.2 for the conditionally mean zero case, it can be shown that the tangent set $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$ for η is the subset of $L_2^0(P_{\beta,\eta})$ which consists of all functions $g(e, Z) \in L_2^0(P_{\beta,\eta})$ which satisfy

$$E[\text{sign}(e)g(e, Z)|Z] = \frac{\int_{\mathbb{R}} \text{sign}(e)g(e, Z)\eta(e, Z)de}{\int_{\mathbb{R}} \eta(e, Z)de} = 0,$$

almost surely. It can also be shown that this set is the orthocomplement in $L_2^0(P_{\beta,\eta})$ of all functions of the form $\text{sign}(e)f(Z)$, where f satisfies $P_{\beta,\eta}f^2(Z) < \infty$. Hence the efficient score $\tilde{\ell}_{\beta,\eta}$ is the projection in $L_2^0(P_{\beta,\eta})$ of $-Z(\dot{\eta}_1/\eta)(e, Z)$ onto $\{\text{sign}(e)f(Z) : P_{\beta,\eta}f^2(Z) < \infty\}$. Hence

$$\tilde{\ell}_{\beta,\eta}(Y, Z) = -Z\text{sign}(e) \int_{\mathbb{R}} \dot{\eta}_1(e, Z)\text{sign}(e)de = Z\text{sign}(e)\eta(0, Z),$$

where the second equality follows from the facts that $\int_{\mathbb{R}} \dot{\eta}_1(e, Z)de = 0$ and $\int_{-\infty}^0 \dot{\eta}_1(e, Z)de = \eta(0, Z)$. When $\eta(0, z)$ is non-constant in z , $\tilde{\ell}_{\beta,\eta}(Y, Z)$ is not proportional to $\text{sign}(Z)(Y - \beta'Z)$, and thus the least-absolute-deviation estimator is not efficient in this instance. Efficient estimation in this situation appears to require estimation of $\eta(0, Z)$, but we will not pursue this further.

We now return our attention to the setting where the median zero residuals are independent of the covariates. We will also assume that $0 < \eta(0) < \infty$. Recall that in this setting, the usual score for β is $\dot{\ell}_{\beta,\eta} \equiv -Z(\dot{\eta}/\eta)(e)$, where $\dot{\eta}$ is the derivative of the density η of e . If we temporarily make the fairly strong assumption that the residuals have a Laplace density (ie., $\eta(e) = (\nu/2)\exp(-\nu|e|)$ for a real parameter $\nu > 0$), then the usual score simplifies to $Z\text{sign}(e)$, and thus the least-absolute deviation estimator is fully efficient for this special case.

Relaxing the assumptions to allow for arbitrary median zero density, we can follow arguments similar to those we used in the previous section to obtain that the tangent set for η , $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$, consists of all functions in $L_2^0(e)$ which are orthogonal to $\text{sign}(e)$. The structure of the tangent set follows from the median zero residual constraint which can be expressed as $\int_{\mathbb{R}} \text{sign}(e)\eta(e)de = 0$. The projection of $\dot{\ell}_{\beta,\eta}$ on $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$ is a function $h \in \dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$ such that $-(\dot{\eta}/\eta)(e)Z - h(e)$ is uncorrelated with $b(e)$ for all $b \in \dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$.

We will now prove that $h(e) = -(\dot{\eta}/\eta)(e)\mu - 2\eta(0)\text{sign}(e)\mu$ satisfies the above constraints. First, it is easy to see that $h(e)$ is square-integrable. Second, since $-\int_{\mathbb{R}} \text{sign}(e)\dot{\eta}(e)de = 2\eta(0)$, we have that $\text{sign}(e)h(e)$ has zero expectation. Thus $h \in \dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$. Thirdly, it is straightforward to verify that $\dot{\ell}_{\beta,\eta}(Y, Z) - h(e) = -(\dot{\eta}/\eta)(e)(Z - \mu) + 2\eta(0)\text{sign}(e)\mu$ is orthogonal to

all elements of $\dot{\mathcal{P}}_{P_{\beta,\eta}}^{(\eta)}$. Hence the efficient score is $\tilde{\ell}_{\beta,\eta} = \dot{\ell}_{\beta,\eta} - h$. Thus the efficient information is

$$\tilde{I}_{\beta,\eta} = P_{\beta,\eta} \left[\tilde{\ell}_{\beta,\eta} \tilde{\ell}'_{\beta,\eta} \right] = P_{\beta,\eta} \left[\left\{ \frac{\dot{\eta}}{\eta}(e) \right\}^2 (Z - \mu)(Z - \mu)' \right] + 4\eta^2(0)\mu\mu'.$$

Note that

$$\begin{aligned} 4\eta^2(0) &= 4 \left[\int_{-\infty}^0 \dot{\eta}(e) de \right]^2 \\ &= 4 \left[\int_{-\infty}^0 \frac{\dot{\eta}}{\eta}(e) \eta(e) de \right]^2 \\ &\leq 4 \int_{-\infty}^0 \left\{ \frac{\dot{\eta}}{\eta}(e) \right\}^2 \eta(e) de \int_{-\infty}^0 \eta(e) de \\ (4.4) \quad &= 2 \int_{-\infty}^0 \left\{ \frac{\dot{\eta}}{\eta}(e) \right\}^2 \eta(e) de. \end{aligned}$$

Similar arguments yield that

$$4\eta^2(0) \leq 2 \int_0^{\infty} \left\{ \frac{\dot{\eta}}{\eta}(e) \right\}^2 \eta(e) de.$$

Combining this last inequality with (4.4), we obtain that

$$4\eta^2(0) \leq \int_{\mathbb{R}} \left\{ \frac{\dot{\eta}}{\eta}(e) \right\}^2 \eta(e) de.$$

Hence the efficient estimator can have strictly lower variance than the least-absolute-deviation estimator in this situation.

While the least-absolute-deviation estimator is only guaranteed to be fully efficient for Laplace distributed residuals, it does have excellent robustness properties. In particular, the *breakdown point* of this estimator is 50%. The breakdown point is the maximum proportion of contamination—from arbitrarily large symmetrically distributed residuals—tolerated by the estimator without resulting in inconsistency.

4.2 Counting Process Regression

We now examine in detail the counting process regression model considered in chapter 1. The observed data are $X = (N, Y, Z)$, where for $t \in [0, \tau]$, $N(t)$ is a counting process and $Y(t) = 1\{V \geq t\}$ is an at-risk process based on a random time $V \geq 0$ which may depend on N , with $PY(0) = 1$, $\inf_Z P[Y(\tau)|Z] > 0$, $PN^2(\tau) < \infty$, and where $Z \in \mathbb{R}^k$ is a regression

covariate. The regression model (1.3) is assumed, implying $E\{dN(t)|Z\} = E\{Y(t)|Z\}e^{\beta'Z}d\Lambda(t)$, for some $\beta \in B \subset \mathbb{R}^k$ and continuous nondecreasing function $\Lambda(t)$ with $\Lambda(0) = 0$ and $0 < \Lambda(\tau) < \infty$. We assume Z is restricted to a bounded set, $\text{var}(Z)$ is positive definite, and that B is open, convex and bounded. We first consider inference for β and Λ in this general model, and then examine the specialization to the Cox proportional hazards model for right-censored failure time data.

4.2.1 The General Case

As described in chapter 1, we estimate β with the estimating equation $U_n(t, \beta) = \mathbb{P}_n \int_0^t [Z - E_n(s, \beta)] dN(s)$, where

$$E_n(t, \beta) = \frac{\mathbb{P}_n ZY(t)e^{\beta'Z}}{\mathbb{P}_n Y(t)e^{\beta'Z}}.$$

Specifically, the estimator $\hat{\beta}$ is a root of $U_n(\tau, \beta) = 0$. The estimator for Λ , is $\hat{\Lambda}(t) = \int_0^t [\mathbb{P}_n Y(s)e^{\hat{\beta}'Z}]^{-1} \mathbb{P}_n dN(s)$. We first show that $\hat{\beta}$ is consistent for β , and that $\sqrt{n}(\hat{\beta} - \beta)$ is asymptotically normal with a consistently estimable covariance matrix. We then derive consistency and weak convergence results for $\hat{\Lambda}$ and suggest a simple method of inference based on the influence function.

We first argue that certain classes of functions are Donsker, and therefore also Glivenko-Cantelli. We first present the following lemma:

LEMMA 4.1 *For $-\infty < a < b < \infty$, Let $\{X(t), t \in [a, b]\}$ be a monotone cadlag or caglad stochastic process with $P[|X(a)| \vee |X(b)|]^2 < \infty$. Then X is P -Donsker.*

The proof will be given later in part II. Note that we usually speak of classes of functions as being Donsker or Glivenko-Cantelli, but in lemma 4.1, we are saying this about a process. Let \mathcal{X} be the sample space for the stochastic process $\{X(t) : t \in T\}$. Then $\sqrt{n}(\mathbb{P}_n - P)X$ converges weakly in $\ell^\infty(T)$ to a tight, mean zero Gaussian process if and only if $\mathcal{F} = \{f_t : t \in T\}$ is P -Donsker, where for any $x \in \mathcal{X}$ and $t \in T$, $f_t(x) = x(t)$. Viewed in this manner, this modified use of the term Donsker is, in fact, not a modification at all.

Since Y and N both satisfy the conditions of lemma 4.1, they are both Donsker as processes in $\ell^\infty([0, \tau])$. Trivially, the classes $\{\beta \in B\}$ and $\{Z\}$ are both Donsker classes, and therefore so is $\{\beta'Z : \beta \in B\}$ since products of bounded Donsker classes are Donsker. Now the class $\{e^{\beta'Z} : \beta \in B\}$ is Donsker since exponentiation is Lipschitz continuous on compacts. Hence $\{Y(t)e^{\beta'Z} : t \in [0, \tau], \beta \in B\}$, $\{ZY(t)e^{\beta'Z} : t \in [0, \tau], \beta \in B\}$, and $\{ZZ'Y(t)e^{\beta'Z} : t \in [0, \tau], \beta \in B\}$, are all Donsker since they are all products of bounded Donsker classes.

Now the derivative of $U_n(\tau, \beta)$ with respect to β can be shown to be $-V_n(\beta)$, where

$$V_n(\beta) = \int_0^\tau \left[\frac{\mathbb{P}_n ZZ'Y(t)e^{\beta'Z}}{\mathbb{P}_n Y(t)e^{\beta'Z}} - \left\{ \frac{\mathbb{P}_n ZY(t)e^{\beta'Z}}{\mathbb{P}_n Y(t)e^{\beta'Z}} \right\}^{\otimes 2} \right] \mathbb{P}_n dN(t),$$

and where superscript $\otimes 2$ denotes outer product. Because all of the classes involved are Glivenko-Cantelli and the limiting values of the denominators are bounded away from zero, we have the $\sup_{\beta \in B} |V_n(\beta) - V(\beta)| \xrightarrow{\text{as}^*} 0$, where

$$(4.5) \quad V(\beta) = \int_0^\tau \left[\frac{PZZ'Y(t)e^{\beta'Z}}{PY(t)e^{\beta'Z}} - \left\{ \frac{PZY(t)e^{\beta'Z}}{PY(t)e^{\beta'Z}} \right\}^{\otimes 2} \right] \times P \left[Y(t)e^{\beta'Z} \right] d\Lambda(t).$$

After some work, it can be shown that there exists a $c > 0$ such that $V(\beta) \geq c \text{var}(Z)$, where for two symmetric matrices A, B , $A \geq B$ means that $A - B$ is positive semidefinite. Thus $U_n(\tau, \beta)$ is almost surely convex for all $n \geq 1$ large enough. Thus $\hat{\beta}$ is almost surely consistent for the true parameter β_0 .

Using algebra, $U_n(\tau, \beta) = \mathbb{P}_n \int_0^\tau [Z - E_n(s, \beta)] dM_\beta(s)$, where $M_\beta(t) = N(t) - \int_0^t Y(s)e^{\beta'Z} d\Lambda_0(s)$ and Λ_0 is the true value of Λ . Let $U(t, \beta) = P \left\{ \int_0^t [Z - E(s, \beta)] dM_\beta(s) \right\}$, where $E(t, \beta) = P \left[ZY(t)e^{\beta'Z} \right] / P \left[Y(t)e^{\beta'Z} \right]$. It is not difficult to verify that

$$\sqrt{n} \left[U_n(\tau, \hat{\beta}) - U(\tau, \hat{\beta}) \right] - \sqrt{n} \left[U_n(\tau, \beta_0) - U(\tau, \beta_0) \right] \xrightarrow{P} 0,$$

since

$$(4.6) \quad \begin{aligned} & \sqrt{n} \mathbb{P}_n \int_0^\tau \left[E_n(s, \hat{\beta}) - E(s, \hat{\beta}) \right] dM_{\hat{\beta}}(s) \\ &= \sqrt{n} \mathbb{P}_n \int_0^\tau \left[E_n(s, \hat{\beta}) - E(s, \hat{\beta}) \right] \\ & \quad \times \left\{ dM_{\beta_0}(s) - Y(s) \left[e^{\hat{\beta}'Z} - e^{\beta_0'Z} \right] d\Lambda_0(s) \right\} \\ & \xrightarrow{P} 0. \end{aligned}$$

This follows from the following lemma (which we prove later in part II), with $[a, b] = [0, \tau]$, $A_n(t) = E_n(t, \hat{\beta})$, and $B_n(t) = \sqrt{n} \mathbb{P}_n M_{\beta_0}(t)$:

LEMMA 4.2 *Let $B_n \in D[a, b]$ and $A_n \in \ell^\infty([a, b])$ be either cadlag or caglad, and assume $\sup_{t \in [a, b]} |A_n(t)| \xrightarrow{P} 0$, A_n has uniformly bounded total variation, and B_n converges weakly to a tight, mean zero process with sample paths in $D[a, b]$. Then $\int_a^b A_n(s) dB_n(s) \xrightarrow{P} 0$.*

Thus the conditions of the Z-estimator master theorem, theorem 2.11, are all satisfied, and $\sqrt{n}(\hat{\beta} - \beta_0)$ converges weakly to a mean zero random vector with covariance $C = V^{-1}(\beta_0)P \left[\int_0^\tau [Z - E(s, \beta_0)] dM_{\beta_0}(s) \right]^{\otimes 2} V^{-1}(\beta_0)$. With a little additional work, it can be verified that C can be consistently estimated with $\hat{C} = V_n^{-1}(\hat{\beta})\mathbb{P}_n \left\{ \int_0^\tau [Z - E_n(s, \hat{\beta})] d\hat{M}(s) \right\}^{\otimes 2} V_n^{-1}(\hat{\beta})$, where $\hat{M}(t) = N(t) - \int_0^t Y(s)e^{\hat{\beta}'Z} d\hat{\Lambda}(s)$, and where

$$\hat{\Lambda}(t) = \int_0^t \frac{\mathbb{P}_n dN(s)}{\mathbb{P}_n Y(s)e^{\hat{\beta}'Z}},$$

as defined in chapter 1.

Let Λ_0 be the true value of Λ . Then

$$\begin{aligned} \hat{\Lambda}(t) - \Lambda_0(t) &= \int_0^t 1\{\mathbb{P}_n Y(s) > 0\} \left\{ \frac{(\mathbb{P}_n - P)dN(s)}{\mathbb{P}_n Y(s)e^{\hat{\beta}'Z}} \right\} \\ &\quad - \int_0^t 1\{\mathbb{P}_n Y(s) = 0\} \left\{ \frac{PdN(s)}{\mathbb{P}_n Y(s)e^{\hat{\beta}'Z}} \right\} \\ &\quad - \int_0^t \frac{(\mathbb{P}_n - P)Y(s)e^{\hat{\beta}'Z}}{PY(s)e^{\hat{\beta}'Z}} \left\{ \frac{PdN(s)}{\mathbb{P}_n Y(s)e^{\hat{\beta}'Z}} \right\} \\ &\quad - \int_0^t \frac{P \left[Y(s) \left(e^{\hat{\beta}'Z} - e^{\beta_0'Z} \right) \right]}{PY(s)e^{\hat{\beta}'Z}} d\Lambda_0(s) \\ &= A_n(t) - B_n(t) - C_n(t) - D_n(t). \end{aligned}$$

By the smoothness of these functions of $\hat{\beta}$ and the almost sure consistency of $\hat{\beta}$, each of the processes A_n , B_n , C_n and D_n converge uniformly to zero, and thus $\sup_{t \in [0, \tau]} |\hat{\Lambda}(t) - \Lambda_0(t)| \xrightarrow{\text{as*}} 0$. Since $P\{\mathbb{P}_n Y(t) = 0\} \leq [P\{V < \tau\}]^n$, $B_n(t) = o_p(n^{-1/2})$, where the $o_p(n^{-1/2})$ term is uniform in t . It is also not hard to verify that $A_n(t) = (\mathbb{P}_n - P)\tilde{A}(t) + o_p(n^{-1/2})$, $C_n(t) = (\mathbb{P}_n - P)\tilde{C}(t) + o_p(n^{-1/2})$, where $\tilde{A}(t) = \int_0^t [PY(s)e^{\beta_0'Z}]^{-1} dN(s)$, $\tilde{C}(t) = \int_0^t [PY(s)e^{\beta_0'Z}]^{-1} Y(s)e^{\beta_0'Z} d\Lambda_0(s)$, and both remainder terms are uniform in t . In addition,

$$D_n(t) = (\hat{\beta} - \beta_0)' \int_0^t \left\{ \frac{PZY(s)e^{\beta_0'Z}}{PY(s)e^{\beta_0'Z}} \right\} d\Lambda_0(t) + o_p(n^{-1/2}),$$

where the remainder term is again uniform in t .

Taking this all together, we obtain the expansion

$$\begin{aligned} \sqrt{n} [\hat{\Lambda}(t) - \Lambda_0(t)] &= \sqrt{n}(\mathbb{P}_n - P) \int_0^t \frac{dM_{\beta_0}(s)}{PY(s)e^{\beta_0'Z}} \\ &\quad - \sqrt{n}(\hat{\beta} - \beta_0)' \int_0^t E(s, \beta_0) d\Lambda_0(s) + o_p(1), \end{aligned}$$

where the remainder term is uniform in t . By previous arguments, $\sqrt{n}(\hat{\beta} - \beta_0) = V^{-1}(\beta_0)\mathbb{P}_n \int_0^\tau [Z - E(s, \beta_0)] dM_{\beta_0}(s) + o_p(1)$, and thus $\hat{\Lambda}$ is asymptotically linear with influence function

$$(4.7) \quad \begin{aligned} \psi(t) &= \int_0^t \frac{dM_{\beta_0}(s)}{PY(s)e^{\beta_0'Z}} \\ &\quad - \left\{ \int_0^\tau [Z - E(s, \beta_0)]' dM_{\beta_0}(s) \right\} V^{-1}(\beta_0) \int_0^t E(s, \beta_0) d\Lambda_0(s). \end{aligned}$$

Since $\{\psi(t) : t \in T\}$ is Donsker, $\sqrt{n}(\hat{\Lambda} - \Lambda)$ converges weakly in $D[0, \tau]$ to a tight, mean zero Gaussian process \mathcal{Z} with covariance $P[\psi(s)\psi(t)]$. Let $\hat{\psi}(t)$ be $\psi(t)$ with $\hat{\beta}$ and $\hat{\Lambda}$ substituted for β_0 and Λ_0 , respectively.

After some additional analysis, it can be verified that

$$\sup_{t \in [0, \tau]} \mathbb{P}_n [\hat{\psi}(t) - \psi(t)]^2 \xrightarrow{P} 0.$$

Since ψ has envelope $kN(\tau)$, for some fixed $k < \infty$, the class $\{\psi(s)\psi(t) : s, t \in [0, \tau]\}$ is Glivenko-Cantelli, and thus $\mathbb{P}_n [\hat{\psi}(s)\hat{\psi}(t)]$ is uniformly consistent (in probability) for $P[\psi(s)\psi(t)]$ by the discussion in the beginning of section 2.2.3. However, this is not particularly helpful for inference, and the following approach is better. Let ξ be standard normal and independent of the data $X = (N, Y, Z)$, and consider the wild bootstrap $\tilde{\Delta}(t) = \mathbb{P}_n \xi \hat{\psi}(t)$. Although some work is needed, it can be shown that $\sqrt{n} \tilde{\Delta} \xrightarrow[\xi]{P} \mathcal{Z}$. This is computationally quite simple, since it requires saving $\hat{\psi}_1(t_j), \dots, \hat{\psi}_n(t_j)$ only at all of the observed jump points t_1, \dots, t_{m_n} , drawing a sample of standard normals ξ_1, \dots, ξ_n , evaluating $\sup_{1 \leq j \leq m_n} |\tilde{\Delta}(t_j)|$, and repeating often enough to obtain a reliable estimate of the $(1 - \alpha)$ -level quantile \hat{c}_α . The confidence band $\{\hat{\Lambda}(t) \pm \hat{c}_\alpha : t \in [0, \tau]\}$ thus has approximate coverage $1 - \alpha$. A number of modifications of this are possible, including a modification where the width of the band at t is roughly proportional to the variance of $\sqrt{n}(\hat{\Lambda}(t) - \Lambda_0(t))$.

4.2.2 The Cox Model

For the Cox regression model applied to right-censored failure time data, we observe $X = (W, \delta, Z)$, where $W = T \wedge C$, $\delta = 1\{W = T\}$, $Z \in \mathbb{R}^k$

is a regression covariate, T is a right-censored failure time with integrated hazard $e^{\beta'Z}\Lambda(t)$ given the covariate, and where C is a censoring time independent of T given Z . We also assume that censoring is uninformative of β or Λ . This is a special case of the general counting process regression model of the previous section, with $N(t) = \delta 1\{W \leq t\}$ and $Y(t) = 1\{W \geq t\}$. The consistency of $\hat{\beta}$, a zero of $U_n(\tau, \beta)$, and $\hat{\Lambda}$ both follow from the previous general results, as does the asymptotic normality and validity of the multiplier bootstrap based on the estimated influence function. There are, however, some interesting special features of the Cox model that are of interest, including the martingale structure, the limiting covariance, and the efficiency of the estimators.

First, it is not difficult to show that $M_{\beta_0}(t)$ and $U_n(t, \beta_0)$ are both continuous-time martingales (Fleming and Harrington, 1991; Andersen, Borgan, Gill and Keiding, 1993). This implies that

$$P \left\{ \int_0^\tau [Z - E(s, \beta_0)] dM_{\beta_0}(s) \right\}^{\otimes 2} = V(\beta_0),$$

and thus the asymptotic limiting variance of $\sqrt{n}(\hat{\beta} - \beta_0)$ is simply $V^{-1}(\beta_0)$. Thus $\hat{\beta}$ is efficient. In addition, the results of this section verify conditions (3.6) and (3.8) of theorem 3.1 for $\hat{\ell}_{\hat{\beta}, n}$ defined in (3.9). Verification of (3.7) is left as an exercise. This provides another proof of the efficiency of $\hat{\beta}$. What remains to be verified is that $\hat{\Lambda}$ is also efficient (in the uniform sense). The influence function for $\hat{\Lambda}$ is ψ given in (4.7). We will now use the methods of section 3.2 to verify that $\psi(u)$ is the efficient influence function for the parameter $\Lambda(u)$, for each $u \in [0, \tau]$. This will imply that $\hat{\Lambda}$ is uniformly efficient for Λ , since $\sqrt{n}(\hat{\Lambda} - \Lambda)$ converges weakly to a tight, mean zero Gaussian process.

The tangent space for the Cox model (for both parameters together) is $\{A(a, b) = \int_0^\tau [Z'a + b(s)] dM_\beta(s) : (a, b) \in H\}$, where $H = \mathbb{R}^k \times L_2(\Lambda)$. $A : H \mapsto L_2^0(P_{\beta, \Lambda})$ is thus a score operator for the full model. The natural inner product for pairs of elements in H is $\langle (a, b), (c, d) \rangle_H = a'b + \int_0^\tau c(s)d(s)d\Lambda(s)$, for $(a, b), (c, d) \in H$; and the natural inner product for $g, h \in L_2^0(P_{\beta, \Lambda})$ is $\langle g, h \rangle_{P_{\beta, \Lambda}} = P_{\beta, \Lambda}[gh]$. Hence the adjoint of A , A^* , satisfies $\langle A(a, b), g \rangle_{P_{\beta, \Lambda}} = \langle (a, b), A^*g \rangle_H$ for all $(a, b) \in H$ and all $g \in L_2^0(P_{\beta, \Lambda})$. It can be shown that $A^*g = (P_{\beta, \Lambda}[\int_0^\tau Z dM(s)g], P_{\beta, \Lambda}[dM(t)g]/d\Lambda_0(t))$ satisfies these equations and is thus the adjoint of A .

Let the one-dimensional submodel $\{P_t : t \in N_\epsilon\}$ be a perturbation in the direction $A(a, b)$, for some $(a, b) \in H$. In section 3.2, we showed that $\Lambda_t(u) = \int_0^u (1 + tb(s))d\Lambda(s) + o(t)$, and thus, for the parameter $\chi(P_{\beta, \Lambda}) = \Lambda(u)$, $\partial\chi(P_t)/(\partial t)|_{t=0} = \int_0^u b(s)d\Lambda(s)$. Thus $\langle \tilde{\chi}, (a, b) \rangle_H = \int_0^u b(s)d\Lambda(s)$, and therefore $\tilde{\chi} = (0, 1\{s \leq v\}) \in H$ (s is the function argument here). Our proof will be complete if we can show that $A^*\psi = \tilde{\chi}$, since this would imply that $\psi(u)$ is the efficient influence function for $\Lambda(u)$. First,

$$\begin{aligned}
A^*\psi(u) &= P_{\beta,\Lambda} \left\{ \int_0^\tau Z dM(s) \int_0^u \frac{dM(s)}{P_{\beta,\Lambda}[Y(s)e^{\beta'Z}]} \right\} \\
&\quad - P_{\beta,\Lambda} \left\{ \int_0^\tau Z dM(s) \int_0^\tau [Z - E(s, \beta)]' dM(s) \right\} \\
&\quad \times V^{-1}(\beta) \int_0^u E(s, \beta) d\Lambda(s) \\
&= \int_0^u E(s, \beta) d\Lambda(s) - V(\beta) V^{-1}(\beta) \int_0^u E(s, \beta) d\Lambda(s) \\
&= 0.
\end{aligned}$$

Second, it is not difficult to verify that $P_{\beta,\Lambda}[dM(s)\psi(u)] = 1\{s \leq u\}d\Lambda(s)$, and thus we obtain the desired result. Therefore, both parameter estimators $\hat{\beta}$ and $\hat{\Lambda}$ are uniformly efficient for estimating β and Λ in the Cox model.

4.3 The Kaplan-Meier Estimator

In this section, we consider the same right-censored failure time data considered in section 4.2.2, except that there is no regression covariate. The precise set-up is described in section 2.2.5. Thus T and C are assumed to be independent, where T has distribution function F and C has distribution function G . Assume also that $F(0) = 0$. We denote P_F to be the probability measure for the observed data $X = (W, \delta)$, and allow both F and G to have jumps. Define $S = 1 - F$, $L = 1 - G$, and $\pi(t) = P_F Y(t)$, and let $\tau \in (0, \infty)$ satisfy $F(\tau) > 0$ and $\pi(\tau) > 0$. Also define $\hat{\Lambda}(t) = \int_0^t [\mathbb{P}_n Y(s)]^{-1} \mathbb{P}_n dN(s)$. The Kaplan-Meier estimator \hat{S} has the product integral form $\hat{S}(t) = \prod_{0 < s \leq t} [1 - d\hat{\Lambda}(s)]$.

We have already established in section 2.2.5, using the self-consistency representation of \hat{S} , that \hat{S} is uniformly consistent for S over $[0, \tau]$ and that $\sqrt{n}[\hat{S} - S]$ converges weakly in $D[0, \tau]$ to a tight, mean zero Gaussian process. In this section, we verify that \hat{S} is also uniformly efficient. We first derive the influence function $\check{\psi}$ for \hat{S} , and then show that this satisfies the appropriate version of the adjoint formula (3.5) for estimating $S(u)$, for each $u \in [0, \tau]$. This pointwise efficiency will then imply uniform efficiency because of the weak convergence of $\sqrt{n}[\hat{S} - S]$.

Standard calculations (Fleming and Harrington, 1991, chapter 3) reveal that

$$\begin{aligned}
\hat{S}(u) - S(u) &= -\hat{S}(u) \int_0^u \frac{S(v-)}{\hat{S}(v)} \{d\hat{\Lambda}(v) - d\Lambda(v)\} \\
&= -\hat{S}(u) \int_0^u 1 \{\mathbb{P}_n Y(v) > 0\} \frac{S(v-)}{\hat{S}(v)} \left\{ \frac{\mathbb{P}_n dM(v)}{\mathbb{P}_n Y(v)} \right\} \\
&\quad + \hat{S}(u) \int_0^u 1 \{\mathbb{P}_n Y(v) = 0\} \frac{S(v-)}{\hat{S}(v)} d\Lambda(v) \\
&= A_n(u) + B_n(u),
\end{aligned}$$

where $M(v) = N(v) - \int_0^v Y(s)d\Lambda(s)$ is a martingale. Since $\pi(\tau) > 0$, $B_n(u) = o_p(n^{-1/2})$. Using martingale methods, it can be shown that $A_n(u) = \mathbb{P}_n \check{\psi}(u) + o_p(n^{1/2})$, where

$$\check{\psi}(u) = -S(u) \int_0^u \frac{1}{1 - \Delta\Lambda(v)} \left\{ \frac{dM(v)}{\pi(v)} \right\}.$$

Since all error terms are uniform for $u \in [0, \tau]$, we obtain that \hat{S} is asymptotically linear, with influence function $\check{\psi}$. As an element of $L_2^0(P_F)$, $\check{\psi}(u)(W, \delta) = \check{g}(W, \delta)$, where

$$\begin{aligned}
\check{g}(W, 1) &= -S(u) \left[\frac{1\{W \leq u\}}{[1 - \Delta\Lambda(W)] \pi(W)} - \int_0^u \frac{1\{W \geq s\} d\Lambda(s)}{[1 - \Delta\Lambda(s)] \pi(s)} \right] \text{ and} \\
\check{g}(W, 0) &= - \int_0^u \frac{1\{W \geq s\} d\Lambda(s)}{[1 - \Delta\Lambda(s)] \pi(s)}.
\end{aligned}$$

For each $h \in L_2^0(F)$, there exists a one-dimensional submodel $\{F_t : t \in N_\epsilon\}$, with $F_t(v) = \int_0^v (1 + th(s))dF(s) + o(t)$. This is clearly the maximal tangent set for F . This collection of submodels can be shown to generate the tangent set for the observed data model $\{P_F : F \in \mathcal{D}\}$, where \mathcal{D} is the collection of all failure time distribution functions with $F(0) = 0$, via the score operator $A : L_2^0(F) \mapsto L_2^0(P_F)$ defined by $(Ah)(W, \delta) = \delta h(W) + (1 - \delta) \int_W^\infty h(v)dF(v)/S(W)$. For $a, b \in L_2^0(F)$, let $\langle a, b \rangle_F = \int_0^\infty a(s)b(s)dF(s)$; and for $j, k \in L_2^0(P_F)$, let $\langle j, k \rangle_{P_F} = P_F[jk]$. Thus the adjoint of A must satisfy $\langle Ah, g \rangle_{P_F} = \langle h, A^*g \rangle_F$. Accordingly,

$$\begin{aligned}
P_F [(Ah)g] &= P_F \left[\delta h(W)g(W, 1) + (1 - \delta) \frac{\int_W^\infty h(v)dF(v)}{S(W)} g(W, 0) \right] \\
&= \int_0^\infty g(v, 1)L(v-)h(v)dF(v) + \int_0^\infty \frac{\int_w^\infty h(v)dF(v)}{S(w)} S(w)dG(w) \\
&= \int_0^\infty g(v, 1)L(v-)h(v)dF(v) - \int_0^\infty \int_{[v, \infty]} g(w, 0)dG(w)h(v)dF(v),
\end{aligned}$$

by the fact that $\int_s^\infty h(v)dF(v) = -\int_0^s h(v)dF(v)$ and by changing the order of integration on the right-hand-side. Thus $A^*g(v) = g(v, 1)L(v-) - \int_{[v, \infty]} g(s, 0)dG(s)$.

With $S_t(u) = 1 - F_t(u)$ based on a submodel perturbed in the direction $h \in L_2^0(F)$, we have that $\partial S_t(u)/(\partial t)|_{t=0} = \int_u^\infty h(v)dF(v) = -\int_0^u h(v)dF(v) = \langle \tilde{\chi}, h \rangle_F$, where $\tilde{\chi} = -1\{v \leq u\}$. We now verify that $A^*[\check{\psi}(u)] = \tilde{\chi}$, and thereby prove that $\hat{S}(u)$ is efficient for estimating the parameter $S(u)$, since it can be shown that $\check{\psi} \in R(A)$. We now have

(4.8)

$$\begin{aligned} (A^*[\check{\psi}(u)])(v) &= \check{\psi}(u)(v, 1)L(v-) - \int_{[v, \infty]} \check{\psi}(u)(s, 0)dG(s) \\ &= -\frac{S(u)}{S(v)}1\{v \leq u\} + S(u)L(v-) \int_0^u \frac{1\{v \geq s\}d\Lambda(s)}{[1 - \Delta\Lambda(s)]\pi(s)} \\ &\quad - \int_{[v, \infty]} \left\{ S(u) \int_0^u \frac{1\{s \geq r\}d\Lambda(r)}{[1 - \Delta\Lambda(r)]\pi(r)} \right\} dG(s). \end{aligned}$$

Since $\int_{[v, \infty]} 1\{s \geq r\}dG(s) = L([v \vee r]-) = 1\{v \geq r\}L(v-) + 1\{v < r\}L(r-)$, we now have that

$$\begin{aligned} (4.8) &= -\frac{S(u)}{S(v)}1\{v \leq u\} - S(u) \int_0^u \frac{1\{v < r\}L(r-)d\Lambda(r)}{[1 - \Delta\Lambda(r)]\pi(r)} \\ &= -1\{v \leq u\} \left[\frac{1}{S(v)} + \int_v^u \frac{dF(r)}{S(r)S(r-)} \right] S(u) \\ &= -1\{v \leq u\}. \end{aligned}$$

Thus (3.5) is satisfied, and we obtain the result that \hat{S} is pointwise and, therefore, uniformly efficient for estimating S .

4.4 Efficient Estimating Equations for Regression

We now consider a generalization of the conditionally mean zero residual linear regression model considered previously. A typical observation is assumed to be $X = (Y, Z)$, where $Y = g_\theta(Z) + e$, $E\{e|Z\} = 0$, $Z, \theta \in \mathbb{R}^k$, and $g_\theta(Z)$ is a known, sufficiently smooth function of θ . In addition to linear regression, generalized linear models—as well as many nonlinear regression models—fall into this structure. We assume that (Z, e) has a density η , and, therefore, that the observation (Y, Z) has density $\eta(y - g_\theta(z), z)$ with the only restriction being that $\int_{\mathbb{R}} e\eta(e, z)de = 0$. As we observed in section 3.2, these conditions force the score functions for η to be all square-integrable functions $a(e, z)$ which satisfy

$$E\{ea(e, Z)|Z\} = \frac{\int_{\mathbb{R}} ea(e, Z)\eta(e, Z)de}{\int_{\mathbb{R}} \eta(e, Z)de} = 0$$

almost surely. As also demonstrated in section 3.2, the above equality implies that the tangent space for η is the orthocomplement in $L_2^0(P_{\theta, \eta})$ of the set \mathcal{H} of all functions of the form $eh(Z)$, where $Eh^2(Z) < \infty$.

Hence, as pointed out previously, the efficient score for θ is obtained by projecting the ordinary score $\dot{\ell}_{\theta,\eta}(e, z) = -[\dot{\eta}_1(e, z)/\eta(e, z)]\dot{g}_\theta(z)$ onto \mathcal{H} , where $\dot{\eta}_1$ is the derivative with respect to the first argument of η and \dot{g}_θ is the derivative of g_θ with respect to θ . Of course, we are assuming that these derivatives exist and are square-integrable. Since the projection of an arbitrary $b(e, z)$ onto \mathcal{H} is $eE\{eb(e, Z)|Z\}/V(Z)$, where $V(z) \equiv E\{e^2|Z = z\}$, the efficient score for θ is

$$(4.9) \quad \tilde{\ell}_{\theta,\eta}(Y, Z) = -\frac{\dot{g}_\theta(Z)e \int_{\mathbb{R}} \dot{\eta}_1(u, Z)udu}{V(Z) \int_{\mathbb{R}} \eta u, Z du} = \frac{\dot{g}_\theta(Z)(Y - g_\theta(Z))}{V(Z)}.$$

This, of course, implies the result in section 3.2 for the special case of linear regression.

In practice, the form of $V(Z)$ is typically not known and needs to be estimated. Let this estimator be denoted $\hat{V}(Z)$. It can be shown that, even if \hat{V} is not consistent but converges to $\tilde{V} \neq V$, the estimating equation

$$\hat{S}_n(\theta) \equiv \mathbb{P}_n \left[\frac{\dot{g}_\theta(Z)(Y - g_\theta(Z))}{\hat{V}(Z)} \right]$$

is approximately equivalent to the estimating equation

$$\tilde{S}_n(\theta) \equiv \mathbb{P}_n \left[\frac{\dot{g}_\theta(Z)(Y - g_\theta(Z))}{\tilde{V}(Z)} \right],$$

both of which can still yield \sqrt{n} consistent estimators of θ . The closer \tilde{V} is to V , the more efficient will be the estimator based on solving $\tilde{S}_n(\theta) = 0$. Another variant of the question of optimality is “for what choice of \tilde{V} will the estimator obtained by solving $\tilde{S}_n(\theta) = 0$ yield the smallest possible variance?” Godambe (1960) showed that, for univariate θ , the answer is $\tilde{V} = V$. This result is not based on semiparametric efficiency analysis, but is obtained from minimizing the limiting variance of the estimator over all “reasonable” choices of \tilde{V} .

For any real, measurable function w of $z \in \mathbb{R}^k$ (measurable on the probability space for Z), let

$$S_{n,w}(\theta) \equiv \mathbb{P}_n [\dot{g}_\theta(Z)(Y - g_\theta(Z))w(Z)/V(Z)],$$

$U(z) \equiv \dot{g}(z)\dot{g}'(z)/V(z)$, and let $\hat{\theta}_{n,w}$ be a solution of $S_{n,w}(\theta) = 0$. Standard methods can be used to show that if both $E[U(Z)]$ and $E[U(Z)w(Z)]$ are positive definite, then the limiting variance of $\sqrt{n}(\hat{\theta}_{n,w} - \theta)$ is

$$\{E[U(Z)w(Z)]\}^{-1} \{E[U(Z)w^2(Z)]\} \{E[U(Z)w(Z)]\}^{-1}.$$

The following proposition yields Godambe’s (1960) result generalized for arbitrary $k \geq 1$:

PROPOSITION 4.3 *Assume $E[U(Z)]$ is positive definite. Then, for any real, Z -measurable function w for which $E[U(Z)w(Z)]$ is positive definite,*

$$C_{0,w} \equiv \{E[U(Z)w(Z)]\}^{-1} E[U(Z)w^2(Z)] \{E[U(Z)w(Z)]\}^{-1} - \{E[U(Z)]\}^{-1}$$

is positive semidefinite.

Proof. Define $B(z) \equiv \{E[U(Z)w(Z)]\}^{-1} w(z) - \{E[U(Z)]\}^{-1}$, and note that $C_w(Z) \equiv B(Z)E[U(Z)]B'(Z)$ must therefore be positive semidefinite. The desired result now follows since $E[C_w(Z)] = C_{0,w}$. \square

Note that when $A - B$ is positive semidefinite, for two $k \times k$ variance matrices A and B , we know that B is the smaller variance. This follows since, for any $v \in \mathbb{R}^k$, $v'Av \geq v'Bv$. Thus the choice $w = 1$ will yield the minimum variance, or, in other words, the choice $\tilde{V} = V$ will yield the lowest possible variance for estimators asymptotically equivalent to the solution of $\tilde{S}_n(\theta) = 0$.

We now verify that estimation based on solving $\hat{S}_n(\theta) = 0$ is asymptotically equivalent to estimation based on solving $\tilde{S}_n(\theta) = 0$, under reasonable regularity conditions. Assume that $\hat{\theta}$ satisfies $\hat{S}_n(\hat{\theta}) = o_p(n^{-1/2})$ and that $\hat{\theta} \xrightarrow{P} \theta$. Assume also that for every $\epsilon > 0$, there exists a P -Donsker class \mathcal{G} such that the inner probability that $\hat{V} \in \mathcal{G}$ is $> 1 - \epsilon$ for all n large enough and all $\epsilon > 0$, and that for some $\tau > 0$, the class

$$\mathcal{F}_1 = \left\{ \frac{\dot{g}_{\theta_1}(Z)(Y - g_{\theta_2}(Z))}{W(Z)} : \|\theta_1 - \theta\| \leq \tau, \|\theta_2 - \theta\| \leq \tau, W \in \mathcal{G} \right\}$$

is P -Donsker, and the class

$$\mathcal{F}_2 = \left\{ \frac{\dot{g}_{\theta_1}(Z)\dot{g}'_{\theta_2}(Z)}{W(Z)} : \|\theta_1 - \theta\| \leq \tau, \|\theta_2 - \theta\| \leq \tau, W \in \mathcal{G} \right\}$$

is P -Glivenko-Cantelli. We also need that

$$(4.10) \quad P \left[\frac{\dot{g}_{\hat{\theta}}(Z)}{\hat{V}(Z)} - \frac{\dot{g}_{\theta}(Z)}{\tilde{V}(Z)} \right]^2 \xrightarrow{P} 0,$$

and, for any $\tilde{\theta} \xrightarrow{P} \theta$,

$$(4.11) \quad \left| P \frac{\dot{g}_{\tilde{\theta}}(Z)\dot{g}'_{\tilde{\theta}}(Z)}{\hat{V}(Z)} - P \frac{\dot{g}_{\theta}(Z)\dot{g}'_{\theta}(Z)}{\tilde{V}(Z)} \right| \xrightarrow{P} 0.$$

We have the following lemma:

LEMMA 4.4 *Assume that $E[Y|Z = z] = g_{\theta}(z)$, that*

$$U_0 \equiv E \left[\frac{\dot{g}_{\theta}(Z)\dot{g}'_{\theta}(Z)}{V(Z)} \right]$$

is positive definite, that $\hat{\theta}$ satisfies $\hat{S}_n(\hat{\theta}) = o_p(n^{-1/2})$, and that $\hat{\theta} \xrightarrow{P} \theta$. Suppose also that \mathcal{F}_1 is P -Donsker, that \mathcal{F}_2 is P -Glivenko-Cantelli, and that both (4.10) and (4.11) hold for any $\hat{\theta} \xrightarrow{P} \theta$. Then $\sqrt{n}(\hat{\theta} - \theta)$ is asymptotically normal with variance

$$U_0^{-1} \mathbb{E} \left[\frac{\dot{g}\theta(Z)\dot{g}'_{\theta}(Z)V(Z)}{\tilde{V}^2(Z)} \right] U_0^{-1}.$$

Moreover, if $\tilde{V} = V$, Z -almost surely, then $\hat{\theta}$ is optimal in the sense of proposition 4.3.

Proof. For any $h \in \mathbb{R}^k$,

$$\begin{aligned} o_p(n^{-1/2}) &= h' \mathbb{P}_n \left[\frac{\dot{g}_{\hat{\theta}}(Z)(Y - g_{\hat{\theta}}(Z))}{\hat{V}(Z)} \right] = h' \mathbb{P}_n \left[\frac{\dot{g}_{\theta}(Z)(Y - g_{\theta}(Z))}{\tilde{V}(Z)} \right] \\ &\quad + h' \mathbb{P}_n \left[\left\{ \frac{\dot{g}_{\hat{\theta}}(Z)}{\hat{V}(Z)} - \frac{\dot{g}_{\theta}(Z)}{\tilde{V}(Z)} \right\} (Y - g_{\theta}(Z)) \right] \\ &\quad - h' \mathbb{P}_n \left[\frac{\dot{g}_{\hat{\theta}}(Z)\dot{g}'_{\hat{\theta}}(Z)}{\hat{V}(Z)} (\hat{\theta} - \theta) \right], \end{aligned}$$

where $\tilde{\theta}$ is on the line segment between θ and $\hat{\theta}$. Now the conditions of the theorem can be seen to imply that

$$\sqrt{n}(\hat{\theta} - \theta) = \sqrt{n} U_0^{-1} \mathbb{P}_n \left\{ \frac{\dot{g}_{\theta}(Z)(Y - g_{\theta}(Z))}{\tilde{V}(Z)} \right\} + o_p(1),$$

and the first conclusion of the lemma follows. The final conclusion now follows directly from proposition 4.3. \square

The conditions of lemma 4.4 are easily satisfied for a variety of regression models and variance estimators \hat{V} . For example, if $g_{\theta}(Z) = (1 + e^{-\theta'Z})^{-1}$ is the conditional expectation of a Bernoulli outcome Y , given Z , and both θ and Z are assumed to be bounded, then all the conditions are easily satisfied with $\hat{\theta}$ being a zero of \hat{S}_n and $\hat{V}(Z) = g_{\hat{\theta}}(Z) [1 - g_{\hat{\theta}}(Z)]$, provided $\mathbb{E}[ZZ']$ is positive definite. Note that for this example, $\hat{\theta}$ is also semiparametric efficient. We now consider two additional examples in some detail. The first example is a special case of the semiparametric model discussed at the beginning of this section. The model is simple linear regression but with an unspecified form for $V(Z)$. Note that for general g_{θ} , if the conditions of lemma 4.4 are satisfied for $\tilde{V} = V$, then the estimator $\hat{\theta}$ is semiparametric efficient by the form of the efficient score given in (4.9). The second example considers estimation of a semiparametric Poisson mixture regression model, where the mixture induces extra-Poisson variation. We will develop an optimal estimating equation procedure in the sense of proposition 4.3. Unfortunately, it is unclear in this instance how to strengthen this result to obtain semiparametric efficiency.

4.4.1 Simple Linear Regression

Consider simple linear regression based on a univariate Z . Let $\theta = (\alpha, \beta)' \in \mathbb{R}^2$ and assume $g_\theta(Z) = \alpha + \beta Z$. We also assume that the support of Z is a known compact interval $[a, b]$. We will use a modified kernel method of estimating $V(z) \equiv E[e^2 | Z = z]$. Let the kernel $L : \mathbb{R} \mapsto [0, 1]$ satisfy $L(x) = 0$ for all $|x| > 1$ and $L(x) = 1 - |x|$ otherwise, and let $h \leq (b - a)/2$ be the bandwidth for this kernel. For the sample $(Y_1, Z_1), \dots, (Y_n, Z_n)$, let $\hat{\theta} = (\hat{\alpha}, \hat{\beta})'$ be the usual least-squares estimator of θ , and let $\hat{e}_i = Y_i - \hat{\alpha} - \hat{\beta}Z_i$ be the estimated residuals. Let $\hat{F}_n(u)$ be the empirical distribution of Z_1, \dots, Z_n , and let $\hat{H}_n(u) \equiv n^{-1} \sum_{i=1}^n \hat{e}_i^2 1\{Z_i \leq u\}$. We will denote F as the true distribution of Z , with density f , and also define $H(z) \equiv \int_a^z V(u) dF(u)$. Now define

$$\hat{V}(z) \equiv \frac{\int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) \hat{H}_n(du)}{\int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) \hat{F}_n(du)},$$

for $z \in [a + h, b - h]$, and let $\hat{V}(z) = \hat{V}(a + h)$ for all $z \in [a, a + h)$ and $\hat{V}(z) = \hat{V}(b - h)$ for all $z \in (b - h, b]$.

We need to assume in addition that both F and V are twice differentiable with second derivatives uniformly bounded on $[a, b]$, that for some $M < \infty$ we have $M^{-1} \leq f(z), V(z) \leq M$ for all $a \leq x \leq b$, and that the possibly data-dependent bandwidth satisfies $h = o_P(1)$ and $h^{-1} = o_P(n^{1/4})$. If we let $U_i \equiv (1, Z_i)'$, $i = 1, \dots, n$, then $\hat{H}_n(z) =$

$$\begin{aligned} n^{-1} \sum_{i=1}^n \hat{e}_i^2 1\{Z_i \leq z\} &= n^{-1} \sum_{i=1}^n \left[e_i - (\hat{\theta} - \theta)' U_i \right]^2 1\{Z_i \leq z\} \\ &= n^{-1} \sum_{i=1}^n e_i^2 1\{Z_i \leq z\} \\ &\quad - 2(\hat{\theta} - \theta)' n^{-1} \sum_{i=1}^n U_i e_i 1\{Z_i \leq z\} \\ &\quad + (\hat{\theta} - \theta)' n^{-1} \sum_{i=1}^n U_i U_i' 1\{Z_i \leq z\} (\hat{\theta} - \theta) \\ &= A_n(z) - B_n(z) + C_n(z). \end{aligned}$$

In chapter 9, we will show in an exercise that $\mathcal{G}_1 \equiv \{Ue \cdot 1\{Z \leq z\}, z \in [a, b]\}$ is Donsker. Hence $\|\mathbb{P}_n - P\|_{\mathcal{G}_1} = O_P(n^{-1/2})$. Since also $E[e|Z] = 0$ and $\|\hat{\theta} - \theta\| = O_P(n^{-1/2})$, we now have that $\sup_{z \in [a, b]} |B_n(z)| = O_P(n^{-1})$. By noting that $\|U\|$ is bounded under our assumptions, we also obtain that $\sup_{z \in [a, b]} |C_n(z)| = O_P(n^{-1})$. In another exercise in chapter 9, we will verify that $\mathcal{G}_2 \equiv \{e^2 \cdot 1\{Z \leq z\}, z \in [a, b]\}$ is also Donsker. Hence $\sup_{z \in [a, b]} |A_n(z) - H(z)| = O_P(n^{-1/2})$, and thus also $\sup_{z \in [a, b]} |\hat{H}_n(z) -$

$H(z) = O_P(n^{-1/2})$. Standard results also verify that $\sup_{z \in [a, b]} |\hat{F}_n(z) - F(z)| = O_P(n^{-1/2})$.

Let $D_n \equiv \hat{H}_n - H$, \dot{L} be the derivative of L , and note that for any $z \in [a + h, b - h]$ we have by integration by parts and by the form of \dot{L} that

$$\int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) D_n(du) = - \int_{z-h}^{z+h} D_n(u) h^{-2} \dot{L}\left(\frac{z-u}{h}\right) du.$$

Thus

$$\left| \int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) D_n(du) \right| \leq h^{-1} \sup_{z \in [a, b]} |\hat{H}_n(z) - H(z)|.$$

Since the right-hand-side does not depend on z , and by the result of the previous paragraph, we obtain that the supremum of the left-hand-side over $z \in [a + h, b - h]$ is $O_P(h^{-1}n^{-1/2})$. Letting \dot{H} be the derivative of H , we save it as an exercise to verify that both

$$\sup_{z \in [a+h, b-h]} \left| \int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) H(du) - \dot{H}(z) \right| = O(h)$$

and $\left(\sup_{z \in [a, a+h]} |\dot{H}(z) - \dot{H}(a+h)| \right) \vee \left(\sup_{z \in (b-h, b]} |\dot{H}(z) - \dot{H}(b-h)| \right) = O(h)$. Hence

$$\hat{R}_n(z) \equiv \begin{cases} \int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) \hat{H}_n(du) & \text{for } z \in [a+h, b-h], \\ \hat{R}_n(a+h) & \text{for } z \in [a, a+h), \\ \hat{R}_n(b-h) & \text{for } z \in (b-h, b] \end{cases}$$

is uniformly consistent for \dot{H} with uniform error $O_P(h + h^{-1}n^{-1/2}) = o_P(1)$.

Similar, but somewhat simpler arguments compared to those in the previous paragraph, can also be used to verify that

$$\hat{Q}_n(z) \equiv \begin{cases} \int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) \hat{F}_n(du) & \text{for } z \in [a+h, b-h], \\ \hat{Q}_n(a+h) & \text{for } z \in [a, a+h), \\ \hat{Q}_n(b-h) & \text{for } z \in (b-h, b] \end{cases}$$

is uniformly consistent for f also with uniform error $O_P(h + h^{-1}n^{-1/2}) = o_P(1)$. Since f is bounded below, we now have that $\sup_{z \in [a, b]} |\hat{V}(z) - V(z)| = o_P(1)$. Since $\dot{g}_\theta(Z) = (1, Z)'$, we have established that both (4.10) and (4.11) hold for this example since V is also bounded below.

We will now show that there exists a $k_0 < \infty$ such that the probability that “ \hat{V} goes below $1/k_0$ or the first derivative of \hat{V} exceeds k_0 ” goes to

zero as $n \rightarrow \infty$. If we let \mathcal{G} be the class of all functions $q : [a, b] \mapsto [k_0^{-1}, k_0]$ such that the first derivative of q , \dot{q} , satisfies $|\dot{q}| \leq k_0$, then our result will imply that the inner probability that $\hat{V} \in \mathcal{G}$ is $> 1 - \epsilon$ for all n large enough and all $\epsilon > 0$. An additional exercise in chapter 9 will then show that, for this simple linear regression example, the class \mathcal{F}_1 defined above is Donsker and \mathcal{F}_2 defined above is Glivenko-Cantelli. Thus lemma 4.4 applies. This means that if we first use least-squares estimators of α and β to construct the estimator \hat{V} , and then compute the “two-stage” estimator

$$\tilde{\theta} \equiv \left[\sum_{i=1}^n \frac{U_i U_i'}{\hat{V}(Z_i)} \right]^{-1} \sum_{i=1}^n \frac{U_i Y_i}{\hat{V}(Z_i)},$$

then this $\tilde{\theta}$ will be efficient for θ .

Since \hat{V} is uniformly consistent, the only thing remaining to show is that the derivative of \hat{V} , denoted \hat{V}_n , is uniformly bounded. Note that the derivative of \hat{R}_n , which we will denote \hat{R}_n , satisfies the following for all $z \in [a+h, b-h]$:

$$\begin{aligned} \hat{R}_n(z) &= - \int_{\mathbb{R}} h^{-2} \dot{L} \left(\frac{z-u}{h} \right) \hat{H}_n(du) \\ &= h^{-2} \left[\hat{H}_n(z) - \hat{H}_n(z-h) - \hat{H}_n(z+h) + \hat{H}_n(z) \right] \\ &= O_P(h^{-2} n^{-1/2}) + h^{-2} [H(z) - H(z-h) - H(z+h) + H(z)], \end{aligned}$$

where the last equality follows from the previously established fact that $\sup_{z \in [a+h, b-h]} |\hat{H}_n(z) - H(z)| = O_P(n^{-1/2})$. Now the uniform boundedness of the second derivative of H ensures that $\sup_{z \in [a+h, b-h]} |\hat{H}_n(z)| = O_P(1)$. Similar arguments can be used to establish that the derivative of \hat{Q}_n , which we will denote \hat{Q}_n , satisfies $\sup_{z \in [a+h, b-h]} |\hat{Q}_n(z)| = O_P(1)$. Now we have uniform boundedness in probability of \hat{V}_n over $[a+h, b-h]$. Since \hat{V}_n does not change over either $[a, a+h]$ or $(b-h, b]$, we have also established that $\sup_{z \in [a, b]} |\hat{V}_n(z)| = O_P(1)$, and the desired results follow.

4.4.2 A Poisson Mixture Regression Model

In this section, we consider a Poisson mixture regression model in which the nuisance parameter is not \sqrt{n} consistent in the uniform norm. Given a regression vector $Z \in \mathbb{R}^k$ and a nonnegative random quantity $W \in \mathbb{R}$, the observation Y is Poisson with parameter $W e^{\beta' Z}$, for some $\beta \in \mathbb{R}^k$. We only observe the pair (Y, Z) . Thus the density of Y given Z is

$$Q_{\beta, G}(y) = \int_0^\infty \frac{e^{-w e^{\beta' Z}} [w e^{\beta' Z}]^y}{y!} dG(w),$$

where G is the unknown distribution function for W . For identifiability, we assume that one of the components of Z is the constant 1 and that the expectation of W is also 1. We denote the joint distribution of the observed data as $P_{\beta,G}$, and assume that Z is bounded, $P_{\beta,G}[ZZ']$ is full rank, and that $P_{\beta,G}Y^2 < \infty$.

In this situation, it is not hard to verify that

$$\begin{aligned} V(z) &= \mathbb{E} \left\{ \left[Y - e^{\beta'Z} \right]^2 \middle| Z = z \right\} \\ &= \mathbb{E} \left[\mathbb{E} \left\{ \left[Y - Ue^{\beta'Z} \right]^2 \middle| U, Z = z \right\} + (U - 1)^2 e^{2\beta'z} \right] \\ &= e^{\beta'z} + \sigma^2 e^{2\beta'z}, \end{aligned}$$

where $\sigma^2 \equiv \mathbb{E}[U - 1]^2$. Let $\tilde{\beta}$ be the estimator obtained by solving

$$\mathbb{P}_n \left[Z(Y - e^{\beta'Z}) \right] = 0.$$

Standard arguments reveal that $\sqrt{n}(\tilde{\beta} - \beta)$ is asymptotically mean zero Gaussian with finite variance matrix. Relatively simple calculations also reveal that $\mathbb{E}[Y(Y - 1)|Z = z] = \int_0^\infty u^2 dG(u)e^{2\beta'Z} = (\sigma^2 + 1)e^{2\beta'Z}$. Hence $\hat{\sigma}^2 \equiv -1 + n^{-1} \sum_{i=1}^n e^{-2\tilde{\beta}'Z_i} Y_i(Y_i - 1)$ will satisfy $\hat{\sigma}^2 = \sigma^2 + O_P(n^{-1/2})$. Now let $\hat{\beta}$ be the solution of

$$\mathbb{P}_n \left[\frac{Z(Y - e^{\beta'Z})}{\hat{V}(Z)} \right] = 0,$$

where $\hat{V}(z) \equiv e^{\tilde{\beta}'z} + \hat{\sigma}^2 e^{2\tilde{\beta}'z}$. It is left as an exercise to verify that $\hat{\beta}$ satisfies the conditions of lemma 4.4 for $\tilde{V} = V$, and thus the desired optimality is achieved.

4.5 Partly Linear Logistic Regression

For the partly linear logistic regression example given in chapter 1, the observed data are n independent realizations of the random triplet (Y, Z, U) , where $Z \in \mathbb{R}^p$ and $U \in \mathbb{R}$ are covariates which are not linearly dependent, Y is a dichotomous outcome with conditional expectation $\nu[\beta'Z + \eta(U)]$, $\beta \in \mathbb{R}^p$, Z is restricted to a bounded set, $U \in [0, 1]$, $\nu(t) = 1/(1 + e^{-t})$, and where η is an unknown smooth function. Hereafter, for simplicity, we will also assume that $p = 1$. We further assume, for some integer $k \geq 1$, that the first $k - 1$ derivatives of η exist and are absolutely continuous with $J^2(\eta) = \int_0^1 [\eta^{(k)}(t)]^2 dt < \infty$. To estimate β and η based on an i.i.d. sample $X_i = (Y_i, Z_i, U_i)$, $i = 1, \dots, n$, we use the following penalized log-likelihood:

$$\tilde{L}_n(\beta, \eta) = n^{-1} \sum_{i=1}^n \log p_{\beta, \eta}(X_i) - \hat{\lambda}_n^2 J^2(\eta),$$

where

$$p_{\beta, \eta}(x) = \{\nu [\beta z + \eta(u)]\}^y \{1 - \nu [\beta z + \eta(u)]\}^{1-y}$$

and $\hat{\lambda}_n$ is chosen to satisfy $\hat{\lambda}_n = o_p(n^{-1/4})$ and $\hat{\lambda}_n^{-1} = O_p(n^{k/(2k+1)})$. Denote $\hat{\beta}_n$ and $\hat{\eta}_n$ to be the maximizers of $\tilde{L}_n(\beta, \eta)$, let $P_{\beta, \eta}$ denote expectation under the model, and let β_0 and η_0 to be the true values of the parameters.

Consistency of $\hat{\beta}_n$ and $\hat{\eta}_n$ and efficiency of $\hat{\beta}_n$ are established for partly linear generalized linear models in Mammen and van de Geer (1997). We now derive the efficient score for β and then sketch a verification that $\hat{\beta}_n$ is asymptotically linear with influence function equal to the efficient influence function. Let \mathcal{H} be the linear space of functions $h : [0, 1] \mapsto \mathbb{R}$ with $J(h) < \infty$. For $t \in [0, \epsilon)$ and ϵ sufficiently small, let $\beta_t = \beta + tv$ and $\eta_t = \eta + th$ for $v \in \mathbb{R}$ and $h \in \mathcal{H}$. If we differentiate the non-penalized log-likelihood, we deduce that the score for β and η , in the direction (v, h) , is $(vZ + h(U))(Y - \mu_{\beta, \eta}(Z, U))$, where $\mu_{\beta, \eta}(Z, U) = \nu[\beta Z + \eta(U)]$. Now let

$$h_1(u) = \frac{\mathbb{E}\{ZV_{\beta, \eta}(Z, U)|U = u\}}{\mathbb{E}\{V_{\beta, \eta}(Z, U)|U = u\}},$$

where $V_{\beta, \eta} = \mu_{\beta, \eta}(1 - \mu_{\beta, \eta})$, and assume that $h_1 \in \mathcal{H}$. It can easily be verified that $Z - h_1(U)$ is uncorrelated with any $h(U)$, $h \in \mathcal{H}$, and thus the efficient score for β is $\tilde{\ell}_{\beta, \eta}(Z, U) = (Z - h_1(U))(Y - \mu(Z, U))$. Hence the efficient information for β is $\tilde{I}_{\beta, \eta} = P_{\beta, \eta} [(Z - h_1(U))^2 V_{\beta, \eta}(Z, U)]$ and the efficient influence function is $\tilde{\psi}_{\beta, \eta} = \tilde{I}_{\beta, \eta}^{-1} \tilde{\ell}_{\beta, \eta}$, provided $\tilde{I}_{\beta, \eta} > 0$, which we assume hereafter to be true for $\beta = \beta_0$ and $\eta = \eta_0$.

In order to prove asymptotic linearity of $\hat{\beta}_n$, we need to also assume that $P_{\beta_0, \eta_0} [Z - \tilde{h}_1(U)]^2 > 0$, where $\tilde{h}_1(u) = \mathbb{E}\{Z|U = u\}$. First, Mammen and van de Geer established that $\hat{\beta}_n$ and $\hat{\eta}_n$ are both uniformly consistent for β_0 and η_0 , respectively, and that

$$(4.12) \quad \mathbb{P}_n \left[(\hat{\beta}_n - \beta_0)Z + \hat{\eta}_n(U) - \eta_0(U) \right]^2 = o_p(n^{-1/2}).$$

Let $\hat{\beta}_{ns} = \hat{\beta}_n + s$ and $\hat{\eta}_{ns}(u) = \hat{\eta}_n(u) - sh_1(u)$. If we now differentiate $\tilde{L}_n(\hat{\beta}_{ns}, \hat{\eta}_{ns})$ and evaluate at $s = 0$, we obtain

$$\begin{aligned} 0 &= \mathbb{P}_n [(Y - \mu_{\beta_0, \eta_0})(Z - h_1(U))] \\ &\quad - \mathbb{P}_n \left[(\mu_{\hat{\beta}_n, \hat{\eta}_n} - \mu_{\beta_0, \eta_0})(Z - h_1(U)) \right] - \lambda_n^2 \{ \partial J^2(\hat{\eta}_{ns}) / (\partial s) |_{s=0} \} \\ &= A_n - B_n - C_n, \end{aligned}$$

since $\tilde{L}_n(\beta, \eta)$ is maximized at $\hat{\beta}_n$ and $\hat{\eta}_n$ by definition.

Using (4.12), we obtain

$$B_n = \mathbb{P}_n \left[V_{\beta_0, \eta_0}(Z, U) \left\{ (\hat{\beta}_n - \beta_0)Z + \hat{\eta}_n(U) - \eta_0(U) \right\} (Z - h_1(U)) \right] + o_p(n^{-1/2}),$$

since $\partial\nu(t)/(\partial t) = t(1-t)$. By definition of h_1 , we have

$$P_{\beta_0, \eta_0} [V_{\beta_0, \eta_0}(Z, U)(\hat{\eta}_n(U) - \eta_0(U))(Z - h_1(U))] = 0,$$

and we also have that $P_{\beta_0, \eta_0} [V_{\beta_0, \eta_0}(Z, U)(\hat{\eta}_n(U) - \eta_0(U))(Z - h_1(U))]^2 \xrightarrow{P} 0$. Thus, if we can establish that for each $\tau > 0$, $\hat{\eta}_n(U) - \eta_0(U)$ lies in a bounded P_{β_0, η_0} -Donsker class with probability $> (1 - \tau)$ for all $n \geq 1$ large enough and all $\tau > 0$, then

$$\mathbb{P}_n [V_{\beta_0, \eta_0}(Z, U)(\hat{\eta}_n(U) - \eta_0(U))(Z - h_1(U))] = o_p(n^{-1/2}),$$

and thus

$$(4.13) \quad B_n = (\hat{\beta}_n - \beta_0)P_{\beta_0, \eta_0} [V_{\beta_0, \eta_0}(Z, U)(Z - h_1(U))^2] + o_p(n^{-1/2}),$$

since products of bounded Donsker classes are Donsker (and therefore also Glivenko-Cantelli), and since

$$P_{\beta_0, \eta_0} [V_{\beta_0, \eta_0}(Z, U)Z(Z - h_1(U))] = P_{\beta_0, \eta_0} [V_{\beta_0, \eta_0}(Z, U)(Z - h_1(U))^2]$$

by definition of h_1 . Let \mathcal{H}_c be the subset of \mathcal{H} with functions h satisfying $J(h) \leq c$. We will show in part II that $\{h(U) : h \in \mathcal{H}_c\}$ is indeed Donsker for each $c < \infty$. Since Mammen and van de Geer verify that $J(\hat{\eta}_n) = O_p(1)$, we have the desired Donsker property, and (4.13) follows.

It is not difficult to verify that $C_n \leq 2\lambda_n^2 J(\hat{\eta}_n)J(h_1) = o_p(n^{-1/2})$, since $\lambda_n = o_p(n^{-1/4})$ by assumption. Hence $\sqrt{n}(\hat{\beta}_n - \beta_0) = \sqrt{n}\mathbb{P}_n \tilde{\psi}_{\beta_0, \eta_0} + o_p(n^{-1/2})$, and we have verified that $\hat{\beta}_n$ is efficient for β_0 .

4.6 Exercises

4.6.1. For the linear regression example, verify the inequality (4.3).

4.6.2. For the general counting process regression model setting, show that $V(\beta) \geq c \text{var}(Z)$, for some $c > 0$, where $V(\beta)$ is given in (4.5).

4.6.3. Show how to use lemma 4.2 to establish (4.6). Hint: Let $A_n(t) = E_n(t, \hat{\beta}) - E(t, \hat{\beta})$ and $B_n(t) = \sqrt{n}\mathbb{P}_n M_{\beta_0}(t)$ for part of it, and let

$$A_n(t) = \mathbb{P}_n \left\{ \int_0^t Y(s) \left[e^{\hat{\beta}'Z} - e^{\beta_0'Z} \right] d\Lambda_0(s) \right\}$$

and $B_n(t) = \sqrt{n} \left[E_n(t, \hat{\beta}) - E(t, \hat{\beta}) \right]$ for the other part.

4.6.4. For the Cox model example (section 4.2.2), verify condition (3.7) of theorem 3.1 for $\hat{\ell}_{\hat{\beta},n}$ defined in (3.9).

4.6.5. For the Kaplan-Meier example of section 4.3, verify that $\check{\psi} \in R(A)$.

4.6.6. For the Poisson mixture model example of section 4.4.2, verify that the conditions of lemma 4.4 are satisfied for the given choice of \hat{V} :

(a) Show that the class $\mathcal{G} \equiv \left\{ e^{t'Z} + s^2 e^{2t'Z} : \|t - \beta\| \leq \epsilon_1, |s^2 - \sigma^2| \leq \epsilon_2 \right\}$ is Donsker for some $\epsilon_1, \epsilon_2 > 0$. Hint: First show that $\{t'Z : \|t - \beta\| \leq \epsilon_1\}$ is Donsker from the fact that the product of two (in this case trivial) bounded Donsker classes is also Donsker. Now complete the proof by using the facts that Lipschitz functions of Donsker classes are Donsker, that products of bounded Donsker classes are Donsker (used earlier), and that sums of Donsker classes are Donsker.

(b) Now complete the verification of lemma 4.4.

4.6.7. For the partly linear logistic regression example of section 4.5, verify that $(Z - h_1(U))(Y - \mu_{\beta_0, \eta_0})$ is uncorrelated with $h(U)(Y - \mu_{\beta_0, \eta_0})$ for all $h \in \mathcal{H}$, where the quantities are as defined in the example.

4.7 Notes

Least absolute deviation regression was studied using equicontinuity arguments in Bassett and Koenker (1978). More succinct results based on empirical processes can be found in Pollard (1991). An excellent discussion of breakdown points and other robustness issues is given in Huber (1981).

Part II

Empirical Processes

5

Introduction to Empirical Processes

The goal of part II is to provide an in depth coverage of the basics of empirical process techniques which are useful in statistics. Chapter 6 presents preliminary mathematical background which provides a foundation for later technical development. The topics covered include metric spaces, outer expectations, linear operators and functional differentiation. The main topics overviewed in chapter 2 of part I will then be covered in greater depth, along with several additional topics, in chapters 7 through 14. Part II finishes in chapter 15 with several case studies. The main approach is to present the mathematical and statistical ideas in a logical, linear progression, and then to illustrate the application and integration of these ideas in the case study examples. The scaffolding provided by the overview, part I, should enable the reader to maintain perspective during the sometimes rigorous developments of this section.

Stochastic convergence is studied in chapter 6. An important aspect of the modes of convergence explored in this book are the notions of outer integrals and outer measure which were mentioned briefly in section 2.2.1. While many of the standard relationships between the modes of stochastic convergence apply when using outer measure, there are a few important differences which we will examine. While these differences may, in some cases, add complexity to an already difficult asymptotic theory, the gain in breadth of applicability to semiparametric statistical estimators is well worth the trouble. For example, convergence based on outer measure permits the use of the uniform topology for studying convergence of empirical processes with complex index sets. This contrasts with more traditional approaches which require special topologies that can be harder to use in

applications, such as the Skorohod topology for cadlag processes (see chapter 3 of Billingsley, 1968).

The main techniques for proving empirical process central limit theorems will be presented in chapter 8. Establishing Glivenko-Cantelli and Donsker theorems requires bounding expectations involving suprema of stochastic processes. Maximal inequalities and symmetrization techniques are important tools for accomplishing this, and careful measurability arguments are also sometimes needed. Symmetrization involves replacing inequalities for the empirical process $f \mapsto (\mathbb{P}_n - P)f$, $f \in \mathcal{F}$, with inequalities for the “symmetrized” process $n^{-1} \sum_{i=1}^n \epsilon_i f(X_i)$, where $\epsilon_1, \dots, \epsilon_n$ are i.i.d. Rademacher random variables (eg., $P\{\epsilon_1 = -1\} = P\{\epsilon_1 = 1\} = 1/2$) independent of X_1, \dots, X_n . Several tools for assessing measurability in statistical applications will also be discussed.

Entropy with bracketing, uniform entropy, and other measures of entropy are essential aspects in all of these results. This is the topic of chapter 9. The associated entropy calculations can be quite challenging, but the work is often greatly simplified by using Donsker preservation results to build larger Donsker classes from smaller ones. Similar preservation results are available for Glivenko-Cantelli classes.

Bootstrapping of empirical processes, based on multinomial or other Monte Carlo weights, is studied in chapter 10. The bootstrap is a valuable way to conduct inference for empirical processes because of its broad applicability. In many semiparametric settings, there are no viable alternatives for inference. A central role in establishing validity of the bootstrap is played by multiplier central limit theorems which establish weak convergence of processes of the form $\sqrt{n}\mathbb{P}_n\xi(f(X) - Pf)$, $f \in \mathcal{F}$, where ξ is independent of X , has mean zero and variance 1, and $\int_0^\infty \sqrt{P(|\xi| > x)}dx < \infty$.

In chapter 11, several extensions of empirical process results are presented for function classes which either consist of sequences or change with the sample size n , as well as results for independent but not identically distributed data. These results are useful in a number of statistical settings, including asymptotic analysis of the Cramér-von Mises statistic and regression settings where the covariates are assumed fixed or when a biased coin study design (see Wei, 1978, for example) is used. Extensions of the bootstrap for conducting inference in these new situations is also discussed.

Many interesting statistical quantities can be expressed as functionals of empirical processes. The functional delta method, discussed in chapter 12, can be used to translate weak convergence and bootstrap results for empirical processes to corresponding inference results for these functionals. Most Z- and M- estimators are functionals of empirical processes. For example, under reasonable regularity conditions, the functional which extracts the zero (root) of a Z-estimating equation is sufficiently smooth to permit the delta method to carry over inference results for the estimating equation to the corresponding Z-estimator. The results also apply to M-estimators which can be expressed as approximate Z-estimators.

Z-estimation is discussed in chapter 13, while M-estimation is discussed in chapter 14. A key challenge with many important M-estimators is to establish the rate of convergence, especially in settings where the estimators are not \sqrt{n} consistent. This issue was only briefly mentioned in section 2.2.6 because of the technical complexity of the problem. There are a number of tools which can be used to establish these rates, and several such tools will be studied in chapter 14. These techniques rely significantly on accurate entropy calculations for the M-estimator empirical process, as indexed by the parameter set, within a small neighborhood of the true parameter.

The case studies presented in chapter 15 demonstrate that the technical power of empirical process methods facilitates valid inference for flexible models in many interesting and important statistical settings.

6

Preliminaries for Empirical Processes

In this chapter, we cover several mathematical topics that play a central role in the empirical process results we present later. Metric spaces are crucial since they provide the descriptive language by which the most important results about stochastic processes are derived and expressed. Outer expectations, or, more correctly, outer integrals are key to defining and utilizing outer modes of convergence for quantities which are not measurable. Since many statistical quantities of interest are not measurable with respect to the uniform topology, which is often the topology of choice for applications, outer modes of convergence will be the primary approach for stochastic convergence throughout this book. Linear operators and functional derivatives also play a major role in empirical process methods and are key tools for the functional delta method and Z-estimator theory discussed in chapters 12 and 13.

6.1 Metric Spaces

We now introduce a number of concepts and results for metric spaces. Before defining metric spaces, we briefly review topological spaces, σ -fields, and measure spaces. A collection \mathcal{O} of subsets of a set X is a *topology in* X if:

- (i) $\emptyset \in \mathcal{O}$ and $X \in \mathcal{O}$, where \emptyset is the empty set;
- (ii) If $U_j \in \mathcal{O}$ for $j = 1, \dots, m$, then $\bigcap_{j=1, m} U_j \in \mathcal{O}$;

- (iii) If $\{U_\alpha\}$ is an arbitrary collection of members of \mathcal{O} (finite, countable or uncountable), then $\bigcup_\alpha U_\alpha \in \mathcal{O}$.

When \mathcal{O} is a topology in X , then X (or the pair (X, \mathcal{O})) is a *topological space*, and the members of \mathcal{O} are called the *open sets* in X . For a subset $A \subset X$, the *relative topology* on A consists of the sets $\{A \cap B : B \in \mathcal{O}\}$.

A map $f : X \mapsto Y$ between topological spaces is *continuous* if $f^{-1}(U)$ is open in X whenever U is open in Y . A set B in X is *closed* if and only if its complement in X , denoted $X - B$, is open. The *closure* of an arbitrary set $E \in X$, denoted \overline{E} , is the smallest closed set containing E ; while the *interior* of an arbitrary set $E \in X$, denoted E° , is the largest open set contained in E . A subset A of a topological space X is *dense* if $\overline{A} = X$. A topological space X is *separable* if it has a countable dense subset.

A *neighborhood* of a point $x \in X$ is any open set that contains x . A topological space is *Hausdorff* if distinct points in X have disjoint neighborhoods. A sequence of points $\{x_n\}$ in a topological space X *converges* to a point $x \in X$ if every neighborhood of x contains all but finitely many of the x_n . This convergence is denoted $x_n \rightarrow x$. Suppose $x_n \rightarrow x$ and $x_n \rightarrow y$. Then x and y share all neighborhoods, and $x = y$ when X is Hausdorff. If a map $f : X \mapsto Y$ between topological spaces is continuous, then $f(x_n) \rightarrow f(x)$ whenever $x_n \rightarrow x$ in X . To see this, let $\{x_n\} \subset X$ be a sequence with $x_n \rightarrow x \in X$. Then for any open $U \subset Y$ containing $f(x)$, all but finitely many $\{x_n\}$ are in $f^{-1}(U)$, and thus all but finitely many $\{f(x_n)\}$ are in U . Since U was arbitrary, we have $f(x_n) \rightarrow f(x)$.

We now review the important concept of *compactness*. A subset K of a topological space is *compact* if for every set $A \supset K$, where A is the union of a collection of open sets \mathcal{S} , K is also contained in some finite union of sets in \mathcal{S} . When the topological space involved is also Hausdorff, then compactness of K is equivalent to the assertion that every sequence in K has a convergent subsequence (converging to a point in K). We omit the proof of this equivalence. This result implies that compact subsets of Hausdorff topological spaces are necessarily closed. Note that a compact set is sometimes called a *compact* for short. A σ -*compact* set is a countable union of compacts.

A collection \mathcal{A} of subsets of a set X is a σ -*field in X* (sometimes called a σ -*algebra*) if:

- (i) $X \in \mathcal{A}$;
- (ii) If $U \in \mathcal{A}$, then $X - U \in \mathcal{A}$;
- (iii) The countable union $\bigcup_{j=1}^\infty U_j \in \mathcal{A}$ whenever $U_j \in \mathcal{A}$ for all $j \geq 1$.

Note that (iii) clearly includes finite unions. When (iii) is only required to hold for finite unions, then \mathcal{A} is called a *field*. When \mathcal{A} is a σ -field in X , then X (or the pair (X, \mathcal{A})) is a *measurable space*, and the members of \mathcal{A} are called the *measurable sets* in X . If X is a measurable space and Y

is a topological space, then a map $f : X \mapsto Y$ is *measurable* if $f^{-1}(U)$ is measurable in X whenever U is open in Y .

If \mathcal{O} is a collection of subsets of X (not necessary open), then there exists a smallest σ -field \mathcal{A}^* in X so that $\mathcal{O} \in \mathcal{A}^*$. This \mathcal{A}^* is called the σ -field *generated* by \mathcal{O} . To see that such an \mathcal{A}^* exists, let \mathcal{S} be the collection of all σ -fields in X which contain \mathcal{O} . Since the collection of all subsets of X is one such σ -field, \mathcal{S} is not empty. Define \mathcal{A}^* to be the intersection of all $\mathcal{A} \in \mathcal{S}$. Clearly, $\mathcal{O} \in \mathcal{A}^*$ and \mathcal{A}^* is in every σ -field containing \mathcal{O} . All that remains is to show that \mathcal{A}^* is itself a σ -field. Assume that $A_j \in \mathcal{A}^*$ for all integers $j \geq 1$. If $\mathcal{A} \in \mathcal{S}$, then $\bigcup_{j \geq 1} A_j \in \mathcal{A}$. Since $\bigcup_{j \geq 1} A_j \in \mathcal{A}$ for every $\mathcal{A} \in \mathcal{S}$, we have $\bigcup_{j \geq 1} A_j \in \mathcal{A}^*$. Also $X \in \mathcal{A}^*$ since $X \in \mathcal{A}$ for all $\mathcal{A} \in \mathcal{S}$; and for any $A \in \mathcal{A}^*$, both A and $X - A$ are in every $\mathcal{A} \in \mathcal{S}$. Thus \mathcal{A}^* is indeed a σ -field.

A σ -field is *separable* if it is generated by a countable collection of subsets. Note that we have already defined “separable” as a characteristic of certain topological spaces. There is a connection between the two definitions which we will point out in a few paragraphs when we discuss metric spaces. When X is a topological space, the smallest σ -field \mathcal{B} generated by the open sets is called the *Borel σ -field* of X . Elements of \mathcal{B} are called *Borel sets*. A function $f : X \mapsto Y$ between topological spaces is *Borel-measurable* if it is measurable with respect to the Borel σ -field of X . Clearly, a continuous function between topological spaces is also Borel-measurable.

For a σ -field \mathcal{A} in a set X , a map $\mu : \mathcal{A} \mapsto \overline{\mathbb{R}}$ is a *measure* if:

- (i) $\mu(A) \in [0, \infty]$ for all $A \in \mathcal{A}$;
- (ii) $\mu(\emptyset) = 0$;
- (iii) For a disjoint sequence $\{A_j\} \in \mathcal{A}$, $\mu\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} \mu(A_j)$ (*countable additivity*).

If $X = A_1 \cup A_2 \cup \dots$ for some finite or countable sequence of sets in \mathcal{A} with $\mu(A_j) < \infty$ for all indices j , then μ is σ -*finite*. The triple (X, \mathcal{A}, μ) is called a *measure space*. If $\mu(X) = 1$, then μ is a *probability measure*. For a probability measure P on a set Ω with σ -field \mathcal{A} , the triple (Ω, \mathcal{A}, P) is called a *probability space*. If the set $[0, \infty]$ in part (i) is extended to $(-\infty, \infty]$ or replaced by $[-\infty, \infty)$ (but not both), then μ is a *signed measure*. For a measure space (X, \mathcal{A}, μ) , let \mathcal{A}^* be the collection of all $E \in \mathcal{A}$ for which there exists $A, B \in \mathcal{A}$ with $A \subset E \subset B$ and $\mu(B - A) = 0$, and define $\mu(E) = \mu(A)$ in this setting. Then \mathcal{A}^* is a σ -field, μ is still a measure, and \mathcal{A}^* is called the μ -*completion* of \mathcal{A} .

A *metric space* is a set \mathbb{D} together with a *metric*. A metric or *distance function* is a map $d : \mathbb{D} \times \mathbb{D} \mapsto [0, \infty)$ where:

- (i) $d(x, y) = d(y, x)$;
- (ii) $d(x, z) \leq d(x, y) + d(y, z)$ (the *triangle inequality*);

(iii) $d(x, y) = 0$ if and only if $x = y$.

A *semimetric* or *pseudometric* satisfies (i) and (ii) but not necessarily (iii). Technically, a metric space consists of the pair (\mathbb{D}, d) , but usually only \mathbb{D} is given and the underlying metric d is implied by the context. This is similar to topological and measurable spaces, where only the set of all points X is given while the remaining components are omitted except where needed to clarify the context. A semimetric space is also a topological space with the open sets generated by applying arbitrary unions to the *open r -balls* $B_r(x) \equiv \{y : d(x, y) < r\}$ for $r \geq 0$ and $x \in \mathbb{D}$ (where $B_0(x) \equiv \emptyset$). A metric space is also Hausdorff, and, in this case, a sequence $\{x_n\} \in \mathbb{D}$ converges to $x \in \mathbb{D}$ if $d(x_n, x) \rightarrow 0$. For a semimetric space, $d(x_n, x) \rightarrow 0$ ensures only that x_n converges to elements in the *equivalence class* of x , where the equivalence class of x consists of all $\{y \in \mathbb{D} : d(x, y) = 0\}$. Accordingly, the closure \overline{A} of a set $A \in \mathbb{D}$ is not only the smallest closed set containing A , as stated earlier, but \overline{A} also equals the set of all points that are limits of sequences $\{x_n\} \in A$. Showing this relationship is saved as an exercise. In addition, two semimetrics d_1 and d_2 on a set \mathbb{D} are considered equivalent (in a topological sense) if they both generate the same open sets. It is left as an exercise to show that equivalent metrics yield the same convergent subsequences.

A map $f : \mathbb{D} \mapsto \mathbb{E}$ between two semimetric spaces is *continuous at a point* x if and only if $f(x_n) \rightarrow f(x)$ for every sequence $x_n \rightarrow x$. The map f is continuous (in the topological sense) if and only if it is continuous at all points $x \in \mathbb{D}$. Verifying this last equivalence is saved as an exercise. The following lemma helps to define *semicontinuity* for real valued maps:

LEMMA 6.1 *Let $f : \mathbb{D} \mapsto \mathbb{R}$ be a function on the metric space \mathbb{D} . Then the following are equivalent:*

(i) *For all $c \in \mathbb{R}$, the set $\{y : f(y) \geq c\}$ is closed.*

(ii) *For all $y_0 \in \mathbb{D}$, $\limsup_{y \rightarrow y_0} f(y) \leq f(y_0)$.*

Proof. Assume (i) holds but that (ii) is untrue from some $y_0 \in \mathbb{D}$. This implies that for some $\delta > 0$, $\limsup_{y \rightarrow y_0} f(y) = f(y_0) + \delta$. Thus $H \cap \{y : d(y, y_0) < \epsilon\}$, where $H \equiv \{y : f(y) \geq f(y_0) + \delta\}$, is nonempty for all $\epsilon > 0$. Since H is closed by (i), we now have that $y_0 \in H$. But this implies that $f(y_0) = f(y_0) + \delta$, which is impossible. Hence (ii) holds. The proof that (ii) implies (i) is saved as an exercise. \square

A function $f : \mathbb{D} \mapsto \mathbb{R}$ satisfying either (i) or (ii) (and hence both) of the conditions in lemma 6.1 is said to be *upper semicontinuous*. A function $f : \mathbb{D} \mapsto \mathbb{R}$ is *lower semicontinuous* if $-f$ is upper semicontinuous. Using condition (ii), it is easy to see that a function which is both upper and lower semicontinuous is also continuous. The set of all continuous and bounded functions $f : \mathbb{D} \mapsto \mathbb{R}$, which we denote $C_b(\mathbb{D})$, plays an important role in weak convergence on the metric space \mathbb{D} which we will explore in chapter 7.

It is not hard to see that the Borel σ -field on a metric space \mathbb{D} is the smallest σ -field generated by the open balls. It turns out that the Borel σ -field \mathcal{B} of \mathbb{D} is also the smallest σ -field \mathcal{A} making all of $C_b(\mathbb{D})$ measurable. To see this, note that any closed $A \subset \mathbb{D}$ is the preimage of the closed set $\{0\}$ for the continuous bounded function $x \mapsto d(x, A) \wedge 1$, where for any set $B \subset \mathbb{D}$, $d(x, B) \equiv \inf\{d(x, y) : y \in B\}$. Thus $\mathcal{B} \subset \mathcal{A}$. Since it is obvious that $\mathcal{A} \subset \mathcal{B}$, we now have $\mathcal{A} = \mathcal{B}$. A Borel-measurable map $X : \Omega \mapsto \mathbb{D}$ defined on a probability space (Ω, \mathcal{A}, P) is called a *random element* or *random map* with values in \mathbb{D} . Borel measurability is, in many ways, the natural concept to use on metric spaces since it connects nicely with the topological structure.

A *Cauchy sequence* is a sequence $\{x_n\}$ in a semimetric space (\mathbb{D}, d) such that $d(x_n, x_m) \rightarrow 0$ as $n, m \rightarrow \infty$. A semimetric space \mathbb{D} is *complete* if every Cauchy sequence has a limit $x \in \mathbb{D}$. Every metric space \mathbb{D} has a completion $\overline{\mathbb{D}}$ which has a dense subset *isometric* with \mathbb{D} . Two metric spaces are isometric if there exists a *bijection* (a one-to-one and onto map) between them which preserves distances.

When a metric space \mathbb{D} is separable, and therefore has a countable dense subset, the Borel σ -field for \mathbb{D} is itself a separable σ -field. To see this, let $A \in \mathbb{D}$ be a countable dense subset and consider the collection of open balls with centers at points in A and with rational radii. Clearly, the set of such balls is countable and generates all open sets in \mathbb{D} . A topological space X is *Polish* if it is separable and if there exists a metric making X into a complete metric space. Hence any complete and separable metric space is Polish. Furthermore, any open subset of a Polish space is also Polish. Examples of Polish spaces include Euclidean spaces and many other interesting spaces which we will explore shortly. A *Suslin set* is the continuous image of a Polish space. If a Suslin set is also a Hausdorff topological space, then it is a *Suslin space*. An *analytic set* is a subset A of a Polish space (X, \mathcal{O}) which is Suslin with respect to the relative topology $\{A \cap B : B \in \mathcal{O}\}$. Since there always exists a continuous and onto map $f : X \mapsto A$ for any Borel subset A of a Polish space (X, \mathcal{O}) , every Borel subset of a Polish space is Suslin and therefore also analytic.

A subset K is *totally bounded* if and only if for every $r > 0$, K can be covered by finitely many open r -balls. Furthermore, it can be shown that a subset K of a complete semimetric space is compact if and only if it is totally bounded and closed. A totally bounded subset K is also called *precompact* because every sequence in K has a Cauchy subsequence. To see this, assume K is totally bounded and choose any sequence $\{x_n\} \in K$. There exists a nested series of subsequence indices $\{N_m\}$ and a nested series of 2^{-m} -balls $\{A_m\} \subset K$, such that for each integer $m \geq 1$, N_m is infinite, $N_{m+1} \subset N_m$, $A_{m+1} \subset A_m$, and $x_j \in A_m$ for all $j \in N_m$. This follows from the total boundedness properties. For each $m \geq 1$, choose a $n_m \in N_m$, and note that the subsequence $\{x_{n_m}\}$ is Cauchy. Now assume every sequence in K has a Cauchy subsequence. It is not difficult to verify that if K were not totally bounded, then it is possible to come up with

a sequence which has no Cauchy subsequences (see exercise 6.5.4). This relationship between compactness and total boundedness implies that a σ -compact set in a metric space is separable. These definitions of compactness agree with the previously given compactness properties for Hausdorff spaces. This happens because a semimetric space \mathbb{D} can be made into a metric—and hence Hausdorff—space \mathbb{D}_H by equating points in \mathbb{D}_H with equivalence classes in \mathbb{D} .

A very important example of a metric space is a *normed space*. A normed space \mathbb{D} is a vector space (also called a linear space) equipped with a *norm*, and a norm is a map $\|\cdot\| : \mathbb{D} \mapsto [0, \infty)$ such that, for all $x, y \in \mathbb{D}$ and $\alpha \in \mathbb{R}$,

$$(i) \|x + y\| \leq \|x\| + \|y\| \text{ (another triangle inequality);}$$

$$(ii) \|\alpha x\| = |\alpha| \times \|x\|;$$

$$(iii) \|x\| = 0 \text{ if and only if } x = 0.$$

A *seminorm* satisfies (i) and (ii) but not necessarily (iii). A normed space is a metric space (and a seminormed space is a semimetric space) with $d(x, y) = \|x - y\|$, for all $x, y \in \mathbb{D}$. A complete normed space is called a *Banach space*. Two seminorms $\|\cdot\|_1$ and $\|\cdot\|_2$ on a set \mathbb{D} are equivalent if the following is true for all $x, \{x_n\} \in \mathbb{D}$: $\|x_n - x\|_1 \rightarrow 0$ if and only if $\|x_n - x\|_2 \rightarrow 0$.

In our definition of a normed space \mathbb{D} , we require the space to also be a vector space (and therefore it contains all linear combinations of elements in \mathbb{D}). However, it is sometimes of interest to apply norms to subsets $K \subset \mathbb{D}$ which may not be linear subspaces. In this setting, let $\text{lin}K$ denote the *linear span of K* (all linear combinations of elements in K), and let $\overline{\text{lin}K}$ the closure of $\text{lin}K$. Note that both $\text{lin}K$ and $\overline{\text{lin}K}$ are now vector spaces and that $\overline{\text{lin}K}$ is also a Banach space.

We now present several specific examples of metric spaces. The Euclidean space \mathbb{R}^d is a Banach space with squared norm $\|x\|^2 = \sum_{j=1}^d x_j^2$. This space is equivalent under several other norms, including $\|x\| = \max_{1 \leq j \leq d} |x_j|$ and $\|x\| = \sum_{j=1}^d |x_j|$. A Euclidean space is separable with a countably dense subset consisting of all vectors with rational coordinates. By the Heine-Borel theorem, a subset in a Euclidean space is compact if and only if it is closed and bounded. The Borel σ -field is generated by the intervals of the type $(-\infty, x]$, for rational x , where the interval is defined as follows: $y \in (-\infty, x]$ if and only if $y_j \in (-\infty, x_j]$ for all coordinates $j = 1, \dots, d$. For one-dimensional Euclidean space, \mathbb{R} , the norm is $\|x\| = |x|$ (absolute value). The extended real line $\overline{\mathbb{R}} = [-\infty, \infty]$ is a metric space with respect to the metric $d(x, y) = |G(x) - G(y)|$, where $G : \overline{\mathbb{R}} \mapsto \mathbb{R}$ is any strictly monotone increasing, continuous and bounded function, such as the arctan function. For any sequence $\{x_n\} \in \overline{\mathbb{R}}$, $|x_n - x| \rightarrow 0$ implies $d(x_n, x) \rightarrow 0$, while divergence of $d(x_n, x)$ implies divergence of $|x_n - x|$. In addition, it

is possible for a sequence to converge, with respect to d , to either $-\infty$ or ∞ . This makes $\bar{\mathbb{R}}$ compact.

Another important example is the set of bounded real functions $f : T \mapsto \mathbb{R}$, where T is an arbitrary set. This is a vector space if sums $z_1 + z_2$ and products with scalars, αz , are defined pointwise for all $z, z_1, z_2 \in \ell^\infty(T)$. Specifically, $(z_1 + z_2)(t) = z_1(t) + z_2(t)$ and $(\alpha z)(t) = \alpha z(t)$, for all $t \in T$. This space is denoted $\ell^\infty(T)$. The *uniform norm* $\|x\|_T \equiv \sup_{t \in T} |x(t)|$ makes $\ell^\infty(T)$ into a Banach space consisting exactly of all functions $z : T \mapsto \mathbb{R}$ satisfying $\|z\|_T < \infty$. It is not hard to show that $\ell^\infty(T)$ is separable if and only if T is countable.

Two useful subspaces of $\ell^\infty([a, b])$, where $a, b \in \bar{\mathbb{R}}$, are $C[a, b]$ and $D[a, b]$. The space $C[a, b]$ consists of continuous functions $z : [a, b] \mapsto \mathbb{R}$, and $D[a, b]$ is the space of *cadlag* functions which are right-continuous with left-hand limits (cadlag is an abbreviation for *continue à droite, limites à gauche*). We usually equip these spaces with the uniform norm $\|\cdot\|_{[a, b]}$ inherited from $\ell^\infty([a, b])$. Note that $C[a, b] \subset D[a, b] \subset \ell^\infty([a, b])$. Relative to the uniform norm, $C[a, b]$ is separable, and thus also Polish by the completeness established in exercise 6.5.5(a), but $D[a, b]$ is not separable. Sometimes, $D[a, b]$ is called the *Skorohod space*, although Skorohod equipped $D[a, b]$ with a special metric—quite different than the uniform metric—resulting in a separable space.

An important subspace of $\ell^\infty(T)$ is the space $UC(T, \rho)$, where ρ is a semimetric on T . $UC(T, \rho)$ consists of all bounded function $f : T \mapsto \mathbb{R}$ which are uniformly ρ -continuous, i.e.,

$$\lim_{\delta \downarrow 0} \sup_{\rho(s, t) < \delta} |f(s) - f(t)| = 0.$$

When (T, ρ) is totally bounded, the boundedness requirement for functions in $UC(T, \rho)$ is superfluous since a uniformly continuous function on a totally bounded set must necessarily be bounded. We denote $C(T, \rho)$ to be the space of ρ -continuous (not necessarily continuous) function on T . It is left as an exercise to show that the spaces $C[a, b]$, $D[a, b]$, $UC(T, \rho)$, $C(T, \rho)$, when (T, ρ) is a totally bounded semimetric space, and $UC(T, \rho)$ and $\ell^\infty(T)$, for an arbitrary set T , are all complete with respect to the uniform metric. When (T, ρ) is a compact semimetric space, T is totally bounded, and a ρ -continuous function in T is automatically uniformly ρ -continuous. Thus, when T is compact, $C(T, \rho) = UC(T, \rho)$. Actually, every space $UC(T, \rho)$ is equivalent to a space $C(\bar{T}, \rho)$, because the completion \bar{T} of a totally bounded space T is compact and, furthermore, every uniformly continuous function on T has a unique continuous extension to \bar{T} . Showing this is saved as an exercise.

The forgoing structure makes it clear that $UC(T, \rho)$ is a Polish space which is made complete by the uniform norm. Hence $UC(T, \rho)$ is also σ -compact. In fact, any σ -compact set in $\ell^\infty(T)$ is contained in $UC(T, \rho)$,

for some totally bounded semimetric space (T, ρ) , and all compact sets in $\ell^\infty(T)$ have a specific form:

THEOREM 6.2 (*Arzelà-Ascoli*)

(a) *The closure of $K \subset UC(T, \rho)$, where (T, ρ) is totally bounded, is compact if and only if*

(i) $\sup_{x \in K} |x(t_0)| < \infty$, for some $t_0 \in T$; and

(ii)

$$\limsup_{\delta \downarrow 0} \sup_{x \in K} \sup_{s, t \in T: \rho(s, t) < \delta} |x(s) - x(t)| = 0.$$

(b) *The closure of $K \subset \ell^\infty(T)$ is σ -compact if and only if $K \subset UC(T, \rho)$ for some semimetric ρ making T totally bounded.*

The proof is given in section 6.4. Since all compact sets are trivially σ -compact, theorem 6.2 implies that any compact set in $\ell^\infty(T)$ is actually contained in $UC(T, \rho)$ for some semimetric ρ making T totally bounded.

Another important class of metric spaces are product spaces. For a pair of metric spaces (\mathbb{D}, d) and (\mathbb{E}, e) , the *Cartesian product* $\mathbb{D} \times \mathbb{E}$ is a metric space with respect to the metric $\rho((x_1, y_1), (x_2, y_2)) \equiv d(x_1, x_2) \vee e(y_1, y_2)$, for $x_1, x_2 \in \mathbb{D}$ and $y_1, y_2 \in \mathbb{E}$. This resulting topology is the *product topology*. In this setting, convergence of $(x_n, y_n) \rightarrow (x, y)$ is equivalent to convergence of both $x_n \rightarrow x$ and $y_n \rightarrow y$. There are two natural σ -fields for $\mathbb{D} \times \mathbb{E}$ which we can consider. The first is the Borel σ -field for $\mathbb{D} \times \mathbb{E}$ generated from the product topology. The second is the product σ -field generated by all sets of the form $A \times B$, where $A \in \mathcal{A}$, $B \in \mathcal{B}$, and \mathcal{A} and \mathcal{B} are the respective σ -fields for \mathbb{D} and \mathbb{E} . These two are equal when \mathbb{D} and \mathbb{E} are separable, but they may be unequal otherwise, with the first σ -field larger than the second. Suppose $X : \Omega \mapsto \mathbb{D}$ and $Y : \Omega \mapsto \mathbb{E}$ are Borel-measurable maps defined on a measurable space (Ω, \mathcal{A}) . Then $(X, Y) : \Omega \mapsto \mathbb{D} \times \mathbb{E}$ is a measurable map for the product of the two σ -fields by the definition of a measurable map. Unfortunately, when the Borel σ -field for $\mathbb{D} \times \mathbb{E}$ is larger than the product σ -field, then it is possible for (X, Y) to not be Borel-measurable.

6.2 Outer Expectation

An excellent overview of outer expectations is given in chapter 1.2 of van der Vaart and Wellner (1996). The concept applies to an arbitrary probability space (Ω, \mathcal{A}, P) and an arbitrary map $T : \Omega \mapsto \bar{\mathbb{R}}$, where $\bar{\mathbb{R}} \equiv [-\infty, \infty]$. As described in chapter 2, the outer expectation of T , denoted E^*T , is the infimum over all EU , where $U : \Omega \mapsto \mathbb{R}$ is measurable, $U \geq T$, and EU exists. For EU to exist, it must not be indeterminate, although it can be $\pm\infty$, provided the sign is clear. Since T is not necessarily a random variable,

the proper term for E^*T is *outer integral*. However, we will use the term outer expectation throughout the remainder of this book in deference to its connection with the classical notion of expectation. We analogously define inner expectation: $E_*T = -E^*[-T]$. The following lemma verifies the existence of a minimal measurable majorant $T^* \geq T$:

LEMMA 6.3 For any $T : \Omega \mapsto \bar{\mathbb{R}}$, there exists a minimal measurable majorant $T^* : \Omega \mapsto \bar{\mathbb{R}}$ with

(i) $T^* \geq T$;

(ii) For every measurable $U : \Omega \mapsto \bar{\mathbb{R}}$ with $U \geq T$ a.s., $T^* \leq U$ a.s.

For any T^* satisfying (i) and (ii), $E^*T = ET^*$, provided ET^* exists. The last statement is true if $E^*T < \infty$.

The proof is given in section 6.4 at the end of this chapter. The following lemma, the proof of which is left as an exercise, is an immediate consequence of lemma 6.3 and verifies the existence of a maximal measurable minorant:

LEMMA 6.4 For any $T : \Omega \mapsto \bar{\mathbb{R}}$, the maximal measurable minorant $T_* \equiv -(-T)^*$ exists and satisfies

(i) $T_* \leq T$;

(ii) For every measurable $U : \Omega \mapsto \bar{\mathbb{R}}$ with $U \leq T$ a.s., $T_* \geq U$ a.s.

For any T_* satisfying (i) and (ii), $E_*T = ET_*$, provided ET_* exists. The last statement is true if $E_*T > -\infty$.

An important special case of outer expectation is outer probability. The outer probability of an arbitrary $B \subset \Omega$, denoted $P^*(B)$, is the infimum over all $P(A)$ such that $A \supset B$ and $A \in \mathcal{A}$. The inner probability of an arbitrary $B \subset \Omega$ is defined to be $P_*(B) = 1 - P^*(\Omega - B)$. The following lemma gives the precise connection between outer/inner expectations and outer/inner probabilities:

LEMMA 6.5 For any $B \subset \Omega$,

(i) $P^*(B) = E^*1\{B\}$ and $P_*(B) = E_*1\{B\}$;

(ii) there exists a measurable set $B^* \supset B$ so that $P(B^*) = P^*(B)$; for any such B^* , $1\{B^*\} = (1\{B\})^*$;

(iii) For $B_* \equiv \Omega - \{\Omega - B\}^*$, $P_*(B) = P(B_*)$;

(iv) $(1\{B\})^* + (1\{\Omega - B\})_* = 1$.

Proof. From the definitions, $P^*(B) = \inf_{\{A \in \mathcal{A}: A \supset B\}} E1\{A\} \geq E^*1\{B\}$. Next, $E^*1\{B\} = E(1\{B\})^* \geq E1\{(1\{B\})^* \geq 1\} = P\{(1\{B\})^* \geq 1\} \geq$

$P^*(B)$, where the last inequality follows from the definition of P^* . Combining the two conclusions yields that all inequalities are actually equalities. This gives the first parts of (i) and (ii), with $B^* = \{(1\{B\})^* \geq 1\}$. The second part of (i) results from $P_*(B) = 1 - P^*(\Omega - B) = 1 - E(1 - 1\{B\})^* = 1 - E(1 - (1\{B\})_*)$. The second part of (ii) follows from $(1\{B\})^* \leq 1\{B^*\} = 1\{(1\{B\})^* \geq 1\} \leq (1\{B\})^*$. The definition of P_* implies (iii) directly. To verify (iv), we have $(1\{\Omega - B\})_* = (1 - 1\{B\})_* = -(1\{B\} - 1)^* = 1 - (1\{B\})^*$. \square

The following three lemmas, lemmas 6.6–6.8, provide several relations which will prove useful later on but which might be skipped on a first reading. The proofs are given in section 6.4 and in the exercises.

LEMMA 6.6 *Let $S, T : \Omega \mapsto \mathbb{R}$ be arbitrary maps. The following statements are true almost surely, provided the statements are well-defined:*

- (i) $S_* + T_* \leq (S + T)_* \leq S_* + T_*$, with all equalities if S is measurable;
- (ii) $S_* + T_* \leq (S + T)_* \leq S_* + T_*$, with all equalities if T is measurable;
- (iii) $(S - T)^* \geq S^* - T^*$;
- (iv) $|S^* - T^*| \leq |S - T|^*$;
- (v) $(1\{T > c\})^* = 1\{T^* > c\}$, for any $c \in \mathbb{R}$;
- (vi) $(1\{T \geq c\})_* = 1\{T_* \geq c\}$, for any $c \in \mathbb{R}$;
- (vii) $(S \vee T)^* = S^* \vee T^*$;
- (viii) $(S \wedge T)^* \leq S^* \wedge T^*$, with equality if S is measurable.

LEMMA 6.7 *For any sets $A, B \subset \Omega$,*

- (i) $(A \cup B)^* = A^* \cup B^*$ and $(A \cap B)_* = A_* \cap B_*$;
- (ii) $(A \cap B)^* \subset A^* \cap B^*$ and $(A \cup B)_* \supset A_* \cup B_*$, with the inclusions replaced by equalities if either A or B is measurable;
- (iii) If $A \cap B = \emptyset$, then $P_*(A) + P_*(B) \leq P_*(A \cup B) \leq P^*(A \cup B) \leq P^*(A) + P^*(B)$.

LEMMA 6.8 *Let $T : \Omega \mapsto \mathbb{R}$ be an arbitrary map and let $\phi : \mathbb{R} \mapsto \mathbb{R}$ be monotone, with an extension to $\bar{\mathbb{R}}$. The following statements are true almost surely, provided the statements are well-defined:*

A. *If ϕ is nondecreasing, then*

- (i) $\phi(T^*) \geq [\phi(T)]^*$, with equality if ϕ is left-continuous on $[-\infty, \infty)$;
- (ii) $\phi(T_*) \leq [\phi(T)]_*$, with equality if ϕ is right-continuous on $(-\infty, \infty]$.

B. *If ϕ is nonincreasing, then*

- (i) $\phi(T^*) \leq [\phi(T)]_*$, with equality if ϕ is left-continuous on $[-\infty, \infty)$;
(ii) $\phi(T_*) \geq [\phi(T)]^*$, with equality if ϕ is right-continuous on $(-\infty, \infty]$.

We next present an outer-expectation version of the unconditional Jensen's inequality:

LEMMA 6.9 (Jensen's inequality) *Let $T : \Omega \mapsto \mathbb{R}$ be an arbitrary map, with $E^*|T| < \infty$, and assume $\phi : \mathbb{R} \mapsto \mathbb{R}$ is convex. Then*

- (i) $E^*\phi(T) \geq \phi(E_*T) \vee \phi(E^*T)$;
(ii) if ϕ is also monotone, $E_*\phi(T) \geq \phi(E_*T) \wedge \phi(E^*T)$.

Proof. Assume first that ϕ is monotone increasing. Since ϕ is also continuous (by convexity), $E^*\phi(T) = E\phi(T^*) \geq \phi(E^*T)$, where the equality follows from A(i) of lemma 6.8 and the inequality from the usual Jensen's inequality. Similar arguments verify that $E_*\phi(T) \geq \phi(E_*T)$ based on A(ii) of the same lemma. Note also that $\phi(E^*T) \geq \phi(E_*T)$. Now assume that ϕ is monotone decreasing. Using B(i) and B(ii) of lemma 6.8 and arguments similar to those used for increasing ϕ , we obtain that both $E^*\phi(T) \geq \phi(E_*T)$ and $E_*\phi(T) \geq \phi(E^*T)$. Note in this case that $\phi(E_*T) \geq \phi(E^*T)$. Thus when ϕ is monotone (either increasing or decreasing), we have that both $E^*\phi(T) \geq \phi(E^*T) \vee \phi(E_*T)$ and $E_*\phi(T) \geq \phi(E^*T) \wedge \phi(E_*T)$. Hence (ii) follows.

We have also proved (i) in the case where ϕ is monotone. Suppose now that ϕ is not monotone. Then there exists a $c \in \mathbb{R}$ so that ϕ is nonincreasing over $(-\infty, c]$ and nondecreasing over (c, ∞) . Let $g_1(t) \equiv \phi(t)1\{t \leq c\} + \phi(c)1\{t > c\}$ and $g_2(t) \equiv \phi(c)1\{t \leq c\} + \phi(t)1\{t > c\}$, and note that $\phi(t) = g_1(t) + g_2(t) - \phi(c)$ and that both g_1 and g_2 are convex and monotone. Now $[\phi(T)]^* = [g_1(T) + g_2(T) - \phi(c)]^* \geq [g_1(T)]_* + [g_2(T)]^*$ by part (i) of lemma 6.6. Now, using the results in the previous paragraph, we have that $E^*\phi(T) \geq g_1(E^*T) + g_2(E^*T) - \phi(c) = \phi(E^*T)$. However, we also have that $[g_1(T) + g_2(T) - \phi(c)]^* \geq [g_1(T)]^* + [g_2(T)]_*$. Thus, again using results from the previous paragraph, we have $E^*\phi(T) \geq g_1(E_*T) + g_2(E_*T) - \phi(c) = \phi(E_*T)$. Hence (i) follows. \square

The following outer-expectation version of Chebyshev's inequality is also useful:

LEMMA 6.10 (Chebyshev's inequality) *Let $T : \Omega \mapsto \mathbb{R}$ be an arbitrary map, with $\phi : [0, \infty) \mapsto [0, \infty)$ nondecreasing and strictly positive on $(0, \infty)$. Then, for every $u > 0$,*

$$P^*(|T| \geq u) \leq \frac{E^*\phi(|T|)}{\phi(u)}.$$

Proof. The result follows from the standard Chebyshev inequality as a result of the following chain of inequalities:

$$(1\{|T| \geq u\})^* \leq (1\{\phi(|T|) \geq \phi(u)\})^* \leq 1\{[\phi(|T|)]^* \geq \phi(u)\}.$$

The first inequality follows from the fact that $|T| \geq u$ implies $\phi(|T|) \geq \phi(u)$, and the second inequality follows from A(i) of lemma 6.8. \square

There are many other analogies between outer and standard versions of expectation and probability, but we only present a few more in this chapter. We next present versions of the monotone and dominated convergence theorems. The proofs are given in section 6.4. Some additional results are given in the exercises.

LEMMA 6.11 (Monotone convergence) *Let $T_n, T : \Omega \mapsto \mathbb{R}$ be arbitrary maps on a probability space, with $T_n \uparrow T$ pointwise on a set of inner probability one. Then $T_n^* \uparrow T^*$ almost surely. Provided $E^*T_n > -\infty$ for some n , then $E^*T_n \uparrow E^*T$.*

LEMMA 6.12 (Dominated convergence) *Let $T_n, T, S : \Omega \mapsto \mathbb{R}$ be maps on a probability space, with $|T_n - T|^* \xrightarrow{\text{as}} 0$, $|T_n| \leq S$ for all n , and $E^*S < \infty$. Then $E^*T_n \rightarrow E^*T$.*

Let $(\Omega, \tilde{\mathcal{A}}, \tilde{P})$ be the P -completion of the probability space (Ω, \mathcal{A}, P) , as defined in the previous section (for general measure spaces). Recall that a completion of a measure space is itself a measure space. One can also show that $\tilde{\mathcal{A}}$ is the σ -field of all sets of the form $A \cup N$, with $A \in \mathcal{A}$ and $N \subset \Omega$ so that $P^*(N) = 0$, and \tilde{P} is the probability measure satisfying $\tilde{P}(A \cup N) = P(A)$. A nice property of $(\Omega, \tilde{\mathcal{A}}, \tilde{P})$ is that for every measurable map $\tilde{S} : (\Omega, \tilde{\mathcal{A}}) \mapsto \mathbb{R}$, there is a measurable map $S : (\Omega, \mathcal{A}) \mapsto \mathbb{R}$ such that $P^*(S \neq \tilde{S}) = 0$. Furthermore, a minimal measurable cover T^* of a map $T : (\Omega, \mathcal{A}, P) \mapsto \bar{\mathbb{R}}$ is a version of a minimal measurable cover \tilde{T}^* for T as a map on the P -completion of (Ω, \mathcal{A}, P) , i.e., $P^*(T^* \neq \tilde{T}^*) = 0$. While it is not difficult to show this, we do not prove it.

We close this section with two results which have application to product probability spaces. The first result involves *perfect* maps, and the second result is a special formulation of Fubini's theorem. Consider composing a map $T : \Omega \mapsto \mathbb{R}$ with a measurable map $\phi : \tilde{\Omega} \mapsto \Omega$, defined on some probability space, to form $T \circ \phi : (\tilde{\Omega}, \tilde{\mathcal{A}}, \tilde{P}) \mapsto \mathbb{R}$, where $\phi : (\tilde{\Omega}, \tilde{\mathcal{A}}, \tilde{P}) \mapsto (\Omega, \mathcal{A}, P)$. Denote T^* as the minimal measurable cover of T for $\tilde{P} \circ \phi^{-1}$. It is easy to see that since $T^* \circ \phi \geq T \circ \phi$, we have $(T \circ \phi)^* \leq T^* \circ \mathbb{R}$. The map ϕ is perfect if $(T \circ \phi)^* = T^* \circ \phi$, for every bounded $T : \Omega \mapsto \mathbb{R}$. This property ensures that $P^*(\phi \in A) = (\tilde{P} \circ \phi^{-1})^*(A)$ for every set $A \subset \Omega$.

An important example of a perfect map is a coordinate projection in a product probability space. Specifically, let T be a real valued map defined on $(\Omega_1 \times \Omega_2, \mathcal{A}_1 \times \mathcal{A}_2, P_1 \times P_2)$ which only depends on the first coordinate of $\omega = (\omega_1, \omega_2)$. T^* can then be computed by just ignoring Ω_2 and thinking of T as a map on Ω_1 . More precisely, suppose $T = T_1 \circ \pi_1$, where π_1 is

the projection on the first coordinate. The following lemma shows that $T^* = T_1^* \circ \pi_1$, and thus coordinate projections such as π_1 are perfect. We will see other examples of perfect maps later on in chapter 7.

LEMMA 6.13 *A coordinate projection on a product probability space with product measure is perfect.*

Proof. Let $\pi_1 : (\Omega_1 \times \Omega_2, \mathcal{A}_1 \times \mathcal{A}_2, P_1 \times P_2) \mapsto \Omega_1$ be the projection on the first coordinate, and let $T : \Omega_1 \mapsto \mathbb{R}$ be bounded, but otherwise arbitrary. Let T^* be the minimal measurable cover of T for $P_1 = (P_1 \times P_2) \circ \pi_1^{-1}$. By definition, $(T \circ \pi_1)^* \leq T^* \circ \pi_1$. Now suppose $U \geq T \circ \pi_1$, $P_1 \times P_2$ -a.s., and is measurable, where $U : \Omega_1 \times \Omega_2 \mapsto \mathbb{R}$. Fubini's theorem yields that for P_2 -almost all ω_2 , we have $U(\omega_1, \omega_2) \geq T(\omega_1)$ for P_2 -almost all ω_2 . But for fixed ω_2 , U is a measurable function of ω_1 . Thus for P_2 -almost all ω_2 , $U(\omega_1, \omega_2) \geq T^*(\omega_1)$ for P_1 -almost all ω_1 . Applying Fubini's theorem again, the jointly measurable set $\{(\omega_1, \omega_2) : U < T^* \circ \pi_1\}$ is $P_1 \times P_2$ -null. Hence $(T \circ \pi_1)^* = T^* \circ \pi_1$ almost surely. \square

Now we consider Fubini's theorem for maps on product spaces which may not be measurable. There is no generally satisfactory version of Fubini's theorem that will work in all nonmeasurable settings of interest, and it is frequently necessary to establish at least some kind of measurability to obtain certain key empirical process results. The version of Fubini's theorem we now present basically states that repeated outer expectations are always less than joint outer expectations. Let T be an arbitrary real map defined on the product space $(\Omega_1 \times \Omega_2, \mathcal{A}_1 \times \mathcal{A}_2, P_1 \times P_2)$. We write $E_1^* E_2^* T$ to mean outer expectations taken in turn. For fixed ω_1 , let $(E_2^* T)(\omega_1)$ be the infimum of $\int_{\Omega_2} U(\omega_2) dP_2(\omega_2)$ taken over all measurable $U : \Omega_2 \mapsto \bar{\mathbb{R}}$ with $U(\omega_2) \geq T(\omega_1, \omega_2)$ for every ω_2 for which $\int_{\Omega_2} U(\omega_2) dP_2(\omega_2)$ exists. Next, $E_1^* E_2^* T$ is the outer integral of $E_2^* T : \Omega_1 \mapsto \mathbb{R}$. Repeated inner expectations are analogously defined. The following version of Fubini's theorem gives bounds for this repeated expectation process. We omit the proof.

LEMMA 6.14 (Fubini's theorem) *Let T be an arbitrary real valued map on a product probability space. Then $E_* T \leq E_{1*} E_{2*} T \leq E_1^* E_2^* T \leq E^* T$.*

6.3 Linear Operators and Functional Differentiation

A *linear operator* is a map $T : \mathbb{D} \mapsto \mathbb{E}$ between normed spaces with the property that $T(ax + by) = aT(x) + bT(y)$ for all scalars a, b and any $x, y \in \mathbb{D}$. When the range space \mathbb{E} is \mathbb{R} , then T is a *linear functional*. When T is linear, we will often use Tx instead of $T(x)$. A linear operator $T : \mathbb{D} \mapsto \mathbb{E}$ is a *bounded linear operator* if

$$(6.1) \quad \|T\| \equiv \sup_{x \in \mathbb{D} : \|x\| \leq 1} \|Tx\| < \infty.$$

Here, the norms $\|\cdot\|$ are defined by the context. We have the following proposition:

PROPOSITION 6.15 *For a linear operator $T : \mathbb{D} \mapsto \mathbb{E}$ between normed spaces, the following are equivalent:*

- (i) T is continuous at a point $x_0 \in \mathbb{D}$;
- (ii) T is continuous on all of \mathbb{D} ;
- (iii) T is bounded.

Proof. We save the implication (i) \Rightarrow (ii) as an exercise. Note that by linearity, boundedness of T is equivalent to there existing some $0 < c < \infty$ for which

$$(6.2) \quad \|Tx\| \leq c\|x\| \text{ for all } x \in \mathbb{D}.$$

Assume T is continuous but that there exists no $0 < c < \infty$ satisfying (6.2). Then there exists a sequence $\{x_n\} \in \mathbb{D}$ so that $\|x_n\| = 1$ and $\|Tx_n\| \geq n$ for all $n \geq 1$. Define $y_n = \|Tx_n\|^{-1}x_n$ and note that $\|Ty_n\| = 1$ by linearity. Now $y_n \rightarrow 0$ and thus $Ty_n \rightarrow 0$, but this is a contradiction. Thus there exists some $0 < c < \infty$ satisfying (6.2), and (iii) follows. Now assume (iii) and let $\{x_n\} \in X$ be any sequence satisfying $x_n \rightarrow 0$. Then by (6.2), $\|Tx_n\| \rightarrow 0$, and thus (i) is satisfied at $x_0 = 0$. \square

For normed spaces \mathbb{D} and \mathbb{E} , let $B(\mathbb{D}, \mathbb{E})$ be the space of all bounded linear operators $T : \mathbb{D} \mapsto \mathbb{E}$. This structure makes the space $B(\mathbb{D}, \mathbb{E})$ into a normed space with norm $\|\cdot\|$ defined in (6.1). When \mathbb{E} is a Banach space, then any convergent sequence $T_n x_n$ will be contained in \mathbb{E} , and thus $B(\mathbb{D}, \mathbb{E})$ is also a Banach space. When \mathbb{D} is not a Banach space, T has a unique continuous extension to $\overline{\mathbb{D}}$. To see this, fix $x \in \overline{\mathbb{D}}$ and let $\{x_n\} \in \mathbb{D}$ be a sequence converging to x . Then, since $\|Tx_n - Tx_m\| \leq c\|x_n - x_m\|$, Tx_n converges to some point in $\overline{\mathbb{E}}$. Next, note that if both sequences $\{x_n\}, \{y_n\} \in \mathbb{D}$ converge to x , then $\|Tx_n - Ty_n\| \leq c\|x_n - y_n\| \rightarrow 0$. Thus we can define an extension $\overline{T} : \overline{\mathbb{D}} \mapsto \overline{\mathbb{E}}$ to be the unique linear operator with $\overline{T}x = \lim_{n \rightarrow \infty} Tx_n$, where x is any point in $\overline{\mathbb{D}}$ and $\{x_n\}$ is any sequence in \mathbb{D} converging to x .

For normed spaces \mathbb{D} and \mathbb{E} , and for any $T \in B(\mathbb{D}, \mathbb{E})$, $N(T) \equiv \{x \in \mathbb{D} : Tx = 0\}$ is the *null space* of T and $R(T) \equiv \{y \in \mathbb{E} : Tx = y \text{ for some } x \in \mathbb{D}\}$ is the *range space* of T . It is clear that T is one-to-one if and only if $N(T) = \{0\}$. We have the following two results for inverses, which we give without proof:

LEMMA 6.16 *Assume \mathbb{D} and \mathbb{E} are normed spaces and that $T \in B(\mathbb{D}, \mathbb{E})$. Then*

- (i) T has a continuous inverse $T^{-1} : R(T) \mapsto \mathbb{D}$ if and only if there exists a $c > 0$ so that $\|Tx\| \geq c\|x\|$ for all $x \in \mathbb{D}$;

(ii) **(Banach's theorem)** *If \mathbb{D} and \mathbb{E} are complete and T is continuous with $N(T) = \{0\}$, then T^{-1} is continuous if and only if $R(T)$ is closed.*

A linear operator $T : \mathbb{D} \mapsto \mathbb{E}$ between normed spaces is a *compact operator* if $T(U)$ is compact in $\overline{\mathbb{E}}$, where $U = \{x \in \mathbb{D} : \|x\| \leq 1\}$ is the unit ball in \mathbb{D} and, for a set $A \in \mathbb{D}$, $T(A) \equiv \{Tx : x \in A\}$. The operator T is *onto* if for every $y \in \mathbb{E}$, there exists an $x \in \mathbb{D}$ so that $Tx = y$. Later on in the book, we will encounter linear operators of the form $T + K$, where T is continuously invertible and onto and K is compact. The following result will be useful:

LEMMA 6.17 *Let $A = T + K : \mathbb{D} \mapsto \mathbb{E}$ be a linear operator between Banach spaces, where T is both continuously invertible and onto and K is compact. Then if $N(A) = \{0\}$, A is also continuously invertible and onto.*

Proof. We only sketch the proof. Since T^{-1} is continuous, the operator $T^{-1}K : \mathbb{E} \mapsto \mathbb{D}$ is compact. Hence $I + T^{-1}K$ is one-to-one and therefore also onto by a result of Riesz for compact operators (see, for example, theorem 3.4 of Kress, 1999). Thus $T + K$ is also onto. We will be done if we can show that $I + T^{-1}K$ is continuously invertible, since that would imply that $(T + K)^{-1} = (I + T^{-1}K)^{-1}T^{-1}$ is bounded. Assume $L \equiv I + T^{-1}K$ is not bounded. Then there exists a sequence $\{x_n\} \in \mathbb{D}$ with $\|x_n\| = 1$ and $\|L^{-1}x_n\| \geq n$ for all integers $n \geq 1$. Let $y_n = (\|L^{-1}x_n\|)^{-1}x_n$ and $\phi_n = (\|L^{-1}x_n\|)^{-1}L^{-1}x_n$, and note that $y_n \rightarrow 0$ while $\|\phi_n\| = 1$ for all $n \geq 1$. Since $T^{-1}K$ is compact, there exists a subsequence $\{n'\}$ so that $T^{-1}K\phi_{n'} \rightarrow \phi \in \mathbb{D}$. Since $\phi_n + T^{-1}K\phi_n = y_n$ for all $n \geq 1$, we have $\phi_{n'} \rightarrow -\phi$. Hence $\phi \in N(L)$, which implies $\phi = 0$ since L is one-to-one. But this contradicts $\|\phi_n\| = 1$ for all $n \geq 1$. Thus L^{-1} is bounded. \square

The following simple inversion result for *contraction operators* is also useful. An operator A is a contraction operator if $\|A\| < 1$.

PROPOSITION 6.18 *Let $A : \mathbb{D} \mapsto \mathbb{D}$ be a linear operator with $\|A\| < 1$. Then $I - A$, where I is the identity, is continuously invertible and onto with inverse $(I - A)^{-1} = \sum_{j=0}^{\infty} A^j$.*

Proof. Let $B \equiv \sum_{j=0}^{\infty} A^j$, and note that $\|B\| \leq \sum_{j=0}^{\infty} \|A\|^j = (1 - \|A\|)^{-1} < \infty$. Thus B is a bounded linear operator on \mathbb{D} . Since $(I - A)B = I$ by simple algebra, we have that $B = (I - A)^{-1}$, and the result follows. \square

We now shift our attention to differentiation. Let \mathbb{D} and \mathbb{E} be two normed spaces, and let $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ be a function. We allow the domain \mathbb{D}_ϕ of the function to be an arbitrary subset of \mathbb{D} . The function $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ is *Gateaux-differentiable* at $\theta \in \mathbb{D}_\phi$, in the direction h , if there exists a quantity $\phi'_\theta(h) \in \mathbb{E}$ so that

$$\frac{\phi(\theta + t_n h) - \phi(\theta)}{t_n} \rightarrow \phi'_\theta(h),$$

as $n \rightarrow \infty$, for any scalar sequence $t_n \rightarrow 0$. Gateaux differentiability is usually not strong enough for the applications of functional derivatives needed for Z-estimators and the delta method. The stronger differentiability we need is *Hadamard* and *Fréchet* differentiability.

A map $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ is Hadamard differentiable at $\theta \in \mathbb{D}_\phi$ if there exists a continuous linear operator $\phi'_\theta : \mathbb{D} \mapsto \mathbb{E}$ such that

$$(6.3) \quad \frac{\phi(\theta + t_n h_n) - \phi(\theta)}{t_n} \rightarrow \phi'_\theta(h),$$

as $n \rightarrow \infty$, for any scalar sequence $t_n \rightarrow 0$ and any $h, \{h_n\} \in \mathbb{D}$, with $h_n \rightarrow h$, and so that $\theta + t_n h_n \in \mathbb{D}_\phi$ for all n . It is left as an exercise to show that Hadamard differentiability is equivalent to *compact differentiability*, where compact differentiability satisfies

$$(6.4) \quad \sup_{h \in K, \theta + th \in \mathbb{D}_\phi} \left\| \frac{\phi(\theta + th) - \phi(\theta)}{t} - \phi'_\theta(h) \right\| \rightarrow 0, \text{ as } t \rightarrow 0,$$

for every compact $K \subset \mathbb{D}$. Hadamard differentiability can be refined by restricting the h values to be in a set $\mathbb{D}_0 \subset \mathbb{D}$. More precisely, if in (6.3) it is required that $h_n \rightarrow h$ only for $h \in \mathbb{D}_0 \subset \mathbb{D}$, we say ϕ is *Hadamard-differentiable tangentially* to the set \mathbb{D}_0 . There appears to be no easy way to refine compact differentiability in an equivalent manner.

A map $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ is *Fréchet-differentiable* if there exists a continuous linear operator $\phi'_\theta : \mathbb{D} \mapsto \mathbb{E}$ so that (6.4) holds uniformly in h on bounded subsets of \mathbb{D} . This is equivalent to $\|\phi(\theta + h) - \phi(\theta) - \phi'_\theta(h)\| = o(\|h\|)$, as $\|h\| \rightarrow 0$. Since compact sets are bounded, Fréchet differentiability implies Hadamard differentiability. Fréchet differentiability will be needed for Z-estimator theory, while Hadamard differentiability is useful in the delta method. The following chain rule for Hadamard differentiability will also prove useful:

LEMMA 6.19 (Chain rule) *Assume $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}_\psi \subset \mathbb{E}$ is Hadamard differentiable at $\theta \in \mathbb{D}_\phi$ tangentially to $\mathbb{D}_0 \subset \mathbb{D}$, and $\psi : \mathbb{E}_\psi \subset \mathbb{E} \mapsto \mathbb{F}$ is Hadamard differentiable at $\phi(\theta)$ tangentially to $\phi'_\theta(\mathbb{D}_0)$. Then $\psi \circ \phi : \mathbb{D}_\phi \mapsto \mathbb{F}$ is Hadamard differentiable at θ tangentially to \mathbb{D}_0 with derivative $\psi'_{\phi(\theta)} \circ \phi'_\theta$.*

Proof. First, $\psi \circ \phi(\theta + th_t) - \psi \circ \phi(\theta) = \psi(\phi(\theta) + tk_t) - \psi(\phi(\theta))$, where $k_t = [\phi(\theta + th_t) - \phi(\theta)]/t$. Note that if $h \in \mathbb{D}_0$, then $k_t \rightarrow k \equiv \phi'_\theta(h) \in \phi'_\theta(\mathbb{D}_0)$, as $t \rightarrow 0$, by the Hadamard differentiability of ϕ . Now, $[\psi(\phi(\theta) + tk_t) - \psi(\phi(\theta))]/t \rightarrow \psi'_{\phi(\theta)}(k)$ by the Hadamard differentiability of ψ . \square

6.4 Proofs

Proof of theorem 6.2. First assume $\overline{K} \subset \ell^\infty(T)$ is compact. Let $\rho(s, t) = \sup_{x \in K} |x(s) - x(t)|$. We will now establish that (T, ρ) is totally bounded.

Fix $\eta > 0$ and cover K with finitely many open balls of radius η , centered at x_1, \dots, x_k . Partition \mathbb{R}^k into cubes with edges of length η . For every such cube for which it is possible, choose at most one $t \in T$ so that $(x_1(t), \dots, x_k(t))$ is in the cube. This results in a finite set $T_\eta \equiv \{t_1, \dots, t_m\}$ in T because z_1, \dots, z_k are uniformly bounded. For each $s \in T_\eta$, we have for all values of $t \in T$ for which $(z_1(t), \dots, z_k(t))$ and $(z_1(s), \dots, z_k(s))$ are in the same cube, that

$$\begin{aligned} \rho(s, t) &= \sup_{z \in K} |z(s) - z(t)| \\ &\leq \sup_{1 \leq j \leq k} |z_j(s) - z_j(t)| + 2 \sup_{z \in K} \inf_{1 \leq j \leq k} \sup_{u \in T} |z_j(u) - z(u)| \\ &< 3\eta. \end{aligned}$$

Thus the balls $\{t : \rho(t, t_1) < 3\eta\}, \dots, \{t : \rho(t, t_m) < 3\eta\}$ completely cover T . Hence (T, ρ) is totally bounded since η was arbitrary. Also, by construction, the condition in part (a.ii) of the theorem is satisfied. Combining this with the total boundedness of (T, ρ) yields condition (a.i). We have now obtained the fairly strong result that compactness of the closure of $K \subset \ell^\infty(T)$ implies that there exists a semimetric ρ which makes (T, ρ) totally bounded and which enables conditions (a.i) and (a.ii) to be satisfied for K .

Now assume that the closure of $K \subset \ell^\infty(T)$ is σ -compact. This implies that there exists a sequence of compact sets $K_1 \subset K_2 \subset \dots$ for which $\overline{K} = \cup_{i \geq 1} K_i$. The previous result yields for each $i \geq 1$, that there exists a semimetric ρ_i making T totally bounded and for which conditions (a.i) and (a.ii) are satisfied for each K_i . Now let $\rho(s, t) = \sum_{i=1}^{\infty} 2^{-i} (\rho_i(s, t) \wedge 1)$. Fix $\eta > 0$, and select a finite integer m so that $2^{-m} < \eta$. Cover T with finitely many open ρ_m balls of radius η , and let $T_\eta = \{t_1, \dots, t_k\}$ be their centers. Because $\rho_1 \leq \rho_2 \leq \dots$, there is for very every $t \in T$ an $s \in T_\eta$ with $\rho(s, t) \leq \sum_{i=1}^m 2^{-i} \rho_i(s, t) + 2^{-m} \leq 2\eta$. Thus (T, ρ) is totally bounded by ρ since η was arbitrary. Now, for any $x \in K$, $x \in K_m$ for some finite $m \geq 1$, and thus x is both bounded and uniformly ρ -continuous since $\rho_m \leq 2^m \rho$. Hence σ -compactness of $K \subset \ell^\infty(T)$ implies $K \subset UC(T, \rho)$ for some semimetric ρ making T totally bounded. In the discussion preceding the statement of theorem 6.2, we argued that $UC(T, \rho)$ is σ -compact whenever (T, ρ) is totally bounded. Hence we have proven part (b).

The only part of the proof which remains is to show that if $K \subset UC(T, \rho)$ satisfies conditions (a.i) and (a.ii), then the closure of K is compact. Assume conditions (a.i) and (a.ii) hold for K . Define

$$m_\delta(x) \equiv \sup_{s, t \in T: \rho(s, t) < \delta} |x(s) - x(t)|$$

and $m_\delta \equiv \sup_{x \in K} m_\delta(x)$, and note that $m_\delta(x)$ is continuous in x and that $m_{1/n}(x)$ is nonincreasing in n , with $\lim_{n \rightarrow \infty} m_{1/n}(x) = 0$. Choose $k < \infty$ large enough so that $m_{1/k} < \infty$. For every $\delta > 0$, let $T_\delta \subset T$ be a finite mesh satisfying $\sup_{t \in T} \inf_{s \in T_\delta} \rho(s, t) < \delta$, and let N_δ be the number of points in

T_δ . Now, for any $t \in T$, $|x(t)| \leq |x(t_0)| + |x(t) - x(t_0)| \leq |x(t_0)| + N_{1/k} m_{1/k}$, and thus $\alpha \equiv \sup_{x \in K} \sup_{t \in T} |x(t)| < \infty$. For each $\epsilon > 0$, pick a $\delta > 0$ so that $m_\delta < \epsilon$ and an integer $n < \infty$ so that $\alpha/n \leq \epsilon$. Let $U \equiv T_{1/(2\delta)}$ and define a “bracket” to be a finite collection $h = \{h_t : t \in U\}$ so that $h_t = -\alpha + j\alpha/n$, for some integer $1 \leq j \leq 2n - 1$, for each $t \in U$. Say that $x \in \ell^\infty(T)$ is “in” the bracket h , denoted $x \in h$, if $x(t) \in [h_t, h_t + \alpha/n]$ for all $t \in U$. Let $B(K)$ be the set of all brackets h for which $x \in h$ for some $x \in K$. For each $h \in B(K)$, choose one and only one $x \in K$ with $x \in h$, discard duplicates, and denote the resulting set $X(K)$. It is not hard to verify that $\sup_{x \in K} \inf_{y \in X(K)} \|x - y\|_T < 2\epsilon$, and thus the union of 2ϵ -balls with centers in $X(K)$ is a finite cover of K . Since ϵ is arbitrary, K is totally bounded. \square

Proof of lemma 6.3. Begin by selecting a measurable sequence $U_m \geq T$ such that $E \arctan U_m \downarrow E^* \arctan T$, and set $T^*(\omega) = \lim_{m \rightarrow \infty} \inf_{1 \leq k \leq m} U_k(\omega)$. This gives a measurable function T^* taking values in \mathbb{R} , with $T^* \geq T$, and $E \arctan T^* = E^* \arctan T$ by monotone convergence. For any measurable $U \geq T$, $\arctan U \wedge T^* \geq \arctan T$, and thus $E \arctan U \wedge T^* \geq E^* \arctan T = E \arctan T^*$. However, $U \wedge T^*$ is trivially smaller than T^* , and since both quantities therefore have the same expectation, $\arctan U \wedge T^* = \arctan T^*$ a.s. Hence $T^* \leq U$ a.s., and (i) and (ii) follow. When ET^* exists, it is larger than E^*T by (i) yet smaller by (ii), and thus $ET^* = E^*T$. When $E^*T < \infty$, there exists a measurable $U \geq T$ with $EU^+ < \infty$, where z^+ is the positive part of z . Hence $E(T^*)^+ \leq EU^+$ and ET^* exists. \square

Proof of lemma 6.6. The second inequality in (i) is obvious. If S and $U \geq S + T$ are both measurable, then $U - S \geq T$ and $U - S \geq T^*$ since $U - S$ is also measurable. Hence $U = (S + T)^* \geq S + T^*$ and the second inequality is an equality. Now $(S + T)^* \geq (S_* + T)^* = S_* + T^*$, and we obtain the first inequality. If S is measurable, then $S_* = S^*$ and thus $S_* + T^* = S^* + T^*$. Part (ii) is left as an exercise. Part (iii) follows from the second inequality in (i) after relabeling and rearranging. Part (iv) follows from $S^* - T^* \leq (S - T)^* \leq |S - T|^*$ and then exchanging the roles of S and T . For part (v), it is clear that $(1\{T > c\})^* \geq 1\{T^* > c\}$. If $U \geq 1\{T > c\}$ is measurable, then $S = T^*1\{U \geq 1\} + (T^* \wedge c)1\{U < 1\} \geq T$ and is measurable. Hence $S \geq T^*$, and thus $T^* \leq c$ whenever $U < 1$. This trivially implies $1\{T^* > c\} = 0$ when $U < 1$, and thus $1\{T^* > c\} \leq U$. Part (vi) is left as an exercise.

For part (vii), $(S \vee T)^* \leq S^* \vee T^*$ trivially. Let $U = (S \vee T)^*$ and note that U is measurable and both $U \geq T$ and $U \geq S$. Hence both $U \geq T^*$ and $U \geq S^*$, and thus $(S \vee T)^* \geq S^* \vee T^*$, yielding the desired inequality. The inequality in (viii) is obvious. Assume S is measurable and let $U = (S \wedge T)^*$. Clearly, $U \leq S^* \wedge T^*$. Define $\tilde{T} \equiv U1\{U < S\} + T^*1\{U \geq S\}$; and note that $\tilde{T} \geq T^*$, since if $U < S$, then $T < S$ and thus $U \geq T$. Fix $\omega \in \Omega$. If $U < S$, then $S \wedge \tilde{T} = U$. If $U \geq S$, then $U = S$ since $U \leq S$, and, furthermore, we will now show that $T^* \geq S$. If it were not true, then $T^* < S$ and $U \leq S \wedge T^* < S$, which is clearly a contradiction. Thus when $U \geq S$,

$U = S = S \wedge \tilde{T} \geq S \wedge T^*$. Hence $U = S \wedge \tilde{T} \geq S \wedge T^*$ a.s., and the desired equality in (viii) follows. \square

Proof of lemma 6.7. The first part of (i) is a consequence of the following chain of equalities: $1\{(A \cup B)^*\} = (1\{A \cup B\})^* = (1\{A\} \vee 1\{B\})^* = (1\{A\})^* \vee (1\{B\})^* = 1\{A^*\} \vee 1\{B^*\} = 1\{A^* \cup B^*\}$. The first and fourth equalities follow from (ii) of lemma 6.5, the second and fifth equalities follow directly, and the third equality follows from (vii) of lemma 6.6. For the second part of (i), we have $(A \cap B)_* = \Omega - [(\Omega - A) \cup (\Omega - B)]^* = \Omega - (\Omega - A)^* \cup (\Omega - B)^* = \Omega - (\Omega - A_*) \cup (\Omega - B_*) = A_* \cap B_*$, where the second equality follows from the first part of (i).

The inclusions in part (ii) are obvious. Assume A is measurable. Then (viii) of lemma 6.6 can be used to validate the following string of equalities: $1\{(A \cap B)^*\} = (1\{A \cap B\})^* = (1\{A\} \wedge 1\{B\})^* = (1\{A\})^* \wedge (1\{B\})^* = 1\{A^*\} \wedge 1\{B^*\} = 1\{A^* \cap B^*\}$. Thus $(A \cap B)^* = A^* \cap B^*$. By symmetry, this works whether A or B is measurable. The proof that $(A \cup B)_* = A_* \cup B_*$ when either A or B is measurable is left as an exercise. The proof of part (iii) is also left as an exercise. \square

Proof of lemma 6.8. All of the inequalities follow from the definitions. For the equality in A(i), assume that ϕ is nondecreasing and left-continuous. Define $\phi^{-1}(u) = \inf\{t : \phi(t) \geq u\}$, and note that $\phi(t) > u$ if and only if $t > \phi^{-1}(u)$. Thus, for any $c \in \mathbb{R}$, $1\{\phi(T^*) > c\} = 1\{T^* > \phi^{-1}(c)\} = (1\{T > \phi^{-1}(c)\})^* = (1\{\phi(T) > c\})^* = 1\{[\phi(T)]^* > c\}$. The second and fourth equalities follow from (v) of lemma 6.6. Hence $\phi(T^*) = [\phi(T)]^*$. For the equality in A(ii), assume that ϕ is nondecreasing and left-continuous; and define $\phi^{-1}(u) = \sup\{t : \phi(t) \leq u\}$. Note that $\phi(t) \geq u$ if and only if $t \geq \phi^{-1}(u)$. The proof proceeds in the same manor as for A(i), only part (vi) in lemma 6.6 is used in place of part (v). We leave the proof of part B as an exercise. \square

Proof of lemma 6.11. Clearly, $\liminf T_n^* \leq \limsup T_n^* \leq T^*$. Conversely, $\liminf T_n^* \geq \liminf T_n = T$ and is measurable, and thus $\liminf T_n^* \geq T^*$. Hence $T_n^* \uparrow T^*$. Now $E^*T_n^* = ET_n^* \uparrow ET^*$ by monotone convergence for measurable maps. Note we are allowing $+\infty$ as a possible value for E^*T_n , for some n , or E^*T . \square

Proof of lemma 6.12. Since $|T| \leq |T_n| + |T - T_n|$ for all n , we have $|T - T_n|^* \leq 2S^*$ a.s. Fix $\epsilon > 0$. Since $E^*S < \infty$, there exists a $0 < k < \infty$ so that $E[S^*1\{S^* > k\}] \leq \epsilon/2$. Thus $E^*|T - T_n| \leq Ek \wedge |T - T_n|^* + 2E[S^*1\{S^* > k\}] \rightarrow \epsilon$. The result now follows since ϵ was arbitrary. \square

6.5 Exercises

6.5.1. Show that part (iii) in the definition of σ -field can be replaced, without really changing the definition, by the following: The countable intersection $\bigcap_{j=1}^{\infty} U_j \in \mathcal{A}$ whenever $U_j \in \mathcal{A}$ for all $j \geq 1$.

6.5.2. Show the following:

- (a) For a metric space \mathbb{D} and set $A \in \mathbb{D}$, the closure \overline{A} consists of all limits of sequences $\{x_n\} \in A$.
- (b) Two metrics d_1 and d_2 on a set \mathbb{D} are equivalent if and only if we have the following for any sequence $\{x_j\} \in \mathbb{D}$, as $n, m \rightarrow \infty$: $d_1(x_n, x_m) \rightarrow 0$ if and only if $d_2(x_n, x_m) \rightarrow 0$.
- (c) A function $f : \mathbb{D} \mapsto \mathbb{E}$ between two metric spaces is continuous (in the topological sense) if and only if, for all $x \in \mathbb{D}$ and all sequences $\{x_n\} \in \mathbb{D}$, $f(x_n) \rightarrow f(x)$ whenever $x_n \rightarrow x$.

6.5.3. Verify the implication (ii) \Rightarrow (i) in lemma 6.1.

6.5.4. Show that if a subset K of a metric space is not totally bounded, then it is possible to construct a sequence $\{x_n\} \in K$ which has no Cauchy subsequences.

6.5.5. Show that the following spaces are complete with respect to the uniform metric:

- (a) $C[a, b]$ and $D[a, b]$;
- (b) $UC(T, \rho)$ and $C(\overline{T}, \rho)$, where (T, ρ) is a totally bounded semimetric space;
- (c) $UC(T, \rho)$ and $\ell^\infty(T)$, where T is an arbitrary set.

6.5.6. Show that a uniformly continuous function $f : T \mapsto \mathbb{R}$, where T is totally bounded, has a unique continuous extension to \overline{T} .

6.5.7. Let $C_L[0, 1] \subset C[0, 1]$ be the space of all Lipschitz-continuous functions on $[0, 1]$, and endow it with the uniform metric:

- (a) Show that $C_L[0, 1]$ is dense in $C[0, 1]$.
- (b) Show that no open ball in $C[0, 1]$ is contained in $C_L[0, 1]$.
- (c) Show that $C_L[0, 1]$ is not complete.
- (d) Show that for $0 < c < \infty$, the set

$$\{f \in C_L[0, 1] : |f(x)| \leq c \text{ and } |f(x) - f(y)| \leq c|x - y|, \forall x, y \in [0, 1]\}$$

is compact.

6.5.8. A collection \mathcal{F} of maps $f : \mathbb{D} \mapsto \mathbb{E}$ between metric spaces, with respective Borel σ -fields \mathcal{D} and \mathcal{E} , can generate a (possibly) new σ -field for \mathbb{D} by considering all inverse images $f^{-1}(A)$, for $f \in \mathcal{F}$ and $A \in \mathcal{E}$. Show that the σ -field σ_p generated by the coordinate projections $x \mapsto x(t)$ on

$C[a, b]$ is equal to the Borel σ -field σ_c generated by the uniform norm. Hint: Show first that continuity of the projection maps implies $\sigma_p \subset \sigma_c$. Second, show that the open balls in σ_c can be created from countable set operations on sets in σ_p .

6.5.9. Show that the following metrics generate the product topology on $\mathbb{D} \times \mathbb{E}$, where d and e are the respective metrics for \mathbb{D} and \mathbb{E} :

- (i) $\rho_1((x_1, y_1), (x_2, y_2)) \equiv d(x_1, x_2) + e(y_1, y_2)$.
- (ii) $\rho_2((x_1, y_1), (x_2, y_2)) \equiv \sqrt{d^2(x_1, x_2) + e^2(y_1, y_2)}$.

6.5.10. For any map $T : \Omega \mapsto \bar{\mathbb{R}}$, show that E_*T is the supremum of all EU , where $U \leq T$, $U : \Omega \mapsto \bar{\mathbb{R}}$ measurable and EU exists. Show also that for any set $B \in \Omega$, $P_*(B)$ is the supremum of all $P(A)$, where $A \subset B$ and $A \in \mathcal{A}$.

6.5.11. Use lemma 6.3 to prove lemma 6.4.

6.5.12. Prove parts (ii) and (vi) of lemma 6.6 using parts (i) and (v), respectively.

6.5.13. Let $S, T : \Omega \mapsto \bar{\mathbb{R}}$ be arbitrary. Show the following:

- (a) $|S^* - T_*| \leq |S - T|^* + (S^* - S_*) \wedge (T^* - T_*)$;
- (b) $|S - T|^* \leq (S^* - T_*) \vee (T^* - S_*) \leq |S - T|^* + (S^* - S_*) \wedge (T^* - T_*)$.

6.5.14. Finish the proof of lemma 6.7:

- (a) Show that $(A \cup B)_* = A_* \cup B_*$ if either A or B is measurable.
- (b) Prove part (iii) of the lemma.

6.5.15. Prove part B of lemma 6.8 using the approach outlined in the proof of part A.

6.5.16. Prove the following “converse” to Jensen’s inequality:

LEMMA 6.20 (converse to Jensen’s inequality) *Let $T : \Omega \mapsto \bar{\mathbb{R}}$ be an arbitrary map, with $E^*|T| < \infty$, and assume $\phi : \bar{\mathbb{R}} \mapsto \bar{\mathbb{R}}$ is concave. Then*

- (a) $E_*\phi(T) \leq \phi(E_*T) \wedge \phi(E^*T)$;
- (b) if ϕ is also monotone, $E^*\phi(T) \leq \phi(E_*T) \vee \phi(E^*T)$.

6.5.17. In the proof of proposition 6.15, show that (i) implies (ii).

6.5.18. Show that Hadamard and compact differentiability are equivalent.

6.6 Notes

Much of the material in section 6.2 is an amalgamation of several concepts and presentation styles found in chapter 2 of Billingsley (1986), sections 1.3 and 1.7 of van der Vaart and Wellner (1996), section 18.1 of van der Vaart (1998), and in chapter 1 of both Rudin (1987) and Rudin (1991). A nice proof of the equivalence of the several definitions of compactness can be found in appendix I of Billingsley (1968).

Many of the ideas in section 6.3 come from chapter 1.2 of van der Vaart and Wellner (1996), abbreviated VW hereafter. Lemma 6.3 corresponds to lemma 1.2.1 of VW, lemma 6.4 is given in exercise 1.2.1 of VW, and parts (i), (ii) and (iii) correspond to lemma 1.2.3 of VW. Most of lemma 6.6 is given in lemma 1.2.2 of VW, although the first inequalities in parts (i) and (ii), as well as part (vi), are new. Lemma 6.7 is given in exercise 1.2.15 in VW. Lemmas 6.11 and 6.12 correspond to exercises 1.2.3 and 1.2.4 of VW, respectively, after some modification. Also, lemmas 6.13 and 6.14 correspond to lemmas 1.2.5 and 1.2.6 of VW, respectively.

Much of the material on linear operators can be found in appendix A.1 of Bickel, Klaassen, Ritov and Wellner (1997), and in chapter 2 of Kress (1999). Lemma 6.16 is proposition 7, parts A and B, in appendix A.1 of Bickel, et al (1997). The presentation on functional differentiation is motivated by the first few pages in chapter 3.9 of VW, and the chain rule (lemma 6.19) is lemma 3.9.3 of VW.

7

Stochastic Convergence

In this chapter, we study concepts and theory useful in understanding the limiting behavior of stochastic processes. We begin with a general discussion of stochastic processes in metric spaces. The focus of this discussion is on measurable stochastic processes since most limits of empirical processes in statistical applications are measurable. We next discuss weak convergence both in general and in the specific case of bounded stochastic processes. One of the interesting aspects of the approach we take to weak convergence is that the processes studied need not be measurable except in the limit. This is useful in applications since many empirical processes in statistics are not measurable with respect to the uniform metric. The final section of this chapter considers other modes of convergence, such as in probability and outer almost surely, and their relationships to weak convergence.

7.1 Stochastic Processes in Metric Spaces

In this section, we introduce several important concepts about stochastic processes in metric spaces. Recall that for a stochastic process $\{X(t), t \in T\}$, $X(t)$ is a measurable real random variable for each $t \in T$ on a probability space (Ω, \mathcal{A}, P) . The sample paths of such a process typically reside in the metric space $\mathbb{D} = \ell^\infty(T)$ with the uniform metric. Often, however, when X is viewed as a map from Ω to \mathbb{D} , it is no longer Borel measurable. A classic example of this duality comes from Billingsley (1968, pages 152–153). The example hinges on the fact that there exists a set $H \subset [0, 1]$ which is not a

Borel set. Define the stochastic process $X(t) = 1\{U \leq t\}$, where $t \in [0, 1]$ and U is uniformly distributed on $[0, 1]$. The probability space is (Ω, \mathcal{B}, P) , where $\Omega = [0, 1]$, \mathcal{B} are the Borel sets on $[0, 1]$, and P is the uniform probability measure on $[0, 1]$. A natural metric space for the sample paths of X is $\ell^\infty([0, 1])$. Define the set $A = \cup_{s \in H} B_s(1/2)$, where $B_s(1/2)$ is the uniform open ball of radius $1/2$ around the function $t \mapsto f_s(t) \equiv 1\{t \leq s\}$. Since A is an open set in $\ell^\infty([0, 1])$, and since the uniform distance between f_{s_1} and f_{s_2} is 1 whenever $s_1 \neq s_2$, the set $\{\omega \in \Omega : X(\omega) \in A\}$ equals H . Since H is not a Borel set, X is not Borel measurable.

This lack of measurability is actually the usual state for most of the empirical processes we are interested in studying, especially since most of the time the uniform metric is the natural choice of metric. Much of the associated technical difficulties can be resolved through the use of outer measure and outer expectation as defined in the previous chapter and which we will utilize in our study of weak convergence. In contrast, most of the limiting processes we will be studying are, in fact, Borel measurable. For this reason, a brief study of Borel measurable processes is valuable. The following lemma, for example, provides two ways of establishing equivalence between Borel probability measures. Recall from section 2.2.3 that $BL_1(\mathbb{D})$ is the set of all functions $f : \mathbb{D} \mapsto \mathbb{R}$ bounded by 1 and with Lipschitz norm bounded by 1, i.e., with $|f(x) - f(y)| \leq d(x, y)$ for all $x, y \in \mathbb{D}$. When the choice of metric space is clear by the context, we will sometimes use the abbreviated notation BL_1 as was done in chapter 2. Define also a *vector lattice* $\mathcal{F} \subset C_b(\mathbb{D})$ to be a vector space for which if $f \in \mathcal{F}$ then $f \vee 0 \in \mathcal{F}$. We also say that a set \mathcal{F} of real functions on \mathbb{D} *separates points* of \mathbb{D} if, for any $x, y \in \mathbb{D}$ with $x \neq y$, there exists $f \in \mathcal{F}$ such that $f(x) \neq f(y)$. We are now ready for the lemma:

LEMMA 7.1 *Let L_1 and L_2 be Borel probability measures on a metric space \mathbb{D} . The following are equivalent:*

(i) $L_1 = L_2$.

(ii) $\int f dL_1 = \int f dL_2$ for all $f \in C_b(\mathbb{D})$.

If L_1 and L_2 are also separable, then (i) and (ii) are both equivalent to

(iii) $\int f dL_1 = \int f dL_2$ for all $f \in BL_1$.

Moreover, if L_1 and L_2 are also tight, then (i)–(iii) are all equivalent to

(iv) $\int f dL_1 = \int f dL_2$ for all f in a vector lattice $\mathcal{F} \subset C_b(\mathbb{D})$ that both contains the constant functions and separates points in \mathbb{D} .

The proof is given in section 7.4. We say that two Borel random maps X and X' , with respective laws L and L' , are *versions* of each other if $L = L'$.

In addition to being Borel measurable, most of the limiting stochastic processes of interest are *tight*. A Borel probability measure L on a metric

space \mathbb{D} is tight if for every $\epsilon > 0$, there exists a compact $K \subset \mathbb{D}$ so that $L(K) \geq 1 - \epsilon$. A Borel random map $X : \Omega \mapsto \mathbb{D}$ is tight if its law L is tight. Tightness is equivalent to there being a σ -compact set that has probability 1 under L or X . L or X is *separable* if there is a measurable and separable set which has probability 1. L or X is *Polish* if there is a measurable Polish set having probability 1. Note that tightness, separability and Polishness are all topological properties and do not depend on the metric. Since both σ -compact and Polish sets are also separable, separability is the weakest of the three properties. Whenever we say X has any one of these three properties, we tacetly imply that X is also Borel measurable.

On a complete metric space, tightness, separability and Polishness are equivalent. This equivalence for Polishness and separability follows from the definitions. To see the remaining equivalence, assume L is separable. By completeness, there is a $\mathbb{D}_0 \subset \mathbb{D}$ having probability 1 which is both separable and closed. Fix any $\epsilon \in (0, 1)$. By separability, there exists a sequence $\{x_k\} \in \mathbb{D}_0$ which is dense in \mathbb{D}_0 . For every $\delta > 0$, the union of the balls of radius δ centered on the $\{x_k\}$ covers \mathbb{D}_0 . Hence for every integer $j \geq 1$, there exists a finite collection of balls of radius $1/j$ whose union D_j has probability $\geq 1 - \epsilon/2^j$. Thus the closure of the intersection $\bigcap_{j \geq 1} D_j$ is totally bounded and has probability $\geq 1 - \epsilon$. Since ϵ is arbitrary, L is tight.

For a stochastic process $\{X(t), t \in T\}$, where (T, ρ) is a separable, semi-metric space, there is another meaning for separable. X is *separable* (as a stochastic process) if there exists a countable subset $S \subset T$ and a null set N so that, for each $\omega \notin N$ and $t \in T$, there exists a sequence $\{s_m\} \in S$ with $\rho(s_m, t) \rightarrow 0$ and $|X(s_m, \omega) - X(t, \omega)| \rightarrow 0$. It turns out that many of the empirical processes we will be studying are separable in this sense, even though they are not Borel measurable and therefore cannot satisfy the other meaning for separable. Throughout the remainder of the book, the distinction between these two definitions will either be explicitly stated or made clear by the context.

Most limiting processes X of interest will reside in $\ell^\infty(T)$, where the index set T is often a class of real functions \mathcal{F} with domain equal to the sample space. When such limiting processes are tight, the following lemma demands that X resides on $UC(T, \rho)$, where ρ is some semimetric making T totally bounded, with probability 1:

LEMMA 7.2 *Let X be a Borel measurable random element in $\ell^\infty(T)$. Then the following are equivalent:*

- (i) X is tight.
- (ii) There exists a semimetric ρ making T totally bounded and for which $X \in UC(T, \rho)$ with probability 1.

Furthermore, if (ii) holds for any ρ , then it also holds for the semimetric $\rho_0(s, t) \equiv E \arctan |X(s) - X(t)|$.

The proof is given in section 7.4. A nice feature of tight processes in $\ell^\infty(T)$ is that the laws of such processes are completely defined by their finite-dimensional marginal distributions $(X(t_1), \dots, X(t_k))$, where $t_1, \dots, t_k \in T$ and $k \geq 1$ is an integer:

LEMMA 7.3 *Let X and Y be tight, Borel measurable stochastic processes in $\ell^\infty(T)$. Then the Borel laws of X and Y are equal if and only if all corresponding finite-dimensional marginal distributions are equal.*

Proof. Consider the collection $\mathcal{F} \subset C_b(\mathbb{D})$ of all functions $f : \ell^\infty(T) \mapsto \mathbb{R}$ of the form $f(x) = g(x(t_1), \dots, x(t_k))$, where $g \in C_b(\mathbb{R}^k)$ and $k \geq 1$ is an integer. We leave it as an exercise to show that \mathcal{F} is a vector lattice, an algebra, and separates points of $\ell^\infty(T)$. The desired result now follows from lemma 7.1. \square

While the semimetric ρ_0 defined in lemma 7.2 is always applicable when X is tight, it is frequently not the most convenient semimetric to work with. The family of semimetrics $\rho_p(s, t) \equiv (E|X(s) - X(t)|^p)^{1/(p \vee 1)}$, for some choice of $p \in (0, \infty)$, is sometimes more useful. There is an interesting link between ρ_p and other semimetrics for which lemma 7.2 holds. For a process X in $\ell^\infty(T)$ and a semimetric ρ on T , we say that X is *uniformly ρ -continuous in p th mean* if $E|X(s_n) - X(t_n)|^p \rightarrow 0$ whenever $\rho(s_n, t_n) \rightarrow 0$. The following lemma is a conclusion from lemma 1.5.9 of van der Vaart and Wellner (1996) (abbreviated VW hereafter), and we omit the proof:

LEMMA 7.4 *Let X be a tight Borel measurable random element in $\ell^\infty(T)$, and let ρ be any semimetric for which (i) of lemma 7.2 holds. If X is ρ -continuous in p th mean for some $p \in (0, \infty)$, then (ii) of lemma 7.2 also holds for the semimetric ρ_p .*

Perhaps the most frequently occurring limiting process in $\ell^\infty(T)$ is a *Gaussian* process. A stochastic process $\{X(t), t \in T\}$ is Gaussian if all finite-dimensional marginals $\{X(t_1), \dots, X(t_k)\}$ are multivariate normal. If a Gaussian process X is tight, then by lemma 7.2, there is a semimetric ρ making T totally bounded and for which the sample paths $t \mapsto X(t)$ are uniformly ρ -continuous. An interesting feature of Gaussian processes is that this result implies that the map $t \mapsto X(t)$ is uniformly ρ -continuous in p th mean for all $p \in (0, \infty)$. To see this, take $p = 2$, and note that $|X(s_n) - X(t_n)| \rightarrow 0$ in probability if and only if $E|X(s_n) - X(t_n)|^2 \rightarrow 0$. Thus whenever $\rho(s_n, t_n) \rightarrow 0$, $E|X(s_n) - X(t_n)|^2 \rightarrow 0$ and hence also $E|X(s_n) - X(t_n)|^p \rightarrow 0$ for any $p \in (0, \infty)$ since $X(s_n) - X(t_n)$ is normally distributed for all $n \geq 1$. Lemma 7.4 now implies that tightness of a Gaussian process is equivalent to T being totally bounded by ρ_p with almost all sample paths of X being uniformly ρ_p -continuous for all $p \in (0, \infty)$.

For a general Banach space \mathbb{D} , a Borel measurable random element X on \mathbb{D} is Gaussian if and only if $f(X)$ is Gaussian for every continuous, linear map $f : \mathbb{D} \mapsto \mathbb{R}$. When $\mathbb{D} = \ell^\infty(T)$ for some set T , this definition appears to contradict the definition of Gaussianity given in the preceding

paragraph, since now we are using all continuous linear functionals instead of just linear combinations of coordinate projections. These two definitions are not really reconcilable in general, and so some care must be taking in reading the literature to determine the appropriate context. However, when the process in question is tight, the two definitions are equivalent, as verified in the following proposition:

PROPOSITION 7.5 *Let X be a tight, Borel measurable map into $\ell^\infty(T)$. Then the following are equivalent:*

- (i) *The vector $(X_{t_1}, \dots, X_{t_k})$ is multivariate normal for every finite set $\{t_1, \dots, t_k\} \subset T$.*
- (ii) *$\phi(X)$ is Gaussian for every continuous, linear map $\phi : \ell^\infty(T) \mapsto \mathbb{R}$.*
- (iii) *$\phi(X)$ is Gaussian for every continuous, linear map $\phi : \ell^\infty(T) \mapsto \mathbb{D}$ into any Banach space \mathbb{D} .*

Proof. The proof that (i) \Rightarrow (ii) is given in the proof of lemma 3.9.8 of VW, and we omit the details here. Now assume (ii), and fix any Banach space \mathbb{D} and any continuous, linear map $\phi : \ell^\infty(T) \mapsto \mathbb{D}$. Now for any continuous, linear map $\psi : \mathbb{D} \mapsto \mathbb{R}$, the composition map $\psi \circ \phi : \ell^\infty(T) \mapsto \mathbb{R}$ is continuous and linear, and thus by (ii) we have that $\psi(\phi(X))$ is Gaussian. Since ψ is arbitrary, we have by the definition of a Gaussian process on a Banach space that $\phi(X)$ is Gaussian. Since both \mathbb{D} and ϕ were also arbitrary, conclusion (iii) follows. Finally, (iii) \Rightarrow (i) since multivariate coordinate projections are special examples of continuous, linear maps into Banach spaces. \square

7.2 Weak Convergence

We first discuss the general theory of weak convergence in metric spaces and then discuss results for the special metric space of uniformly bounded functions, $\ell^\infty(T)$, for an arbitrary index set T . This last space is where most—if not all—of the action occurs for statistical applications of empirical processes.

7.2.1 General Theory

The extremely important concept of weak convergence of sequences arises in many areas of statistics. To be as flexible as possible, we allow the probability spaces associated with the sequences to change with n . Let $(\Omega_n, \mathcal{A}_n, P_n)$ be a sequence of probability spaces and $X_n : \Omega_n \mapsto \mathbb{D}$ a sequence of maps. We say that X_n *converges weakly* to a Borel measurable $X : \Omega \mapsto \mathbb{D}$ if

$$(7.1) \quad E^* f(X_n) \rightarrow E f(X), \text{ for every } f \in C_b(\mathbb{D}).$$

If L is the law of X , (7.1) can be reexpressed as

$$E^* f(X_n) \rightarrow \int_{\Omega} f(x) dL(x), \text{ for every } f \in C_b(\mathbb{D}).$$

This weak convergence is denoted $X_n \rightsquigarrow X$ or, equivalently, $X_n \rightsquigarrow L$. Weak convergence is equivalent to “convergence in distribution” and “convergence in law.” By lemma 7.1, this definition of weak convergence ensures that the limiting distributions are unique. Note that the choice of probability spaces $(\Omega_n, \mathcal{A}_n, P_n)$ is important since these dictate the outer expectation used in the definition of weak convergence. In most of the settings discussed in this book, $\Omega_n = \Omega$ for all $n \geq 1$. Some important exceptions to this rule will be discussed in chapter 11. Fortunately, even in those settings where Ω_n does change with n , one can frequently readjust the probability spaces so that the sample spaces are all the same. It is also possible to generalize the concept of weak convergence of sequences to weak convergence of nets as done in VW, but we will restrict ourselves to sequences throughout this book.

The forgoing definition of weak convergence does not obviously appear to be related to convergence of probabilities, but this is in fact true for the probabilities of sets $B \subset \Omega$ which have boundaries δB satisfying $L(\delta B) = 0$. Here and elsewhere, we define the boundary δB of a set B in a topological space to be the closure of B minus the interior of B . Several interesting equivalent formulations of weak convergence on a metric space \mathbb{D} are given in the following portmanteau theorem:

THEOREM 7.6 (*Portmanteau*) *The following are equivalent:*

- (i) $X_n \rightsquigarrow L$;
- (ii) $\liminf P_*(X_n \in G) \geq L(G)$ for every open G ;
- (iii) $\limsup P^*(X_n \in F) \leq L(F)$ for every closed F ;
- (iv) $\liminf E_* f(X_n) \geq \int_{\Omega} f(x) dL(x)$ for every lower semicontinuous f bounded below;
- (v) $\limsup E^* f(X_n) \leq \int_{\Omega} f(x) dL(x)$ for every upper semicontinuous f bounded above;
- (vi) $\lim P^*(X_n \in B) = \lim P_*(X_n \in B) = L(B)$ for every Borel B with $L(\delta B) = 0$;
- (vii) $\liminf E_* f(X_n) \geq \int_{\Omega} f(x) dL(x)$ for every bounded, Lipschitz continuous, nonnegative f .

Furthermore, if L is separable, then (i)–(vii) are also equivalent to

(viii) $\sup_{f \in BL_1} |E^* f(X_n) - Ef(X)| \rightarrow 0$.

The proof is given in section 7.4. Depending on the setting, one or more of these alternative definitions will prove more useful than the others. For example, definition (vi) is probably the most intuitive from a statistical point of view, while definition (viii) is convenient for studying certain properties of the bootstrap.

Another very useful result is the continuous mapping theorem:

THEOREM 7.7 (Continuous mapping) *Let $g : \mathbb{D} \mapsto \mathbb{E}$ be continuous at all points in $\mathbb{D}_0 \subset \mathbb{D}$, where \mathbb{D} and \mathbb{E} are metric spaces. Then if $X_n \rightsquigarrow X$ in \mathbb{D} , with $P_*(X \in \mathbb{D}_0) = 1$, then $g(X_n) \rightsquigarrow g(X)$.*

The proof is given in section 7.4. As mentioned in chapter 2, a common application of this theorem is in the construction of confidence bands based on the supremum distance.

A potential issue is that there may sometimes be more than one choice of metric space \mathbb{D} to work with in a given weak convergence setting. For example, if we are studying weak convergence of the usual empirical process $\sqrt{n}(\hat{F}_n(t) - F(t))$ based on data in $[0, 1]$, we could let \mathbb{D} be either $\ell^\infty([0, 1])$ or $D[0, 1]$. The following lemma tells us that the choice of metric space is generally not a problem. Recall from chapter 6 that for a topological space (X, \mathcal{O}) , the relative topology on $A \subset X$ consists of the open sets $\{A \cap B : B \in \mathcal{O}\}$.

LEMMA 7.8 *Let the metric spaces $\mathbb{D}_0 \subset \mathbb{D}$ have the same metric, and assume X and X_n reside in \mathbb{D}_0 . Then $X_n \rightsquigarrow X$ in \mathbb{D}_0 if and only if $X_n \rightsquigarrow X$ in \mathbb{D} .*

Proof. Since any set $B_0 \in \mathbb{D}_0$ is open if and only if it is of the form $B \cap \mathbb{D}_0$ for some open B in \mathbb{D} , the result follows from part (ii) of the portmanteau theorem. \square

Recall from chapter 2 that a sequence X_n is asymptotically measurable if and only if

$$(7.2) \quad E^* f(X_n) - E_* f(X_n) \rightarrow 0,$$

for all $f \in C_b(\mathbb{D})$. An important, related concept is that of *asymptotic tightness*. A sequence X_n is asymptotically tight if for every $\epsilon > 0$, there is a compact K so that $\liminf P_*(X_n \in K^\delta) \geq 1 - \epsilon$, for every $\delta > 0$, where for a set $A \subset \mathbb{D}$, $A^\delta = \{x \in \mathbb{D} : d(x, A) < \delta\}$ is the “ δ -enlargement” around A . The following lemma tells us that when X_n is asymptotically tight, we can determine asymptotic measurability by verifying (7.2) only for a subset of functions in $C_b(\mathbb{D})$. For the purposes of this lemma, an *algebra* $\mathcal{F} \subset C_b(\mathbb{D})$ is a vector space for which if $f, g \in \mathcal{F}$ then $fg \in \mathcal{F}$.

LEMMA 7.9 *Assume the sequence X_n is asymptotically tight and that (7.2) holds for all f in a subalgebra $\mathcal{F} \subset C_b(\mathbb{D})$ that separates points of \mathbb{D} . Then X_n is asymptotically measurable.*

We omit the proof of this lemma, but it can be found in chapter 1.3 of VW.

When \mathbb{D} is a Polish space and X_n and X are both Borel measurable, tightness of X_n for each $n \geq 1$ plus asymptotic tightness is equivalent to the concept of *uniform tightness* used in the classical theory of weak convergence (see p. 37, Billingsley, 1968). More precisely, a Borel measurable sequence $\{X_n\}$ is uniformly tight if for every $\epsilon > 0$, there is a compact K so that $P(X_n \in K) \geq 1 - \epsilon$ for all $n \geq 1$. The following is a more formal statement of the equivalence we are describing:

LEMMA 7.10 *Assume \mathbb{D} is a Polish space and that the maps X_n and X are Borel measurable. Then $\{X_n\}$ is uniformly tight if and only if X_n is tight for each $n \geq 1$ and $\{X_n\}$ is asymptotically tight.*

The proof is given in section 7.4. Because X_n will typically not be measurable in many of the applications of interest to us, uniform tightness will not prove as useful a concept as asymptotic tightness.

Two good properties of asymptotic tightness are that it does not depend on the metric chosen—only on the topology—and that weak convergence often implies asymptotic tightness. The first of these two properties are verified in the following lemma:

LEMMA 7.11 *X_n is asymptotically tight if and only if for every $\epsilon > 0$ there exists a compact K so that $\liminf P_*(X_n \in G) \geq 1 - \epsilon$ for every open $G \supset K$.*

Proof. Assume first that X_n is asymptotically tight. Fix $\epsilon > 0$, and let the compact set K satisfy $\liminf P_*(X_n \in K^\delta) \geq 1 - \epsilon$ for every $\delta > 0$. If $G \supset K$ is open, then there exists a $\delta_0 > 0$ so that $G \supset K^{\delta_0}$. If this were not true, then there would exist a sequence $\{x_n\} \notin G$ so that $d(x_n, K) \rightarrow 0$. This implies the existence of a sequence $\{y_n\} \in K$ so that $d(x_n, y_n) \rightarrow 0$. Thus, since K is compact and the complement of G is closed, there is a subsequence n' and a point $y \notin G$ so that $d(y_{n'}, y) \rightarrow 0$, but this is impossible. Hence $\liminf P_*(X_n \in G) \geq 1 - \epsilon$. Now assume that X_n satisfies the alternative definition. Fix $\epsilon > 0$, and let the compact set K satisfy $\liminf P_*(X_n \in G) \geq 1 - \epsilon$ for every open $G \supset K$. For every $\delta > 0$, K^δ is an open set. Thus $\liminf P_*(X_n \in K^\delta) \geq 1 - \epsilon$ for every $\delta > 0$. \square

The second good property of asymptotic tightness is given in the second part of the following lemma, the first part of which gives the necessity of asymptotic measurability for weakly convergent sequences:

LEMMA 7.12 *Assume $X_n \rightsquigarrow X$. Then*

- (i) X_n is asymptotically measurable.
- (ii) X_n is asymptotically tight if and only if X is tight.

Proof. For part (i), fix $f \in C_b(\mathbb{D})$. Note that weak convergence implies both $E^*f(X_n) \rightarrow Ef(X)$ and $E_*f(X_n) = -E^*[-f(X_n)] \rightarrow -E[-f(X)] =$

$Ef(X)$, and the desired result follows since f is arbitrary. For part (ii), fix $\epsilon > 0$. Assume X is tight, and choose a compact K so that $P(X \in K) \geq 1 - \epsilon$. By part (ii) of the portmanteau theorem, $\liminf P_*(X_n \in K^\delta) \geq P(X \in K^\delta) \geq 1 - \epsilon$ for every $\delta > 0$. Hence X_n is asymptotically tight. Now assume that X_n is asymptotically tight, fix $\epsilon > 0$, and choose a compact K so that $\liminf P_*(X_n \in K^\delta) \geq 1 - \epsilon$ for every $\delta > 0$. By part (iii) of the portmanteau theorem, $P(X \in \overline{K^\delta}) \geq \limsup P^*(X_n \in \overline{K^\delta}) \geq \liminf P_*(X_n \in \overline{K^\delta}) \geq 1 - \epsilon$. By letting $\delta \downarrow 0$, we obtain that X is tight. \square

Prohorov's theorem (given below) tells us that asymptotic measurability and asymptotic tightness together almost gives us weak convergence. This "almost-weak-convergence" is *relative compactness*. A sequence X_n is relatively compact if every subsequence $X_{n'}$ has a further subsequence $X_{n''}$ which converges weakly to a tight Borel law. Weak convergence happens when all of the limiting Borel laws are the same. Note that when all of the limiting laws assign probability one to a fixed Polish space, there is a converse of Prohorov's theorem, that relative compactness of X_n implies asymptotic tightness of X_n (and hence also uniform tightness). Details of this result are discussed in chapter 1.12 of VW, but we do not pursue it further here.

THEOREM 7.13 (Prohorov's theorem) *If the sequence X_n is asymptotically measurable and asymptotically tight, then it has a subsequence $X_{n'}$ that converges weakly to a tight Borel law.*

The proof, which we omit, is given in chapter 1.3 of VW. Note that the conclusion of Prohorov's theorem does not state that X_n is relatively compact, and thus it appears as if we have broken our earlier promise. However, if X_n is asymptotically measurable and asymptotically tight, then every subsequence $X_{n'}$ is also asymptotically measurable and asymptotically tight. Thus repeated application of Prohorov's theorem does indeed imply relative compactness of X_n .

A natural question to ask at this juncture is: under what circumstances does asymptotic measurability and/or tightness of the marginal sequences X_n and Y_n imply asymptotic measurability and/or tightness of the joint sequence (X_n, Y_n) ? This question is answered in the following lemma:

LEMMA 7.14 *Let $X_n : \Omega_n \mapsto \mathbb{D}$ and $Y_n : \Omega_n \mapsto \mathbb{E}$ be sequences of maps. Then the following are true:*

- (i) *X_n and Y_n are both asymptotically tight if and only if the same is true for the joint sequence $(X_n, Y_n) : \Omega_n \mapsto \mathbb{D} \times \mathbb{E}$.*
- (ii) *Asymptotically tight sequences X_n and Y_n are both asymptotically measurable if and only if $(X_n, Y_n) : \Omega_n \mapsto \mathbb{D} \times \mathbb{E}$ is asymptotically measurable.*

Proof. Let d and e be the metrics for \mathbb{D} and \mathbb{E} , respectively, and let $\mathbb{D} \times \mathbb{E}$ be endowed with the product topology. Now note that a set in $\mathbb{D} \times \mathbb{E}$ of the form $K_1 \times K_2$ is compact if and only if K_1 and K_2 are both compact. Let $\pi_j K$, $j = 1, 2$, be the projections of K onto \mathbb{D} and \mathbb{E} , respectively. To be precise, the projection $\pi_1 K$ consists of all $x \in \mathbb{D}$ such that $(x, y) \in K$ for some $y \in \mathbb{E}$, and $\pi_2 K$ is analogously defined for \mathbb{E} . We leave it as an exercise to show that $\pi_1 K$ and $\pi_2 K$ are both compact. It is easy to see that K is contained in $\pi_1 K \times \pi_2 K$. Using the product space metric $\rho((x_1, y_1), (x_2, y_2)) = d(x_1, x_2) \vee e(y_1, y_2)$ (one of several metrics generating the product topology), we now have for $K_1 \in \mathbb{D}$ and $K_2 \in \mathbb{E}$ that $(K_1 \times K_2)^\delta = K_1^\delta \times K_2^\delta$. Part (i) now follows from the definition of asymptotic tightness.

Let $\pi_1 : \mathbb{D} \times \mathbb{E} \mapsto \mathbb{D}$ be the projection onto the first coordinate, and note that π_1 is continuous. Thus for any $f \in C_b(\mathbb{D})$, $f \circ \pi_1 \in C_b(\mathbb{D} \times \mathbb{E})$. Hence joint asymptotic measurability of (X_n, Y_n) implies asymptotic measurability of X_n by the definition of asymptotic measurability. The same argument holds for Y_n . The difficult part of proving part (ii) is the implication that asymptotic tightness plus asymptotic measurability of both marginal sequences yields asymptotic measurability of the joint sequence. We omit this part of the proof, but it can be found in chapter 1.4 of VW. \square

A very useful consequence of lemma 7.14 is Slutsky's theorem. Note that the proof of Slutsky's theorem (given below) also utilizes both Prohorov's theorem and the continuous mapping theorem.

THEOREM 7.15 (Slutsky's theorem) *Suppose $X_n \rightsquigarrow X$ and $Y_n \rightsquigarrow c$, where X is separable and c is a fixed constant. Then the following are true:*

(i) $(X_n, Y_n) \rightsquigarrow (X, c)$.

(ii) *If X_n and Y_n are in the same metric space, then $X_n + Y_n \rightsquigarrow X + c$.*

(iii) *Assume in addition that the Y_n are scalars. Then whenever $c \in \mathbb{R}$, $Y_n X_n \rightsquigarrow cX$. Also, whenever $c \neq 0$, $X_n/Y_n \rightsquigarrow X/c$.*

Proof. By completing the metric space for X , we can without loss of generality assume that X is tight. Thus by lemma 7.14, (X_n, Y_n) is asymptotically tight and asymptotically measurable. Thus by Prohorov's theorem, all subsequences of (X_n, Y_n) have further subsequences which converge to tight limits. Since these limit points have marginals X and c , and since the marginals in this case completely determine the joint distribution, we have that all limiting distributions are uniquely determined as (X, c) . This proves part (i). Parts (ii) and (iii) now follow from the continuous mapping theorem. \square

7.2.2 Spaces of Bounded Functions

Now we consider the setting where the X_n are stochastic processes with index set T . The natural metric space for weak convergence in this setting is $\ell^\infty(T)$. A nice feature of this setting is the fact that asymptotic measurability of X_n follows from asymptotic measurability of $X_n(t)$ for each $t \in T$:

LEMMA 7.16 *Let the sequence of maps X_n in $\ell^\infty(T)$ be asymptotically tight. Then X_n is asymptotically measurable if and only if $X_n(t)$ is asymptotically measurable for each $t \in T$.*

Proof. Let $f_t : \ell^\infty(T) \mapsto \mathbb{R}$ be the marginal projection at $t \in T$, i.e., $f_t(x) = x(t)$ for any $x \in \ell^\infty(T)$. Since each f_t is continuous, asymptotic measurability of X_n implies asymptotic measurability of $X_n(t)$ for each $t \in T$. Now assume that $X_n(t)$ is asymptotically measurable for each $t \in T$. Then lemma 7.14 implies asymptotic measurability for all finite-dimensional marginals $(X_n(t_1), \dots, X_n(t_k))$. Consequently, $f(X_n)$ is asymptotically measurable for all $f \in \mathcal{F}$, for the subset of $C_b(\mathbb{D})$ defined in the proof of lemma 7.3 given above in section 7.2, where $\mathbb{D} = \ell^\infty(T)$. Since \mathcal{F} is an algebra that separates points in $\ell^\infty(T)$, asymptotic measurability of X_n follows from lemma 7.9. \square

We now verify that convergence of finite dimensional distributions plus asymptotic tightness is equivalent to weak convergence in $\ell^\infty(T)$:

THEOREM 7.17 *The sequence X_n converges to a tight limit in $\ell^\infty(T)$ if and only if X_n is asymptotically tight and all finite-dimensional marginals converge weakly to limits. Moreover, if X_n is asymptotically tight and all of its finite-dimensional marginals $(X_n(t_1), \dots, X_n(t_k))$ converge weakly to the marginals $(X(t_1), \dots, X(t_k))$ of a stochastic process X , then there is a version of X such that $X_n \rightsquigarrow X$ and X resides in $UC(T, \rho)$ for some semimetric ρ making T totally bounded.*

Proof. The result that “asymptotic tightness plus convergence of finite-dimensional distributions implies weak convergence” follows from Prohorov’s theorem and lemmas 7.9 and 7.1, using the vector lattice and subalgebra $\mathcal{F} \subset C_b(\mathbb{D})$ defined above in the proof of lemma 7.3. The implication in the opposite direction follows easily from lemma 7.12 and the continuous mapping theorem. Now assume X_n is asymptotically tight and that all finite-dimensional distributions of X_n converge to those of a stochastic process X . By asymptotic tightness of X_n , the probability that a version of X lies in some σ -compact $K \subset \ell^\infty(T)$ is one. By theorem 6.2, $K \subset UC(T, \rho)$ for some semimetric ρ making T totally bounded. \square

Recall theorem 2.1 and the condition (2.6) from chapter 2. When condition (2.6) holds for every $\epsilon > 0$, then we say that the sequence X_n is *asymptotically uniformly ρ -equicontinuous in probability*. We are now in a position to prove theorem 2.1. Note that the statement of this theorem is slightly in-

formal: conditions (i) and (ii) of the theorem actually imply that $X_n \rightsquigarrow X'$ in $\ell^\infty(T)$ for some tight version X' of X . Recall that $\|x\|_T \equiv \sup_{t \in T} |x(t)|$.

Proof of theorem 2.1. First assume $X_n \rightsquigarrow X$ in $\ell^\infty(T)$, where X is tight. Convergence of all finite-dimensional distributions follows from the continuous mapping theorem. Now by theorem 6.2, $P(X \in UC(T, \rho))$ for some semimetric ρ making T totally bounded. Hence for every $\eta > 0$, there exists some compact $K \subset UC(T, \rho)$ so that

$$(7.3) \quad \liminf_{n \rightarrow \infty} P_*(X_n \in K^\delta) \geq 1 - \eta, \text{ for all } \delta > 0.$$

Fix $\epsilon, \eta > 0$, and let the compact set K satisfy (7.3). By theorem 6.2, there exists a $\delta_0 > 0$ so that $\sup_{x \in K} \sup_{s, t: \rho(s, t) < \delta_0} |x(s) - x(t)| \leq \epsilon/3$. Now

$$\begin{aligned} & P^* \left[\sup_{s, t \in T: \rho(s, t) < \delta_0} |X_n(s) - X_n(t)| > \epsilon \right] \\ & \leq P^* \left[\sup_{s, t \in T: \rho(s, t) < \delta_0} |X_n(s) - X_t(t)| > \epsilon, X_n \in K^{\epsilon/3} \right] + P^*(X_n \notin K^{\epsilon/3}) \\ & \equiv E_n \end{aligned}$$

satisfies $\limsup_{n \rightarrow \infty} E_n \leq \eta$, since if $x \in K^{\epsilon/3}$ then $\sup_{s, t \in T: \rho(s, t) < \delta_0} |x(s) - x(t)| < \epsilon$. Thus X_n is asymptotically uniformly ρ -continuous in probability since ϵ and η were arbitrary.

Now assume that conditions (i) and (ii) of the theorem hold. Lemma 7.18 below, the proof of which is given in section 7.4, yields that X_n is asymptotically tight. Thus the desired weak converge of X_n follows from theorem 7.17 above. \square

LEMMA 7.18 *Assume conditions (i) and (ii) of theorem 2.1 hold. Then X_n is asymptotically tight.*

The proof of theorem 2.1 verifies that whenever $X_n \rightsquigarrow X$ and X is tight, any semimetric ρ defining a σ -compact set $UC(T, \rho)$ such that $P(X \in UC(T, \rho)) = 1$ will also result in X_n being uniformly ρ -equicontinuous in probability. What is not clear at this point is the converse, that any semimetric ρ_* which enables uniform asymptotic equicontinuity of X_n will also define a σ -compact set $UC(T, \rho_*)$ wherein X resides with probability 1. The following theorem shows that, in fact, any semimetric which works for one of these implications will work for the other:

THEOREM 7.19 *Assume $X_n \rightsquigarrow X$ in $\ell^\infty(T)$, and let ρ be a semimetric making (T, ρ) totally bounded. Then the following are equivalent:*

- (i) X_n is asymptotically uniformly ρ -equicontinuous in probability.
- (ii) $P(X \in UC(T, \rho)) = 1$.

Proof. If we assume (ii), then (i) will follow by arguments given in the proof of theorem 2.1 above. Now assume (i). For $x \in \ell^\infty(T)$, define $M_\delta(x) \equiv \sup_{s,t \in T: \rho(s,t) < \delta} |x(s) - x(t)|$. Note that if we restrict δ to $(0, 1)$ then $x \mapsto M_{(\cdot)}(x)$, as a map from $\ell^\infty(T)$ to $\ell^\infty((0, 1))$, is continuous since $|M_\delta(x) - M_\delta(y)| \leq 2\|x - y\|_T$ for all $\delta \in (0, 1)$. Hence $M_{(\cdot)}(X_n) \rightsquigarrow M_{(\cdot)}(X)$ in $\ell^\infty((0, 1))$. Condition (i) now implies that there exists a positive sequence $\delta_n \downarrow 0$ so that $\mathbb{P}^*(M_{\delta_n}(X_n) > \epsilon) \rightarrow 0$ for every $\epsilon > 0$. Hence $M_{\delta_n}(X) \rightsquigarrow 0$. This implies (ii) since X is tight by theorem 2.1. \square

An interesting consequence of theorems 2.1 and 7.19, in conjunction with lemma 7.4, happens when $X_n \rightsquigarrow X$ in $\ell^\infty(T)$ and X is a tight Gaussian process. Recall from section 7.1 the semimetric $\rho_p(s, t) \equiv (E|X(s) - X(t)|^p)^{1/(p \vee 1)}$, for any $p \in (0, \infty)$. Then for any $p \in (0, \infty)$, (T, ρ_p) is totally bounded, the sample paths of X are ρ_p -continuous, and, furthermore, X_n is asymptotically uniformly ρ_p -equicontinuous in probability. While any value of $p \in (0, \infty)$ will work, the choice $p = 2$ (the “standard deviation” metric) is often the most convenient to work with.

We now point out an equivalent condition for X_n to be asymptotically uniformly ρ -equicontinuous in probability. This new condition, which is expressed in the following lemma, is sometimes easier to verify for certain settings (one of which occurs in the next chapter):

LEMMA 7.20 *Let X_n be a sequence of stochastic processes indexed by T . Then the following are equivalent:*

- (i) *There exists a semimetric ρ making T totally bounded and for which X_n is uniformly ρ -equicontinuous in probability.*
- (ii) *For every $\epsilon, \eta > 0$, there exists a finite partition $T = \cup_{i=1}^k T_i$ such that*

$$(7.4) \quad \limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\sup_{1 \leq i \leq k} \sup_{s, t \in T_i} |X_n(s) - X_n(t)| > \epsilon \right) < \eta.$$

The proof is given in section 4 below.

7.3 Other Modes of Convergence

Recall the definitions of convergence in probability and outer almost surely, for arbitrary maps $X_n : \Omega \mapsto \mathbb{D}$, as defined in chapter 2. We now introduce two additional modes of convergence which can be useful in some settings. X_n *converges almost uniformly* to X if, for every $\epsilon > 0$, there exists a measurable set A such that $\mathbb{P}(A) \geq 1 - \epsilon$ and $d(X_n, X) \rightarrow 0$ uniformly on A . X_n *converges almost surely* to X if $\mathbb{P}_*(\lim_{n \rightarrow \infty} d(X_n, X) = 0) = 1$. Note that an important distinction between almost sure and outer almost sure convergence is that, in the latter mode, there must exist a measurable majorant of $d(X_n, X)$ which goes to zero. This distinction is quite important

because it can be shown that almost sure convergence does not in general imply convergence in probability when $d(X_n, X)$ is not measurable. For this reason, we do not use the almost sure convergence mode in this book except rarely. One of those rare times is in exercise 7.5.7, another is in proposition 7.22 below which will be used in chapter 10. The following lemma characterizes the relationships among the three remaining modes:

LEMMA 7.21 *Let $X_n, X : \Omega \mapsto \mathbb{D}$ be maps with X Borel measurable. Then*

- (i) $X_n \xrightarrow{\text{as}^*} X$ implies $X_n \xrightarrow{\text{P}} X$.
- (ii) $X_n \xrightarrow{\text{P}} X$ if and only if every subsequence $X_{n'}$ has a further subsequence $X_{n''}$ such that $X_{n''} \xrightarrow{\text{as}^*} X$.
- (iii) $X_n \xrightarrow{\text{as}^*} X$ if and only if X_n converges almost uniformly to X if and only if $\sup_{m \geq n} d(X_m, X) \xrightarrow{\text{P}} 0$.

The proof is given section 7.4. Since almost uniform convergence and outer almost sure convergence are equivalent for sequences, we will not use the almost uniform mode very much.

The following proposition gives a connection between almost sure convergence and convergence in probability. We need this proposition for a continuous mapping result for bootstrapped processes presented in chapter 10:

PROPOSITION 7.22 *Let $X_n, Y_n : \Omega \mapsto \mathbb{D}$ be maps with Y_n measurable. Suppose every subsequence n' has a further subsequence n'' such that $X_{n''} \rightarrow 0$ almost surely. Suppose also that $d(X_n, Y_n) \xrightarrow{\text{P}} 0$. Then $X_n \xrightarrow{\text{P}} 0$.*

Proof. For every subsequence n' there exists a further subsequence n'' such that both $X_{n''} \rightarrow 0$ and $d(X_{n''}, Y_{n''})^* \rightarrow 0$ almost surely for some versions $d(X_{n''}, Y_{n''})^*$. Since $d(Y_n, 0) \leq d(X_n, 0) + d(X_n, Y_n)^*$, we have that $Y_{n''} \rightarrow 0$ almost surely. But this implies $Y_{n''} \xrightarrow{\text{as}^*} 0$ since the Y_n are measurable. Since the subsequence n' was arbitrary, we now have that $Y_n \xrightarrow{\text{P}} 0$. Thus $X_n \xrightarrow{\text{P}} 0$ since $d(X_n, 0) \leq d(Y_n, 0) + d(X_n, Y_n)$. \square

The next lemma describes several important relationships between weak convergence and convergence in probability. Before presenting it, we need to extend the definition of convergence in probability, in the setting where the limit is a constant, to allow the probability spaces involved to change with n as is already permitted for weak convergence. We denote this modified convergence $X_n \xrightarrow{\text{P}} c$, and distinguish it from the previous form of convergence in probability only by context.

LEMMA 7.23 *Let $X_n, Y_n : \Omega_n \mapsto \mathbb{D}$ be maps, $X : \Omega \mapsto \mathbb{D}$ be Borel measurable, and $c \in \mathbb{D}$ be a constant. Then*

- (i) If $X_n \rightsquigarrow X$ and $d(X_n, Y_n) \xrightarrow{\text{P}} 0$, then $Y_n \rightsquigarrow X$.

(ii) $X_n \xrightarrow{P} X$ implies $X_n \rightsquigarrow X$.

(iii) $X_n \xrightarrow{P} c$ if and only if $X_n \rightsquigarrow c$.

Proof. We first prove (i). Let $F \subset \mathbb{D}$ be closed, and fix $\epsilon > 0$. Then $\limsup_{n \rightarrow \infty} P^*(Y_n \in F) = \limsup_{n \rightarrow \infty} P^*(Y_n \in F, d(X_n, Y_n)^* \leq \epsilon) \leq \limsup_{n \rightarrow \infty} P^*(X_n \in \overline{F^\epsilon}) \leq P(X \in \overline{F^\epsilon})$. The result follows by letting $\epsilon \downarrow 0$. Now assume $X_n \xrightarrow{P} X$. Since $X \rightsquigarrow X$, $d(X, X_n) \xrightarrow{P} 0$ implies $X_n \rightsquigarrow X$ by (i), thus (ii) follows. We now prove (iii). $X_n \xrightarrow{P} c$ implies $X_n \rightsquigarrow c$ by (ii). Now assume $X_n \rightsquigarrow c$, and fix $\epsilon > 0$. Note that $P^*(d(X_n, c) \geq \epsilon) = P^*(X_n \notin B(c, \epsilon))$, where $B(c, \epsilon)$ is the open ϵ -ball around $c \in \mathbb{D}$. By the portmanteau theorem, $\limsup_{n \rightarrow \infty} P^*(X_n \notin B(c, \epsilon)) \leq P(X \notin B(c, \epsilon)) = 0$. Thus $X_n \xrightarrow{P} c$ since ϵ is arbitrary, and (iii) follows. \square

We now present a generalized continuous mapping theorem that allows for sequences of maps g_n which converge to g in a fairly general sense. In the exercises below, we consider an instance of this where one is interested in maximizing a stochastic process $\{X_n(t), t \in T\}$ over an ‘‘approximation’’ T_n of a subset $T_0 \subset T$. As a specific motivation, suppose T is high dimensional. The computational burden of computing the supremum of $X_n(t)$ over T may be reduced by choosing a finite mesh T_n which closely approximates T .

THEOREM 7.24 (Extended continuous mapping). *Let $\mathbb{D}_n \subset \mathbb{D}$ and $g_n : \mathbb{D}_n \mapsto \mathbb{E}$ satisfy the following: if $x_n \rightarrow x$ with $x_n \in \mathbb{D}_n$ for all $n \geq 1$ and $x \in \mathbb{D}_0$, then $g_n(x_n) \rightarrow g(x)$, where $\mathbb{D}_0 \subset \mathbb{D}$ and $g : \mathbb{D}_0 \mapsto \mathbb{E}$. Let X_n be maps taking values in \mathbb{D}_n , and let X be Borel measurable and separable. Then*

(i) $X_n \rightsquigarrow X$ implies $g_n(X_n) \rightsquigarrow g(X)$.

(ii) $X_n \xrightarrow{P} X$ implies $g_n(X_n) \xrightarrow{P} g(X)$.

(iii) $X_n \xrightarrow{\text{as}^*} X$ implies $g_n(X_n) \xrightarrow{\text{as}^*} g(X)$.

The proof can be found in chapter 1.11 of VW, and we omit it here. It is easy to see that if $g : \mathbb{D} \mapsto \mathbb{E}$ is continuous at all points in \mathbb{D}_0 , and if we set $g_n = g$ and $\mathbb{D}_n = \mathbb{D}$ for all $n \geq 1$, then the standard continuous mapping theorem (theorem 7.7), specialized to the setting where X is separable, is a corollary of part (i) of the above theorem.

The following theorem gives another kind of continuous mapping result for sequences which converge in probability and outer almost surely. When X is separable, the conclusions of this theorem are a simple corollary of theorem 7.24.

THEOREM 7.25 *Let $g : \mathbb{D} \mapsto \mathbb{E}$ be continuous at all points in $\mathbb{D}_0 \subset \mathbb{D}$, and let X be Borel measurable with $P_*(X \in \mathbb{D}_0) = 1$. Then*

(i) $X_n \xrightarrow{P} X$ implies $g(X_n) \xrightarrow{P} g(X)$.

(ii) $X_n \xrightarrow{\text{as}^*} X$ implies $g(X_n) \xrightarrow{\text{as}^*} g(X)$.

Proof. Assume $X_n \xrightarrow{P} X$, and fix $\epsilon > 0$. Define B_k to be all $x \in \mathbb{D}$ such that the $1/k$ -ball around x contains points y and z with $e(g(y), g(z)) > \epsilon$. Part of the proof of theorem 7.7 in section 7.4 verifies that B_k is open. It is clear that B_k decreases as k increases. Furthermore, $P(X \in B_k) \downarrow 0$, since every point in $\bigcap_{k=1}^{\infty} B_k$ is a point of discontinuity of g . Now the outer probability that $e(g(X_n), g(X)) > \epsilon$ is bounded above by the outer probability that either $X \in B_k$ or $d(X_n, X) \geq 1/k$. But this last outer probability converges to $P^*(X \in B_k)$ since $d(X_n, X) \xrightarrow{P} 0$. Part (i) now follows by letting $k \downarrow 0$ and noting that ϵ was arbitrary. Assume that $X_n \xrightarrow{\text{as}^*} X$. Note that a minor modification of the proof of part (i) verifies that $\sup_{m \geq n} d(X_m, X) \xrightarrow{P} 0$ implies $\sup_{m \geq n} e(g(X_n), g(X)) \xrightarrow{P} 0$. Now part (iii) of lemma 7.21 yields that $X_n \xrightarrow{\text{as}^*} X$ implies $g(X_n) \xrightarrow{\text{as}^*} g(X)$. \square

We now present a useful outer almost sure representation result for weak convergence. Such representations allow the conversion of certain weak convergence problems into problems about convergence of fixed sequences. We give an illustration of this approach in the proof of proposition 7.27 below.

THEOREM 7.26 *Let $X_n : \Omega_n \mapsto \mathbb{D}$ be a sequence of maps, and let X_∞ be Borel measurable and separable. If $X_n \rightsquigarrow X_\infty$, then there exists a probability space $(\tilde{\Omega}, \tilde{\mathcal{A}}, \tilde{P})$ and maps $\tilde{X}_n : \tilde{\Omega} \mapsto \mathbb{D}$ with*

(i) $\tilde{X}_n \xrightarrow{\text{as}^*} \tilde{X}_\infty$;

(ii) $E^* f(\tilde{X}_n) = E^* f(X_n)$, for every bounded $f : \mathbb{D} \mapsto \mathbb{R}$ and all $1 \leq n \leq \infty$.

Moreover, \tilde{X}_n can be chosen to be equal to $X_n \circ \phi_n$, for all $1 \leq n \leq \infty$, where the $\phi_n : \tilde{\Omega} \mapsto \Omega_n$ are measurable and perfect maps and $P_n = \tilde{P} \circ \phi_n$.

The proof can be found in chapter 1.10 of VW, and we omit it here. Recall the definition of perfect maps from chapter 6. In the setting of the above theorem, if the \tilde{X}_n are constructed from the perfect maps ϕ_n , then $[f(\tilde{X}_n)]^* = [f(X_n)]^* \circ \phi_n$ for all bounded $f : \mathbb{D} \mapsto \mathbb{R}$. Thus the equivalence between \tilde{X}_n and X_n can be made much stronger than simply equivalence in law.

The following proposition can be useful in studying weak convergence of certain statistics which can be expressed as stochastic integrals. For example, the Wilcoxon statistic can be expressed in this way. The proof of the proposition provides the illustration of theorem 7.26 promised above.

PROPOSITION 7.27 *Let $X_n, G_n \in D[a, b]$ be stochastic processes with $X_n \rightsquigarrow X$ and $G_n \xrightarrow{P} G$ in $D[a, b]$, where X is bounded with continuous*

sample paths, G is fixed, and G_n and G have total variation bounded by $K < \infty$. Then $\int_a^{(\cdot)} X_n(s) dG_n(s) \rightsquigarrow \int_a^{(\cdot)} X(s) dG(s)$ in $D[a, b]$.

Proof. First, Slutsky's theorem and lemma 7.23 establish that $(X_n, G_n) \rightsquigarrow (X, G)$. Next, theorem 7.26 tells us that there exists a new probability space and processes $\tilde{X}_n, \tilde{X}, \tilde{G}_n$ and \tilde{G} which have the same outer integrals for bounded functions as X_n, X, G_n and G , respectively, but which also satisfy $(\tilde{X}_n, \tilde{G}_n) \xrightarrow{\text{as}^*} (\tilde{X}, \tilde{G})$. For each integer $m \geq 1$, define $t_j = a + (b - a)j/m$, $j = 0, \dots, m$; let

$$M_m \equiv \max_{1 \leq j \leq m} \sup_{s, t \in (t_{j-1}, t_j]} |\tilde{X}(s) - \tilde{X}(t)|;$$

and define $\tilde{X}_m \in D[a, b]$ such that $\tilde{X}_m(a) = \tilde{X}(a)$ and $\tilde{X}_m(t) \equiv \sum_{j=1}^m 1\{t_{j-1} < t \leq t_j\} \tilde{X}(t_j)$, for $t \in (a, b]$. Note that for integrals over the range $(a, t]$, for $t \in [a, b]$, we define the value of the integral to be zero when $t = a$ since $(a, a]$ is the null set. We now have, for any $t \in [a, b]$, that

$$\begin{aligned} & \left| \int_a^t \tilde{X}_n(s) d\tilde{G}_n(s) - \int_a^t \tilde{X}(s) d\tilde{G}(s) \right| \\ & \leq \int_a^b \left| \tilde{X}_n(s) - \tilde{X}(s) \right| \times |d\tilde{G}_n(s)| + \int_a^b \left| \tilde{X}_m(s) - \tilde{X}(s) \right| \times |d\tilde{G}_n(s)| \\ & \quad + \left| \int_a^t \tilde{X}_m(s) \left\{ d\tilde{G}_n(s) - d\tilde{G}(s) \right\} \right| \\ & \leq K \left(\|\tilde{X}_n - \tilde{X}\|_{[a, b]} + M_m \right) \\ & \quad + \left| \sum_{j=1}^m \tilde{X}(t_j) \int_{(t_{j-1}, t_j] \cap (a, t]} \left\{ d\tilde{G}_n(s) - d\tilde{G}(s) \right\} \right| \\ & \leq K \left(\|\tilde{X}_n - \tilde{X}\|_{[a, b]}^* + M_m \right) + m \left(\|\tilde{X}\|_{[a, b]} \times \|\tilde{G}_n - \tilde{G}\|_{[a, b]}^* \right) \\ & \equiv E_n(m). \end{aligned}$$

Note that $E_n(m)$ is measurable and $\rightarrow 0$ almost surely. Define D_n to be the infimum of $E_n(m)$ over all integers $m \geq 1$. Since $D_n \xrightarrow{\text{as}^*} 0$ and D_n is measurable, we have that $\int_a^{(\cdot)} \tilde{X}_n(s) d\tilde{G}_n(s) \xrightarrow{\text{as}^*} \int_a^{(\cdot)} \tilde{X}(s) d\tilde{G}(s)$. Choose any $f \in C_b(D[a, b])$, and note that the map $(x, y) \mapsto f\left(\int_a^{(\cdot)} x(s) dy(s)\right)$, for $x, y \in D[a, b]$ with the total variation of y bounded, is bounded. Thus

$$\begin{aligned}
\mathbb{E}^* f \left(\int_a^{(\cdot)} X_n(s) dG_n(s) \right) &= \mathbb{E}^* f \left(\int_a^{(\cdot)} \tilde{X}_n(s) d\tilde{G}_n(s) \right) \\
&\rightarrow \mathbb{E} f \left(\int_a^{(\cdot)} \tilde{X}(s) d\tilde{G}(s) \right) \\
&= \mathbb{E} f \left(\int_a^{(\cdot)} X(s) dG(s) \right).
\end{aligned}$$

Since this convergence holds for all $f \in C_b(D[a, b])$, the desired result now follows. \square

We give one more result before closing this chapter. The result applies to certain weak convergence settings involving questions that are easier to answer for measurable maps. The following lemma shows that a nonmeasurable, weakly convergent sequence X_n is usually quite close to a measurable sequence Y_n :

LEMMA 7.28 *Let $X_n : \Omega_n \mapsto \mathbb{D}$ be a sequence of maps. If $X_n \rightsquigarrow X$, where X is Borel measurable and separable, then there exists a Borel measurable sequence $Y_n : \Omega_n \mapsto \mathbb{D}$ with $d(X_n, Y_n) \xrightarrow{P} 0$.*

The proof can be found in chapter 1.10 of VW, and we omit it here.

7.4 Proofs

Proof of lemma 7.1. Clearly (i) implies (ii). Now assume (ii). For every open $G \subset \mathbb{D}$, define the sequence of functions $f_m(x) = [md(x, \mathbb{D} - G)] \wedge 1$, for integers $m \geq 1$, and note that each f_m is bounded and Lipschitz continuous and $f_m \uparrow 1\{G\}$ as $m \rightarrow \infty$. By monotone convergence, $L_1(G) = L_2(G)$. Since this is true for every open $G \subset \mathbb{D}$, including $G = \mathbb{D}$, the collection of Borel sets for which $L_1(B) = L_2(B)$ is a σ -field and is at least as large as the Borel σ -field. Hence (ii) implies (i). The equivalence of (i) and (iii) under separability follows from theorem 1.12.2 of VW and we omit the details here.

The fact that (i) implies (iv) is obvious. Now assume L_1 and L_2 are tight and that (iv) holds. Fix $\epsilon > 0$, and choose a compact $K \subset \mathbb{D}$ such that $L_1(K) \wedge L_2(K) \geq 1 - \epsilon$. According to a version of the Stone-Weierstrass theorem given in Jameson (1974, p. 263), a vector lattice $\mathcal{F} \subset C_b(K)$ that includes the constants and separates points of K is uniformly dense in $C_b(K)$. Choose a $g \in C_b(\mathbb{D})$ for which $0 \leq g \leq 1$, and select an $f \in \mathcal{F}$ such that $\sup_{x \in K} |g(x) - f(x)| \leq \epsilon$. Now we have $|\int g dL_1 - \int g dL_2| \leq |\int_K g dL_1 - \int_K g dL_2| + 2\epsilon \leq |\int_K (f \wedge 1)^+ dL_1 - \int_K (f \wedge 1)^+ dL_2| + 4\epsilon = 4\epsilon$. The last equality follows since $(f \wedge 1)^+ \in \mathcal{F}$. Thus $\int g dL_1 = \int g dL_2$ since ϵ is arbitrary. By adding and subtracting scalars, we can verify that the same result holds for all $g \in C_b(\mathbb{D})$. Hence (i) holds. \square

Proof of lemma 7.2. The equivalence of (i) and (ii) is an immediate consequence of theorem 6.2. Now assume (ii) holds. Then $|X|$ is bounded almost surely. Hence for any pair of sequences $s_n, t_n \in T$ such that $\rho_0(s_n, t_n) \rightarrow 0$, $X(s_n) - X(t_n) \xrightarrow{P} 0$. Thus $X \in UC(T, \rho_0)$ with probability 1. It remains to show that (T, ρ_0) is totally bounded. Let the pair of sequences $s_n, t_n \in T$ satisfy $\rho(s_n, t_n) \rightarrow 0$. Then $X(s_n) - X(t_n) \xrightarrow{P} 0$ and thus $\rho_0(s_n, t_n) \rightarrow 0$. This means that since (T, ρ) is totally bounded, we have for every $\epsilon > 0$ that there exists a finite $T_\epsilon \subset T$ such that $\sup_{t \in T} \inf_{s \in T_\epsilon} \rho_0(s, t) < \epsilon$. Thus (T, ρ_0) is also totally bounded, and the desired result follows. \square

Proof of theorem 7.6. Assume (i), and note that (i) implies (vii) trivially. Now assume (vii), and fix an open $G \subset \mathbb{D}$. As in the proof of lemma 7.1 above, there exists a sequence of nonnegative, Lipschitz continuous functions f_m with $0 \leq f_m \uparrow 1\{G\}$. Now for each integer $m \geq 1$, $\liminf P_*(X_n \in G) \geq \liminf E_* f_m(X_n) = E f_m(X)$. Taking the limit as $m \rightarrow \infty$ yields (ii). Thus (vii) \Rightarrow (ii). The equivalence of (ii) and (iii) follows by taking complements.

Assume (ii) and let f be lower semicontinuous with $f \geq 0$. Define the sequence of functions $f_m = \sum_{i=1}^{m^2} (1/m) 1\{G_i\}$, where $G_i = \{x : f(x) > i/m\}$, $i = 1, \dots, m^2$. Thus f_m “rounds” f down to i/m if $f(x) \in (i/m, (i+1)/m]$ for any $i = 0, \dots, m^2 - 1$ and $f_m = m$ when $f(x) > m$. Hence $0 \leq f_m \leq f \wedge m$ and $|f_m - f|(x) \leq 1/m$ whenever $f(x) \leq m$. Fix m . Note that each G_i is open by the definition of lower semicontinuity. Thus

$$\begin{aligned} \liminf_{n \rightarrow \infty} E_* f(X_n) &\geq \liminf_{n \rightarrow \infty} E_* f_m(X_n) \\ &\geq \sum_{i=1}^{m^2} (1/m) \left[\liminf_{n \rightarrow \infty} P_*(X_n \in G_i) \right] \\ &\geq \sum_{i=1}^{m^2} (1/m) P(X \in G_i) \\ &= E f_m(X). \end{aligned}$$

Thus (ii) implies (iv) after letting $m \rightarrow \infty$ and adding then subtracting a constant as needed to compensate for the lower bound of f . The equivalence of (iv) and (v) follows by replacing f with $-f$. Assume (v) (and thus also (iv)). Since a continuous function is both upper and lower semicontinuous, we have for any $f \in C_b(\mathbb{D})$ that $E f(X) \geq \limsup E_* f(X_n) \geq \liminf E_* f(X_n) \geq E f(X)$. Hence (v) implies (i).

Assume (ii) (and hence also (iii)). For any Borel set $B \subset \mathbb{D}$, $L(B^\circ) \leq \liminf_{n \rightarrow \infty} P_*(X_n \in B^\circ) \leq \limsup_{n \rightarrow \infty} P^*(X_n \in \overline{B}) \leq L(\overline{B})$; however, the forgoing inequalities all become equalities when $L(\delta B) = 0$. Thus (ii) implies (vi). Assume (vi), and let F be closed. For each $\epsilon > 0$ define $F^\epsilon = \{x : d(x, F) < \epsilon\}$. Since the sets δF^ϵ are disjoint, $L(\delta F^\epsilon) > 0$

for at most countably many ϵ . Hence we can choose a sequence $\epsilon_m \downarrow 0$ so that $L(\delta F^{\epsilon_m}) = 0$ for each integer $m \geq 1$. Note that for fixed m , $\limsup_{n \rightarrow \infty} P^*(X_n \in F) \leq \limsup_{n \rightarrow \infty} P^*(X_n \in \overline{F^{\epsilon_m}}) = L(\overline{F^{\epsilon_m}})$. By letting $m \rightarrow \infty$, we obtain that (vi) implies (ii). Thus conditions (i)–(vii) are all equivalent.

The equivalence of (i)–(vii) to (viii) when L is separable follows from theorem 1.12.2 of VW, and we omit the details. \square

Proof of theorem 7.7. The set of all points at which g is not continuous can be expressed as $D_g \equiv \bigcup_{m=1}^{\infty} \bigcap_{k=1}^{\infty} G_k^m$, where G_k^m consists of all $x \in \mathbb{D}$ so that $e(g(y), g(z)) > 1/m$ for some $y, z \in B_{1/k}(x)$, where e is the metric for \mathbb{E} . Note that the complement of G_k^m , $(G_k^m)^c$, consists of all x for which $e(g(y), g(z)) \leq 1/m$ for all $y, z \in B_{1/k}(x)$. Now if $x_n \rightarrow x$, then for any $y, z \in B_{1/k}(x)$, we have that $y, z \in B_{1/k}(x_n)$ for all n large enough. Hence $(G_k^m)^c$ is closed and thus G_k^m is open. This means that D_g is a Borel set. Let $F \subset \mathbb{E}$ be closed, and let $\{x_n\} \in g^{-1}(F)$ be a sequence for which $x_n \rightarrow x$. If x is a continuity point of g , then $x \in g^{-1}(F)$. Otherwise $x \in D_g$. Hence $\overline{g^{-1}(F)} \subset g^{-1}(F) \cup D_g$. Since g is continuous on the range of X , there is a version of $g(X)$ that is Borel Measurable. By the portmanteau theorem, $\limsup P^*(g(X_n) \in F) \leq \limsup P^*(X_n \in \overline{g^{-1}(F)}) \leq P(X \in \overline{g^{-1}(F)})$. Since D_g has probability zero under the law of X , $P(X \in \overline{g^{-1}(F)}) = P(g(X) \in F)$. Reapplying the portmanteau theorem, we obtain the desired result. \square

Proof of lemma 7.10. It is easy to see that uniform tightness implies asymptotic tightness. To verify equivalence going the other direction, assume that X_n is tight for each $n \geq 1$ and that the sequence is asymptotically tight. Fix $\epsilon > 0$, and choose a compact K_0 for which $\liminf P(X_n \in K_0^\delta) \geq 1 - \epsilon$ for all $\delta > 0$. For each integer $m \geq 1$, choose an $n_m < \infty$ so that $P(X_n \in K_0^{1/m}) \geq 1 - 2\epsilon$ for all $n \geq n_m$. For each integer $n \in (n_m, n_{m+1}]$, choose a compact \tilde{K}_n so that $P(X_n \in \tilde{K}_n) \geq 1 - \epsilon/2$ and an $\eta_n \in (0, 1/m)$ so that $P(X_n \in K_0^{1/m} - \overline{K_0^{\eta_n}}) < \epsilon/2$. Let $K_n = (\tilde{K}_n \cup K_0) \cap \overline{K_0^{\eta_n}}$, and note that $K_0 \subset K_n \subset K_0^{1/m}$ and that K_n is compact. We leave it as an exercise to show that $K \equiv \bigcup_{n=1}^{\infty} K_n$ is also compact. Now $P(X_n \in K_n) \geq 1 - 3\epsilon$ for all $n \geq 1$, and thus $P(X_n \in K) \geq 1 - 3\epsilon$ for all $n \geq 1$. Uniform tightness follows since ϵ was arbitrary. \square

Proof of lemma 7.18. Fix $\zeta > 0$. The conditions imply that $P^*(\|X_n\|_T^* > M) < \zeta$ for some $M < \infty$. Let ϵ_m be a positive sequence converging down to zero and let $\eta_m \equiv 2^{-m}\zeta$. By condition (ii), there exists a positive sequence $\delta_m \downarrow 0$ so that

$$\limsup_{n \rightarrow \infty} P^* \left(\sup_{s, t \in T: \rho(s, t) < \delta_m} |X_n(s) - X_n(t)| > \epsilon_m \right) < \eta_m.$$

Now fix m . By the total boundedness of T , there exists a finite set of disjoint partitions T_1, \dots, T_k so that $T = \bigcup_{i=1}^k T_i$ and so that

$$\mathbf{P}^* \left(\max_{1 \leq i \leq k} \sup_{s, t \in T_i} |X_n(s) - X_n(t)| > \epsilon_m \right) < \eta_m.$$

Let z_1, \dots, z_p be the set of all functions on $\ell^\infty(T)$ which are constant on each T_i and which take only the values $\pm i\epsilon_m$, for $i = 0, \dots, K$, where K is the largest integer $\leq M/\epsilon_m$. Now let K_m be the union of the closed balls of radius ϵ_m around the z_i , $i = 1, \dots, p$. This construction ensures that if $\|x\|_T \leq M$ and $\max_{1 \leq i \leq k} \sup_{s, t \in T_i} |x(s) - x(t)| \leq \epsilon_m$, then $x \in K_m$. This construction can be repeated for each $m \geq 1$.

Let $K = \bigcap_{m=1}^{\infty} K_m$, and note that K is totally bounded and closed. Total boundedness follows since each K_m is a union of finite ϵ_m balls which cover K and $\epsilon_m \downarrow 0$. We leave the proof that K is closed as an exercise. Now we show that for every $\delta > 0$, there is an $m < \infty$ so that $K^\delta \supset \bigcap_{i=1}^m K_i$. If this were not true, then there would be a sequence $\{z_m\} \notin K^\delta$ with $z_m \in \bigcap_{i=1}^m K_i$ for every $m \geq 1$. This sequence has a subsequence $\{z_{m_1(k)}, k \geq 1\}$ contained in one of the open balls making up K_1 . This subsequence has a further subsequence $\{z_{m_2(k)}, k \geq 1\}$ contained in one of the open balls making up K_2 . We can continue with this process to generate, for each integer $j \geq 1$, a subsequence $\{z_{m_j(k)}, k \geq 1\}$ contained in the intersection $\bigcap_{i=1}^j B_i$, where each B_i is one of the open balls making up K_i . Define a new subsequence $\tilde{z}_k = z_{m_k(k)}$, and note that \tilde{z}_k is Cauchy with limit in K since K is closed. However, this contradicts $d(z_m, K) \geq \delta$ for all m . Hence the complement of K^δ is contained in the complement of $\bigcap_{i=1}^m K_i$ for some $m < \infty$. Thus $\limsup_{n \rightarrow \infty} \mathbf{P}^*(X_n \notin K^\delta) \leq \limsup_{n \rightarrow \infty} \mathbf{P}^*(X_n \notin \bigcap_{i=1}^m K_i) \leq \zeta + \sum_{i=1}^m \eta_m \leq 2\zeta$. Since this result holds for all $\delta > 0$, we now have that $\liminf_{n \rightarrow \infty} \mathbf{P}^*(X_n \in K^\delta) \geq 1 - 2\zeta$ for all $\delta > 0$. Asymptotic tightness follows since ζ is arbitrary. \square

Proof of lemma 7.20. The fact that (i) implies (ii) we leave as an exercise. Assume (ii). Then there exists a sequence of finite partitions $\mathcal{T}_1, \mathcal{T}_2, \dots$ such that

$$\limsup_{n \rightarrow \infty} \mathbf{P}^* \left(\sup_{U \in \mathcal{T}_k} \sup_{s, t \in U} |X_n(s) - X_n(t)| > 2^{-k} \right) < 2^{-k},$$

for all integers $k \geq 1$. Here, each partition \mathcal{T}_k is a collection of disjoint sets U_1, \dots, U_{m_k} with $T = \bigcup_{i=1}^{m_k} U_i$. Without loss of generality, we can insist that the \mathcal{T}_k are nested in the sense that any $U \in \mathcal{T}_k$ is a union of sets in \mathcal{T}_{k+1} . Such a nested sequence of partitions is easy to construct from any other sequence of partitions \mathcal{T}_k^* by letting \mathcal{T}_k consist of all nontrivial intersections of all sets in \mathcal{T}_1^* up through and including \mathcal{T}_k^* . Let \mathcal{T}_0 denote the partition consisting of the single set T .

For any $s, t \in T$, define $K(s, t) \equiv \sup\{k : s, t \in U \text{ for some } U \in \mathcal{T}_k\}$ and $\rho(s, t) \equiv 2^{-K(s, t)}$. Also, for any $\delta > 0$, let $J(\delta) \equiv \inf\{k : 2^{-k} < \delta\}$. It is not hard to verify that for any $s, t \in T$ and $\delta > 0$, $\rho(s, t) < \delta$ if and only if $s, t \in U$ for some $U \in \mathcal{T}_{J(\delta)}$. Thus

$$\sup_{s,t \in T: \rho(s,t) < \delta} |X_n(s) - X_n(t)| = \sup_{U \in \mathcal{T}_{J(\delta)}} \sup_{s,t \in U} |X_n(s) - X_n(t)|$$

for all $0 < \delta \leq 1$. Since $J(\delta) \rightarrow \infty$ as $\delta \downarrow 0$, we now have that X_n is asymptotically ρ -equicontinuous in probability, as long as ρ is a pseudometric. The only difficulty here is to verify that the triangle inequality holds for ρ , which we leave as an exercise. It is easy to see that T is totally bounded with respect to ρ , and (i) follows. \square

Proof of lemma 7.21. We first prove (iii). Assume that $X_n \xrightarrow{\text{as}^*} X$, and define $A_n^k \equiv \{\sup_{m \geq n} d(X_m, X)^* > 1/k\}$. For each integer $k \geq 1$, we have $P(A_n^k) \downarrow 0$ as $n \rightarrow \infty$. Now fix $\epsilon > 0$, and note that for each $k \geq 1$ we can choose an n_k so that $P(A_n^k) \leq \epsilon/2^k$. Let $A = \Omega - \cup_{k=1}^{\infty} A_{n_k}^k$, and observe that, by this construction, $P(A) \geq 1 - \epsilon$ and $d(X_n, X)^* \leq 1/k$, for all $n \geq n_k$ and all $\omega \in A$. Thus X_n converges to X almost uniformly since ϵ is arbitrary. Assume now that X_n converges to X almost uniformly. Fix $\epsilon > 0$, and let A be measurable with $P(A) \geq 1 - \epsilon$ and $d(X_n, X) \rightarrow 0$ uniformly over $\omega \in A$. Fix $\eta > 0$, and note that $\eta \geq (d(X_n, X)1\{A\})^*$ for all sufficiently large n , since η is measurable and satisfies $\eta \geq d(X_n, X)1\{A\}$ for sufficiently large n . Now let $S, T : \Omega \mapsto [0, \infty)$ be maps with S measurable and T^* bounded. Then for any $c > 0$, $[(S+c)T]^* \leq (S+c)T^*$, and $(S+c)T^* \leq [(S+c)T]^*$ since $T^* \leq [(S+c)T]^*/(S+c)$. Hence $[(S+c)T]^* = (S+c)T^*$. By letting $c \downarrow 0$, we obtain that $(ST)^* = ST^*$. Hence $d(X_n, X)^*1\{A\} = (d(X_n, X) \inf\{A\})^* \leq \eta$ for all n large enough, and thus $d(X_n, X)^* \rightarrow 0$ for almost all $\omega \in A$. Since ϵ is arbitrary, $X_n \xrightarrow{\text{as}^*} X$.

Now assume that $X_n \xrightarrow{\text{as}^*} X$. This clearly implies that $\sup_{m \geq n} d(X_m, X) \xrightarrow{P} 0$. Fix $\epsilon > 0$. Generate a subsequence n_k , by finding, for each integer $k \geq 1$, an integer $n_k \geq n_{k-1} \geq 1$ which satisfies $P^*(\sup_{m \geq n_k} d(X_m, X) > 1/k) \leq \epsilon/2^k$. Call the set inside this outer probability statement A_k , and define $A = \Omega - \cap_{k=1}^{\infty} A_k^*$. Now $P(A) \geq 1 - \epsilon$, and for each $\omega \in A$ and all $m \geq n_k$, $d(X_m, X) \leq 1/k$ for all $k \geq 1$. Hence $\sup_{\omega \in A} d(X_n, X)(\omega) \rightarrow 0$, as $n \rightarrow \infty$. Thus X_n converges almost uniformly to X , since ϵ is arbitrary, and therefore $X_n \xrightarrow{\text{as}^*} X$. Thus we have proven (iii).

We now prove (i). It is easy to see that X_n converging to X almost uniformly will imply $X_n \xrightarrow{P} X$. Thus (i) follows from (iii). We next prove (ii). Assume $X_n \xrightarrow{P} X$. Construct a subsequence $1 \leq n_1 < n_2 < \dots$ so that $P(d(X_{n_j}, X)^* > 1/j) < 2^{-j}$ for all integers $j \geq 1$. Then

$$P(d(X_{n_j}, X)^* > 1/j, \text{ for infinitely many } j) = 0$$

by the Borel-Cantelli lemma. Hence $X_{n_j} \xrightarrow{\text{as}^*} X$ as a sequence in j . This now implies that every sequence has an outer almost surely convergent subsequence. Assume now that every subsequence has an outer almost surely convergent subsequence. By (i), the almost surely convergent subsequences also converge in probability. Hence $X_n \xrightarrow{P} X$. \square

7.5 Exercises

7.5.1. Show that \mathcal{F} defined in the proof of lemma 7.3 is a vector lattice, an algebra, and separates points in \mathbb{D} .

7.5.2. Show that the set K defined in the proof of lemma 7.10 in section 7.2 is compact.

7.5.3. In the setting of the proof of lemma 7.14, show that $\pi_1 K$ and $\pi_2 K$ are both compact whenever $K \in \mathbb{D} \times \mathbb{E}$ is compact.

7.5.4. In the proof of lemma 7.18 (given in section 7.4), show that the set $K = \bigcap_{m=1}^{\infty} K_m$ is closed.

7.5.5. Suppose that T_n and T_0 are subsets of a semimetric space (T, ρ) such that $T_n \rightarrow T_0$ in the sense that

- (i) Every $t \in T_0$ is the limit of a sequence $t_n \in T_n$;
- (ii) For every closed $S \subset T - T_0$, $S \cap T_n = \emptyset$ for all n large enough.

Suppose that X_n and X are stochastic processes indexed by T for which $X_n \rightsquigarrow X$ in $\ell^\infty(T)$ and X is Borel measurable with $P(X \in UC(T, \rho)) = 1$, where (T, ρ) is not necessarily totally bounded. Show that $\sup_{t \in T_n} X_n(t) \rightsquigarrow \sup_{t \in T_0} X(t)$. Hint: show first that for any $x_n \rightarrow x$, where $\{x_n\} \in \ell^\infty(T)$ and $x \in UC(T, \rho)$, we have $\lim_{n \rightarrow \infty} \sup_{t \in T_n} x_n(t) = \sup_{t \in T_0} x(t)$.

7.5.6. Complete the proof of lemma 7.20:

1. Show that (i) implies (ii).
2. Show that for any $s, t \in T$ and $\delta > 0$, $\rho(s, t) < \delta$ if and only if $s, t \in U$ for some $U \in \mathcal{T}_{J(\delta)}$, where $\rho(s, t) \equiv 2^{-K(s, t)}$ and K is as defined in the proof.
3. Verify that the triangle inequality holds for ρ .

7.5.7. Let X_n and X be maps into \mathbb{R} with X Borel measurable. Show the following:

- (i) $X_n \xrightarrow{\text{as}^*} X$ if and only if both X_n^* and $X_{n*} \equiv (X_n)_*$ converge almost surely to X .
- (ii) $X_n \rightsquigarrow X$ if and only if $X_n^* \rightsquigarrow X$ and $X_{n*} \rightsquigarrow X$.

Hints: For (i), first show $|X_n - X|^* = |X_n^* - X| \vee |X_{n*} - X|$. For (ii), assume $X_n \rightsquigarrow X$ and apply the extended almost sure representation theorem to find a new probability space and a perfect sequence of maps ϕ_n such that $\tilde{X}_n = X_n \circ \phi_n \xrightarrow{\text{as}^*} \tilde{X}$. By (i), $\tilde{X}_n^* \rightarrow \tilde{X}$ almost surely, and thus $\tilde{X}_n^* \rightsquigarrow \tilde{X}$. Since ϕ_n is perfect, $\tilde{X}_n^* = X_n^* \circ \phi_n$; and thus $\text{E}f(\tilde{X}_n^*) = \text{E}f(X_n^*)$ for every

measurable f . Hence $X_n^* \rightsquigarrow X$. Now show $X_{n^*} \rightsquigarrow X$. For the converse, use the facts that $P(X_n^* \leq x) \leq P^*(X_n \leq x) \leq P(X_{n^*} \leq x)$ and that distributions for real random variable are completely determined by their cumulative distribution functions.

7.5.8. Using the ideas in the proof of proposition 7.27, prove lemma 4.2.

7.6 Notes

Many of the ideas and results of this chapter come from chapters 1.3–1.5 and 1.9–1.12 of VW specialized to sequences (rather than nets). Lemma 7.1 is a composite of lemmas 1.3.12 and theorem 1.12.2 of VW, while Lemma 7.4 is lemma 1.5.3 of VW. Components (i)–(vii) of the portmanteau theorem are a specialization to sequences of the portmanteau theorem in chapter 1.3 of VW. Theorem 7.7, lemmas 7.8, 7.9, and 7.12, and theorem 7.13 correspond to VW theorems 1.3.6 and 1.3.10, lemmas 1.3.13 and 1.3.8, and theorem 1.3.9, respectively. Lemma 7.14 is a composite of lemmas 1.4.3 and 1.4.4 of VW. Lemma 7.16 and theorem 7.17 are essentially VW lemma 1.5.2 and theorem 1.5.4. Lemmas 7.21 and 7.23 and theorem 7.24 are specializations to sequences of VW lemmas 1.9.2 and 1.10.2 and theorem 1.11.1. Theorem 7.26 is a composition of VW theorem 1.10.4 and addendum 1.10.5 applied to sequences, while lemma 7.28 is essentially proposition 1.10.12 of VW.

Proposition 7.5 on Gaussian processes is essentially a variation of lemma 3.9.8 of VW. Proposition 7.27 is a modification of lemma A.3 of Billias, Gu and Ying (1997) who use this result to obtain weak convergence of the proportional hazards regression parameter in a continuously monitored sequential clinical trial.

8

Empirical Process Methods

Recall from section 2.2.2 the concepts of bracketing and uniform entropy along with the corresponding Glivenko-Cantelli and Donsker theorems. We now briefly review the set-up. Given a probability space (Ω, \mathcal{A}, P) , the data of interest consist of n independent copies X_1, \dots, X_n of a map $X : \Omega \mapsto \mathcal{X}$, where \mathcal{X} is the sample space. We are interested in studying the limiting behavior of empirical processes indexed by classes \mathcal{F} of functions $f : \mathcal{X} \mapsto \mathbb{R}$ which are measurable in the sense that each composite map $f(X) : \Omega \mapsto \mathcal{X} \mapsto \mathbb{R}$ is \mathcal{A} -measurable. With \mathbb{P}_n denoting the empirical measure based on the data X_1, \dots, X_n , the empirical process of interest is $\mathbb{P}_n f$ viewed as a stochastic process indexed by $f \in \mathcal{F}$.

A class \mathcal{F} is P -Glivenko-Cantelli if $\|\mathbb{P}_n - P\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$, where for any $u \in \ell^\infty(T)$, $\|u\|_T \equiv \sup_{t \in T} |u(t)|$. We say \mathcal{F} is *weak P -Glivenko-Cantelli*, if the outer almost sure convergence is replaced by convergence in probability. Sometimes, for clarification, we will call a P -Glivenko-Cantelli class a *strong P -Glivenko-Cantelli* class to remind ourselves that the convergence is outer almost sure. A class \mathcal{F} is P -Donsker if $\mathbb{G}_n \rightsquigarrow \mathbb{G}$ weakly in $\ell^\infty(\mathcal{F})$, where \mathbb{G} is a tight Brownian bridge. Of course, the P prefix can be dropped if the context is clear.

As mentioned in section 4.2.1, these ideas also apply directly to i.i.d. samples of stochastic processes. In this setting, X has the form $\{X(t), t \in T\}$, where $X(t)$ is measurable for each $t \in T$, and \mathcal{X} is typically $\ell^\infty(T)$. We say that X is *P -Glivenko-Cantelli* if $\sup_{t \in T} |(\mathbb{P}_n - P)X(t)| \xrightarrow{\text{as}^*} 0$ and that X is *P -Donsker* if $\mathbb{G}_n X$ converges weakly to a tight Gaussian process. This is exactly equivalent to considering whether the class $\mathcal{F}_T \equiv \{f_t, t \in$

$T\}$, where $f_t(x) \equiv x(t)$ for all $x \in \mathcal{X}$ and $t \in T$, is Glivenko-Cantelli or Donsker. The limiting Brownian bridge \mathbb{G} for the class \mathcal{F}_T has covariance $P[(\mathbb{G}f_s)(\mathbb{G}f_t)] = P[(X(s) - PX(s))(X(t) - PX(t))]$. This duality between the function class and stochastic process viewpoint will prove useful from time to time, and which approach we take will depend on the setting.

The main goal of this chapter is to present the empirical process techniques needed to prove the Glivenko-Cantelli and Donsker theorems of section 2.2.2. The approach we take is guided by chapters 2.2–2.5 of VW, although we leave out many technical details. The most difficult step in these proofs is going from point-wise convergence to uniform convergence. Maximal inequalities are very useful tools for accomplishing this step. For uniform entropy results, an additional tool, symmetrization, is also needed. To use symmetrization, several measurability conditions are required on the class of function \mathcal{F} beyond the usual requirement that each $f \in \mathcal{F}$ be measurable. In the sections presenting the Glivenko-Cantelli and Donsker theorem proofs, results for bracketing entropy are presented before the uniform entropy results.

8.1 Maximal Inequalities

We first present several results about Orlicz norms which are useful for controlling the size of the maximum of a finite collection of random variables. Several maximal inequalities for stochastic processes will be given next. These inequalities include a general maximal inequality for separable stochastic processes and a maximal inequality for sub-Gaussian processes. The results will utilize Orlicz norms combined with a method known as *chaining*. The results of this section will play a key role in the proofs of the Donsker theorems developed later on in this chapter.

8.1.1 Orlicz Norms and Maxima

A very useful class of norms for random variables used in maximal inequalities are the *Orlicz norms*. For a nondecreasing, nonzero convex function $\psi : [0, \infty] \mapsto [0, \infty]$, with $\psi(0) = 0$, the Orlicz norm $\|X\|_\psi$ of a real random variable X is defined as

$$\|X\|_\psi \equiv \inf \left\{ c > 0 : E\psi \left(\frac{|X|}{c} \right) \leq 1 \right\},$$

where the norm takes the value ∞ if no finite c exists for which $E\psi(|X|/c) \leq 1$. Exercise 8.5.1 below verifies that $\|\cdot\|_\psi$ is indeed a norm on the space of random variables with $\|X\|_\psi < \infty$. The Orlicz norm $\|\cdot\|_\psi$ is also called the *ψ -norm*, in order to specify the choice of ψ . When ψ is of the form $x \mapsto x^p$, where $p \geq 1$, the corresponding Orlicz norm is just the L_p -norm

$\|X\|_p \equiv (\mathbb{E}|X|^p)^{1/p}$. For maximal inequalities, Orlicz norms defined with $\psi_p(x) \equiv e^{x^p} - 1$, for $p \geq 1$, are of greater interest because of their sensitivity to behavior in the tails. Clearly, since $x^p \leq \psi_p(x)$, we have $\|X\|_p \leq \|X\|_{\psi_p}$. Also, by the series representation for exponentiation, $\|X\|_p \leq p! \|X\|_{\psi_1}$ for all $p \geq 1$. The following result shows how Orlicz norms based on ψ_p relate fairly precisely to the tail probabilities:

LEMMA 8.1 *For a real random variable X and any $p \in [1, \infty)$, the following are equivalent:*

(i) $\|X\|_{\psi_p} < \infty$.

(ii) *There exist constants $0 < C, K < \infty$ such that*

$$(8.1) \quad \mathbb{P}(|X| > x) \leq K e^{-Cx^p}, \text{ for all } x > 0.$$

Moreover, if either condition holds, then $K = 2$ and $C = \|X\|_{\psi_p}^{-p}$ satisfies (8.1), and, for any $C, K \in (0, \infty)$ satisfying (8.1), $\|X\|_{\psi_p} \leq ((1 + K)/C)^{1/p}$.

Proof. Assume (i). Then $\mathbb{P}(|X| > x)$ equals

$$(8.2) \quad \mathbb{P}\{\psi_p(|X|/\|X\|_{\psi_p}) \geq \psi_p(x/\|X\|_{\psi_p})\} \leq 1 \wedge \left(\frac{1}{\psi_p(x/\|X\|_{\psi_p})} \right),$$

by Markov's inequality. By exercise 8.5.2 below, $1 \wedge (e^u - 1)^{-1} \leq 2e^{-u}$ for all $u > 0$. Thus the right-hand-side of (8.2) is bounded above by (8.1) with $K = 2$ and $C = \|X\|_{\psi_p}^{-p}$. Hence (ii) and the first half of the last sentence of the lemma follow. Now assume (ii). For any $c \in (0, C)$, Fubini's theorem gives us

$$\begin{aligned} \mathbb{E}(e^{c|X|^p} - 1) &= \mathbb{E} \int_0^{|X|^p} c e^{cs} ds \\ &= \int_0^\infty \mathbb{P}(|X| > s^{1/p}) c e^{cs} ds \\ &\leq \int_0^\infty K e^{-Cs} c e^{cs} ds \\ &= Kc/(C - c), \end{aligned}$$

where the inequality follows from the assumption. Now $Kc/(C - c) \leq 1$ whenever $c \leq C/(1 + K)$ or, in other words, whenever $c^{-1/p} \geq ((1 + K)/C)^{1/p}$. This implies (i) and the rest of the lemma. \square

An important use for Orlicz norms is to control the behavior of maxima. This control is somewhat of an extension of the following simple result for L_p -norms: For any random variables X_1, \dots, X_m , $\|\max_{1 \leq i \leq m} X_i\|_p$

$$\leq \left(\mathbb{E} \max_{1 \leq i \leq m} |X_i|^p \right)^{1/p} \leq \left(\mathbb{E} \sum_{i=1}^m |X_i|^p \right)^{1/p} \leq m^{1/p} \max_{1 \leq i \leq m} \|X_i\|_p.$$

The following lemma shows that a similar result holds for certain Orlicz norms but with the $m^{1/p}$ replaced with a constant times $\psi^{-1}(m)$:

LEMMA 8.2 *Let $\psi : [0, \infty) \mapsto [0, \infty)$ be convex, nondecreasing and nonzero, with $\psi(0) = 0$ and $\limsup_{x, y \rightarrow \infty} \psi(x)\psi(y)/\psi(cxy) < \infty$ for some constant $c < \infty$. Then, for any random variables X_1, \dots, X_m ,*

$$\left\| \max_{1 \leq i \leq m} X_i \right\|_{\psi} \leq K \psi^{-1}(m) \max_{1 \leq i \leq m} \|X_i\|_{\psi},$$

where the constant K depends only on ψ .

Proof. We first make the stronger assumption that $\psi(1) \leq 1/2$ and that $\psi(x)\psi(y) \leq \psi(cxy)$ for all $x, y \geq 1$. Under this stronger assumption, we also have $\psi(x/y) \leq \psi(cx)/\psi(y)$ for all $x \geq y \geq 1$. Hence, for any $y \geq 1$ and $k > 0$,

$$\begin{aligned} \max_{1 \leq i \leq m} \psi \left(\frac{|X_i|}{ky} \right) &\leq \max_i \left[\frac{\psi(c|X_i|/k)}{\psi(y)} \mathbf{1} \left\{ \frac{|X_i|}{ky} \geq 1 \right\} \right. \\ &\quad \left. + \psi \left(\frac{|X_i|}{ky} \right) \mathbf{1} \left\{ \frac{|X_i|}{ky} < 1 \right\} \right] \\ &\leq \sum_{i=1}^m \left[\frac{\psi(c|X_i|/k)}{\psi(y)} \right] + \psi(1). \end{aligned}$$

In the summation, set $k = c \max_i \|X_i\|_{\psi}$ and take expectations of both sides to obtain

$$\mathbb{E} \psi \left(\frac{\max_i |X_i|}{ky} \right) \leq \frac{m}{\psi(y)} + \frac{1}{2}.$$

With $y = \psi^{-1}(2m)$, the right-hand-side is ≤ 1 . Thus $\|\max_i |X_i|\|_{\psi} \leq c\psi^{-1}(2m) \max_i \|X_i\|_{\psi}$. Since ψ is convex and $\psi(0) = 0$, $x \mapsto \psi^{-1}(x)$ is concave and one-to-one for $x > 0$. Thus $\psi^{-1}(2m) \leq 2\psi^{-1}(m)$, and the result follows with $K = 2c$ for the special ψ functions specified at the beginning of the proof.

By exercise 8.5.3 below, we have for any ψ satisfying the conditions of the lemma, that there exists constants $0 < \sigma \leq 1$ and $\tau > 0$ such that $\phi(x) \equiv \sigma\psi(\tau x)$ satisfies $\phi(1) \leq 1/2$ and $\phi(x)\phi(y) \leq \phi(cxy)$ for all $x, y \geq 1$. Furthermore, for this ϕ , $\phi^{-1}(u) \leq \psi^{-1}(u)/(\sigma\tau)$, for all $u > 0$, and, for any random variable X , $\|X\|_{\psi} \leq \|X\|_{\phi}/(\sigma\tau) \leq \|X\|_{\psi}/\sigma$. Hence

$$\begin{aligned} \sigma\tau \left\| \max_i X_i \right\|_{\psi} &\leq \left\| \max_i X_i \right\|_{\phi} \\ &\leq 2c\phi^{-1}(m) \max_i \|X_i\|_{\phi} \leq \frac{2c}{\sigma} \psi^{-1}(m) \max_i \|X_i\|_{\psi}, \end{aligned}$$

and the desired result follows with $K = 2c/(\sigma^2\tau)$. \square

An important consequence of lemma 8.2 is that maximums of random variables with bounded ψ -norm grow at the rate $\psi^{-1}(m)$. Based on exercise 8.5.4, ψ_p satisfies the conditions of lemma 8.2 with $c = 1$, for any $p \in [1, \infty)$. The implication in this situation is that the growth of maxima is at most logarithmic, since $\psi_p^{-1}(m) = (\log(m+1))^{1/p}$. These results will prove quite useful in the next section.

We now present an inequality for collections X_1, \dots, X_m of random variables which satisfy

$$(8.3) \quad P(|X_i| > x) \leq 2e^{-\frac{1}{2} \frac{x^2}{b+ax}}, \text{ for all } x > 0,$$

for $i = 1, \dots, m$ and some $a, b \geq 0$. This setting will arise later in the development of a Donsker theorem based on bracketing entropy.

LEMMA 8.3 *Let X_1, \dots, X_m be random variables that satisfy the tail bound (8.3) for $1 \leq i \leq m$ and some $a, b \geq 0$. Then*

$$\left\| \max_{1 \leq i \leq m} |X_i| \right\|_{\psi_1} \leq K \left\{ a \log(1+m) + \sqrt{b} \sqrt{\log(1+m)} \right\},$$

where the constant K is universal, in the sense that it does not depend on a, b , or on the random variables.

Proof. Assume for now that $a, b > 0$. The condition implies for all $x \leq b/a$ the upper bound $2 \exp(-x^2/(4b))$ for $P(|X_i| > x)$, since in this case $b+ax \leq 2b$. For all $x > b/a$, the condition implies an upper bound of $2 \exp(-x/(4a))$, since $b/x+a \leq 2a$ in this case. This implies that $P(|X_i|1\{|X_i| \leq b/a\} > x) \leq 2 \exp(-x^2/(4b))$ and $P(|X_i|1\{|X_i| > b/a\} > x) \leq 2 \exp(-x/(4a))$ for all $x > 0$. Hence, by lemma 8.1, the Orlicz norms $\| |X_i|1\{|X_i| \leq b/a\} \|_{\psi_2}$ and $\| |X_i|1\{|X_i| > b/a\} \|_{\psi_1}$ are bounded by $\sqrt{12b}$ and $12a$, respectively. The result now follows by lemma 8.2 combined with the inequality

$$\| \max_i |X_i| \|_{\psi_1} \leq \| \max_i [|X_i|1\{|X_i| \leq b/a\}] \|_{\psi_1} + \| \max_i [|X_i|1\{|X_i| > b/a\}] \|_{\psi_2},$$

where the replacement of ψ_1 with ψ_2 in the last term follows since ψ_p norms increase in p .

Suppose now that $a > 0$ but $b = 0$. Then the tail bound (8.3) holds for all $b > 0$, and the result of the lemma is thus true for all $b > 0$. The desired result now follows by letting $b \downarrow 0$. A similar argument will verify that the result holds when $a = 0$ and $b > 0$. Finally, the result is trivially true when $a = b = 0$ since, in this case, $X_i = 0$ almost surely for $i = 1, \dots, m$. \square

8.1.2 Maximal Inequalities for Processes

The goals of this section are to first establish a general maximal inequality for *separable* stochastic processes and then specialize the result to *sub-Gaussian* processes. A stochastic process $\{X(t), t \in T\}$ is separable when

there exists a countable subset $T_* \subset T$ such that $\sup_{t \in T} \inf_{s \in T_*} |X(t) - X(s)| = 0$ almost surely. For example, any cadlag process indexed by a closed interval in \mathbb{R} is separable because the rationals are a separable subset of \mathbb{R} . The need for separability of certain processes in the Glivenko-Cantelli and Donsker theorems is hidden in other conditions of the involved theorems, and direct verification of separability is seldom required in statistical applications.

A stochastic process is sub-Gaussian when

$$(8.4) \quad P(|X(t) - X(s)| > x) \leq 2e^{-\frac{1}{2}x^2/d^2(s,t)}, \text{ for all } s, t \in T, x > 0,$$

for a semimetric d on T . In this case, we say that X is sub-Gaussian with respect to d . An important example of a separable sub-Gaussian stochastic process, the Rademacher process, will be presented at the end of this section. These processes will be utilized later in this chapter in the development of a Donsker theorem based on uniform entropy. Another example of a sub-Gaussian process is Brownian motion on $[0, 1]$, which can easily be shown to be sub-Gaussian with respect to $d(s, t) = |s - t|^{1/2}$. Because the sample paths are continuous, Brownian motion is also separable.

The conclusion of lemma 8.2 above is not immediately useful for maximizing $X(t)$ over $t \in T$ since a potentially infinite number of random variables is involved. However, a method called *chaining*, which involves linking up increasingly refined finite subsets of T and repeatedly applying lemma 8.2, does make such maximization possible in some settings. The technique depends on the *metric entropy* of the index set T based on the semimetric $d(s, t) = \|X(s) - X(t)\|_\psi$.

For an arbitrary semimetric space (T, d) , the *covering number* $N(\epsilon, T, d)$ is the minimal number of closed d -balls of radius ϵ required to cover T . The *packing number* $D(\epsilon, T, d)$ is the maximal number of points that can fit in T while maintaining a distance greater than ϵ between all points. When the choice of index set T is clear by context, the notation for covering and packing numbers will be abbreviated as $N(\epsilon, d)$ and $D(\epsilon, d)$, respectively. The associated *entropy numbers* are the respective logarithms of the covering and packing numbers. Taken together, these concepts define metric entropy.

For a semimetric space (T, d) and each $\epsilon > 0$,

$$N(\epsilon, d) \leq D(\epsilon, d) \leq N(\epsilon/2, d).$$

To see this, note that there exists a minimal subset $T_\epsilon \subset T$ such that the cardinality of $T_\epsilon = D(\epsilon, d)$ and the minimum distance between distinct points in T_ϵ is $> \epsilon$. If we now place closed ϵ -balls around each point in T_ϵ , we have a covering of T . If this were not true, there would exist a point $t \in T$ which has distance $> \epsilon$ from all the points in T_ϵ , but this would mean that $D(\epsilon, d) + 1$ points can fit into T while still maintaining a separation $> \epsilon$ between all points. But this contradicts the maximality of $D(\epsilon, d)$. Thus

$N(\epsilon, d) \leq D(\epsilon, d)$. Now note that no ball of radius $\leq \epsilon/2$ can cover more than one point in T_ϵ , and thus at least $D(\epsilon, d)$ closed $\epsilon/2$ -balls are needed to cover T_ϵ . Hence $D(\epsilon, d) \leq N(\epsilon/2, d)$.

This forgoing discussion reveals that covering and packing numbers are essentially equivalent in behavior as $\epsilon \downarrow 0$. However, it turns out to be slightly more convenient for our purposes to focus on packing numbers in this section. Note that T is totally bounded if and only if $D(\epsilon, d)$ is finite for each $\epsilon > 0$. The success of the following maximal inequality depends on how fast $D(\epsilon, d)$ increases as $\epsilon \downarrow 0$:

THEOREM 8.4 (General maximal inequality) *Let ψ satisfy the conditions of lemma 8.2, and let $\{X(t), t \in T\}$ be a separable stochastic process with $\|X(s) - X(t)\|_\psi \leq rd(s, t)$, for all $s, t \in T$, some semimetric d on T , and a constant $r < \infty$. Then for any $\eta, \delta > 0$,*

$$\left\| \sup_{s, t \in T: d(s, t) \leq \delta} |X(s) - X(t)| \right\|_\psi \leq K \left[\int_0^\eta \psi^{-1}(D(\epsilon, d)) d\epsilon + \delta \psi^{-1}(D^2(\eta, d)) \right],$$

for a constant $K < \infty$ which depends only on ψ and r . Moreover,

$$\left\| \sup_{s, t \in T} |X(s) - X(t)| \right\|_\psi \leq 2K \int_0^{\text{diam } T} \psi^{-1}(D(\epsilon, d)) d\epsilon,$$

where $\text{diam } T \equiv \sup_{s, t \in T} d(s, t)$ is the diameter of T .

Before we present the proof of this theorem, recall in the discussion following the proof of lemma 8.2 that ψ_p -norms, for any $p \in [1, \infty)$, satisfy the conditions of this lemma. The case $p = 2$, which applies to sub-Gaussian processes, is of most interest to us and is explicitly evaluated in corollary 8.5 below. This corollary plays a key role in the proof of the Donsker theorem for uniform entropy (see theorem 8.19 in section 8.4).

Proof of theorem 8.4. Note that if the first integral were infinite, the inequalities would be trivially true. Hence we can, without loss of generality assume that the packing numbers and associated integral are bounded. Construct a sequence of finite nested sets $T_0 \subset T_1 \subset \dots \subset T$ such that for each T_j , $d(s, t) > \eta 2^{-j}$ for every distinct $s, t \in T_j$, and that each T_j is “maximal” in the sense that no additional points can be added to T_j without violating the inequality. Note that by the definition of packing numbers, the number of points in T_j is bounded above by $D(\eta 2^{-j}, d)$.

Now we will do the chaining part of the proof. Begin by “linking” every point $t_{j+1} \in T_{j+1}$ to one and only one $t_j \in T_j$ such that $d(t_j, t_{j+1}) \leq \eta 2^{-j}$, for all points in T_{j+1} . Continue this process to link all points in T_j with points in T_{j-1} , and so on, to obtain for every $t_{j+1} (\in T_{j+1})$ a chain $t_{j+1}, t_j, t_{j-1}, \dots, t_0$ that connects to a point in T_0 . For any integer $k \geq 0$ and arbitrary points $s_{k+1}, t_{k+1} \in T_{k+1}$, the difference in increments along their respective chains connecting to s_0, t_0 can be bounded as follows:

$$\begin{aligned}
& |\{X(s_{k+1}) - X(t_{k+1})\} - \{X(s_0) - X(t_0)\}| \\
&= \left| \sum_{j=0}^k \{X(s_{j+1}) - X(s_j)\} - \sum_{j=0}^k \{X(t_{j+1}) - X(t_j)\} \right| \\
&\leq 2 \sum_{j=0}^k \max |X(u) - X(v)|,
\end{aligned}$$

where for fixed j the maximum is taken over all links (u, v) from T_{j+1} to T_j . Hence the j th maximum is taken over at most the cardinality of T_{j+1} links, with each link having $\|X(u) - X(v)\|_\psi$ bounded by $rd(u, v) \leq r\eta 2^{-j}$. By lemma 8.2, we have for a constant $K_0 < \infty$ depending only on ψ and r ,

$$\begin{aligned}
(8.5) \quad & \left\| \max_{s, t \in T_{k+1}} |\{X(s) - X(s_0)\} - \{X(t) - X(t_0)\}| \right\|_\psi \\
& \leq K_0 \sum_{j=0}^k \psi^{-1}(D(\eta 2^{-j-1}, d)) \eta 2^{-j} \\
& = 4K_0 \sum_{j=0}^k \psi^{-1}(D(\eta 2^{-k+j-1}, d)) \eta 2^{-k+j-2} \\
& \leq 4\eta K_0 \int_0^1 \psi^{-1}(D(\eta u, d)) du \\
& = 4K_0 \int_0^\eta \psi^{-1}(D(\epsilon, d)) d\epsilon.
\end{aligned}$$

In this bound, s_0 and t_0 depend on s and t in that they are the endpoints of the chains starting at s and t , respectively.

The maximum of the increments $|X(s_{k+1}) - X(t_{k+1})|$, over all s_{k+1} and t_{k+1} in T_{k+1} with $d(s_{k+1}, t_{k+1}) < \delta$, is bounded by the left-hand-side of (8.5) plus the maximum of the discrepancies at the ends of the chains $|X(s_0) - X(t_0)|$ for those points in T_{k+1} which are less than δ apart. For every such pair of endpoints s_0, t_0 of chains starting at two points in T_{k+1} within in distance δ of each other, choose one and only one pair s_{k+1}, t_{k+1} in T_{k+1} , with $d(s_{k+1}, t_{k+1}) < \delta$, whose chains end at s_0, t_0 . By definition of T_0 , this results in at most $D^2(\eta, d)$ pairs. Now,

$$\begin{aligned}
(8.6) \quad |X(s_0) - X(t_0)| & \leq |\{X(s_0) - X(s_{k+1})\} - \{X(t_0) - X(t_{k+1})\}| \\
& \quad + |X(s_{k+1}) - X(t_{k+1})|.
\end{aligned}$$

Take the maximum of (8.6) over all pairs of endpoints s_0, t_0 . The maximum of the first term of the right-hand-side of (8.6) is bounded by the left-hand-side of (8.5). The maximum of the second term of the right-hand-side of (8.6) is the maximum of $D^2(\eta, d)$ terms with ψ -norm bounded by

$r\delta$. By lemma 8.2, this maximum is bounded by some constant C times $\delta\psi^{-1}(D^2(\eta, d))$. Combining this with (8.5), we obtain

$$\begin{aligned} & \left\| \max_{s,t \in T_{k+1}: d(s,t) < \delta} |X(s) - X(t)| \right\|_{\psi} \\ & \leq 8K_0 \int_0^{\eta} \psi^{-1}(D(\epsilon, d))d\epsilon + C\delta\psi^{-1}(D^2(\eta, d)). \end{aligned}$$

By the fact that the right-hand-side does not depend on k , we can replace T_{k+1} with $T_{\infty} = \cup_{j=0}^{\infty} T_j$ by the monotone convergence theorem. If we can verify that taking the supremum over T_{∞} is equivalent to taking the supremum over T , then the first conclusion of the theorem follows with $K = (8K_0) \vee C$.

Since X is separable, there exists a countable subset $T_* \subset T$ such that $\sup_{t \in T} \inf_{s \in T_*} |X(t) - X(s)| = 0$ almost surely. Let Ω_* denote the subset of the sample space of X for which this supremum is zero. Accordingly $P(\Omega_*) = 1$. Now, for any point t and sequence $\{t_n\}$ in T , it is easy to see that $d(t, t_n) \rightarrow 0$ implies $|X(t) - X(t_n)| \rightarrow 0$ almost surely (see exercise 8.5.5 below). For each $t \in T_*$, let Ω_t be the subset of the sample space of X for which $\inf_{s \in T_{\infty}} |X(s) - X(t)| = 0$. Since T_{∞} is a dense subset of the semimetric space (T, d) , $P(\Omega_t) = 1$. Letting $\tilde{\Omega} \equiv \Omega_* \cap (\cap_{t \in T_*} \Omega_t)$, we now have $P(\tilde{\Omega}) = 1$. This, combined with the fact that

$$\begin{aligned} \sup_{t \in T} \inf_{s \in T_{\infty}} |X(t) - X(s)| & \leq \sup_{t \in T} \inf_{s \in T_*} |X(t) - X(s)| \\ & \quad + \sup_{t \in T_*} \inf_{s \in T_{\infty}} |X(s) - X(t)|, \end{aligned}$$

implies that $\sup_{t \in T} \inf_{s \in T_{\infty}} |X(t) - X(s)| = 0$ almost surely. Thus taking the supremum over T is equivalent to taking the supremum over T_{∞} .

The second conclusion of the theorem follows from the previous result by setting $\delta = \eta = \text{diam } T$ and noting that, in this case, $D(\eta, d) = 1$. Now we have

$$\begin{aligned} \delta\psi^{-1}(D^2(\eta, d)) & = \eta\psi^{-1}(D(\eta, d)) \\ & = \int_0^{\eta} \psi^{-1}(D(\eta, d))d\epsilon \\ & \leq \int_0^{\eta} \psi^{-1}(D(\epsilon, d))d\epsilon, \end{aligned}$$

and the second conclusion follows. \square

As a consequence of exercise 8.5.5 below, the conclusions of theorem 8.4 show that X has d -continuous sample paths almost surely whenever the integral $\int_0^{\eta} \psi^{-1}(D(\epsilon, d))d\epsilon$ is bounded for some $\eta > 0$. It is also easy to verify that the maximum of the process of X is bounded, since $\|\sup_{t \in T} X(t)\|_{\psi} \leq \|X(t_0)\|_{\psi} + \|\sup_{s,t \in T} |X(t) - X(s)|\|_{\psi}$, for any choice of $t_0 \in T$. Thus X

is tight and takes its values in $UC(T, d)$ almost surely. These results will prove quite useful in later developments.

An important application of theorem 8.4 is to *sub-Gaussian* processes:

COROLLARY 8.5 *Let $\{X(t), t \in T\}$ be a separable sub-Gaussian process with respect to d . Then for all $\delta > 0$,*

$$\mathbb{E} \left(\sup_{s, t \in T: d(s, t) \leq \delta} |X(s) - X(t)| \right) \leq K \int_0^\delta \sqrt{\log D(\epsilon, d)} d\epsilon,$$

where K is a universal constant. Also, for any $t_0 \in T$,

$$\mathbb{E} \left(\sup_{t \in T} |X(t)| \right) \leq \mathbb{E}|X(t_0)| + K \int_0^{\text{diam } T} \sqrt{\log D(\epsilon, d)} d\epsilon.$$

Proof. Apply theorem 8.4 with $\psi = \psi_2$ and $\eta = \delta$. Because $\psi_2^{-1}(m) = \sqrt{\log(1+m)}$, $\psi_2^{-1}(D^2(\delta, d)) \leq \sqrt{2}\psi^{-1}(D(\delta, d))$. Hence the second term of the general maximal inequality can be replaced by

$$\sqrt{2}\delta\psi^{-1}(D(\delta, d)) \leq \sqrt{2} \int_0^\delta \psi^{-1}(D(\epsilon, d)) d\epsilon,$$

and we obtain

$$\left\| \sup_{d(s, t) \leq \delta} |X(s) - X(t)| \right\|_{\psi_2} \leq K \int_0^\delta \sqrt{\log(1 + D(\epsilon, d))} d\epsilon,$$

for an enlarged universal constant K . Note that $D(\epsilon, d) \geq 2$ for all ϵ strictly less than $\text{diam } T$. Since $(1+m) \leq m^2$ for all $m \geq 2$, the 1 inside of the logarithm can be removed at the cost of increasing K again, whenever $\delta < \text{diam } T$. Thus it is also true for all $\delta \leq \text{diam } T$. We are done with the first conclusion since $d(s, t) \leq \text{diam } T$ for all $s, t \in T$. Since the second conclusion is an easy consequence of the first, the proof is complete. \square

The next corollary shows how to use the previous corollary to establish bounds on the *modulus of continuity* of certain sub-Gaussian processes. Here the modulus of continuity for a stochastic process $\{X(t) : t \in T\}$, where (T, d) is a semimetric space, is defined as

$$m_X(\delta) \equiv \sup_{s, t \in T: d(s, t) \leq \delta} |X(s) - X(t)|.$$

COROLLARY 8.6 *Assume the conditions of corollary 8.5. Also assume there exists a differentiable function $\delta \mapsto h(\delta)$, with derivative $\dot{h}(\delta)$, satisfying $h(\delta) \geq \sqrt{\log D(\delta, d)}$ for all $\delta > 0$ small enough and $\lim_{\delta \downarrow 0} [\delta \dot{h}(\delta)/h(\delta)] = 0$. Then*

$$\lim_{M \rightarrow \infty} \limsup_{\delta \downarrow 0} \mathbb{P} \left(\frac{m_X(\delta)}{\delta h(\delta)} > M \right) = 0.$$

Proof. Using L'Hospital's rule and the assumptions of the theorem, we obtain that

$$\frac{\int_0^\delta \sqrt{\log D(\epsilon, d)} d\epsilon}{\delta h(\delta)} \leq \frac{\int_0^\delta h(\epsilon) d\epsilon}{\delta h(\delta)} \rightarrow 1,$$

as $\delta \downarrow 0$. The result now follows from the first assertion of corollary 8.5. \square

In the situation where $D(\epsilon, d) \leq K(1/\epsilon)^r$, for constants $0 < r, K < \infty$ and all $\epsilon > 0$ small enough, the above corollary works for $h(\delta) = c\sqrt{\log(1/\delta)}$, for some constant $0 < c < \infty$. This follows from simple calculations. This situation applies, for example, when X is either a standard Brownian motion or a Brownian bridge on $T = [0, 1]$. Both of these processes are sub-Gaussian with respect to the metric $d(s, t) = |s - t|^{1/2}$, and if we let $\eta = \delta^2$, we obtain from the corollary that

$$\lim_{M \rightarrow \infty} \limsup_{\eta \downarrow 0} \mathbb{P} \left(\frac{m_X(\eta)}{\sqrt{\eta \log(1/\eta)}} > M \right) = 0.$$

The rate in the denominator is quite precise in this instance since the Lévy modulus theorem (see theorem 9.25 of Karatzas and Shreve, 1991) yields

$$\mathbb{P} \left(\limsup_{\eta \downarrow 0} \frac{m_X(\eta)}{\sqrt{\eta \log(1/\eta)}} = \sqrt{2} \right) = 1.$$

The above discussion is also applicable to the modulus of continuity of certain empirical processes, and we will examine this briefly in chapter 11.

We now consider an important example of a sub-Gaussian process useful for studying empirical processes. This is the *Rademacher process*

$$X(a) = \sum_{i=1}^n \epsilon_i a_i, \quad a \in \mathbb{R}^n,$$

where $\epsilon_1, \dots, \epsilon_n$ are i.i.d. *Rademacher random variables* satisfying $P(\epsilon = -1) = P(\epsilon = 1) = 1/2$. We will verify shortly that this is indeed a sub-Gaussian process with respect to the Euclidean distance $d(a, b) = \|a - b\|$ (which obviously makes $T = \mathbb{R}^n$ into a metric space). This process will emerge in our development of Donsker results based on uniform entropy. The following lemma, also known as Hoeffding's inequality, verifies that Rademacher processes are sub-Gaussian:

LEMMA 8.7 (Hoeffding's inequality) *Let $a = (a_1, \dots, a_n) \in \mathbb{R}^n$ and $\epsilon_1, \dots, \epsilon_n$ be independent Rademacher random variables. Then*

$$\mathbb{P} \left(\left| \sum_{i=1}^n \epsilon_i a_i \right| > x \right) \leq 2e^{-\frac{1}{2}x^2/\|a\|^2},$$

for the Euclidean norm $\|\cdot\|$. Hence $\|\sum \epsilon a\|_{\psi_2} \leq \sqrt{6}\|a\|$.

Proof. For any λ and Rademacher variable ϵ , one has $Ee^{\lambda\epsilon} = (e^\lambda + e^{-\lambda})/2 = \sum_{i=0}^{\infty} \lambda^{2i}/(2i)! \leq e^{\lambda^2/2}$, where the last inequality follows from the relation $(2i)! \geq 2^i i!$ for all nonnegative integers. Hence Markov's inequality gives for any $\lambda > 0$

$$\mathbb{P}\left(\sum_{i=1}^n \epsilon_i a_i > x\right) \leq e^{-\lambda x} \mathbb{E} \exp\left\{\lambda \sum_{i=1}^n \epsilon_i a_i\right\} \leq \exp\{(\lambda^2/2)\|a\|^2 - \lambda x\}.$$

Setting $\lambda = x/\|a\|^2$ yields the desired upper bound. Since multiplying $\epsilon_1, \dots, \epsilon_n$ by -1 does not change the joint distribution, we obtain

$$\mathbb{P}\left(-\sum_{i=1}^n \epsilon_i a_i > x\right) = \mathbb{P}\left(\sum_{i=1}^n \epsilon_i a_i > x\right),$$

and the desired upper bound for the absolute value of the sum follows. The bound on the ψ_2 -norm follows directly from lemma 8.1. \square

8.2 The Symmetrization Inequality and Measurability

We now discuss a powerful technique for empirical processes called *symmetrization*. We begin by defining the “symmetrized” empirical process $f \mapsto \mathbb{P}_n^\circ f \equiv n^{-1} \sum_{i=1}^n \epsilon_i f(X_i)$, where $\epsilon_1, \dots, \epsilon_n$ are independent Rademacher random variables which are also independent of X_1, \dots, X_n . The basic idea behind symmetrization is to replace supremums of the form $\|(\mathbb{P}_n - P)f\|_{\mathcal{F}}$ with supremums of the form $\|\mathbb{P}_n^\circ f\|_{\mathcal{F}}$. This replacement is very useful in Glivenko-Cantelli and Donsker theorems based on uniform entropy, and a proof of the validity of this replacement is the primary goal of this section. Note that the processes $(\mathbb{P}_n - P)f$ and $\mathbb{P}_n^\circ f$ both have mean zero. A deeper connection between these two processes is that a Donsker theorem or Glivenko-Cantelli theorem holds for one of these processes if and only if it holds for the other.

One potentially troublesome difficulty is that the supremums involved may not be measurable, and we need to be clear about the underlying product probability spaces so that the outer expectations are well defined. In this setting, we will assume that X_1, \dots, X_n are the coordinate projections of the product space $(\mathcal{X}^n, \mathcal{A}^n, P^n)$, where $(\mathcal{X}, \mathcal{A}, P)$ is the product space for a single observation and \mathcal{A}^n is shorthand for the product σ -field generated from sets of the form $A_1 \times \dots \times A_n$, where $A_1, \dots, A_n \in \mathcal{A}$. In many of the settings of interest to us, the σ -field \mathcal{A}^n will be strictly smaller than the Borel σ -field generated from the product topology, as discussed in section 6.1, but the results we obtain using \mathcal{A}^n will be sufficient for our purposes. In some settings, an additional source of ran-

domness, independent of X_1, \dots, X_n , will be involved which we will denote Z . If we let the probability space for Z be $(\mathcal{Z}, \mathcal{D}, Q)$, we will assume that the resulting underlying joint probability space has the form $(\mathcal{X}^n, \mathcal{A}^n, P^n) \times (\mathcal{Z}, \mathcal{D}, Q) = (\mathcal{X}^n \times \mathcal{Z}, \mathcal{A}^n \times \mathcal{D}, P^n \times Q)$, where we define the product σ -field $\mathcal{A}^n \times \mathcal{D}$ in the same manner as before. Now X_1, \dots, X_n are equal to the coordinate projections on the first n coordinates, while Z is equal to the coordinate projection on the $(n+1)$ st coordinate.

We now present the symmetrization theorem. After its proof, we will discuss a few additional important measurability issues.

THEOREM 8.8 (Symmetrization) *For every nondecreasing, convex $\phi : \mathbb{R} \mapsto \mathbb{R}$ and class of measurable functions \mathcal{F} ,*

$$E^* \phi \left(\frac{1}{2} \|\mathbb{P}_n - P\|_{\mathcal{F}} \right) \leq E^* \phi (\|\mathbb{P}_n^{\circ}\|_{\mathcal{F}}) \leq E^* \phi (2\|\mathbb{P}_n - P\|_{\mathcal{F}} + |R_n| \cdot \|P\|_{\mathcal{F}}),$$

where $R_n \equiv \mathbb{P}_n^{\circ} 1 = n^{-1} \sum_{i=1}^n \epsilon_i$ and the outer expectations are computed based on the product σ -field described in the previous paragraph.

Before giving the proof of this theorem, we make a few observations. Firstly, the constants $1/2$, 1 and 2 appearing in front of the three respective supremum norms in the chain of inequalities can all be replaced by $c/2$, c and $2c$, respectively, for any positive constant c . This follows trivially since, for any positive c , $x \mapsto \phi(cx)$ is nondecreasing and convex whenever $x \mapsto \phi(x)$ is nondecreasing and convex. Secondly, we note that most of our applications of this theorem will be for the setting $\phi(x) = x$. Thirdly, we note that the first inequality in the chain of inequalities will be of greatest use to us. However, the second inequality in the chain can be used to establish the following Glivenko-Cantelli result, the complete proof of which will be given later on, at the tail end of section 8.3:

PROPOSITION 8.9 *For any class of measurable functions \mathcal{F} , the following are equivalent:*

- (i) \mathcal{F} is P -Glivenko-Cantelli and $\|P\|_{\mathcal{F}} < \infty$.
- (ii) $\|\mathbb{P}_n^{\circ}\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$.

As mentioned previously, there is also a similar equivalence involving Donsker results, but we will postpone further discussion of this until we encounter multiplier central limit theorems in chapter 10.

Proof of theorem 8.8. Let Y_1, \dots, Y_n be independent copies of X_1, \dots, X_n . Formally, Y_1, \dots, Y_n are the coordinate projections on the last n coordinates in the product space $(\mathcal{X}^n, \mathcal{A}^n, P^n) \times (\mathcal{Z}, \mathcal{D}, Q) \times (\mathcal{X}^n, \mathcal{A}^n, P^n)$. Here, $(\mathcal{Z}, \mathcal{D}, Q)$ is the probability space for the n -vector of independent Rademacher random variables $\epsilon_1, \dots, \epsilon_n$ used in \mathbb{P}_n° . Since, by lemma 6.13, coordinate projections are perfect maps, the outer expectations in the theorem are unaffected by the enlarged product probability space. For fixed X_1, \dots, X_n , $\|\mathbb{P}_n - P\|_{\mathcal{F}} =$

$$\sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n [f(X_i) - \mathbb{E}f(Y_i)] \right| \leq \mathbb{E}_Y^* \sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n [f(X_i) - f(Y_i)] \right|,$$

where \mathbb{E}_Y^* is the outer expectation with respect to Y_1, \dots, Y_n computed by treating the X_1, \dots, X_n as constants and using the probability space $(\mathcal{X}^n, \mathcal{A}^n, P^n)$. Applying Jensen's inequality, we obtain

$$\phi(\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_Y \phi \left(\left\| \frac{1}{n} \sum_{i=1}^n [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}}^{*Y} \right),$$

where $*Y$ denotes the minimal measurable majorant of the supremum with respect to Y_1, \dots, Y_n and holding X_1, \dots, X_n fixed. Because ϕ is nondecreasing and continuous, the $*Y$ inside of the ϕ in the forgoing expression can be removed after replacing \mathbb{E}_Y with \mathbb{E}_Y^* , as a consequence of lemma 6.8. Now take the expectation of both sides with respect to X_1, \dots, X_n to obtain

$$\mathbb{E}^* \phi(\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_X^* \mathbb{E}_Y^* \phi \left(\frac{1}{n} \left\| \sum_{i=1}^n [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}} \right).$$

The repeated outer expectation can now be bounded above by the joint outer expectation \mathbb{E}^* by lemma 6.14 (Fubini's theorem for outer expectations).

By the product space structure of the underlying probability space, the outer expectation of any function $g(X_1, \dots, X_n, Y_1, \dots, Y_n)$ remains unchanged under permutations of its $2n$ arguments. Since $-[f(X_i) - f(Y_i)] = [f(Y_i) - f(X_i)]$, we have for any n -vector $(e_1, \dots, e_n) \in \{-1, 1\}^n$, that $\|n^{-1} \sum_{i=1}^n e_i [f(X_i) - f(Y_i)]\|_{\mathcal{F}}$ is just a permutation of

$$h(X_1, \dots, X_n, Y_1, \dots, Y_n) \equiv \left\| n^{-1} \sum_{i=1}^n [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}}.$$

Hence

$$\mathbb{E}^* \phi(\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_\epsilon \mathbb{E}_{X,Y}^* \phi \left\| \frac{1}{n} \sum_{i=1}^n e_i [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}}.$$

Now the triangle inequality combined with the convexity of ϕ yields

$$\mathbb{E}^* \phi(\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_\epsilon \mathbb{E}_{X,Y}^* \phi(2\|\mathbb{P}_n^\circ\|_{\mathcal{F}}).$$

By the perfectness of coordinate projections, $\mathbb{E}_{X,Y}^*$ can be replaced by $\mathbb{E}_X^* \mathbb{E}_Y^*$. Now $\mathbb{E}_\epsilon \mathbb{E}_X^* \mathbb{E}_Y^*$ is bounded above by the joint expectation \mathbb{E}^* by reapplication of lemma 6.14. This proves the first inequality.

For the second inequality, let Y_1, \dots, Y_n be independent copies of X_1, \dots, X_n as before. Holding X_1, \dots, X_n and $\epsilon_1, \dots, \epsilon_n$ fixed, we have $\|\mathbb{P}_n^\circ f\|_{\mathcal{F}} = \|\mathbb{P}_n^\circ(f - Pf) + \mathbb{P}_n^\circ Pf\|_{\mathcal{F}} =$

$$\|\mathbb{P}_n^\circ(f - \mathbb{E}f(Y)) + R_n Pf\|_{\mathcal{F}} \leq \mathbb{E}_Y^* \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}} + R_n \|P\|_{\mathcal{F}}.$$

Applying Jensen's inequality, we now have

$$\phi(\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_Y^* \phi \left(\left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}} + R_n \|P\|_{\mathcal{F}} \right).$$

Using the permutation argument we used for the first part of the proof, we can replace the $\epsilon_1, \dots, \epsilon_n$ in the summation with all 1's, and take expectations with respect to X_1, \dots, X_n and $\epsilon_1, \dots, \epsilon_n$ (which are still present in R_n). This gives us

$$\mathbb{E}^* \phi(\|\mathbb{P}_n - P\|_{\mathcal{F}}) \leq \mathbb{E}_\epsilon \mathbb{E}_X^* \mathbb{E}_Y^* \phi \left(\left\| \frac{1}{n} \sum_{i=1}^n [f(X_i) - f(Y_i)] \right\|_{\mathcal{F}} + R_n \|P\|_{\mathcal{F}} \right).$$

After adding and subtracting Pf in the summation and applying the convexity of ϕ , we can bound the right-hand-side by

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_\epsilon \mathbb{E}_X^* \mathbb{E}_Y^* \phi \left(2 \left\| \frac{1}{n} \sum_{i=1}^n [f(X_i) - Pf] \right\|_{\mathcal{F}} + R_n \|P\|_{\mathcal{F}} \right) \\ & + \frac{1}{2} \mathbb{E}_\epsilon \mathbb{E}_X^* \mathbb{E}_Y^* \phi \left(2 \left\| \frac{1}{n} \sum_{i=1}^n [f(Y_i) - Pf] \right\|_{\mathcal{F}} + R_n \|P\|_{\mathcal{F}} \right). \end{aligned}$$

By reapplication of the permutation argument and lemma 6.14, we obtain the desired upper bound. \square

The above symmetrization results will be most useful when the supremum $\|\mathbb{P}_n^\circ\|_{\mathcal{F}}$ is measurable and Fubini's theorem permits taking the expectation first with respect to $\epsilon_1, \dots, \epsilon_n$ given X_1, \dots, X_n and secondly with respect to X_1, \dots, X_n . Without this measurability, only the weaker version of Fubini's theorem for outer expectations applies (theorem 6.14), and thus the desired reordering of expectations may not be valid. To overcome this difficulty, we will assume that the class \mathcal{F} is a *P-measurable class*. A class \mathcal{F} of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$, on the probability space $(\mathcal{X}, \mathcal{A}, P)$, is *P-measurable* if $(X_1, \dots, X_n) \mapsto \|\sum_{i=1}^n e_i f(X_i)\|_{\mathcal{F}}$ is measurable on the completion of $(\mathcal{X}^n, \mathcal{A}^n, P^n)$ for every constant vector $(e_1, \dots, e_n) \in \mathbb{R}^n$. It is possible to weaken this condition, but at least some measurability assumptions will usually be needed. In the Donsker theorem for uniform entropy, it will be necessary to assume that several related classes of \mathcal{F} are also *P-measurable*. These additional classes are $\mathcal{F}_\delta \equiv \{f - g : f, g \in \mathcal{F}, \|f - g\|_{P,2} < \delta\}$, for all $\delta > 0$, and $\mathcal{F}_\infty^2 \equiv \{(f - g)^2 : f, g \in \mathcal{F}\}$ (recall that $\|f\|_{P,2} \equiv (Pf^2)^{1/2}$).

Another assumption on \mathcal{F} which is stronger than *P-measurability* and often easier to verify in statistical applications is *pointwise measurability*.

A class \mathcal{F} of measurable functions is pointwise measurable if there exists a countable subset $\mathcal{G} \subset \mathcal{F}$ such that for every $f \in \mathcal{F}$, there exists a sequence $\{g_m\} \in \mathcal{G}$ with $g_m(x) \rightarrow f(x)$ for every $x \in \mathcal{X}$. Since, by exercise 8.5.6 below, $\|\sum e_i f(X_i)\|_{\mathcal{F}} = \|\sum e_i f(X_i)\|_{\mathcal{G}}$ for all $(e_1, \dots, e_n) \in \mathbb{R}^n$, pointwise measurable classes are P -measurable for all P . Consider, for example, the class $\mathcal{F} = \{1\{x \leq t\} : t \in \mathbb{R}\}$ where the sample space $\mathcal{X} = \mathbb{R}$. Let $\mathcal{G} = \{1\{x \leq t\} : t \in \mathbb{Q}\}$, and fix the function $x \mapsto f(x) = 1\{x \leq t_0\}$ for some $t_0 \in \mathbb{R}$. Note that \mathcal{G} is countable. Let $\{t_m\}$ be a sequence of rationals with $t_m \geq t_0$, for all $m \geq 1$, and with $t_m \downarrow t_0$. Then $x \mapsto g_m(x) = 1\{x \leq t_m\}$ satisfies $g_m \in \mathcal{G}$, for all $m \geq 1$, and $g_m(x) \rightarrow f(x)$ for all $x \in \mathbb{R}$. Since t_0 was arbitrary, we have just proven that \mathcal{F} is pointwise measurable (and hence also P -measurable for all P). Hereafter, we will use the abbreviation PM as a shorthand for denoting pointwise measurable classes.

Another nice feature of PM classes is that they have a number of useful preservation features. An obvious example is that when \mathcal{F}_1 and \mathcal{F}_2 are PM classes, then so is $\mathcal{F}_1 \cup \mathcal{F}_2$. The following lemma provides a number of additional preservation results:

LEMMA 8.10 *Let $\mathcal{F}_1, \dots, \mathcal{F}_k$ be PM classes of real functions on \mathcal{X} , and let $\phi : \mathbb{R}^k \mapsto \mathbb{R}$ be continuous. Then the class $\phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k)$ is PM, where $\phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k)$ denotes the class $\{\phi(f_1, \dots, f_k) : (f_1, \dots, f_k) \in \mathcal{F}_1 \times \dots \times \mathcal{F}_k\}$.*

Proof. Denote $\mathcal{H} \equiv \phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k)$. Fix an arbitrary $h = \phi(f_1, \dots, f_k) \in \mathcal{H}$. By assumption, each \mathcal{F}_j has a countable subset $\mathcal{G}_j \subset \mathcal{F}_j$ such that there exists a subsequence $\{g_m^j\} \in \mathcal{G}_j$ with $g_m^j(x) \rightarrow f_j(x)$, as $m \rightarrow \infty$, for all $x \in \mathcal{X}$ and $j = 1, \dots, k$. By continuity of ϕ , we thus have that $\phi(g_m^1(x), \dots, g_m^k(x)) \rightarrow \phi(f_1(x), \dots, f_k(x)) = h(x)$, as $m \rightarrow \infty$, for all $x \in \mathcal{X}$. Since the choice of h was arbitrary, we therefore have that the set $\phi(\mathcal{G}_1, \dots, \mathcal{G}_k)$ is a countable subset of \mathcal{H} making \mathcal{H} pointwise measurable. \square

Lemma 8.10 automatically yields many other useful PM preservation results, including the following for PM classes \mathcal{F}_1 and \mathcal{F}_2 :

- $\mathcal{F}_1 \wedge \mathcal{F}_2$ (all possible pairwise minimums) is PM.
- $\mathcal{F}_1 \vee \mathcal{F}_2$ (all possible pairwise maximums) is PM.
- $\mathcal{F}_1 + \mathcal{F}_2$ is PM.
- $\mathcal{F}_1 \cdot \mathcal{F}_2 \equiv \{f_1 f_2 : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}$ is PM.

We will use these properties of PM classes to establish Donsker properties for some specific statistical examples later on in the case studies presented in chapter 15. The following proposition shows an additional property of PM classes that potentially simplifies the measurability requirements of the Donsker theorem for uniform entropy, theorem 8.19, given in section 8.4 below:

PROPOSITION 8.11 *Let \mathcal{F} be a class of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$ on the probability space $(\mathcal{X}, \mathcal{A}, P)$. Provided \mathcal{F} is PM with envelope F such that $P^*F^2 < \infty$, then \mathcal{F} , \mathcal{F}_δ and \mathcal{F}_∞^2 are PM for all $\delta > 0$.*

Proof. The fact that both \mathcal{F}_∞ and \mathcal{F}_∞^2 are PM follows easily from lemma 8.10. Assume, without loss of generality, that the envelope F is measurable (if not, simply replace F with F^*). Next, let $\mathcal{H} \subset \mathcal{F}_\infty$ be a countable subset for which there exists for each $g \in \mathcal{F}_\infty$ a sequence $\{h_m\} \in \mathcal{H}$ such that $h_m(x) \rightarrow g(x)$ for all $x \in \mathcal{X}$. Fix $\delta > 0$ and $h \in \mathcal{F}_\delta$. Then there exists an $\epsilon > 0$ such that $Ph^2 = \delta^2 - \epsilon$. Let $\{g_m\} \in \mathcal{H}$ be a sequence for which $g_m(x) \rightarrow h(x)$ for all $x \in \mathcal{X}$, and assume that $Pg_m^2 \geq \delta^2$ infinitely often. Then there exists another sequence $\{\tilde{g}_m\} \in \mathcal{H}$ such that $P\tilde{g}_m^2 \geq \delta^2$ for all $m \geq 1$ and also $\tilde{g}_m(x) \rightarrow h(x)$ for all $x \in \mathcal{X}$. Since $|\tilde{g}_m| \leq F$, for all $m \geq 1$, we have by the dominated convergence theorem that $\delta^2 \leq \liminf_{m \rightarrow \infty} P\tilde{g}_m^2 = Ph^2 = \delta^2 - \epsilon$, which is impossible. Hence, returning to the original sequence $\{g_m\}$, $\|g_m\|_{P,2}$ cannot be $\geq \delta$ infinitely often. Thus there exists a sequence $\{\check{g}_m\} \in \mathcal{H}_\delta \equiv \{g \in \mathcal{H} : \|g\|_{P,2} < \delta\}$ such that $\check{g}_m(x) \rightarrow h(x)$ for all $x \in \mathcal{X}$. Thus \mathcal{F}_δ is PM since h was arbitrary and \mathcal{H}_δ does not depend on h . Since δ was also arbitrary, the proof is complete. \square

We next consider establishing P -measurability for the class

$$\{1\{Y - \beta^T Z \leq t\} : \beta \in \mathbb{R}^k, t \in \mathbb{R}\},$$

where $X \equiv (Y, Z) \in \mathcal{X} \equiv \mathbb{R} \times \mathbb{R}^k$ has distribution P , for arbitrary P . This class was considered in the linear regression example of section 4.1. The desired measurability result is stated in the following lemma:

LEMMA 8.12 *Let $\mathcal{F} \equiv \{1\{Y - \beta^T Z \leq t\} : \beta \in \mathbb{R}^k, t \in \mathbb{R}\}$. Then the classes \mathcal{F} , $\mathcal{F}_\delta \equiv \{f - g : f, g \in \mathcal{F}, \|f - g\|_{P,2} < \delta\}$, and $\mathcal{F}_\infty^2 \equiv \{(f - g)^2 : f, g \in \mathcal{F}\}$ are all P -measurable for any probability measure on \mathcal{X} .*

Proof. We first assume that $\|Z\| \leq M$ for some fixed $M < \infty$. Hence the sample space is $\mathcal{X}_M \equiv \{(y, z) : (y, z) \in \mathcal{X}, \|z\| \leq M\}$. Consider the countable set $\mathcal{G} = \{1\{Y - \beta^T Z \leq t\} : \beta \in \mathbb{Q}^k, t \in \mathbb{Q}\}$, where \mathbb{Q} are the rationals. Fix $\beta \in \mathbb{R}^k$ and $t \in \mathbb{R}$, and construct a sequence $\{(\beta_m, t_m)\}$ as follows: for each $m \geq 1$, pick $\beta_m \in \mathbb{Q}^k$ so that $\|\beta_m - \beta\| < 1/(2mM)$ and pick $t_m \in \mathbb{Q}$ so that $t_m \in (t + 1/(2m), t + 1/m]$. Now, for any (y, z) with $y \in \mathbb{R}$ and $z \in \mathbb{R}^k$ with $\|z\| \leq M$, we have that $1\{y - \beta_m^T z \leq t_m\} = 1\{y - \beta^T z \leq t_m + (\beta_m - \beta)^T z\}$. Since $|(\beta_m - \beta)^T z| < 1/(2m)$ by design, we have that $r_m \equiv t_m + (\beta_m - \beta)^T z - t > 0$ for all m and that $r_m \rightarrow 0$ as $m \rightarrow \infty$. Since the function $t \mapsto \{u \leq t\}$ is right-continuous and since (y, z) was arbitrary, we have just proven that $1\{y - \beta_m^T z \leq t_m\} \rightarrow 1\{y - \beta^T z \leq t\}$ for all $(y, z) \in \mathcal{X}_M$. Thus \mathcal{F} is pointwise measurable with respect to the countable subset \mathcal{G} .

We can also verify that \mathcal{F}_δ and \mathcal{F}_∞^2 are likewise PM classes, under the constraint that the random variable Z satisfies $\|Z\| \leq M$. To see this for

\mathcal{F}_δ , let $f_1, f_2 \in \mathcal{F}$ satisfy $\|f_1 - f_2\|_{P,2} < \delta$ and let $\{g_{m,1}\}, \{g_{m,2}\} \in \mathcal{G}$ be such that $g_{m,1} \rightarrow f_1$ and $g_{m,2} \rightarrow f_2$ pointwise in \mathcal{X}_M . Then, by dominated convergence, $\|g_{m,1} - g_{m,2}\|_{P,2} \rightarrow \|f_1 - f_2\|$, and thus $\|g_{m,1} - g_{m,2}\|_{P,2} < \delta$ for all m large enough. Hence

$$g_{m,1} - g_{m,2} \in \mathcal{G}_\delta \equiv \{f - g : f, g \in \mathcal{G}, \|f - g\|_{P,2} < \delta\}$$

for all m large enough, and thus \mathcal{G}_δ is a separable subset of \mathcal{F}_δ making \mathcal{F}_δ into a PM class. The proof that \mathcal{F}_∞^2 is also PM follows directly from lemma 8.10.

Now let $J_M(x_1, \dots, x_n) = 1\{\max_{1 \leq i \leq n} \|z_i\| \leq M\}$, where $x_i = (y_i, z_i)$, $1 \leq i \leq n$. Since M was arbitrary, the previous two paragraphs have established that

$$(8.7) \quad (x_1, \dots, x_n) \mapsto \left\| \sum_{i=1}^n e_i f(x_i) \right\|_{\mathcal{H}} J_M(x_1, \dots, x_n)$$

is measurable for every n -tuple $(e_1, \dots, e_n) \in \mathbb{R}^n$, every $M < \infty$, and with \mathcal{H} being replaced by \mathcal{F} , \mathcal{F}_δ or \mathcal{F}_∞^2 . Now, for any $(x_1, \dots, x_n) \in \mathcal{X}^n$, $J_M(x_1, \dots, x_n) = 1$ for all M large enough. Thus the map (8.7) is also measurable after replacing J_M with its pointwise limit $1 = \lim_{M \rightarrow \infty} J_M$. Hence \mathcal{F} , \mathcal{F}_δ and \mathcal{F}_∞^2 are all P -measurable classes for any measure P on \mathcal{X} . \square

Another example of a P -measurable class occurs when \mathcal{F} is a Suslin topological space (for an arbitrary topology \mathcal{O}), and the map $(x, f) \mapsto f(x)$ is jointly measurable on $\mathcal{X} \times \mathcal{F}$ for the product σ -field of \mathcal{A} and the Borel σ -field generated from \mathcal{O} . Further insights and results on this *measurable Suslin condition* can be found in example 2.3.5 and chapter 1.7 of VW. While this approach to establishing measurability can be useful in some settings, a genuine need for it does not often occur in statistical applications, and we will not pursue it further here.

8.3 Glivenko-Cantelli Results

We now present several Glivenko-Cantelli (G-C) results. First, we discuss an interesting necessary condition for a class \mathcal{F} to be P -G-C. Next, we present the proofs of G-C theorems for bracketing (theorem 2.2) and uniform (theorem 2.4) entropy. Part of the proof in the uniform entropy case will include the presentation of a new G-C theorem, theorem 8.15 below. Finally, we give the proof of proposition 8.9 which was promised in the previous section.

The following lemma shows that the existence of an integrable envelope of the centered functions of a class \mathcal{F} is a necessary condition for \mathcal{F} to be P -G-C:

LEMMA 8.13 *If the class of functions \mathcal{F} is strong P -G-C, then $P\|f - Pf\|_{\mathcal{F}}^* < \infty$. If in addition $\|P\|_{\mathcal{F}} < \infty$, then also $P\|f\|_{\mathcal{F}}^* < \infty$.*

Proof. Since $f(X_n) - Pf = n(\mathbb{P}_n - P)f - (n-1)(\mathbb{P}_{n-1} - P)f$, we have $n^{-1}\|f(X_n) - Pf\|_{\mathcal{F}} \leq \|\mathbb{P}_n - P\|_{\mathcal{F}} + (1 - n^{-1})\|\mathbb{P}_{n-1} - P\|_{\mathcal{F}}$. Since \mathcal{F} is strong P -G-C, we now have that $P(\|f(X_n) - Pf\|_{\mathcal{F}}^* \geq n \text{ infinitely often}) = 0$. The Borel-Cantelli lemma now yields that $\sum_{n=1}^{\infty} P(\|f(X_n) - Pf\|_{\mathcal{F}}^* \geq n) < \infty$. Since the X_n are i.i.d., the $f(X_n)$ in the summands can be replaced with $f(X_1)$ for all $n \geq 1$. Now we have

$$\begin{aligned} P^*\|f - Pf\|_{\mathcal{F}} &\leq \int_0^{\infty} P(\|f(X_1) - Pf\|_{\mathcal{F}}^* > x) dx \\ &\leq 1 + \sum_{n=1}^{\infty} P(\|f(X_1) - Pf\|_{\mathcal{F}}^* \geq n) \\ &< \infty. \square \end{aligned}$$

Proof of theorem 2.2. Fix $\epsilon > 0$. Since the L_1 -bracketing entropy is bounded, it is possible to choose finitely many ϵ -brackets $[l_i, u_i]$ so that their union contains \mathcal{F} and $P(u_i - l_i) < \epsilon$ for every i . Now, for every $f \in \mathcal{F}$, there is a bracket $[l_i, u_i]$ containing f with $(\mathbb{P}_n - P)f \leq (\mathbb{P}_n - P)u_i + P(u_i - f) \leq (\mathbb{P}_n - P)u_i + \epsilon$. Hence

$$\begin{aligned} \sup_{f \in \mathcal{F}} (\mathbb{P}_n - P)f &\leq \max_i (\mathbb{P}_n - P)u_i + \epsilon \\ &\xrightarrow{\text{as}^*} \epsilon. \end{aligned}$$

Similar arguments can be used to verify that

$$\begin{aligned} \inf_{f \in \mathcal{F}} (\mathbb{P}_n - P)f &\geq \min_i (\mathbb{P}_n - P)l_i - \epsilon \\ &\xrightarrow{\text{as}^*} -\epsilon. \end{aligned}$$

The desired result now follows since ϵ was arbitrary. \square

To prove theorem 2.4, we first restate the theorem to clarify the meaning of “appropriately measurable” in the original statement of the theorem, and then prove a more general version (theorem 8.15 below):

THEOREM 8.14 (*Restated theorem 2.4*) *Let \mathcal{F} be a P -measurable class of measurable functions with envelope F and $\sup_Q N(\epsilon\|F\|_{Q,1}, \mathcal{F}, L_1(Q)) < \infty$, for every $\epsilon > 0$, where the supremum is taken over all finite probability measures Q with $\|F\|_{Q,1} > 0$. If $P^*F < \infty$, the \mathcal{F} is P -G-C.*

Proof. The result is trivial if $P^*F = 0$. Hence we will assume without loss of generality that $P^*F > 0$. Thus there exists an $\eta > 0$ such that, with probability 1, $\mathbb{P}_n F > \eta$ for all n large enough. Fix $\epsilon > 0$. By assumption, there is a $K < \infty$ such that $1\{\mathbb{P}_n F > 0\} \log N(\epsilon\mathbb{P}_n F, \mathcal{F}, L_1(\mathbb{P}_n)) \leq K$

almost surely, since \mathbb{P}_n is a finite probability measure. Hence, with probability 1, $\log N(\epsilon\eta, \mathcal{F}, L_1(\mathbb{P}_n)) \leq K$ for all n large enough. Since ϵ was arbitrary, we now have that $\log N(\epsilon, \mathcal{F}, L_1(\mathbb{P}_n)) = O_P^*(1)$ for all $\epsilon > 0$. Now fix $\epsilon > 0$ (again) and $M < \infty$, and define $\mathcal{F}_M \equiv \{f1\{F \leq M\} : f \in \mathcal{F}\}$. Since, $\|(f-g)1\{F \leq M\}\|_{1, \mathbb{P}_n} \leq \|f-g\|_{1, \mathbb{P}_n}$ for any $f, g \in \mathcal{F}$, $N(\epsilon, \mathcal{F}_M, L_1(\mathbb{P}_n)) \leq N(\epsilon, \mathcal{F}, L_1(\mathbb{P}_n))$. Hence $\log N(\epsilon, \mathcal{F}_M, L_1(\mathbb{P}_n)) = O_P^*(1)$. Finally, since ϵ and M are both arbitrary, the desired result follows from theorem 8.15 below. \square

THEOREM 8.15 *Let \mathcal{F} be a P -measurable class of measurable functions with envelope F such that $P^*F < \infty$. Let \mathcal{F}_M be as defined in the above proof. If $\log N(\epsilon, \mathcal{F}_M, L_1(\mathbb{P}_n)) = o_P^*(n)$ for every $\epsilon > 0$ and $M < \infty$, then $P\|\mathbb{P}_n - P\|_{\mathcal{F}}^* \rightarrow 0$ and \mathcal{F} is strong P -G-C.*

Before giving the proof of theorem 8.15, we give the following lemma which will be needed. This is lemma 2.4.5 of VW, and we omit the proof:

LEMMA 8.16 *Let \mathcal{F} be a class of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$ with integrable envelope. Define a filtration by letting Σ_n be the σ -field generated by all measurable functions $h : \mathcal{X}^\infty \mapsto \mathbb{R}$ that are permutation-symmetric in their first n arguments. Then $\mathbb{E}(\|\mathbb{P}_n - P\|_{\mathcal{F}}^* | \Sigma_{n+1}) \geq \|\mathbb{P}_{n+1} - P\|_{\mathcal{F}}^*$, almost surely. Furthermore, there exist versions of the measurable cover functions $\|\mathbb{P}_n - P\|_{\mathcal{F}}^*$ that are adapted to the filtration. Any such versions form a reverse submartingale and converge almost surely to a constant.*

Proof of theorem 8.15. By the symmetrization theorem 8.8, P -measurability of \mathcal{F} , and by Fubini's theorem, we have for all $M > 0$ that

$$\begin{aligned} \mathbb{E}^* \|\mathbb{P}_n - P\|_{\mathcal{F}} &\leq 2\mathbb{E}_X \mathbb{E}_\epsilon \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}} \\ &\leq 2\mathbb{E}_X \mathbb{E}_\epsilon \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) 1\{F \leq M\} \right\|_{\mathcal{F}} \\ &\quad + 2\mathbb{E}_X \mathbb{E}_\epsilon \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) 1\{F > M\} \right\|_{\mathcal{F}} \\ &\leq 2\mathbb{E}_X \mathbb{E}_\epsilon \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}_M} + 2P^*[F1\{F > M\}]. \end{aligned}$$

The last term can be made arbitrarily small by making M large enough. Thus, for convergence in mean, it is enough to show that the first term goes to zero for each M . Accordingly, fix $M < \infty$. Fix also X_1, \dots, X_n , and let \mathcal{G} be a finite δ -mesh in $L_1(\mathbb{P}_n)$ over \mathcal{F}_M . Thus

$$(8.8) \quad \mathbb{E}_\epsilon \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}_M} \leq \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{G}} + \delta.$$

By definition of the entropy number, the cardinality of \mathcal{G} can be chosen equal to $N(\delta, \mathcal{F}_M, L_1(\mathbb{P}_n))$. Now, we can bound the L_1 -norm on the right-hand-side of (8.8) by the Orlicz-norm for $\psi_2(x) = \exp(x^2) - 1$, and apply the maximal inequality lemma 8.2 to find that the left-hand-side of (8.8) is bounded by a universal constant times

$$\sqrt{1 + \log N(\delta, \mathcal{F}_M, L_1(\mathbb{P}_n))} \sup_{f \in \mathcal{F}} \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\psi_2|X} + \delta,$$

where the Orlicz norms $\|\cdot\|_{\psi_2|X}$ are taken over $\epsilon_1, \dots, \epsilon_n$ with X_1, \dots, X_n still fixed. From exercise 8.5.7 below, we have—by Hoeffding’s inequality (lemma 8.1) combined with lemma 8.1—that the Orlicz norms are all bounded by $\sqrt{6/n} (\mathbb{P}_n f^2)^{1/2}$, which is bounded by $\sqrt{6/n} M$. The last displayed expression is thus bounded by

$$\sqrt{\frac{6\{1 + \log N(\delta, \mathcal{F}_M, L_1(\mathbb{P}_n))\}}{n}} M + \delta \xrightarrow{P} \delta.$$

Thus the left-hand-side of (8.8) goes to zero in probability. Since it is also bounded by M , the bounded convergence theorem implies that its expectation also goes to zero. Since M was arbitrary, we now have that $\mathbb{E}^* \|\mathbb{P}_n - P\|_{\mathcal{F}} \rightarrow 0$. This now implies that \mathcal{F} is weak P -G-C.

From lemma 9.13, we know that there exists a version of $\|\mathbb{P}_n - P\|_{\mathcal{F}}^*$ that converges almost surely to a constant. Since we already know that $\|\mathbb{P}_n - P\|_{\mathcal{F}}^* \xrightarrow{P} 0$, this constant must be zero. The desired result now follows. \square

Proof of proposition 8.9. Assume (i). By the second inequality of the symmetrization theorem (theorem 8.8), $\|\mathbb{P}_n^\circ\|_{\mathcal{F}} \xrightarrow{P} 0$. This convergence can be strengthened to outer almost surely, since $\|\mathbb{P}_n^\circ\|_{\mathcal{F}}$ for a reverse submartingale as in the previous proof. Now assume (ii). By lemma 8.13 and the fact that $P[\epsilon f(X)] = 0$ for a Rademacher ϵ independent of X , we obtain that $P\|f\|_{\mathcal{F}}^* = P\|\epsilon f(X)\|_{\mathcal{F}}^* < \infty$. Now, the fact that \mathcal{F} is weak P -G-C follows from the first inequality in the symmetrization theorem. The convergence can be strengthened to outer almost sure by the reverse martingale argument used previously. Thus (ii) follows. \square

8.4 Donsker Results

We now present several Donsker results. We begin with several interesting necessary and sufficient conditions for a class to be P -Donsker. We next present the proofs of Donsker theorems for bracketing (theorem 2.3) and uniform (theorem 2.5) entropy. Before proceeding, let $\mathcal{F}_\delta \equiv \{f - g : f, g \in \mathcal{F}, \|f - g\|_{P,2} < \delta\}$, for any $0 < \delta \leq \infty$.

The following lemma outlines several properties of Donsker classes and shows that Donsker classes are automatically strong Glivenko–Cantelli classes:

LEMMA 8.17 Let \mathcal{F} be a class of measurable functions, with envelope $F \equiv \|f\|_{\mathcal{F}}$. For any $f, g \in \mathcal{F}$, define $\rho(f, g) \equiv P(f - Pf - g + Pg)^2$; and, for any $\delta > 0$, let $\mathcal{F}_\delta \equiv \{f - g : \rho(f, g) < \delta\}$. Then the following are equivalent:

(i) \mathcal{F} is P -Donsker;

(ii) (\mathcal{F}, ρ) is totally bounded and $\|\mathbb{G}_n\|_{\mathcal{F}_{\delta_n}} \xrightarrow{P} 0$ for every $\delta_n \downarrow 0$;

(iii) (\mathcal{F}, ρ) is totally bounded and $E^*\|\mathbb{G}_n\|_{\mathcal{F}_{\delta_n}} \rightarrow 0$ for every $\delta_n \downarrow 0$.

These conditions imply that $E^*\|\mathbb{G}_n\|_{\mathcal{F}}^r \rightarrow \|\mathbb{G}\|_{\mathcal{F}}^r < \infty$, for every $0 < r < 2$; that $P(\|f - Pf\|_{\mathcal{F}}^* > x) = o(x^{-2})$ as $x \rightarrow \infty$; and that \mathcal{F} is strong P -G-C. If in addition $\|P\|_{\mathcal{F}} < \infty$, then also $P(F^* > x) = o(x^{-2})$ as $x \rightarrow \infty$.

Proof. The equivalence of (i)–(iii) and the first assertion is lemma 2.3.11 of VW, and we omit the equivalence part of the proof. Now assume conditions (i)–(iii) hold. Lemma 2.3.9 of VW states that if \mathcal{F} is Donsker, then

$$(8.9) \quad \lim_{x \rightarrow \infty} x^2 \sup_{n \geq 1} P^*(\|\mathbb{G}_n\|_{\mathcal{F}} > x) = 0.$$

This immediately implies that the r th moment of $\|\mathbb{G}_n\|_{\mathcal{F}}$ is uniformly bounded in n and that $E\|\mathbb{G}\|_{\mathcal{F}}^r < \infty$ for all $0 < r < 2$. Thus the first assertion follows and, therefore, \mathcal{F} is weak P -G-C. Lemma 8.16 now implies \mathcal{F} is strong G-C. Letting $n = 1$ in (8.9) yields that $P(\|f - Pf\|_{\mathcal{F}}^* > x) = o(x^{-2})$ as $x \rightarrow \infty$, and the remaining assertions follow. \square

Proof of theorem 2.3 (Donsker with bracketing entropy). With a given set of $\epsilon/2$ -brackets $[l_i, u_i]$ covering \mathcal{F} , we can construct a set of ϵ -brackets covering \mathcal{F}_∞ by taking differences $[l_i - u_j, u_i - l_j]$ of upper and lower bounds, i.e., if $f \in [l_i, u_i]$ and $g \in [l_j, u_j]$, then $f - g \in [l_i - u_j, u_i - l_j]$. Thus $N_{[]}(\epsilon, \mathcal{F}_\infty, L_2(P)) \leq N_{[]}^2(\epsilon/2, \mathcal{F}, L_2(P))$. From exercise 8.5.8 below, we now have that $J_{[]}(\delta, \mathcal{F}_\infty, L_2(P)) \leq \sqrt{8}J_{[]}(\delta, \mathcal{F}, L_2(P))$.

Now, for a given, small $\delta > 0$, select a minimal number of δ -brackets that cover \mathcal{F} , and use them to construct a finite partition $\mathcal{F} = \cup_{i=1}^m \mathcal{F}_i$ consisting of sets of $L_2(P)$ -diameter δ . For any $i \in \{1, \dots, m\}$, the subset of \mathcal{F}_∞ consisting of all $f - g$ with $f, g \in \mathcal{F}_i$ consists of functions with $L_2(P)$ norm strictly smaller than δ . Hence by lemma 8.18 below, there exists a number $a(\delta) > 0$ satisfying

$$(8.10) \quad E^* \left[\sup_{1 \leq i \leq m} \sup_{f, g \in \mathcal{F}_i} |\mathbb{G}_n(f - g)| \right] \lesssim J_{[]}(\delta, \mathcal{F}, L_2(P)) + \sqrt{n}P^* [G1\{G > a(\delta)\sqrt{n}\}],$$

where G is an envelope for \mathcal{F}_∞ and the relation \lesssim means that the left-hand-side is bounded above by a universal positive constant times the right-hand-side. In this setting, “universal” means that the constant does not depend

on n or δ . If $[l, u]$ is a minimal bracket for covering all of \mathcal{F} , then G can be taken to be $u - l$. Boundedness of the entropy integral implies that there exists some $k < \infty$ so that only one $L_2(P)$ bracket of size k is needed to cover \mathcal{F} . This implies $PG^2 < \infty$.

By exercise 8.5.9 below, the second term on the right-hand-side of (8.10) is bounded above by $[a(\delta)]^{-1}P^* [G^2 \mathbf{1}\{G > a(\delta)\sqrt{n}\}]$ and thus goes to zero as $n \rightarrow \infty$. Since $J_{[]}(\delta, \mathcal{F}, L_2(P)) = o(\delta)$, as $\delta \downarrow 0$, we now have that $\lim_{\delta \downarrow 0} \limsup_{n \rightarrow \infty}$ of the left-hand-side of (8.10) goes to zero. In view of Markov's inequality for outer probability (which follows from Chebyshev's inequality for outer probability as given in lemma 6.10), condition (ii) in lemma 7.20 is satisfied for the stochastic process $X_n(f) = \mathbb{G}_n(f)$ with index set $T = \mathcal{F}$. Now, theorem 2.1 yields the desired result. \square

LEMMA 8.18 *For any class \mathcal{F} of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$ with $Pf^2 < \delta^2$ for every f , we have, with*

$$a(\delta) \equiv \frac{\delta}{\sqrt{1 \vee \log N_{[]}(\delta, \mathcal{F}, L_2(P))}},$$

and F an envelope function for \mathcal{F} , that

$$\mathbb{E}^* \|\mathbb{G}_n\|_{\mathcal{F}} \lesssim J_{[]}(\delta, \mathcal{F}, L_2(P)) + \sqrt{n}P^* [F \mathbf{1}\{F > \sqrt{na}(\delta)\}].$$

This is lemma 19.34 of van der Vaart (1998), who gives a nice proof which utilizes the maximal inequality result given in lemma 8.3 above. The arguments are lengthy, and we omit the proof.

To prove theorem 2.5 (Donsker with uniform entropy), we first restate the theorem with a clarification of the measurability assumption, as done in the previous section for theorem 2.4:

THEOREM 8.19 (Restated theorem 2.5) *Let \mathcal{F} be a class of measurable functions with envelope F and $J(1, \mathcal{F}, L_2) < \infty$. Let the classes \mathcal{F}_δ and $\mathcal{F}_\infty^2 \equiv \{h^2 : h \in \mathcal{F}_\infty\}$ be P -measurable for every $\delta > 0$. If $P^*F^2 < \infty$, then \mathcal{F} is P -Donsker.*

We note here that by proposition 8.11, if \mathcal{F} is PM, then so are \mathcal{F}_δ and \mathcal{F}_∞^2 , for all $\delta > 0$, provided \mathcal{F} has envelope F such that $P^*F^2 < \infty$. Since PM implies P -measurability, all measurability requirements for theorem 8.19 are thus satisfied whenever \mathcal{F} is PM.

Proof of theorem 8.19. Let the positive, decreasing sequence $\delta_n \downarrow 0$ be arbitrary. By Markov's inequality for outer probability (see lemma 6.10) and the symmetrization theorem 8.8,

$$P^* (\|\mathbb{G}_n\|_{\mathcal{F}_{\delta_n}} > x) \leq \frac{2}{x} \mathbb{E}^* \left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}_{\delta_n}},$$

for i.i.d. Rademachers $\epsilon_1, \dots, \epsilon_n$ independent of X_1, \dots, X_n . By the P -measurability assumption for \mathcal{F}_δ , for all $\delta > 0$, the standard version of Fu-

bini's theorem applies, and the outer expectation is just a standard expectation and can be calculated in the order $E_X E_\epsilon$. Accordingly, fix X_1, \dots, X_n . By Hoeffding's inequality (lemma 8.1), the stochastic process $f \mapsto n^{-1/2} \times \sum_{i=1}^n \epsilon_i f(X_i)$ is sub-Gaussian for the $L_2(\mathbb{P}_n)$ -seminorm

$$\|f\|_n \equiv \sqrt{\frac{1}{n} \sum_{i=1}^n f^2(X_i)}.$$

This stochastic process is also separable since, for any measure Q and $\epsilon > 0$, $N(\epsilon, \mathcal{F}_{\delta_n}, L_2(Q)) \leq N(\epsilon, \mathcal{F}_\infty, L_2(Q)) \leq N^2(\epsilon/2, \mathcal{F}, L_2(Q))$, and the latter is finite for any finite dimensional probability measure Q and any $\epsilon > 0$. Thus the second conclusion of corollary 8.5 holds with

$$(8.11) \quad E_\epsilon \left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}_{\delta_n}} \lesssim \int_0^\infty \sqrt{\log N(\epsilon, \mathcal{F}_{\delta_n}, L_2(\mathbb{P}_n))} d\epsilon.$$

Note that we can omit the term $E |n^{-1/2} \sum_{i=1}^n \epsilon_i f_0(X_i)|$ from the conclusion of the corollary because it is also bounded by the right-hand-side of (8.11).

For sufficiently large ϵ , the set \mathcal{F}_{δ_n} fits in a single ball of $L_2(\mathbb{P}_n)$ -radius ϵ around the origin, in which case the integrand on the right-hand-side of (8.11) is zero. This will definitely happen when ϵ is larger than θ_n , where

$$\theta_n^2 \equiv \sup_{f \in \mathcal{F}_{\delta_n}} \|f\|_n^2 = \left\| \frac{1}{n} \sum_{i=1}^n f^2(X_i) \right\|_{\mathcal{F}_{\delta_n}}.$$

Thus the right-hand-side of (8.11) is bounded by

$$\begin{aligned} & \int_0^{\theta_n} \sqrt{\log N(\epsilon, \mathcal{F}_{\delta_n}, L_2(\mathbb{P}_n))} d\epsilon \\ & \lesssim \int_0^{\theta_n} \sqrt{\log N^2(\epsilon/2, \mathcal{F}, L_2(\mathbb{P}_n))} d\epsilon \\ & \lesssim \int_0^{\theta_n/(2\|F\|_n)} \sqrt{\log N(\epsilon\|F\|_n, \mathcal{F}, L_2(\mathbb{P}_n))} d\epsilon \|F\|_n \\ & \lesssim \|F\|_n J(\theta_n, \mathcal{F}, L_2). \end{aligned}$$

The second inequality follows from the change of variables $u = \epsilon/(2\|F\|_n)$ (and then renaming u to ϵ). For the third inequality, note that we can add $1/2$ to the envelope function F without changing the existence of its second moment. Hence $\|F\|_n \geq 1/2$ without loss of generality, and thus $\theta_n/(2\|F\|_n) \leq \theta_n$. Because $\|F\|_n = O_p(1)$, we can now conclude that the left-hand-side of (8.11) goes to zero in probability, provided we can verify that $\theta_n \xrightarrow{P} 0$. This would then imply asymptotic $L_2(P)$ -equicontinuity in probability.

Since $\|Pf^2\|_{\mathcal{F}_{\delta_n}} \rightarrow 0$ and $\mathcal{F}_{\delta_n} \subset \mathcal{F}_\infty$, establishing that $\|\mathbb{P}_n - P\|_{\mathcal{F}_\infty} \xrightarrow{P} 0$ would prove that $\theta_n \xrightarrow{P} 0$. The class \mathcal{F}_∞^2 has integrable envelope $(2F)^2$ and is P -measurable by assumption. Since also, for any $f, g \in \mathcal{F}_\infty$, $|\mathbb{P}_n(f^2 - g^2)| \leq \mathbb{P}_n(|f - g|4F) \leq \|f - g\|_n \|4F\|_n$, we have that the covering number $N(\epsilon \|2F\|_n^2, \mathcal{F}_\infty^2, L_1(\mathbb{P}_n))$ is bounded by $N(\epsilon \|F\|_n, \mathcal{F}_\infty, L_2(\mathbb{P}_n))$. Since this last covering number is bounded by $\sup_Q N^2(\epsilon \|F\|_{Q,2}/2, \mathcal{F}, L_2(Q)) < \infty$, where the supremum is taken over all finitely discrete probability measures with $\|F\|_{Q,2} > 0$, we have by theorem 8.14 that \mathcal{F}_∞^2 is P -Glivenko-Cantelli. Thus $\hat{\theta}_n \xrightarrow{P} 0$. This completes the proof of asymptotic equicontinuity.

The last thing we need to prove is that \mathcal{F} is totally bounded in $L_2(P)$. By the result of the last paragraph, there exists a sequence of discrete probability measures P_n with $\|(P_n - P)f^2\|_{\mathcal{F}_\infty} \rightarrow 0$. Fix $\epsilon > 0$ and take n large enough so that $\|(P_n - P)f^2\|_{\mathcal{F}_\infty} < \epsilon^2$. Note that $N(\epsilon, \mathcal{F}, L_2(P_n))$ is finite by assumption, and, for any $f, g \in \mathcal{F}$ with $\|f - g\|_{P_n,2} < \epsilon$, $P(f - g)^2 \leq P_n(f - g)^2 + |(P_n - P)(f - g)^2| \leq 2\epsilon^2$. Thus any ϵ -net in $L_2(P_n)$ is also a $\sqrt{2}\epsilon$ -net in $L_2(P)$. Hence \mathcal{F} is totally bounded in $L_2(P)$ since ϵ was arbitrary. \square

8.5 Exercises

8.5.1. For any ψ valid for defining an Orlicz norm $\|\cdot\|_\psi$, show that the space H_ψ of real random variables X with $\|X\|_\psi < \infty$ defines a Banach space, where we equate a random variable X with zero if $X = 0$ almost surely:

- (a) Show first that $\|\cdot\|_\psi$ defines a norm on H_ψ . Hint: Use the convexity of ψ to establish that for any $X, Y \in H_\psi$ and any $c_1, c_2 > 0$,

$$\mathbb{E}\psi\left(\frac{|X+Y|}{c_1+c_2}\right) \leq \frac{c_1}{c_1+c_2}\mathbb{E}\psi\left(\frac{|X|}{c_1}\right) + \frac{c_2}{c_1+c_2}\mathbb{E}\psi\left(\frac{|Y|}{c_2}\right).$$

- (b) Now show that H_ψ is complete. Hint: Show that for any Cauchy sequence of random variables $\{X_n\} \in H_\psi$, $\limsup_{n \rightarrow \infty} \|X_n\|_\psi < \infty$. Use Prohorov's theorem to show that every such Cauchy sequence converges to an almost surely unique element of H_ψ .

8.5.2. Show that $1 \wedge (e^u - 1)^{-1} \leq 2e^{-u}$, for any $u > 0$.

8.5.3. For a function ψ meeting the conditions of lemma 8.2, show that there exists constants $0 < \sigma \leq 1$ and $\tau > 0$ such that $\phi(x) \equiv \sigma\psi(\tau x)$ satisfies $\phi(x)\phi(y) \leq \phi(cxy)$ for all $x, y \geq 1$ and $\phi(1) \leq 1/2$. Show that this ϕ also satisfies the following

- (a) For all $u > 0$, $\phi^{-1}(u) \leq \psi^{-1}(u)/(\sigma\tau)$.

(b) For any random variable X , $\|X\|_\psi \leq \|X\|_\phi / (\sigma\tau) \leq \|X\|_\psi / \sigma$.

8.5.4. Show that for any $p \in [1, \infty)$, ψ_p satisfies the conditions of lemma 8.2 with $c = 1$.

8.5.5. Let ψ satisfy the conditions of lemma 8.2. Show that for any sequence of random variables $\{X_n\}$, $\|X_n\|_\psi \rightarrow 0$ implies $X_n \xrightarrow{P} 0$. Hint: Show that $\liminf_{x \rightarrow \infty} \psi(x)/x > 0$, and hence $\|X_n\|_\psi \rightarrow 0$ implies $E|X_n| \rightarrow 0$.

8.5.6. Show that if the class of functions \mathcal{F} has a countable subset $\mathcal{G} \subset \mathcal{F}$ such that for each $f \in \mathcal{F}$ there exists a sequence $\{g_m\} \in \mathcal{G}$ with $g_m(x) \rightarrow f(x)$ for every $x \in \mathcal{X}$, then $\|\sum e_i f(X_i)\|_{\mathcal{F}} = \|\sum e_i f(X_i)\|_{\mathcal{G}}$ for all $(e_1, \dots, e_n) \in \mathbb{R}^n$.

8.5.7. In the context of the proof of theorem 8.15, use Hoeffding's inequality (lemma 8.7) combined with lemma 8.1 to show that

$$\left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\psi_2|X} \leq \sqrt{6/n} (\mathbb{P}_n f^2)^{1/2},$$

where the meaning of the subscript $\psi_2|X$ is given in the body of the proof.

8.5.8. In the context of the proof of theorem 2.3 above, show that by taking logarithms followed by square roots,

$$J_{[]}(\delta, \mathcal{F}_\infty, L_2(P)) \leq \sqrt{8} J_{[]}(\delta, \mathcal{F}, L_2(P)).$$

8.5.9. Show, for any map $X : \Omega \mapsto \mathbb{R}$ and constants $\alpha \in [0, \infty)$ and $m \in (0, \infty)$, that

$$E^* [|X|^\alpha 1\{|X| > m\}] \leq m^{-1} E^* [|X|^{\alpha+1} 1\{|X| > m\}].$$

8.6 Notes

Lemma 8.2, corollary 8.5, and theorems 8.15 and 8.19 correspond to lemma 2.2.1, corollary 2.2.8, and theorems 2.4.3 and 2.5.2 of VW, respectively. The first inequality in theorem 8.8 corresponds to lemma 2.3.1 of VW. Lemma 8.3 is a modification of lemma 2.2.10 of VW, and theorem 8.4 is a modification and combination of both theorem 2.2.4 and corollary 2.2.5 of VW. The version of Hoeffding's inequality we use (lemma 8.7) is lemma 2.2.7 of VW, and lemma 8.13 was inspired by exercise 2.4.1 of VW. The proof of theorem 2.3 follows closely van der Vaart's (1998) proof of his theorem 19.5.

9

Entropy Calculations

The focus of this chapter is on computing entropy for empirical processes. An important use of such entropy calculations is in evaluating whether a class of functions \mathcal{F} is Glivenko-Cantelli and/or Donsker or neither. Several additional uses of entropy bounds will be discussed in chapter 11. Some of these uses will be very helpful in chapter 14 for establishing rates of convergence for M-estimators. An additional use of entropy bounds, one which will receive only limited discussion in this book, is in precisely assessing the asymptotic smoothness of certain empirical processes. Such smoothness results play a role in the asymptotic analysis of a number of statistical applications, including confidence bands for kernel density estimation (eg., Bickel and Rosenblatt, 1973) and certain hypothesis tests for multimodality (Polonik, 1995).

We begin the chapter by describing methods to evaluate uniform entropy. Provided the uniform entropy for a class \mathcal{F} is not too large, \mathcal{F} might be G-C or Donsker, as long as sufficient measurability holds. Since many of the techniques we will describe for building bounded uniform entropy integral (BUEI) classes (which include VC classes) closely parallel the methods for building pointwise measurable (PM) classes described in the previous chapter, we will include a discussion on joint BUEI and PM preservation towards the end of section 9.1.2. We then present methods based on bracketing entropy. Several important results for building G-C classes from other G-C classes (*G-C preservation*), are presented next. Finally, we discuss several useful Donsker preservation results.

One can think of this chapter as a handbag of tools for establishing weak convergence properties of empirical processes. Illustrations of how to

use these tools will be given in various applications scattered throughout later chapters. To help anchor the context for these tools in practice, it might be worthwhile rereading the counting process regression example of section 4.2.1, in the first case studies chapter of this book. In the second case studies chapter of this book (chapter 15), we will provide additional—and more complicated—illustrations of these tools, with special emphasis on Donsker preservation techniques.

9.1 Uniform Entropy

We first discuss the very powerful concept of VC-classes of sets and functions. Such classes are extremely important tools in assessing and controlling bounded uniform entropy. We then discuss several useful and powerful preservation results for bounded uniform entropy integral (BUEI) classes.

9.1.1 VC-Classes

In this section, we introduce Vapnik-Červonenkis (VC) classes of sets, VC-classes of functions, and several related function classes. We then present several examples of VC-classes.

Consider an arbitrary collection $\{x_1, \dots, x_n\}$ of points in a set \mathcal{X} and a collection \mathcal{C} of subsets of \mathcal{X} . We say that \mathcal{C} *picks out* a certain subset A of $\{x_1, \dots, x_n\}$ if $A = C \cap \{x_1, \dots, x_n\}$ for some $C \in \mathcal{C}$. We say that \mathcal{C} *shatters* $\{x_1, \dots, x_n\}$ if all of the 2^n possible subsets of $\{x_1, \dots, x_n\}$ are picked out by the sets in \mathcal{C} . The *VC-index* $V(\mathcal{C})$ of the class \mathcal{C} is the smallest n for which no set of size n $\{x_1, \dots, x_n\} \subset \mathcal{X}$ is shattered by \mathcal{C} . If \mathcal{C} shatters all sets $\{x_1, \dots, x_n\}$ for all $n \geq 1$, we set $V(\mathcal{C}) = \infty$. Clearly, the more refined \mathcal{C} is, the higher the VC-index. We say that \mathcal{C} is a *VC-class* if $V(\mathcal{C}) < \infty$.

For example, let $\mathcal{X} = \mathbb{R}$ and define the collection of sets $\mathcal{C} = \{(-\infty, c] : c \in \mathbb{R}\}$. Consider any two point set $\{x_1, x_2\} \subset \mathbb{R}$ and assume, without loss of generality, that $x_1 < x_2$. It is easy to verify that \mathcal{C} can pick out the null set $\{\}$ and the sets $\{x_1\}$ and $\{x_1, x_2\}$ but can't pick out $\{x_2\}$. Thus $V(\mathcal{C}) = 2$ and \mathcal{C} is a VC-class. As another example, let $\mathcal{C} = \{(a, b] : -\infty \leq a < b \leq \infty\}$. The collection can shatter any two point set, but consider what happens with a three point set $\{x_1, x_2, x_3\}$. Without loss of generality, assume $x_1 < x_2 < x_3$, and note that the set $\{x_1, x_3\}$ cannot be picked out with \mathcal{C} . Thus $V(\mathcal{C}) = 3$ in this instance.

For any class of sets \mathcal{C} and any collection $\{x_1, \dots, x_n\} \subset \mathcal{X}$, let $\Delta_n(\mathcal{C}, x_1, \dots, x_n)$ be the number of subsets of $\{x_1, \dots, x_n\}$ which can be picked out by \mathcal{C} . A surprising combinatorial result is that if $V(\mathcal{C}) < \infty$, then $\Delta_n(\mathcal{C}, x_1, \dots, x_n)$ can increase in n no faster than $O(n^{V(\mathcal{C})-1})$. This is more precisely stated in the following lemma:

LEMMA 9.1 *For a VC-class of sets \mathcal{C} ,*

$$\max_{x_1, \dots, x_n \in \mathcal{X}} \Delta_n(\mathcal{C}, x_1, \dots, x_n) \leq \sum_{j=1}^{V(\mathcal{C})-1} \binom{n}{j}.$$

Since the right-hand-side is bounded by $V(\mathcal{C})n^{V(\mathcal{C})-1}$, the left-hand-side grows polynomially of order at most $O(n^{V(\mathcal{C})-1})$.

This is a corollary of lemma 2.6.2 of VW, and we omit the proof.

Let $1\{\mathcal{C}\}$ denote the collection of all indicator functions of sets in the class \mathcal{C} . The following theorem gives a bound on the L_r covering numbers of $1\{\mathcal{C}\}$:

THEOREM 9.2 *There exists a universal constant $K < \infty$ such that for any VC-class of sets \mathcal{C} , any $r \geq 1$, and any $0 < \epsilon < 1$,*

$$N(\epsilon, 1\{\mathcal{C}\}, L_r(Q)) \leq KV(\mathcal{C})(4e)^{V(\mathcal{C})} \left(\frac{1}{\epsilon}\right)^{r(V(\mathcal{C})-1)}.$$

This is theorem 2.6.4 of VW, and we omit the proof. Since $F = 1$ serves as an envelope for $1\{\mathcal{C}\}$, we have as an immediate corollary that, for $\mathcal{F} = 1\{\mathcal{C}\}$, $\sup_Q N(\epsilon \|F\|_{1,Q}, \mathcal{F}, L_1(Q)) < \infty$ and

$$J(1, \mathcal{F}, L_2) \lesssim \int_0^1 \sqrt{\log(1/\epsilon)} d\epsilon = \int_0^\infty u^{1/2} e^{-u} du \leq 1,$$

where the supremum is over all finite probability measures Q with $\|F\|_{Q,2} > 0$. Thus the uniform entropy conditions required in the G-C and Donsker theorems of the previous chapter are satisfied for indicators of VC-classes of sets. Since the constant 1 serves as a universally applicable envelope function, these classes of indicator functions are therefore G-C and Donsker, provided the requisite measurability conditions hold.

For a function $f : \mathcal{X} \mapsto \mathbb{R}$, the subset of $\mathcal{X} \times \mathbb{R}$ given by $\{(x, t) : t < f(x)\}$ is the *subgraph* of f . A collection \mathcal{F} of measurable real functions on the sample space \mathcal{X} is a *VC-subgraph class* or *VC-class* (for short), if the collection of all subgraphs of functions in \mathcal{F} forms a VC-class of sets (as sets in $\mathcal{X} \times \mathbb{R}$). Let $V(\mathcal{F})$ denote the VC-index of the set of subgraphs of \mathcal{F} . The following theorem, the proof of which is given in section 9.5, shows that covering numbers of VC-classes of functions grow at a polynomial rate just like VC-classes of sets:

THEOREM 9.3 *There exists a universal constant $K < \infty$ such that, for any VC-class of measurable functions \mathcal{F} with integrable envelope F , any $r \geq 1$, any probability measure Q with $\|F\|_{Q,r} > 0$, and any $0 < \epsilon < 1$,*

$$N(\epsilon \|F\|_{Q,r}, \mathcal{F}, L_r(Q)) \leq KV(\mathcal{F})(4e)^{V(\mathcal{F})} \left(\frac{2}{\epsilon}\right)^{r(V(\mathcal{F})-1)}.$$

Thus VC-classes of functions easily satisfy the uniform entropy requirements of the G-C and Donsker theorems of the previous chapter. A related kind of function class is the *VC-hull* class. A class of measurable functions \mathcal{G} is a VC-hull class if there exists a VC-class \mathcal{F} such that each $f \in \mathcal{G}$ is the pointwise limit of a sequence of functions $\{f_m\}$ in the *symmetric convex hull* of \mathcal{F} (denoted $\text{sconv}\mathcal{F}$). A function f is in $\text{sconv}\mathcal{F}$ if $f = \sum_{i=1}^m \alpha_i f_i$, where the α_i s are real numbers satisfying $\sum_{i=1}^m |\alpha_i| \leq 1$ and the f_i s are in \mathcal{F} . The *convex hull* of a class of functions \mathcal{F} , denoted $\text{conv}\mathcal{F}$, is similarly defined but with the requirement that the α_i 's are positive. We use $\overline{\text{conv}}\mathcal{F}$ to denote pointwise closure of $\text{conv}\mathcal{F}$ and $\overline{\text{sconv}}\mathcal{F}$ to denote the pointwise closure of $\text{sconv}\mathcal{F}$. Thus the class of functions \mathcal{F} is a VC-hull class if $\mathcal{F} = \overline{\text{sconv}}\mathcal{G}$ for some VC-class \mathcal{G} . The following theorem provides a useful relationship between the L_2 covering numbers of a class \mathcal{F} (not necessarily a VC-class) and the L_2 covering numbers of $\overline{\text{conv}}\mathcal{F}$ when the covering numbers for \mathcal{F} are polynomial in $1/\epsilon$:

THEOREM 9.4 *Let Q be a probability measure on $(\mathcal{X}, \mathcal{A})$, and let \mathcal{F} be a class of measurable functions with measurable envelope F , such that $QF^2 < \infty$ and, for $0 < \epsilon < 1$,*

$$N(\epsilon \|F\|_{Q,2}, \mathcal{F}, L_2(Q)) \leq C \left(\frac{1}{\epsilon}\right)^V,$$

for constants $C, V < \infty$ (possibly dependent on Q). Then there exist a constant K depending only on V and C such that

$$\log N(\epsilon \|F\|_{Q,2}, \overline{\text{conv}}\mathcal{F}, L_2(Q)) \leq K \left(\frac{1}{\epsilon}\right)^{2V/(V+2)}.$$

This is theorem 2.6.9 of VW, and we omit the proof.

It is not hard to verify that $\text{sconv}\mathcal{F}$ is a subset of the convex hull of $\mathcal{F} \cup \{-\mathcal{F}\} \cup \{0\}$, where $-\mathcal{F} \equiv \{-f : f \in \mathcal{F}\}$ (see exercise 9.6.1 below). Since the covering numbers of $\mathcal{F} \cup \{-\mathcal{F}\} \cup \{0\}$ are at most one plus twice the covering numbers of \mathcal{F} , the conclusion of theorem 9.4 also holds if $\overline{\text{conv}}\mathcal{F}$ is replaced with $\overline{\text{sconv}}\mathcal{F}$. This leads to the following easy corollary for VC-hull classes, the proof of which we save as an exercise:

COROLLARY 9.5 *For any VC-hull class \mathcal{F} of measurable functions and all $0 < \epsilon < 1$,*

$$\sup_Q \log N(\epsilon \|F\|_{Q,2}, \mathcal{F}, L_2(Q)) \leq K \left(\frac{1}{\epsilon}\right)^{2-2/V}, \quad 0 < \epsilon < 1,$$

where the supremum is taken over all probability measures Q with $\|F\|_{Q,2} > 0$, V is the VC-index of the VC-subgraph class associated with \mathcal{F} , and the constant $K < \infty$ depends only on V .

We now present several important examples and results about VC-classes of sets and both VC-subgraph and VC-hull classes of functions. The first lemma, lemma 9.6, applies to vector spaces of functions, a frequently occurring function class in statistical applications.

LEMMA 9.6 *Let \mathcal{F} be a finite-dimensional vector space of measurable functions $f : \mathcal{X} \mapsto \mathbb{R}$. Then \mathcal{F} is VC-subgraph with $V(\mathcal{F}) \leq \dim(\mathcal{F}) + 2$.*

Proof. Fix any collection G of $k = \dim(\mathcal{F}) + 2$ points $(x_1, t_1), \dots, (x_k, t_k)$ in $\mathcal{X} \times \mathbb{R}$. By assumption, the vectors $(f(x_1) - t_1, \dots, f(x_k) - t_k)^T$, as f ranges over \mathcal{F} , are restricted to a $\dim(\mathcal{F}) + 1 = k - 1$ -dimensional subspace H of \mathbb{R}^k . Any vector $0 \neq a \in \mathbb{R}^k$ orthogonal to H satisfies

$$(9.1) \quad \sum_{j:a_j>0} a_j(f(x_j) - t_j) = \sum_{j:a_j<0} (-a_j)(f(x_j) - t_j),$$

for all $f \in \mathcal{F}$, where we define sums over empty sets to be zero. There exists such a vector a with at least one strictly positive coordinate. For this a , the subset of G of the form $\{(x_j, t_j) : a_j > 0\}$ cannot also be of the form $\{(x_j, t_j) : t_j < f(x_j)\}$ for any $f \in \mathcal{F}$, since otherwise the left side of the equation (9.1) would be strictly positive while the right side would be nonpositive. Conclude that the subgraphs of \mathcal{F} cannot pick out the subset $\{(x_j, t_j) : a_j > 0\}$. Since G was arbitrary, we have just shown that the subgraphs of \mathcal{F} cannot shatter any set of k points in $\mathcal{X} \times \mathbb{R}$. The desired result now follows. \square .

The next three lemmas, lemmas 9.7 through 9.9, consist of useful tools for building VC-classes from other VC-classes. One of these lemmas, lemma 9.8, is the important identity that classes of sets are VC if and only if the corresponding classes of indicator functions are VC-subgraph. The proof of lemma 9.9 is relegated to section 9.5.

LEMMA 9.7 *Let \mathcal{C} and \mathcal{D} be VC-classes of sets in a set \mathcal{X} , with respective VC-indices $V_{\mathcal{C}}$ and $V_{\mathcal{D}}$; and let \mathcal{E} be a VC-class of sets in \mathcal{W} , with VC-index $V_{\mathcal{E}}$. Also let $\phi : \mathcal{X} \mapsto \mathcal{Y}$ and $\psi : \mathcal{Z} \mapsto \mathcal{X}$ be fixed functions. Then*

- (i) $\mathcal{C}^c \equiv \{C^c : C \in \mathcal{C}\}$ is VC with $V(\mathcal{C}^c) = V_{\mathcal{C}}$;
- (ii) $\mathcal{C} \cap \mathcal{D} \equiv \{C \cap D : C \in \mathcal{C}, D \in \mathcal{D}\}$ is VC with index $\leq V_{\mathcal{C}} + V_{\mathcal{D}} - 1$;
- (iii) $\mathcal{C} \sqcup \mathcal{D} \equiv \{C \cup D : C \in \mathcal{C}, D \in \mathcal{D}\}$ is VC with index $\leq V_{\mathcal{C}} + V_{\mathcal{D}} - 1$;
- (iv) $\mathcal{D} \times \mathcal{E}$ is VC in $\mathcal{X} \times \mathcal{W}$ with VC index $\leq V_{\mathcal{D}} + V_{\mathcal{E}} - 1$;
- (v) $\phi(\mathcal{C})$ is VC with index $V_{\mathcal{C}}$ if ϕ is one-to-one;
- (vi) $\psi^{-1}(\mathcal{C})$ is VC with index $\leq V_{\mathcal{C}}$.

Proof. For any $C \in \mathcal{C}$, the set C^c picks out the points of a given set x_1, \dots, x_n that C does not pick out. Thus if \mathcal{C} shatters a given set of points, so does \mathcal{C}^c , and vice versa. This proves (i). From n points, \mathcal{C} can pick out $O(n^{V_{\mathcal{C}}-1})$ subsets. From each of these subsets, \mathcal{D} can pick out $O(n^{V_{\mathcal{D}}-1})$ further subsets. Hence $\mathcal{C} \cap \mathcal{D}$ can pick out at most $O(n^{V_{\mathcal{C}}+V_{\mathcal{D}}-2})$ subsets, and thus (ii) follows from the definition of a VC-class. We save it as an exercise to show that (i) and (ii) imply (iii). For (iv), note that $\mathcal{D} \times \mathcal{W}$ and $\mathcal{X} \times \mathcal{E}$ are both VC-classes with respective VC-indices $V_{\mathcal{D}}$ and $V_{\mathcal{E}}$. The desired conclusion now follows from part (ii), since $\mathcal{D} \times \mathcal{E} = (\mathcal{D} \times \mathcal{W}) \cap (\mathcal{X} \times \mathcal{E})$.

For part (v), if $\phi(\mathcal{C})$ shatters $\{y_1, \dots, y_n\}$, then each y_i must be in the range of ϕ and there must therefore exist x_1, \dots, x_n such that ϕ is a bijection between x_1, \dots, x_n and y_1, \dots, y_n . Hence \mathcal{C} must shatter x_1, \dots, x_n , and thus $V(\phi(\mathcal{C})) \leq V(\mathcal{C})$. Conversely, it is obvious that if \mathcal{C} shatters x_1, \dots, x_n , then $\phi(\mathcal{C})$ shatters $\phi(x_1), \dots, \phi(x_n)$. Hence $V(\mathcal{C}) \leq V(\phi(\mathcal{C}))$. For (vi), if $\psi^{-1}(\mathcal{C})$ shatters z_1, \dots, z_n , then all $\psi(z_i)$ must be distinct and the restriction of ψ to z_1, \dots, z_n is a bijection onto its range. Thus \mathcal{C} shatters $\psi(z_1), \dots, \psi(z_n)$, and hence $V(\psi^{-1}(\mathcal{C})) \leq V(\mathcal{C})$. \square

LEMMA 9.8 *For any class \mathcal{C} of sets in a set \mathcal{X} , the class $\mathcal{F}_{\mathcal{C}}$ of indicator functions of sets in \mathcal{C} is VC-subgraph if and only if \mathcal{C} is a VC-class. Moreover, whenever at least one of \mathcal{C} or $\mathcal{F}_{\mathcal{C}}$ is VC, the respective VC-indices are equal.*

Proof. Let \mathcal{D} be the collection of sets of the form $\{(x, t) : t < 1\{x \in C\}\}$ for all $C \in \mathcal{C}$. Suppose that \mathcal{D} is VC, and let $k = V(\mathcal{D})$. Then no set of the form $\{(x_1, 0), \dots, (x_k, 0)\}$ can be shattered by \mathcal{D} , and hence $V(\mathcal{C}) \leq V(\mathcal{D})$. Now suppose that \mathcal{C} is VC with VC-index k . Since for any $t < 0$, $1\{x \in C\} > t$ for all x and all C , we have that no collection $\{(x_1, t_1), \dots, (x_k, t_k)\}$ can be shattered by \mathcal{D} if any of the t_j s are < 0 . It is similarly true that no collection $\{(x_1, t_2), \dots, (x_k, t_k)\}$ can be shattered by \mathcal{D} if any of the t_j s are ≥ 1 , since $1\{x \in C\} > t$ is never true when $t \geq 1$. It can now be deduced that $\{(x_1, t_1), \dots, (x_k, t_k)\}$ can only be shattered if $\{(x_1, 0), \dots, (x_k, 0)\}$ can be shattered. But this can only happen if $\{x_1, \dots, x_k\}$ can be shattered by \mathcal{C} . Thus $V(\mathcal{D}) \leq V(\mathcal{C})$. \square

LEMMA 9.9 *Let \mathcal{F} and \mathcal{G} be VC-subgraph classes of functions on a set \mathcal{X} , with respective VC indices $V_{\mathcal{F}}$ and $V_{\mathcal{G}}$. Let $g : \mathcal{X} \mapsto \mathbb{R}$, $\phi : \mathbb{R} \mapsto \mathbb{R}$, and $\psi : \mathbb{Z} \mapsto \mathcal{X}$ be fixed functions. Then*

- (i) $\mathcal{F} \wedge \mathcal{G} \equiv \{f \wedge g : f \in \mathcal{F}, g \in \mathcal{G}\}$ is VC-subgraph with index $\leq V_{\mathcal{F}} + V_{\mathcal{G}} - 1$;
- (ii) $\mathcal{F} \vee \mathcal{G}$ is VC with index $\leq V_{\mathcal{F}} + V_{\mathcal{G}} - 1$;
- (iii) $\{\mathcal{F} > 0\} \equiv \{\{f > 0\} : f \in \mathcal{F}\}$ is a VC-class of sets with index $V_{\mathcal{F}}$;
- (iv) $-\mathcal{F}$ is VC-subgraph with index $V_{\mathcal{F}}$;

- (v) $\mathcal{F} + g \equiv \{f + g : f \in \mathcal{F}\}$ is VC with index $V_{\mathcal{F}}$;
- (vi) $\mathcal{F} \cdot g \equiv \{fg : f \in \mathcal{F}\}$ is VC with index $\leq 2V_{\mathcal{F}} - 1$;
- (vii) $\mathcal{F} \circ \psi \equiv \{f(\psi) : f \in \mathcal{F}\}$ is VC with index $\leq V_{\mathcal{F}}$;
- (viii) $\phi \circ \mathcal{F}$ is VC with index $\leq V_{\mathcal{F}}$ for monotone ϕ .

The next two lemmas, lemmas 9.10 and 9.11, refer to properties of monotone processes and classes of monotone functions. The proof of lemma 9.11 is relegated to section 9.5.

LEMMA 9.10 *Let $\{X(t), t \in T\}$ be a monotone increasing stochastic process, where $T \subset \mathbb{R}$. Then X is VC-subgraph with index $V(X) = 2$.*

Proof. Let \mathcal{X} be the set of all monotone increasing functions $g : T \mapsto \mathbb{R}$; and for any $s \in T$ and $x \in \mathcal{X}$, define $(s, x) \mapsto f_s(x) = x(s)$. Thus the proof is complete if we can show that the class of functions $\mathcal{F} \equiv \{f_s : s \in T\}$ is VC-subgraph with VC index 2. Now let $(x_1, t_1), (x_2, t_2)$ be any two points in $\mathcal{X} \times \mathbb{R}$. \mathcal{F} shatters $(x_1, t_1), (x_2, t_2)$ if the graph \mathcal{G} of $(f_s(x_1), f_s(x_2))$ in \mathbb{R}^2 “surrounds” the point (t_1, t_2) as s ranges over T . By surrounding a point $(a, b) \in \mathbb{R}^2$, we mean that the graph must pass through all four of the sets $\{(u, v) : u \leq a, v \leq b\}$, $\{(u, v) : u > a, v \leq b\}$, $\{(u, v) : u \leq a, v > b\}$ and $\{(u, v) : u > a, v > b\}$. By the assumed monotonicity of x_1 and x_2 , the graph \mathcal{G} forms a monotone curve in \mathbb{R}^2 , and it is thus impossible for it to surround any point in \mathbb{R}^2 . Thus $(x_1, t_1), (x_2, t_2)$ cannot be shattered by \mathcal{F} , and the desired result follows. \square

LEMMA 9.11 *The set \mathcal{F} of all monotone functions $f : \mathbb{R} \mapsto [0, 1]$ satisfies*

$$\sup_Q \log N(\epsilon, \mathcal{F}, L_2(Q)) \leq \frac{K}{\epsilon}, \quad 0 < \epsilon < 1,$$

where the supremum is taken over all probability measures Q , and the constant $K < \infty$ is universal.

The final lemma of this section, lemma 9.12 below, addresses the claim raised in section 4.1 that the class of functions $\mathcal{F} \equiv \{1\{Y - b'Z \leq v\} : b \in \mathbb{R}^k, v \in \mathbb{R}\}$ is Donsker. Because the indicator functions in \mathcal{F} are a subset of the indicator functions for half-spaces in \mathbb{R}^{k+1} , part (i) of the lemma implies that \mathcal{F} is VC with index $k + 3$. Since lemma 8.12 from chapter 8 verifies that \mathcal{F} , \mathcal{F}_δ and \mathcal{F}_∞^2 are all P -measurable, for any probability measure P , theorem 9.3 combined with theorem 8.19 and the fact that indicator functions are bounded, establishes that \mathcal{F} is P -Donsker for any P . Lemma 9.12 also gives a related result on closed balls in \mathbb{R}^d . In the lemma, $\langle a, b \rangle$ denotes the Euclidean inner product.

LEMMA 9.12 *The following are true:*

(i) The collection of all half-spaces in \mathbb{R}^d , consisting of the sets $\{x \in \mathbb{R}^d : \langle x, u \rangle \leq c\}$ with u ranging over \mathbb{R}^d and c ranging over \mathbb{R} , is VC with index $d + 2$.

(ii) The collection of all closed balls in \mathbb{R}^d is VC with index $\leq d + 3$.

Proof. The class \mathcal{A}^+ of sets $\{x : \langle x, u \rangle \leq c\}$, with u ranging over \mathbb{R}^d and c ranging over $(0, \infty)$, is equivalent to the class of sets $\{x : \langle x, u \rangle - 1 \leq 0\}$ with u ranging over \mathbb{R}^d . In this last class, since $\langle x, u \rangle$ spans a d -dimensional vector space, lemma 9.6 and part (v) of lemma 9.9 yield that the class of functions spanned by $\langle x, u \rangle - 1$ is VC with index $d + 2$. Part (iii) of lemma 9.9 combined with part (i) of lemma 9.7 now yields that the class \mathcal{A}^+ is VC with index $d + 2$. Similar arguments verify that both the class \mathcal{A}^- , with c restricted to $(-\infty, 0)$, and the class \mathcal{A}^0 , with $c = 0$, are VC with index $d + 2$. It is easy to verify that the union of finite VC classes has VC index equal to the maximum of the respective VC indices. This concludes the proof of (i).

Closed balls in \mathbb{R}^d are sets of the form $\{x : \langle x, x \rangle - 2\langle x, u \rangle + \langle u, u \rangle \leq c\}$, where u ranges over \mathbb{R}^d and c ranges over $[0, \infty)$. It is straightforward to check that the class \mathcal{G} all functions of the form $x \mapsto -2\langle x, u \rangle + \langle u, u \rangle - c$ are contained in a $d + 1$ dimensional vector space, and thus \mathcal{G} is VC with index $\leq d + 3$. Combining this with part (v) of lemma 9.9 yields that the class $\mathcal{F} = \mathcal{G} + \langle x, x \rangle$ is also VC with index $d + 3$. Now the desired conclusion follows from part (iii) of lemma 9.9 combined with part (i) of lemma 9.7. \square

9.1.2 BUEI-Classes

Recall for a class of measurable functions \mathcal{F} , with envelope F , the uniform entropy integral

$$J(\delta, \mathcal{F}, L_2) \equiv \int_0^\delta \sqrt{\sup_Q \log N(\epsilon \|F\|_{Q,2}, \mathcal{F}, L_2(Q))} d\epsilon,$$

where the supremum is taken over all finitely discrete probability measures Q with $\|F\|_{Q,2} > 0$. Note the dependence on choice of envelope F . This is crucial since there are many random functions which can serve as an envelope. For example, if F is an envelope, then so is $F + 1$ and $2F$. One must allow that different envelopes may be needed in different settings. We say that the class \mathcal{F} has *bounded uniform entropy integral* (BUEI) with envelope F —or is *BUEI* with envelope F —if $J(1, \mathcal{F}, L_2) < \infty$ for that particular choice of envelope.

Theorem 9.3 tells us that a VC-class \mathcal{F} is automatically BUEI with any envelope. We leave it as an exercise to show that if \mathcal{F} and \mathcal{G} are BUEI with respective envelopes F and G , then $\mathcal{F} \sqcup \mathcal{G}$ is BUEI with envelope $F \vee G$. The following lemma, which is closely related to an important Donsker

preservation theorem in section 9.4 below, is also useful for building BUEI classes from other BUEI classes:

LEMMA 9.13 *Let $\mathcal{F}_1, \dots, \mathcal{F}_k$ be BUEI classes with respective envelopes F_1, \dots, F_k , and let $\phi : \mathbb{R}^k \mapsto \mathbb{R}$ satisfy*

$$(9.2) \quad |\phi \circ f(x) - \phi \circ g(x)|^2 \leq c^2 \sum_{j=1}^k (f_j(x) - g_j(x))^2,$$

for every $f, g \in \mathcal{F}_1 \times \dots \times \mathcal{F}_k$ and x for a constant $0 < c < \infty$. Then the class $\phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k)$ is BUEI with envelope $H \equiv |\phi(f_0)| + c \sum_{j=1}^k (|f_{0j}| + F_j)$, where $f_0 \equiv (f_{01}, \dots, f_{0k})$ is any function in $\mathcal{F}_1 \times \dots \times \mathcal{F}_k$, and where $\phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k)$ is as defined in lemma 8.10.

Proof. Fix $\epsilon > 0$ and a finitely discrete probability measure Q , and let $f, g \in \mathcal{F}_1 \times \dots \times \mathcal{F}_k$ satisfy $\|f_j - g_j\|_{Q,2} \leq \epsilon \|F_j\|_{Q,2}$ for $1 \leq j \leq k$. Now (9.2) implies that

$$\begin{aligned} \|\phi \circ f - \phi \circ g\|_{Q,2} &\leq c \sqrt{\sum_{j=1}^k \|f_j - g_j\|_{Q,2}^2} \\ &\leq \epsilon c \sum_{j=1}^k \|F_j\|_{Q,2} \\ &\leq \epsilon H. \end{aligned}$$

Hence

$$N(\epsilon H, \phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k), L_2(Q)) \leq \prod_{j=1}^k N(\epsilon \|F_j\|_{Q,2}, \mathcal{F}_j, L_2(Q)),$$

and the desired result follows since ϵ and Q were arbitrary. \square

Some useful consequences of lemma 9.13 are given in the following lemma:

LEMMA 9.14 *Let \mathcal{F} and \mathcal{G} be BUEI with respective envelopes F and G , and let $\phi : \mathbb{R} \mapsto \mathbb{R}$ be a Lipschitz continuous function with Lipschitz constant $0 < c < \infty$. Then*

- (i) $\mathcal{F} \wedge \mathcal{G}$ is BUEI with envelope $F + G$;
- (ii) $\mathcal{F} \vee \mathcal{G}$ is BUEI with envelope $F + G$;
- (iii) $\mathcal{F} + \mathcal{G}$ is BUEI with envelope $F + G$;
- (iv) $\phi(\mathcal{F})$ is BUEI with envelope $|\phi(f_0)| + c(|f_0| + F)$, provided $f_0 \in \mathcal{F}$.

The proof, which we omit, is straightforward.

As mentioned earlier, lemma 9.13 is very similar to a Donsker preservation result we will present later in this chapter. In fact, most of the BUEI

preservation results we give in this section have parallel Donsker preservation properties. An important exception, and one which is perhaps the primary justification for the use of BUEI preservation techniques, applies to products of Donsker classes. As verified in the following theorem, the product of two BUEI classes is BUEI, whether or not the two classes involved are bounded (compare with corollary 9.15 below):

THEOREM 9.15 *Let \mathcal{F} and \mathcal{G} be BUEI classes with respective envelopes F and G . Then $\mathcal{F} \cdot \mathcal{G} \equiv \{fg : f \in \mathcal{F}, g \in \mathcal{G}\}$ is BUEI with envelope FG .*

Proof. Fix $\epsilon > 0$ and a finitely discrete probability measure \tilde{Q} with $\|FG\|_{\tilde{Q},2} > 0$, and let $dQ^* \equiv G^2 d\tilde{Q} / \|G\|_{\tilde{Q},2}^2$. Clearly, Q^* is a finitely discrete probability measure with $\|F\|_{Q^*,2} > 0$. Let $f_1, f_2 \in \mathcal{F}$ satisfy $\|f_1 - f_2\|_{Q^*,2} \leq \epsilon \|F\|_{Q^*,2}$. Then

$$\epsilon \geq \frac{\|f_1 - f_2\|_{Q^*,2}}{\|F\|_{Q^*,2}} = \frac{\|(f_1 - f_2)G\|_{\tilde{Q},2}}{\|FG\|_{\tilde{Q},2}},$$

and thus, if we let $\mathcal{F} \cdot G \equiv \{fG : f \in \mathcal{F}\}$,

$$(9.3) \quad \begin{aligned} N(\epsilon \|FG\|_{\tilde{Q},2}, \mathcal{F} \cdot G, L_2(\tilde{Q})) &\leq N(\epsilon \|F\|_{Q^*,2}, \mathcal{F}, L_2(Q^*)) \\ &\leq \sup_Q N(\epsilon \|F\|_{Q,2}, \mathcal{F}, L_2(Q)), \end{aligned}$$

where the supremum is taken over all finitely discrete probability measures Q for which $\|F\|_{Q,2} > 0$. Since the right-hand-side of (9.3) does not depend on \tilde{Q} , and since \tilde{Q} satisfies $\|FG\|_{\tilde{Q},2} > 0$ but is otherwise arbitrary, we have that

$$\sup_Q N(\epsilon \|FG\|_{Q,2}, \mathcal{F} \cdot G, L_2(Q)) \leq \sup_Q N(\epsilon \|F\|_{Q,2}, \mathcal{F}, L_2(Q)),$$

where the supremums are taken over all finitely discrete probability measures Q but with the left side taken over the subset for which $\|FG\|_{Q,2} > 0$ while the right side is taken over the subset for which $\|F\|_{Q,2} > 0$.

We can similarly show that the uniform entropy numbers for the class $\mathcal{G} \cdot F$ with envelope FG is bounded by the uniform entropy numbers for \mathcal{G} with envelope G . Since $|f_1 g_1 - f_2 g_2| \leq |f_1 - f_2|G + |g_1 - g_2|F$ for all $f_1, f_2 \in \mathcal{F}$ and $g_1, g_2 \in \mathcal{G}$, the forgoing results imply that

$$\begin{aligned} \sup_Q N(\epsilon \|FG\|_{Q,2}, \mathcal{F} \cdot \mathcal{G}, L_2(Q)) &\leq \sup_Q N(\epsilon \|F\|_{Q,2}/2, \mathcal{F}, L_2(Q)) \\ &\quad \times \sup_Q N(\epsilon \|G\|_{Q,2}/2, \mathcal{G}, L_2(Q)), \end{aligned}$$

where the supremums are all taken over the appropriate subsets of all finitely discrete probability measures. After taking logs, square roots, and then integrating both sides with respect to ϵ , the desired conclusion follows. \square

In order for BUEI results to be useful for obtaining Donsker results, it is necessary that sufficient measurability be established so that theorem 8.19 can be used. As shown in proposition 8.11 and the comments following theorem 8.19, pointwise measurability (PM) is sufficient measurability for this purpose. Since there are significant similarities between PM preservation and BUEI preservation results, one can construct useful joint PM and BUEI preservation results. Here is one such result:

LEMMA 9.16 *Let the classes $\mathcal{F}_1, \dots, \mathcal{F}_k$ be both BUEI and PM with respective envelopes F_1, \dots, F_k , and let $\phi : \mathbb{R}^k \mapsto \mathbb{R}$ satisfy (9.2) for every $f, g \in \mathcal{F}_1 \times \dots \times \mathcal{F}_k$ and x for a constant $0 < c < \infty$. Then the class $\phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k)$ is both BUEI and PM with envelope $H \equiv |\phi(f_0)| + c \sum_{j=1}^k (|f_{0j}| + F_j)$, where f_0 is any function in $\mathcal{F}_1 \times \dots \times \mathcal{F}_k$.*

Proof. Since a function satisfying (9.2) as specified is also continuous, the desired result is a direct consequence of lemmas 8.10 and 9.13. \square

The following lemma contains some additional joint preservation results:

LEMMA 9.17 *Let the classes \mathcal{F} and \mathcal{G} be both BUEI and PM with respective envelopes F and G , and let $\phi : \mathbb{R} \mapsto \mathbb{R}$ be a Lipschitz continuous function with Lipschitz constant $0 < c < \infty$. Then*

- (i) $\mathcal{F} \cup \mathcal{G}$ is both BUEI and PM with envelope $F \vee G$;
- (ii) $\mathcal{F} \wedge \mathcal{G}$ is both BUEI and PM with envelope $F + G$;
- (iii) $\mathcal{F} \vee \mathcal{G}$ is both BUEI and PM with envelope $F + G$;
- (iv) $\mathcal{F} + \mathcal{G}$ is both BUEI and PM with envelope $F + G$;
- (v) $\mathcal{F} \cdot \mathcal{G}$ is both BUEI and PM with envelope FG ;
- (vi) $\phi(\mathcal{F})$ is both BUEI and PM with envelope $|\phi(f_0)| + c(|f_0| + F)$, where $f_0 \in \mathcal{F}$.

Proof. Verifying (i) is straightforward. Results (ii), (iii), (iv) and (vi) are consequences of lemma 9.16. Result (v) is a consequence of lemma 8.10 and theorem 9.15. \square

If a class of measurable functions \mathcal{F} is both BUEI and PM with envelope F , then theorem 8.19 implies that \mathcal{F} is P -Donsker whenever $P^*F^2 < \infty$. Note that we have somehow avoided discussing preservation for subsets of classes. This is because it is unclear whether a subset of a PM class \mathcal{F} is itself a PM class. The difficulty is that while \mathcal{F} may have a countable dense subset \mathcal{G} (dense in terms of pointwise convergence), it is unclear whether any arbitrary subset $\mathcal{H} \subset \mathcal{F}$ also has a suitable countable dense subset. An easy way around this problem is to use various preservation results to establish that \mathcal{F} is P -Donsker, and then it follows directly that any $\mathcal{H} \subset \mathcal{F}$ is also P -Donsker by the definition of weak convergence. We will explore several additional preservation results as well as several practical examples later in this chapter and in the case studies of chapter 15.

9.2 Bracketing Entropy

We now present several useful bracketing entropy results for certain function classes as well as a few preservation results. We first mention that bracketing numbers are in general larger than covering numbers, as verified in the following lemma:

LEMMA 9.18 *Let \mathcal{F} be any class of real function on \mathcal{X} and $\|\cdot\|$ any norm on \mathcal{F} . Then*

$$N(\epsilon, \mathcal{F}, \|\cdot\|) \leq N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|)$$

for all $\epsilon > 0$.

Proof. Fix $\epsilon > 0$, and let \mathcal{B} be collection of ϵ -brackets that covers \mathcal{F} . From each bracket $B \in \mathcal{B}$, take a function $g(B) \in B \cap \mathcal{F}$ to form a finite collection of functions $\mathcal{G} \subset \mathcal{F}$ of the same cardinality as \mathcal{B} consisting of one function from each bracket in \mathcal{B} . Now every $f \in \mathcal{F}$ lies in a bracket $B \in \mathcal{B}$ such that $\|f - g(B)\| \leq \epsilon$ by the definition of an ϵ -bracket. Thus \mathcal{G} is an ϵ cover of \mathcal{F} of the same cardinality as \mathcal{B} . The desired conclusion now follows. \square

The first substantive bracketing entropy result we present considers classes of smooth functions on a bounded set $\mathcal{X} \subset \mathbb{R}^d$. For any vector $K = (k_1, \dots, k_d)$ of nonnegative integers define the differential operator $D^K \equiv \partial k_{\bullet} / (\partial x_1^{k_1}, \dots, \partial x_d^{k_d})$, where $k_{\bullet} \equiv k_1 + \dots + k_d$. As defined previously, let $\lfloor x \rfloor$ be the largest integer $j \leq x$, for any $x \in \mathbb{R}$. For any function $f : \mathcal{X} \mapsto \mathbb{R}$ and $\alpha > 0$, define the norm

$$\|f\|_{\alpha} \equiv \max_{k_{\bullet} \leq \lfloor \alpha \rfloor} \sup_x |D^k f(x)| + \max_{k: k_{\bullet} = \lfloor \alpha \rfloor} \sup_{x, y} \frac{|D^k f(x) - D^k f(y)|}{\|x - y\|^{\alpha - \lfloor \alpha \rfloor}},$$

where the suprema are taken over $x \neq y$ in the interior of \mathcal{X} . Now let $C_M^{\alpha}(\mathcal{X})$ be the set of all continuous functions $f : \mathcal{X} \mapsto \mathbb{R}$ with $\|f\|_{\alpha} \leq M$. Recall that for a set A in a metric space, $\text{diam } A = \sup_{x, y \in A} d(x, y)$. We have the following theorem:

THEOREM 9.19 *Let $\mathcal{X} \subset \mathbb{R}^d$ be bounded and convex with nonempty interior. There exists a constant $K < \infty$ depending only on α , $\text{diam } \mathcal{X}$, and d such that*

$$\log N_{[]}(\epsilon, C_1^{\alpha}(\mathcal{X}), L_r(Q)) \leq K \left(\frac{1}{\epsilon} \right)^{d/\alpha},$$

for every $r \geq 1$, $\epsilon > 0$, and any probability measure Q on \mathbb{R}^d .

This is corollary 2.7.2 of VW, and we omit the proof.

We now consider several results for Lipschitz and Sobolev function classes. We first present the results for covering numbers based on the uniform norm and then present the relationship to bracketing entropy.

THEOREM 9.20 *For a compact, convex subset $C \subset \mathbb{R}^d$, let \mathcal{F} be the class of all convex functions $f : C \mapsto [0, 1]$ with $|f(x) - f(y)| \leq L\|x - y\|$ for every x, y . For some integer $m \geq 1$, let \mathcal{G} be the class of all functions $g : [0, 1] \mapsto [0, 1]$ with $\int_0^1 [g^{(m)}(x)]^2 dx \leq 1$, where superscript (m) denotes the m 'th derivative. Then*

$$\begin{aligned} \log N(\epsilon, \mathcal{F}, \|\cdot\|_\infty) &\leq K(1+L)^{d/2} \left(\frac{1}{\epsilon}\right)^{d/2}, \text{ and} \\ \log N(\epsilon, \mathcal{G}, \|\cdot\|_\infty) &\leq M \left(\frac{1}{\epsilon}\right)^{1/m}, \end{aligned}$$

where $\|\cdot\|_\infty$ is the uniform norm and the constant $K < \infty$ depends only on d and C and the constant M depends only on m .

The first displayed result is corollary 2.7.10 of VW, while the second displayed result is theorem 2.4 of van de Geer (2000). We omit the proofs.

The following lemma shows how theorem 9.20 applies to bracketing entropy:

LEMMA 9.21 *For any norm $\|\cdot\|$ dominated by $\|\cdot\|_\infty$ and any class of functions \mathcal{F} ,*

$$\log N_{[]} (2\epsilon, \mathcal{F}, \|\cdot\|) \leq \log N(\epsilon, \mathcal{F}, \|\cdot\|_\infty),$$

for all $\epsilon > 0$.

Proof. Let f_1, \dots, f_m be a uniform ϵ -cover of \mathcal{F} . Since the 2ϵ -brackets $[f_i - \epsilon, f_i + \epsilon]$ now cover \mathcal{F} , the result follows. \square

We now present a second Lipschitz continuity result which is in fact a generalization of lemma 9.21. The result applies to function classes of the form $\mathcal{F} = \{f_t : t \in T\}$, where

$$(9.4) \quad |f_s(x) - f_t(x)| \leq d(s, t)F(x)$$

for some metric d on T , some real function F on the sample space \mathcal{X} , and for all $x \in \mathcal{X}$. This special Lipschitz structure arises in a number of settings, including parametric Z- and M- estimation. For example, consider the least absolute deviation regression setting of section 2.2.6, under the assumption that the random covariate U and regression parameter θ are constrained to known compact subsets $\mathcal{U}, \Theta \subset \mathbb{R}^p$. Recall that, in this setting, the outcome given U is modeled as $Y = \theta'U + e$, where the residual error e has median zero. Estimation of the true parameter value θ_0 is accomplished by minimizing $\theta \mapsto \mathbb{P}_n m_\theta$, where $m_\theta(X) \equiv |e - (\theta - \theta_0)'U| - |e|$, $X \equiv (Y, U)$ and $e = Y - \theta_0'U$. From (2.20) in section 2.2.6, we know that the class $\mathcal{F} = \{m_\theta : \theta \in \Theta\}$ satisfies (9.4) with $T = \Theta$, $d(s, t) = \|s - t\|$ and $F(x) = \|u\|$, where $x = (y, u)$ is a realization of X .

The following theorem shows that the bracketing numbers for a general \mathcal{F} satisfying (9.4) are bounded by the covering numbers for the associated index set T .

THEOREM 9.22 *Suppose the class of functions $\mathcal{F} = \{f_t : t \in T\}$ satisfies (9.4) for every $s, t \in T$ and some fixed function F . Then, for any norm $\|\cdot\|$,*

$$N_{[]} (2\epsilon\|F\|, \mathcal{F}, \|\cdot\|) \leq N(\epsilon, T, d).$$

Proof. Note that for any ϵ -net t_1, \dots, t_k that covers T with respect to d , the brackets $[f_{t_j} - \epsilon F, f_{t_j} + \epsilon F]$ cover \mathcal{F} . Since these brackets are all of size $2\epsilon\|F\|$, the proof is complete. \square

Note that when $\|\cdot\|$ is any norm dominated by $\|\cdot\|_\infty$, theorem 9.22 simplifies to lemma 9.21 when $T = \mathcal{F}$ and $d = \|\cdot\|_\infty$ (and thus automatically $F = 1$).

We move now from continuous functions to monotone functions. As was done in lemma 9.11 above for uniform entropy, we can study bracketing entropy of the class of all monotone functions mapping into $[0, 1]$:

THEOREM 9.23 *For each integer $r \geq 1$, there exists a constant $K < \infty$ such that the class \mathcal{F} of monotone functions $f : \mathbb{R} \mapsto [0, 1]$ satisfies*

$$\log N_{[]}(\epsilon, \mathcal{F}, L_r(Q)) \leq \frac{K}{\epsilon},$$

for all $\epsilon > 0$ and every probability measure Q .

The lengthy proof, which we omit, is given in chapter 2.7 of VW.

We now briefly discuss preservation results. Unfortunately, it appears that there are not as many useful preservation results for bracketing entropy as there are for uniform entropy, but the following lemma contains two such results which are easily verified:

LEMMA 9.24 *Let \mathcal{F} and \mathcal{G} be classes of measurable function. Then for any probability measure Q and any $1 \leq r \leq \infty$,*

$$(i) \quad N_{[]} (2\epsilon, \mathcal{F} + \mathcal{G}, L_r(Q)) \leq N_{[]}(\epsilon, \mathcal{F}, L_r(Q)) N_{[]}(\epsilon, \mathcal{G}, L_r(Q));$$

(ii) *Provided \mathcal{F} and \mathcal{G} are bounded by 1,*

$$N_{[]} (2\epsilon, \mathcal{F} \cdot \mathcal{G}, L_r(Q)) \leq N_{[]}(\epsilon, \mathcal{F}, L_r(Q)) N_{[]}(\epsilon, \mathcal{G}, L_r(Q)).$$

The straightforward proof is saved as an exercise.

9.3 Glivenko-Cantelli Preservation

In this section, we discuss methods which are useful for building up Glivenko-Cantelli (G-C) classes from other G-C classes. Such results can be useful for establishing consistency for Z- and M- estimators and their bootstrapped versions. It is clear from the definition of P-G-C classes, that if \mathcal{F} and \mathcal{G} are P-G-C, then $\mathcal{F} \cup \mathcal{G}$ and any subset thereof is also P-G-C. The purpose

of the remainder of this section is to discuss more substantive preservation results. The main tool for this is the following theorem, which is a minor modification of theorem 3 of van der Vaart and Wellner (2000) and which we give without proof:

THEOREM 9.25 *Suppose that $\mathcal{F}_1, \dots, \mathcal{F}_k$ are strong P - G - C classes of functions with $\max_{1 \leq j \leq k} \|P\|_{\mathcal{F}_j} < \infty$, and that $\phi : \mathbb{R}^k \mapsto \mathbb{R}$ is continuous. Then the class $\mathcal{H} \equiv \phi(\mathcal{F}_1, \dots, \mathcal{F}_k)$ is strong P - G - C provided it has an integrable envelope.*

The following corollary lists some obvious consequences of this theorem:

COROLLARY 9.26 *Let \mathcal{F} and \mathcal{G} be P - G - C classes with respective integrable envelopes F and G . Then the following are true:*

- (i) $\mathcal{F} + \mathcal{G}$ is P - G - C .
- (ii) $\mathcal{F} \cdot \mathcal{G}$ is P - G - C provided $P[FG] < \infty$.
- (iii) Let R be the union of the ranges of functions in \mathcal{F} , and let $\psi : \bar{R} \mapsto \mathbb{R}$ be continuous. Then $\psi(\mathcal{F})$ is P - G - C provided it has an integrable envelope.

Proof. The statement (i) is obvious. Since $(x, y) \mapsto xy$ is continuous in \mathbb{R}^2 , statement (ii) follows from theorem 9.25. Statement (iii) also follows from the theorem since ψ has a continuous extension to \mathbb{R} , $\tilde{\psi}$, such that $\|P\tilde{\psi}(f)\|_{\mathcal{F}} = \|P\psi(f)\|_{\mathcal{F}}$. \square

It is interesting to note that the “preservation of products” result in part (ii) of the above corollary does not hold in general for Donsker classes (although, as was shown in section 9.1.2, it does hold for BUEI classes). This preservation result for G - C classes can be useful in formulating master theorems for bootstrapped Z - and M - estimators. Consider, for example, verifying the validity of the bootstrap for a parametric Z -estimator $\hat{\theta}_n$ which is a zero of $\theta \mapsto \mathbb{P}_n \psi_\theta$, for $\theta \in \Theta$, where ψ_θ is a suitable random function. Let $\Psi(\theta) = P\psi_\theta$, where we assume that for any sequence $\{\theta_n\} \in \Theta$, $\Psi(\theta_n) \rightarrow 0$ implies $\theta_n \rightarrow \theta_0 \in \Theta$ (ie., the parameter is identifiable). Usually, to obtain consistency, it is reasonable to assume that the class $\{\psi_\theta, \theta \in \Theta\}$ is P - G - C . Clearly, this condition is sufficient to ensure that $\hat{\theta}_n \xrightarrow{\text{as}^*} \theta_0$.

Now, under a few additional assumptions, the Z -estimator master theorem, theorem 2.11 can be applied, to obtain asymptotic normality of $\sqrt{n}(\hat{\theta}_n - \theta_0)$. In section 2.2.5, we made the claim that if Ψ is appropriately differentiable and the parameter is identifiable (as defined in the previous paragraph), sufficient additional conditions for this asymptotic normality to hold and for the bootstrap to be valid are that the $\{\psi_\theta : \theta \in \Theta\}$ is strong P - G - C with $\sup_{\theta \in \Theta} P|\psi_\theta| < \infty$, that $\{\psi_\theta : \theta \in \Theta, \|\theta - \theta_0\| \leq \delta\}$ is P -Donsker for some $\delta > 0$, and that $P\|\psi_\theta - \psi_{\theta_0}\|^2 \rightarrow 0$ as $\theta \rightarrow \theta_0$. As we will see in chapter 13, where we present the arguments for this result in detail, an important step in the proof of bootstrap validity is to show that

the bootstrap estimate $\hat{\theta}_n^\circ$ is unconditionally consistent for θ_0 . If we use a weighted bootstrap with i.i.d. non-negative weights ξ_1, \dots, ξ_n , which are independent of the data and which satisfy $E\xi_1 = 1$, then result (ii) from the above corollary tells us that $\mathcal{F} \equiv \{\xi\psi_\theta : \theta \in \Theta\}$ is P -G-C. This follows since both classes of functions $\{\xi\}$ (a trivial class with one member) and $\{\psi_\theta : \theta \in \Theta\}$ are P -G-C and since the product class \mathcal{F} has an integral envelope by lemma 8.13. Note here that we are tacetly augmenting P to be the product probability measure of both the data and the independent bootstrap weights. We will expand on these ideas in section 10.3 of the next chapter for the special case where $\Theta \subset \mathbb{R}^p$ and in chapter 13 for the more general case.

Another result which can be useful for inference is the following lemma on covariance estimation. We mentioned this result in the first paragraph of section 2.2.3 in the context of conducting uniform inference for Pf as f ranges over a class of functions \mathcal{F} . The lemma answers the question of when the limiting covariance of \mathbb{G}_n , indexed by \mathcal{F} , can be consistently estimated. Recall that this covariance is $\sigma(f, g) \equiv Pfg - PfPg$, and its estimator is $\hat{\sigma}(f, g) \equiv \mathbb{P}_nfg - \mathbb{P}_nf\mathbb{P}_ng$. Although knowledge of this covariance matrix is usually not sufficient in itself to obtain inference on $\{Pf : f \in \mathcal{F}\}$, it still provides useful information.

LEMMA 9.27 *Let \mathcal{F} be Donsker. Then $\|\hat{\sigma}(f, g) - \sigma(f, g)\|_{\mathcal{F}, \mathcal{F}} \xrightarrow{\text{as}^*} 0$ if and only if $P^*\|f - Pf\|_{\mathcal{F}}^2 < \infty$.*

Proof. Note that since \mathcal{F} is Donsker, \mathcal{F} is also G-C. Hence $\dot{\mathcal{F}} \equiv \{\dot{f} : f \in \mathcal{F}\}$ is G-C, where for any $f \in \mathcal{F}$, $\dot{f} = f - Pf$. Now we first assume that $P^*\|f - Pf\|_{\mathcal{F}}^2 < \infty$. By theorem 9.25, $\dot{\mathcal{F}} \cdot \dot{\mathcal{F}}$ is also G-C. Uniform consistency of $\hat{\sigma}$ now follows since, for any $f, g \in \mathcal{F}$, $\hat{\sigma}(f, g) - \sigma(f, g) = (\mathbb{P}_n - P)\dot{f}\dot{g} - \mathbb{P}_n\dot{f}\mathbb{P}_n\dot{g}$. Assume next that $\|\hat{\sigma}(f, g) - \sigma(f, g)\|_{\mathcal{F}, \mathcal{F}} \xrightarrow{\text{as}^*} 0$. This implies that $\dot{\mathcal{F}} \cdot \dot{\mathcal{F}}$ is G-C. Now lemma 8.13 implies that $P^*\|f - Pf\|_{\mathcal{F}}^2 = P^*\|\dot{f}\dot{g}\|_{\dot{\mathcal{F}}, \dot{\mathcal{F}}} < \infty$. \square

We close this section with the following theorem which provided several interesting necessary and sufficient conditions for \mathcal{F} to be strong G-C:

THEOREM 9.28 *Let \mathcal{F} be a class of measurable functions. Then the following are equivalent:*

- (i) \mathcal{F} is strong P -G-C;
- (ii) $E^*\|\mathbb{P}_n - P\|_{\mathcal{F}} \rightarrow 0$ and $E^*\|f - Pf\|_{\mathcal{F}} < \infty$;
- (iii) $\|\mathbb{P}_n - P\|_{\mathcal{F}} \xrightarrow{P} 0$ and $E^*\|f - Pf\|_{\mathcal{F}} < \infty$.

Proof. Since $\mathbb{P}_n - P$ does not change when the class \mathcal{F} is replaced by $\{f - Pf : f \in \mathcal{F}\}$, we will assume hereafter that $\|P\|_{\mathcal{F}} = 0$ without loss of generality.

(i) \Rightarrow (ii): That (i) implies $E^*\|f\|_{\mathcal{F}} < \infty$ follows from lemma 8.13. Fix $0 < M < \infty$, and note that

$$(9.5) \quad \begin{aligned} \mathbb{E}^* \|\mathbb{P}_n - P\|_{\mathcal{F}} &\leq \mathbb{E}^* \|(\mathbb{P}_n - P)f \times 1\{F \leq M\}\|_{\mathcal{F}} \\ &\quad + 2\mathbb{E}^* [F \times 1\{F > M\}]. \end{aligned}$$

By assertion (ii) of corollary 9.26, $\mathcal{F} \cdot 1\{F \leq M\}$ is strong P -G-C, and thus the first term on the right of (9.5) $\rightarrow 0$ by the bounded convergence theorem. Since the second term on the right of (9.5) can be made arbitrarily small by increasing M , the left side of (9.5) $\rightarrow 0$, and the desired result follows.

(ii) \Rightarrow (iii): This is obvious.

(iii) \Rightarrow (i): By the assumed integrability of the envelope F , lemma 8.16 can be employed to verify that there is a version of $\|\mathbb{P}_n - P\|_{\mathcal{F}}^*$ that converges outer almost surely to a constant. Condition (iii) implies that this constant must be zero. \square

9.4 Donsker Preservation

In this section, we describe several techniques for building Donsker from other Donsker classes. The first theorem, theorem 9.29, gives results for subsets, pointwise closures and symmetric convex hulls of Donsker classes. The second theorem, theorem 9.30, presents a very powerful result for Lipschitz functions of Donsker classes. The corollary that follows present consequences of this theorem that are quite useful in statistical applications.

For a class \mathcal{F} of real-valued, measurable functions on the sample space \mathcal{X} , let $\overline{\mathcal{F}}^{(P,2)}$ be the set of all $f : \mathcal{X} \mapsto \mathbb{R}$ for which there exists a sequence $\{f_m\} \in \mathcal{F}$ such that $f_m \rightarrow f$ both pointwise (ie., for every argument $x \in \mathcal{X}$) and in $L_2(P)$. Similarly, let $\overline{\text{sconv}}^{(P,2)} \mathcal{F}$ be the pointwise and $L_2(P)$ closure of $\text{sconv} \mathcal{F}$ defined in section 9.1.1.

THEOREM 9.29 *Let \mathcal{F} be a P -Donsker class. Then*

(i) *For any $\mathcal{G} \subset \mathcal{F}$, \mathcal{G} is P -Donsker.*

(ii) *$\overline{\mathcal{F}}^{(P,2)}$ is P -Donsker.*

(iii) *$\overline{\text{sconv}}^{(P,2)} \mathcal{F}$ is P -Donsker.*

Proof. The proof of (i) is obvious by the facts that weak convergence consists of marginal convergence plus asymptotic equicontinuity and that the maximum modulus of continuity does not increase when maximizing over a smaller set. For (ii), one can without loss of generality assume that both \mathcal{F} and $\overline{\mathcal{F}}^{(P,2)}$ are mean zero classes. For a class of measurable functions \mathcal{G} , denote the modulus of continuity

$$M_{\mathcal{G}}(\delta) \equiv \sup_{f, g \in \mathcal{G}: \|f - g\|_{P,2} < \delta} |\mathbb{G}_n(f - g)|.$$

Fix $\delta > 0$. We can choose $f, g \in \overline{\mathcal{F}}^{(P,2)}$ such that $|\mathbb{G}_n(f - g)|$ is arbitrarily close to $M_{\overline{\mathcal{F}}^{(P,2)}}(\delta)$ and $\|f - g\|_{P,2} < \delta$. We can also choose $f_*, g_* \in \mathcal{F}$ such that $\|f - f_*\|_{P,2}$ and $\|g - g_*\|_{P,2}$ are arbitrarily small (for fixed data). Thus $M_{\overline{\mathcal{F}}^{(P,2)}}(\delta) \leq M_{\mathcal{F}}(2\delta)$. Since $\delta > 0$ was arbitrary, we obtain that asymptotic equicontinuity in probability for $\overline{\mathcal{F}}^{(P,2)}$ follows from asymptotic equicontinuity in probability of $\{\mathbb{G}_n(f) : f \in \mathcal{F}\}$. Part (iii) is theorem 2.10.3 of VW, and we omit its proof. \square

The following theorem, theorem 2.10.6 of VW, is one of the most useful Donsker preservation results for statistical applications. We omit the proof.

THEOREM 9.30 *Let $\mathcal{F}_1, \dots, \mathcal{F}_k$ be Donsker classes with $\max_{1 \leq i \leq k} \|P\|_{\mathcal{F}_i} < \infty$. Let $\phi : \mathbb{R}^k \mapsto \mathbb{R}$ satisfy*

$$|\phi \circ f(x) - \phi \circ g(x)|^2 \leq c^2 \sum_{i=1}^k (f_i(x) - g_i(x))^2,$$

for every $f, g \in \mathcal{F}_1 \times \dots \times \mathcal{F}_k$ and $x \in \mathcal{X}$ and for some constant $c < \infty$. Then $\phi \circ (\mathcal{F}_1, \dots, \mathcal{F}_k)$ is Donsker provided $\phi \circ f$ is square integrable for at least one $f \in \mathcal{F}_1 \times \dots \times \mathcal{F}_k$.

The following corollary is a useful consequence of this theorem:

COROLLARY 9.31 *Let \mathcal{F} and \mathcal{G} be Donsker classes. Then:*

- (i) $\mathcal{F} \cup \mathcal{G}$ and $\mathcal{F} + \mathcal{G}$ are Donsker.
- (ii) If $\|P\|_{\mathcal{F} \cup \mathcal{G}} < \infty$, then the classes of pairwise infima, $\mathcal{F} \wedge \mathcal{G}$, and pairwise suprema, $\mathcal{F} \vee \mathcal{G}$, are both Donsker.
- (iii) If \mathcal{F} and \mathcal{G} are both uniformly bounded, $\mathcal{F} \cdot \mathcal{G}$ is Donsker.
- (iv) If $\psi : \overline{R} \mapsto \mathbb{R}$ is Lipschitz continuous, where R is the range of functions in \mathcal{F} , and $\|\psi(f)\|_{P,2} < \infty$ for at least one $f \in \mathcal{F}$, then $\psi(\mathcal{F})$ is Donsker.
- (v) If $\|P\|_{\mathcal{F}} < \infty$ and g is a uniformly bounded, measurable function, then $\mathcal{F} \cdot g$ is Donsker.

Proof. For any measurable function f , let $\dot{f} \equiv f - Pf$. Also define $\dot{\mathcal{F}} \equiv \{\dot{f} : f \in \mathcal{F}\}$ and $\dot{\mathcal{G}} \equiv \{\dot{g} : g \in \mathcal{G}\}$. Note that for any $f \in \mathcal{F}$ and $g \in \mathcal{G}$, $\mathbb{G}_n f = \mathbb{G}_n \dot{f}$, $\mathbb{G}_n g = \mathbb{G}_n \dot{g}$ and $\mathbb{G}_n(f + g) = \mathbb{G}_n(\dot{f} + \dot{g})$. Hence $\mathcal{F} \cup \mathcal{G}$ is Donsker if and only if $\dot{\mathcal{F}} \cup \dot{\mathcal{G}}$ is Donsker; and, similarly, $\mathcal{F} + \mathcal{G}$ is Donsker if and only if $\dot{\mathcal{F}} + \dot{\mathcal{G}}$ is Donsker. Clearly, $\|P\|_{\dot{\mathcal{F}} \cup \dot{\mathcal{G}}} = 0$. Hence Lipschitz continuity of the map $(x, y) \mapsto x + y$ on \mathbb{R}^2 yields that $\dot{\mathcal{F}} + \dot{\mathcal{G}}$ is Donsker, via theorem 9.30. Hence also $\mathcal{F} + \mathcal{G}$ is Donsker. Since $\dot{\mathcal{F}} \cup \dot{\mathcal{G}}$ is contained in the union of $\dot{\mathcal{F}} \cup \{0\}$ and $\dot{\mathcal{G}} \cup \{0\}$, $\dot{\mathcal{F}} \cup \dot{\mathcal{G}}$ is Donsker and hence so is $\mathcal{F} \cup \mathcal{G}$. Thus part (i) follows. \square

Proving parts (ii) and (iv) is saved as an exercise. Part (iii) follows since the map $(x, y) \mapsto xy$ is Lipschitz continuous on bounded subsets of \mathbb{R}^2 . For part (v), note that for any $f_1, f_2 \in \mathcal{F}$, $|f_1(x)g(x) - f_2(x)g(x)| \leq \|g\|_\infty |f_1(x) - f_2(x)|$. Hence $\phi \circ \{\mathcal{F}, \{g\}\}$, where $\phi(x, y) = xy$, is Lipschitz continuous in the required manner. \square

9.5 Proofs

Proof of theorem 9.3. Let \mathcal{C} denote the set of all subgraphs C_f of functions $f \in \mathcal{F}$. Note that for any probability measure Q on \mathcal{X} and any $f, g \in \mathcal{F}$,

$$\begin{aligned} Q|f - g| &= \int_{\mathcal{X}} \int_{\mathbb{R}} |1\{t < f(x)\} - 1\{t < g(x)\}| dt Q(dx) \\ &= (Q \times \lambda)(C_f \Delta C_g), \end{aligned}$$

where λ is Lebesgue measure, $A \Delta B \equiv A \cup B - A \cap B$ for any two sets A, B , and the second equality follows from Fubini's theorem. Construct a probability measure P on $\mathcal{X} \times \mathbb{R}$ by restricting $Q \times \lambda$ to the set $\{(x, t) : |t| \leq F(x)\}$ and letting $P = (Q \times \lambda) / (2\|F\|_{Q,1})$. Now $Q|f - g| = 2\|F\|_{Q,1} P|1\{C_f\} - 1\{C_g\}|$. Thus, by theorem 9.2 above,

$$\begin{aligned} (9.6) \quad N(\epsilon 2\|F\|_{Q,1}, \mathcal{F}, L_1(Q)) &= N(\epsilon, \mathcal{C}, L_1(P)) \\ &\leq KV(\mathcal{C})(4e)^{V(\mathcal{C})} \left(\frac{1}{\epsilon}\right)^{V(\mathcal{C})-1}. \end{aligned}$$

For $r > 1$, we have

$$Q|f - g|^r \leq Q\{|f - g|(2F)^{r-1}\} = 2^{r-1} R|f - g| QF^{r-1},$$

where R is the probability measure with density F^{r-1}/QF^{r-1} with respect to Q . Thus

$$\begin{aligned} \|f - g\|_{Q,r} &\leq 2^{1-1/r} \|f - g\|_{R,1}^{1/r} (QF^{r-1})^{1/r} \\ &= 2^{1-1/r} \|f - g\|_{R,1}^{1/r} \|F\|_{Q,r} \left(\frac{QF^{r-1}}{QF^r}\right)^{1/r}, \end{aligned}$$

which implies

$$\frac{\|f - g\|_{Q,r}}{2\|F\|_{Q,r}} \leq \left(\frac{\|f - g\|_{R,1}}{2\|F\|_{R,1}}\right)^{1/r}.$$

Hence $N(\epsilon 2\|F\|_{Q,r}, \mathcal{F}, L_r(Q)) \leq N(\epsilon^r 2\|F\|_{R,1}, \mathcal{F}, L_1(R))$. Since (9.6) applies equally well with Q replaced by R , we now have that

$$N(\epsilon 2 \|F\|_{Q,r}, \mathcal{F}, L_r(Q)) \leq KV(\mathcal{C})(4e)^{V(\mathcal{C})} \left(\frac{1}{\epsilon}\right)^{r(V(\mathcal{C})-1)}.$$

The desired result now follows by bringing the factor 2 in the left-hand-side over to the numerator of $1/\epsilon$ in the right-hand-side. \square

Proof of lemma 9.9. We leave it as an exercise to show that parts (i) and (ii) follow from parts (ii) and (iii) of lemma 9.7. To prove (iii), note first that the sets $\{x : f(x) > 0\}$ are (obviously) one-to-one images of the sets $\{(x, 0) : f(x) > 0\}$ which are the intersections of the open subgraphs $\{(x, t) : f(x) > t\}$ with the set $\mathcal{X} \times \{0\}$. Since a single set has VC index 1, the result now follows by applying (ii) and then (iv) of lemma 9.7. The subgraphs of $-\mathcal{F}$ are the images of the *open supergraphs* $\{(x, t) : t > f(x)\}$ under the one-to-one map $(x, t) \mapsto (x, -t)$. Since the open supergraphs are the complements of the *closed subgraphs* $\{(x, t) : t \geq f(x)\}$, they have the same VC-index as \mathcal{F} by lemma 9.32 below and by part (i) of lemma 9.7.

Part (v) follows from the fact that $\mathcal{F} + g$ shatters a given set of points $(x_1, t_1), \dots, (x_n, t_n)$ if and only if \mathcal{F} shatters $(x_1, t'_1), \dots, (x_n, t'_n)$, where $t'_i = t_i - g(x_i)$, $1 \leq i \leq n$. For part (vi), note that for any $f \in \mathcal{F}$ the subgraph of fg is the union of the sets $C_f^+ \equiv \{(x, t) : t < f(x)g(x), g(x) > 0\}$, $C_f^- \equiv \{(x, t) : t < f(x)g(x), g(x) < 0\}$ and $C_f^0 \equiv \{(x, t) : t < 0, g(x) = 0\}$. Define $\mathcal{C}^+ \equiv \{C_f^+ : f \in \mathcal{F}\}$, $\mathcal{C}^- \equiv \{C_f^- : f \in \mathcal{F}\}$ and $\mathcal{C}^0 \equiv \{C_f^0 : f \in \mathcal{F}\}$. By exercise 9.6.6 below, it suffices to show that these three classes are VC on the respective disjoint sets $\mathcal{X} \cap \{x : g(x) > 0\} \times \mathbb{R}$, $\mathcal{X} \cap \{x : g(x) < 0\} \times \mathbb{R}$ and $\mathcal{X} \cap \{x : g(x) = 0\} \times \mathbb{R}$, with respective VC indices bounded by $V_{\mathcal{F}}$, $V_{\mathcal{F}}$ and 1. Consider first \mathcal{C}^+ on $\mathcal{X} \cap \{x : g(x) < 0\}$. Note that the subset $(x_1, t_1), \dots, (x_m, t_m)$ is shattered by \mathcal{C}^+ if and only if the subset $(x_1, t_1/g(x_1)), \dots, (x_m, t_m/g(x_m))$ is shattered by the subgraphs of \mathcal{F} . Thus the VC-index of \mathcal{C}^+ on the relevant subset of $\mathcal{X} \times \mathbb{R}$ is $V_{\mathcal{F}}$. The same VC-index occurs for \mathcal{C}^- , but the VC-index for \mathcal{C}^0 is clearly 1. This concludes the proof of (vi).

For (vii), the result follows from part (vi) of lemma 9.7 since the subgraphs of the class $\mathcal{F} \circ \psi$ are the inverse images of the subgraphs of \mathcal{F} under the map $(z, t) \mapsto (\psi(z), t)$. For part (viii), suppose that the subgraphs of $\phi \circ \mathcal{F}$ shatter the set of points $(x_1, t_1), \dots, (x_n, t_n)$. Now choose f_1, \dots, f_m from \mathcal{F} so that the subgraphs of the functions $\phi \circ f_j$ pick out all $m = 2^n$ subsets. For each $1 \leq i \leq n$, define s_i to be the largest value of $f_j(x_i)$ over those $j \in \{1, \dots, m\}$ for which $\phi(f_j(x_i)) \leq t_i$. Now note that $f_j(x_i) \leq s_i$ if and only if $\phi(f_j(x_i)) \leq t_i$, for all $1 \leq i \leq n$ and $1 \leq j \leq m$. Hence the subgraphs of f_1, \dots, f_m shatter the points $(x_1, s_1), \dots, (x_n, s_n)$. \square

Proof of lemma 9.11. First consider the class $\mathcal{H}_{+,r}$ of monotone increasing, right-continuous functions $h : \mathbb{R} \mapsto [0, 1]$. For each $h \in \mathcal{H}_{+,r}$, define $h^{-1}(t) \equiv \inf\{x : h(x) \geq t\}$, and note that for any $x, t \in \mathbb{R}$, $h(x) \geq t$ if and only if $x \geq h^{-1}(t)$. Thus the class of indicator functions $\{1\{h(x) \geq t\} : h \in \mathcal{H}_{+,r}\} = \{1\{x \geq h^{-1}(t)\} : h \in \mathcal{H}_{+,r}\} \subset \{1\{x \geq t\} : t \in \mathbb{R}\}$. Since the last class of sets has VC index 2, the first class is also VC

with index 2. Since each function $h \in \mathcal{H}_{+,r}$ is the pointwise limit of the sequence

$$h_m = \sum_{j=1}^m \frac{1}{m} 1 \left\{ h(x) \geq \frac{j}{m} \right\},$$

we have that $\mathcal{H}_{+,r}$ is contained in the closed convex hull of a VC-subgraph class with VC index 2. Thus, by corollary 9.5, we have for all $0 < \epsilon < 1$,

$$\sup_Q \log N(\epsilon, \mathcal{H}_{+,r}, L_2(Q)) \leq \frac{K_0}{\epsilon},$$

where the supremum is taken over all probability measures Q and the constant K_0 is universal. Now consider the class $\mathcal{H}_{+,l}$ of monotone increasing, left-continuous functions, and define $\tilde{h}^{-1}(x) \equiv \sup\{t : h(t) \leq x\}$. Now note that for any $x, t \in \mathbb{R}$, $h(x) > t$ if and only if $x > \tilde{h}^{-1}(t)$. Arguing as before, we deduce that $\{1\{h(x) > t\} : h \in \mathcal{H}_{+,l}\}$ is a VC-class with index 2. Since each $h \in \mathcal{H}_{+,l}$ is the pointwise limit of the sequence

$$h_m = \sum_{j=1}^m \frac{1}{m} 1 \left\{ h(x) > \frac{j}{m} \right\},$$

we can again apply corollary 9.5 to arrive at the same uniform entropy bound we arrived at for $\mathcal{H}_{+,r}$.

Now let \mathcal{H}_+ be the class of all monotone increasing functions $h : \mathbb{R} \mapsto [0, 1]$, and note that each $h \in \mathcal{H}_+$ can be written as $h_r + h_l$, where $h_r \in \mathcal{H}_{+,r}$ and $h_l \in \mathcal{H}_{+,l}$. Hence for any probability measure Q and any $h^{(1)}, h^{(2)} \in \mathcal{H}_+$, $\|h^{(1)} - h^{(2)}\|_{Q,2} \leq \|h_r^{(1)} - h_r^{(2)}\|_{Q,r} + \|h_l^{(1)} - h_l^{(2)}\|_{Q,2}$, where $h_r^{(1)}, h_r^{(2)} \in \mathcal{H}_{+,r}$ and $h_l^{(1)}, h_l^{(2)} \in \mathcal{H}_{+,l}$. Thus $N(\epsilon, \mathcal{H}_+, L_2(Q)) \leq N(\epsilon/2, \mathcal{H}_{+,r}, L_2(Q)) \times N(\epsilon/2, \mathcal{H}_{+,l}, L_2(Q))$, and hence

$$\sup_Q \log N(\epsilon, \mathcal{H}_+, L_2(Q)) \leq \frac{K_1}{\epsilon},$$

where $K_1 = 4K_0$. Since any monotone decreasing function $g : \mathbb{R} \mapsto [0, 1]$ can be written as $1 - h$, where $h \in \mathcal{H}_+$, the uniform entropy numbers for the class of all monotone functions $f : \mathbb{R} \mapsto [0, 1]$, which we denote \mathcal{F} , is $\log(2)$ plus the uniform entropy numbers for \mathcal{H}_+ . Since $0 < \epsilon < 1$, we obtain the desired conclusion given in the statement of the lemma, with $K = \sqrt{2}K_1 = \sqrt{32}K_0$. \square

LEMMA 9.32 *Let \mathcal{F} be a set of measurable functions on \mathcal{X} . Then the closed subgraphs have the same VC-index as the open subgraphs.*

Proof. Suppose the closed subgraphs (the subgraphs of the form $\{(x, t) : t \leq f(x)\}$) shatter the set of points $(x_1, t_1), \dots, (x_n, t_n)$. Now select out of \mathcal{F} functions f_1, \dots, f_m whose closed subgraphs shatter all $m = 2^n$ subsets.

Let $\epsilon = (1/2) \inf\{t_i - f_j(x_i) : t_i - f_j(x_i) > 0\}$, and note that the open subgraphs (the subgraphs of the form $\{(x, t), t < f(x)\}$) of the f_1, \dots, f_m shatter the set of points $(x_1, t_1 - \epsilon), \dots, (x_n, t_n - \epsilon)$. This follows since, by construction, $t_i - \epsilon \geq f_j(x_i)$ if and only if $t_i > f_j(x_i)$, for all $1 \leq i \leq n$ and $1 \leq j \leq m$. Now suppose the open subgraphs shatter the set of points $(x_1, t_1), \dots, (x_n, t_n)$. Select out of \mathcal{F} functions f_1, \dots, f_m whose open subgraphs shatter all $m = 2^n$ subsets, and let $\epsilon = (1/2) \inf\{f_j(x_i) - t_i : f_j(x_i) - t_i > 0\}$. Note now that the closed subgraphs of f_1, \dots, f_m shatter the set of points $(x_1, t_1 + \epsilon), \dots, (x_n, t_n + \epsilon)$, since, by construction, $t_i < f_j(x_i)$ if and only if $t_i + \epsilon \leq f_j(x_i)$. Thus the VC-indices of open and closed subgraphs are the same. \square

9.6 Exercises

9.6.1. Show that $\text{sconv}\mathcal{F} \subset \text{conv}\mathcal{G}$, where $\mathcal{G} = \mathcal{F} \cup \{-\mathcal{F}\} \cup \{0\}$.

9.6.2. Show that the expression $N(\epsilon \|aF\|_{bQ, r}, a\mathcal{F}, L_r(bQ))$ does not depend on the constants $0 < a, b < \infty$, where $1 \leq r < \infty$.

9.6.3. Prove corollary 9.5.

9.6.4. In the proof of lemma 9.7, verify that part (iii) follows from parts (i) and (ii).

9.6.5. Show that parts (i) and (ii) of lemma 9.9 follow from parts (ii) and (iii) of lemma 9.7.

9.6.6. Let $\mathcal{X} = \cup_{i=1}^m \mathcal{X}_i$, where the \mathcal{X}_i are disjoint; and assume \mathcal{C}_i is a VC-class of subsets of \mathcal{X}_i , with VC-index V_i , $1 \leq i \leq m$. Show that $\sqcup_{i=1}^m \mathcal{C}_i$ is a VC-class in \mathcal{X} with VC-index $V_1 + \dots + V_m - m + 1$. Hint: Note that $\mathcal{C}_1 \cup \mathcal{C}_2$ and $\mathcal{X}_1 \cup \mathcal{C}_2$ are VC on $\mathcal{X}_1 \cup \mathcal{X}_2$ with respective indices V_1 and V_2 . Now use part (ii)—not part (iii)—of lemma 9.7 to show that $\mathcal{C}_1 \sqcup \mathcal{C}_2$ is VC on $\mathcal{X}_1 \cup \mathcal{X}_2$ with VC index $V_1 + V_2 - 1$.

9.6.7. Show that if \mathcal{F} and \mathcal{G} are BUEI with respective envelopes F and G , then $\mathcal{F} \sqcup \mathcal{G}$ is BUEI with envelope $F \vee G$.

9.6.8. In the context of the simple linear regression example of section 4.4.1, verify the following:

- (a) Show that both \mathcal{G}_1 and \mathcal{G}_2 are Donsker even though neither U nor e are bounded. Hint: Use BUEI preservation results.
- (b) Verify that both

$$\sup_{z \in [a+h, b-h]} \left| \int_{\mathbb{R}} h^{-1} L\left(\frac{z-u}{h}\right) H(du) - \dot{H}(z) \right| = O(h)$$

and

$$\begin{aligned} & \left(\sup_{z \in [a, a+h]} \left| \dot{H}(z) - \dot{H}(a+h) \right| \right) \vee \left(\sup_{z \in [b-h, b]} \left| \dot{H}(z) - \dot{H}(b-h) \right| \right) \\ & = O(h). \end{aligned}$$

(c) Show that \mathcal{F}_1 is Donsker and \mathcal{F}_2 is Glivenko-Cantelli.

9.6.9. Prove lemma 9.24.

9.6.10. Consider the class \mathcal{F} of all functions $f : [0, 1] \mapsto [0, 1]$ such that $|f(x) - f(y)| \leq |x - y|$. Show that a set of ϵ -brackets can be constructed to cover \mathcal{F} with cardinality bounded by $\exp(C/\epsilon)$ for some $0 < C < \infty$. Hint: Fix $\epsilon > 0$, and let n be the smallest integer $\geq 3/\epsilon$. For any $p = (k_0, \dots, k_n) \in P \equiv \{1, \dots, n\}^{n+1}$, define the path \bar{p} to be the collection of all function in \mathcal{F} such that $f \in \bar{p}$ only if $f(i/n) \in [(k_i - 1)/n, k_i/n]$ for all $i = 0 \dots n$. Show that for all $f \in \mathcal{F}$, if $f(i/n) \in [j/n, (j+1)/n]$, then

$$f \left[\frac{i+1}{n} \right] \in \left[\frac{(j-1) \vee 0}{n}, \frac{(j+2) \wedge n}{n} \right]$$

for $i, j = 0, \dots, n-1$. Show that this implies that the number of paths of the form \bar{p} for $p \in P$ needed to “capture” all elements of \mathcal{F} is bounded by $n3^n$. Now show that for each $p \in P$, there exists a pair of right-continuous “bracketing” functions $L_p, U_p : [0, 1] \mapsto [0, 1]$ such that $\forall x \in [0, 1]$, $L_p(x) < U_p(x)$, $U_p(x) - L_p(x) \leq 3/n \leq \epsilon$, and $L_p(x) \leq f(x) \leq U_p(x)$ for all $f \in \bar{p}$. Now complete the proof.

9.6.11. Show that if \mathcal{F} is Donsker with $\|P\|_{\mathcal{F}} < \infty$ and $f \geq \delta$ for all $f \in \mathcal{F}$ and some constant $\delta > 0$, then $1/\mathcal{F} \equiv \{1/f : f \in \mathcal{F}\}$ is Donsker.

9.6.12. Complete the proof of corollary 9.31:

1. Prove part (ii). Hint: show first that for any real numbers a_1, a_2, b_1, b_2 , $|a_1 \wedge b_1 - a_2 \wedge b_2| \leq |a_1 - a_2| + |b_1 - b_2|$.
2. Prove part (iv).

9.7 Notes

Theorem 9.3 is a minor modification of theorem 2.6.7 of VW. Corollary 9.5, lemma 9.6 and theorem 9.22 are corollary 2.6.12, lemma 2.6.15 and theorem 2.7.11, respectively, of VW. Lemmas 9.7 and 9.9 are modification of lemmas 2.6.17 and 2.6.18, respectively, of VW. Lemma 9.11 was suggested by example 2.6.21 of VW, and lemma 9.12 is a modification of exercise 2.6.14 of VW. Parts (i) and (ii) of theorem 9.29 are theorems 2.10.1

and 2.10.2, respectively, of VW. Corollary 9.31 includes some modifications of examples 2.10.7, 2.10.8 and 2.10.10 of VW. Lemma 9.32 was suggested by exercise 2.6.10 of VW. Exercise 9.6.10 is a modification of exercise 19.5 of van der Vaart (1998).

The bounded uniform entropy integral (BUEI) preservation techniques presented here grew out of the author's work on estimating equations for functional data described in Fine, Yan and Kosorok (2004).

10

Bootstrapping Empirical Processes

The purpose of this chapter is to obtain consistency results for bootstrapped empirical processes. These results can then be applied to many kinds of bootstrapped estimators since most estimators can be expressed as functionals of empirical processes. Much of the bootstrap results for such estimators will be deferred to later chapters where we discuss the functional delta method, Z-estimation and M-estimation. We do, however, present one specialized result for parametric Z-estimators in section 3 of this chapter as a practical illustration of bootstrap techniques.

We note that both conditional and unconditional bootstrap consistency results can be useful depending on the application. For the conditional bootstrap, the goal is to establish convergence of the conditional law given the data to an unconditional limit law. This convergence can be either in probability or outer almost sure. While the later convergence is certainly stronger, convergence in probability is usually sufficient for statistical applications.

The best choice of bootstrap weights for a given statistical application is also an important question, and the answer depends on the application. While the multinomial bootstrap is conceptually simple, its use in survival analysis applications may result in too much tied data. In the presence of censoring, it is even possible that a bootstrap sample could be drawn which consists of only censored observations. To avoid complications of this kind, it may be better to use the Bayesian bootstrap (Rubin, 1981). The weights for the Bayesian bootstrap are $\xi_1/\bar{\xi}, \dots, \xi_n/\bar{\xi}$, where ξ_1, \dots, ξ_n are i.i.d. standard exponential (mean and variance 1), independent of the data X_1, \dots, X_n , and where $\bar{\xi} \equiv n^{-1} \sum_{i=1}^n \xi_i$. Since these weights are strictly

positive, all observations are represented in each bootstrap realization, and the aforementioned problem with tied data won't happen unless the original data has ties. Both the multinomial and Bayesian bootstraps are included in the bootstrap weights we discuss in this chapter.

The multinomial weighted bootstrap is sometimes called the *nonparametric bootstrap* since it amounts to sampling from the empirical distribution which is a nonparametric estimate of the true distribution. In contrast, the *parametric bootstrap* is obtained by sampling from a parametric estimate $P_{\hat{\theta}_n}$ of the true distribution, where $\hat{\theta}_n$ is a consistent estimate of the true value of θ (see, for example, chapter 1 of Shao and Tu, 1995). A detailed discussion of the parametric bootstrap is beyond the scope of this chapter. Another kind of bootstrap is the exchangeable weighted bootstrap, which we only mention briefly in lemma 10.18 below. This lemma is needed for the proof of theorem 10.15.

We also note that the asymptotic results of this chapter are all first order, and in this situation the limiting results do not vary among those schemes that satisfy the stated conditions. A more refined analysis of differences between weighting schemes is beyond the scope of this chapter, but such differences may be important in small samples. A good reference for higher order properties of the bootstrap is Hall (1992).

The first section of this chapter considers unconditional and conditional convergence of bootstrapped empirical processes to limiting laws when the class of functions involved is Donsker. The main result of this section is a proof of theorems 2.6 and 2.7 of section 2.2.3. At the end of the section, we present several special continuous mapping results for bootstrapped processes. The second section considers parallel results when the function class involved is Glivenko-Cantelli. In this case, the limiting laws are degenerate, i.e., constant with probability 1. Such results are helpful for establishing consistency of bootstrapped estimators. The third section presents the simple Z-estimator illustration promised above. Throughout this chapter, we will sometimes for simplicity omit the subscript when referring to a representative of an i.i.d. sample. For example, we may use $E|\xi|$ to refer to $E|\xi_1|$, where ξ_1 is the first member of the sample ξ_1, \dots, ξ_n . The context will make the meaning clear.

10.1 The Bootstrap for Donsker Classes

The overall goal of this section is to prove the validity of the bootstrap central limit theorems given in theorems 2.6 and 2.7 of chapter 2. Both unconditional and conditional multiplier central limit theorems play a pivotal role in this development and will be presented first. At the end of the section, we also present several special continuous mapping results which apply to bootstrapped processes. These results allow the construction of

asymptotically uniformly valid confidence bands for $\{Pf : f \in \mathcal{F}\}$ when \mathcal{F} is Donsker.

10.1.1 An Unconditional Multiplier Central Limit Theorem

In this section, we present a multiplier central limit theorem which forms the basis for proving the unconditional central limit theorems of the next section. We also present an interesting corollary. For a real random variable ξ , recall from section 2.2.3 the quantity $\|\xi\|_{2,1} \equiv \int_0^\infty \sqrt{P(|\xi| > x)} dx$. Exercise 10.5.1 below verifies this is a norm which is slightly larger than $\|\cdot\|_2$. Also recall that δ_{X_i} is the probability measure that assigns a mass of 1 to X_i so that $\mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$ and $\mathbb{G}_n = n^{-1/2} \sum_{i=1}^n (\delta_{X_i} - P)$.

THEOREM 10.1 (Multiplier central limit theorem) *Let \mathcal{F} be a class of measurable functions, and let ξ_1, \dots, ξ_n be i.i.d. random variables with mean zero, variance 1, and with $\|\xi\|_{2,1} < \infty$, independent of the sample data X_1, \dots, X_n . Let $\mathbb{G}'_n \equiv n^{-1/2} \sum_{i=1}^n \xi_i (\delta_{X_i} - P)$ and $\mathbb{G}''_n \equiv n^{-1/2} \sum_{i=1}^n (\xi_i - \bar{\xi}) \delta_{X_i}$, where $\bar{\xi} \equiv n^{-1} \sum_{i=1}^n \xi_i$. Then the following are equivalent:*

- (i) \mathcal{F} is P -Donsker;
- (ii) \mathbb{G}'_n converges weakly to a tight process in $\ell^\infty(\mathcal{F})$;
- (iii) $\mathbb{G}'_n \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$;
- (iv) $\mathbb{G}''_n \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$.

Before giving the proof of this theorem, we will need the following tool. This lemma is lemma 2.9.1 of VW, and we give it without proof:

LEMMA 10.2 (Multiplier inequalities) *Let Z_1, \dots, Z_n be i.i.d. stochastic processes, with index \mathcal{F} such that $E^* \|Z\|_{\mathcal{F}} < \infty$, independent of the i.i.d. Rademacher variables $\epsilon_1, \dots, \epsilon_n$. Then for every i.i.d. sample ξ_1, \dots, ξ_n of real, mean-zero random variables independent of Z_1, \dots, Z_n , and any $1 \leq n_0 \leq n$,*

$$\begin{aligned} \frac{1}{2} \|\xi\|_1 E^* \left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i Z_i \right\|_{\mathcal{F}} &\leq E^* \left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i Z_i \right\|_{\mathcal{F}} \\ &\leq 2(n_0 - 1) E^* \|Z\|_{\mathcal{F}} E \max_{1 \leq i \leq n} \frac{|\xi_i|}{\sqrt{n}} \\ &\quad + 2\sqrt{2} \|\xi\|_{2,1} \max_{n_0 \leq k \leq n} E^* \left\| \frac{1}{\sqrt{k}} \sum_{i=n_0}^k \epsilon_i Z_i \right\|_{\mathcal{F}}. \end{aligned}$$

When the ξ_i are symmetrically distributed, the constants $1/2$, 2 and $2\sqrt{2}$ can all be replaced by 1.

Proof of theorem 10.1. Note that the processes $\mathbb{G}, \mathbb{G}_n, \mathbb{G}'_n$ and \mathbb{G}''_n do not change if they are indexed by $\dot{\mathcal{F}} \equiv \{f - Pf : f \in \mathcal{F}\}$ rather than \mathcal{F} . Thus we can assume throughout the proof that $\|P\|_{\mathcal{F}} = 0$ without loss of generality.

(i) \Leftrightarrow (ii): Convergence of the finite-dimensional marginal distributions of \mathbb{G}_n and \mathbb{G}'_n is equivalent to $\mathcal{F} \subset L_2(P)$, and thus it suffices to show that the asymptotic equicontinuity conditions of both processes are equivalent. By lemma 8.17, if \mathcal{F} is Donsker, then $P^*(F > x) = o(x^{-2})$ as $x \rightarrow \infty$. Similarly, if $\xi \cdot \mathcal{F}$ is Donsker, then $P^*(|\xi| \times F > x) = o(x^{-2})$ as $x \rightarrow \infty$. In both cases, $P^*F < \infty$. Since the variance of ξ is finite, we have by exercise 10.5.2 below that $E^* \max_{1 \leq i \leq n} |\xi_i|/\sqrt{n} \rightarrow 0$. Combining this with the multiplier inequality (lemma 10.2), we have

$$\begin{aligned} \frac{1}{2} \|\xi\|_1 \limsup_{n \rightarrow \infty} E^* \left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}_\delta} &\leq \limsup_{n \rightarrow \infty} E^* \left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i f(X_i) \right\|_{\mathcal{F}_\delta} \\ &\leq 2\sqrt{2} \|\xi\|_{2,1} \sup_{k \geq n_0} E^* \left\| \frac{1}{\sqrt{k}} \sum_{i=1}^k \epsilon_i f(X_i) \right\|_{\mathcal{F}_\delta}, \end{aligned}$$

for every $\delta > 0$ and $n_0 \leq n$. By the symmetrization theorem (theorem 8.8), we can remove the Rademacher variables $\epsilon_1, \dots, \epsilon_n$ at the cost of changing the constants. Hence, for any sequence $\delta_n \downarrow 0$, $E^* \|n^{-1/2} \sum_{i=1}^n (\delta_{X_i} - P)\|_{\mathcal{F}_{\delta_n}} \rightarrow 0$ if and only if $E^* \|n^{-1/2} \sum_{i=1}^n \xi_i (\delta_{X_i} - P)\|_{\mathcal{F}_{\delta_n}} \rightarrow 0$. By lemma 8.17, these mean versions of the asymptotic equicontinuity conditions imply the probability versions, and the desired results follow. We have actually proved that the first three assertions are equivalent.

(iii) \Rightarrow (iv): Note that by the equivalence of (i) and (iii), \mathcal{F} is Glivenko-Cantelli. Since $\mathbb{G}'_n - \mathbb{G}''_n = \sqrt{n} \bar{\xi} \mathbb{P}_n$, we now have that $\|\mathbb{G}'_n - \mathbb{G}''_n\|_{\mathcal{F}} \xrightarrow{P} 0$. Thus (iv) follows.

(iv) \Rightarrow (i): Let (Y_1, \dots, Y_n) be an independent copy of (X_1, \dots, X_n) , and let $(\tilde{\xi}_1, \dots, \tilde{\xi}_n)$ be an independent copy of (ξ_1, \dots, ξ_n) , so that $(\xi_1, \dots, \xi_n, \tilde{\xi}_1, \dots, \tilde{\xi}_n)$ is independent of $(X_1, \dots, X_n, Y_1, \dots, Y_n)$. Let $\bar{\xi}$ be the pooled mean of the ξ_i s and $\tilde{\xi}_i$ s; set

$$\mathbb{G}''_{2n} = (2n)^{-1/2} \left(\sum_{i=1}^n (\xi_i - \bar{\xi}) \delta_{X_i} + \sum_{i=1}^n (\tilde{\xi}_i - \bar{\xi}) \delta_{Y_i} \right)$$

and define

$$\tilde{\mathbb{G}}''_n \equiv (2n)^{-1/2} \left(\sum_{i=1}^n (\tilde{\xi}_i - \bar{\xi}) \delta_{X_i} + \sum_{i=1}^n (\xi_i - \bar{\xi}) \delta_{Y_i} \right).$$

We now have that both $\mathbb{G}''_{2n} \rightsquigarrow \mathbb{G}$ and $\tilde{\mathbb{G}}''_{2n} \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$.

Thus, by the definition of weak convergence, we have that (\mathcal{F}, ρ_P) is totally bounded and that for any sequence $\delta_n \downarrow 0$ both $\|\mathbb{G}''_{2n}\|_{\mathcal{F}_{\delta_n}} \xrightarrow{P} 0$ and

$\|\tilde{\mathbb{G}}''_{2n}\|_{\mathcal{F}_{\delta_n}} \xrightarrow{P} 0$. Hence also $\|\mathbb{G}''_{2n} - \tilde{\mathbb{G}}''_{2n}\|_{\mathcal{F}_{\delta_n}} \xrightarrow{P} 0$. However, since

$$\mathbb{G}''_{2n} - \tilde{\mathbb{G}}''_{2n} = n^{-1/2} \sum_{i=1}^n \frac{(\xi_i - \tilde{\xi}_i)}{\sqrt{2}} (f(X_i) - f(Y_i)),$$

and since the weights $\check{\xi}_i \equiv (\xi_i - \tilde{\xi}_i)/\sqrt{2}$ satisfy the moment conditions for the theorem we are proving, we now have the $\check{\mathbb{G}}_n \equiv n^{-1/2} \sum_{i=1}^n (f(X_i) - f(Y_i)) \rightsquigarrow \sqrt{2}\mathbb{G}$ in $\ell^\infty(\mathcal{F})$ by the already proved equivalence between (iii) and (i). Thus, for any sequence $\delta_n \downarrow 0$, $E^*\|\check{\mathbb{G}}_n\|_{\mathcal{F}_{\delta_n}} \rightarrow 0$. Since also

$$E_Y \left| \sum_{i=1}^n f(X_i) - f(Y_i) \right| \geq \left| \sum_{i=1}^n f(X_i) - E f(Y_i) \right| = \left| \sum_{i=1}^n f(X_i) \right|,$$

we can invoke Fubini's theorem (lemma 6.14) to yield

$$E^*\|\check{\mathbb{G}}_n\|_{\mathcal{F}_{\delta_n}} \geq E^*\|\mathbb{G}_n\|_{\mathcal{F}_{\delta_n}} \rightarrow 0.$$

Hence \mathcal{F} is Donsker. \square

We now present the following interesting corollary which shows the possibly unexpected result that the multiplier empirical process is asymptotically independent of the usual empirical process, even though the same data X_1, \dots, X_n are used in both processes:

COROLLARY 10.3 *Assume the conditions of theorem 10.1 hold and that \mathcal{F} is Donsker. Then $(\mathbb{G}_n, \mathbb{G}'_n, \mathbb{G}''_n) \rightsquigarrow (\mathbb{G}, \mathbb{G}', \mathbb{G}'')$ in $[\ell^\infty(\mathcal{F})]^3$, where \mathbb{G} and \mathbb{G}' are independent P -Brownian bridges.*

Proof. By the preceding theorem, the three processes are asymptotically tight marginally and hence asymptotically tight jointly. Since the first process is uncorrelated with the second process, the limiting distribution of the first process is independent of the limiting distribution of the second process. As argued in the proof of the multiplier central limit theorem, the uniform difference between \mathbb{G}'_n and \mathbb{G}''_n goes to zero in probability, and thus the remainder of the corollary follows. \square

10.1.2 Conditional Multiplier Central Limit Theorems

In this section, the convergence properties of the multiplier processes in the previous section are studied conditional on the data. This yields in-probability and outer-almost-sure conditional multiplier central limit theorems. These results are one step closer to the bootstrap validity results of the next section. For a metric space (\mathbb{D}, d) , define $BL_1(\mathbb{D})$ to be the space of all functions $f : \mathbb{D} \mapsto \mathbb{R}$ with Lipschitz norm bounded by 1, i.e., $\|f\|_\infty \leq 1$ and $|f(x) - f(y)| \leq d(x, y)$ for all $x, y \in \mathbb{D}$. In the current set-up, $\mathbb{D} = \ell^\infty(\mathcal{F})$, for some class of measurable functions \mathcal{F} , and d is the corresponding uniform metric. As we did in section 2.2.3, we will use BL_1 as

shorthand for $BL_1(\ell^\infty(\mathcal{F}))$. The conditional weak convergence arrows we use in theorems 10.4 and 10.6 below were also defined in section 2.2.3.

We now present the in-probability conditional multiplier central limit theorem:

THEOREM 10.4 *Let \mathcal{F} be a class of measurable functions, and let ξ_1, \dots, ξ_n be i.i.d. random variables with mean zero, variance 1, and $\|\xi\|_{2,1} < \infty$, independent of the sample data X_1, \dots, X_n . Let $\mathbb{G}'_n, \mathbb{G}''_n$ and $\bar{\xi}$ be as defined in theorem 10.1. Then the following are equivalent:*

(i) \mathcal{F} is Donsker;

(ii) $\mathbb{G}'_n \xrightarrow[\xi]{P} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ and \mathbb{G}'_n is asymptotically measurable.

(iii) $\mathbb{G}''_n \xrightarrow[\xi]{P} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ and \mathbb{G}''_n is asymptotically measurable.

Before giving the proof of this theorem, we make a few points and present lemma 10.5 below to aid in the proof. In the above theorem, E_ξ denotes taking the expectation conditional on X_1, \dots, X_n . Note that for a continuous function $h : \ell^\infty(\mathcal{F}) \mapsto \mathbb{R}$, if we fix X_1, \dots, X_n , then $(a_1, \dots, a_n) \mapsto h(n^{-1/2} \sum_{i=1}^n a_i(\delta_{X_i} - P))$ is a measurable map from \mathbb{R}^n to \mathbb{R} , provided $\|f(X) - Pf\|_{\mathcal{F}}^* < \infty$ almost surely. This last inequality is tacetly assumed so that the empirical processes under investigation reside in $\ell^\infty(\mathcal{F})$. Thus the expectation E_ξ in conclusions (ii) and (iii) is proper. The following lemma is a conditional multiplier central limit theorem for i.i.d. Euclidean data:

LEMMA 10.5 *Let Z_1, \dots, Z_n be i.i.d. Euclidean random vectors, with $EZ = 0$ and $E\|Z\|^2 < \infty$, independent of the i.i.d. sequence of real random variables ξ_1, \dots, ξ_n with $E\xi = 0$ and $E\xi^2 = 1$. Then, conditionally on Z_1, Z_2, \dots , $n^{-1/2} \sum_{i=1}^n \xi_i Z_i \rightsquigarrow N(0, \text{cov}Z)$, for almost all sequences Z_1, Z_2, \dots*

Proof. By the Lindeberg central limit theorem, convergence to the given normal limit will occur for every sequence Z_1, Z_2, \dots for which

$$n^{-1} \sum_{i=1}^n Z_i Z_i^T \rightarrow \text{cov}Z$$

and

$$(10.1) \quad \frac{1}{n} \sum_{i=1}^n \|Z_i\|^2 E_\xi \xi_i^2 1\{|\xi_i| \times \|Z_i\| > \epsilon\sqrt{n}\} \rightarrow 0,$$

for all $\epsilon > 0$, where E_ξ is the conditional expectation given the Z_1, Z_2, \dots . The first condition is true for almost all sequences by the strong law of

large numbers. We now evaluate the second condition. Fix $\epsilon > 0$. Now, for any $\tau > 0$, the sum in (10.1) is bounded by

$$\frac{1}{n} \sum_{i=1}^n \|Z_i\|^2 \mathbb{E}[\xi^2 1\{|\xi| > \epsilon/\tau\}] + \frac{1}{n} \sum_{i=1}^n \|Z_i\|^2 1\{\|Z_i\| > \sqrt{n}\tau\}.$$

The first sum has an arbitrarily small upper bound in the limit if we choose sufficiently small τ . Since $\mathbb{E}\|Z\|^2 < \infty$, the second sum will go to zero for almost all sequences Z_1, Z_2, \dots . Thus, for almost all sequences Z_1, Z_2, \dots , (10.1) will hold for any $\epsilon > 0$. For, the intersection of the two sets of sequences, all required conditions hold, and the desired result follows. \square

Proof of theorem 10.4. Since the processes \mathbb{G} , \mathbb{G}_n , \mathbb{G}'_n and \mathbb{G}''_n are unaffected if the class \mathcal{F} is replaced with $\{f - Pf : f \in \mathcal{F}\}$, we will assume $\|P\|_{\mathcal{F}} = 0$ throughout the proof, without loss of generality.

(i) \Rightarrow (ii): If \mathcal{F} is Donsker, the sequence \mathbb{G}'_n converges in distribution to a Brownian bridge process by the unconditional multiplier central limit theorem (theorem 10.1). Thus \mathbb{G}'_n is asymptotically measurable. Now, by lemma 8.17, a Donsker class is totally bounded by the semimetric $\rho_P(f, g) \equiv (P[f - g]^2)^{1/2}$. For each fixed $\delta > 0$ and $f \in \mathcal{F}$, denote $\Pi_\delta f$ to be the closest element in a given, finite δ -net (with respect to the metric ρ_P) for \mathcal{F} . We have by continuity of the limit process \mathbb{G} , that $\mathbb{G} \circ \Pi_\delta \rightarrow \mathbb{G}$, almost surely, as $\delta \downarrow 0$. Hence, for any sequence $\delta_n \downarrow 0$,

$$(10.2) \quad \sup_{h \in BL_1} |\mathbb{E}h(\mathbb{G} \circ \Pi_{\delta_n}) - \mathbb{E}h(\mathbb{G})| \rightarrow 0.$$

By lemma 10.5 above, we also have for any fixed $\delta > 0$ that

$$(10.3) \quad \sup_{h \in BL_1} |\mathbb{E}_\xi h(\mathbb{G}'_n \circ \Pi_\delta) - \mathbb{E}h(\mathbb{G} \circ \Pi_\delta)| \rightarrow 0,$$

as $n \rightarrow \infty$, for almost all sequences X_1, X_2, \dots . To see this, let f_1, \dots, f_m be the δ -mesh of \mathcal{F} that defines Π_δ . Now define the map $A : \mathbb{R}^m \mapsto \ell^\infty(\mathcal{F})$ by $(A(y))(f) = y_k$, where $y = (y_1, \dots, y_m)$ and the integer k satisfies $\Pi_\delta f = f_k$. Now $h(\mathbb{G} \circ \Pi_\delta) = g(\mathbb{G}(f_1), \dots, \mathbb{G}(f_m))$ for the function $g : \mathbb{R}^m \mapsto \mathbb{R}$ defined by $g(y) = h(A(y))$. It is not hard to see that if h is bounded Lipschitz on $\ell^\infty(\mathcal{F})$, then g is also bounded Lipschitz on \mathbb{R}^m with a Lipschitz norm no larger than the Lipschitz norm for h . Now (10.3) follows from lemma 10.5. Note also that $BL_1(\mathbb{R}^m)$ is separable with respect to the metric $\rho_{(m)}(f, g) \equiv \sum_{i=1}^\infty 2^{-i} \sup_{x \in K_i} |f(x) - g(x)|$, where $K_1 \subset K_2 \subset \dots$ are compact sets satisfying $\cup_{i=1}^\infty K_i = \mathbb{R}^m$. Hence, since $\mathbb{G}'_n \circ \Pi_\delta$ and $\mathbb{G} \circ \Pi_\delta$ are both tight, the supremum in 10.3 can be replaced by a countable supremum. Thus the displayed quantity is measurable, since $h(\mathbb{G}'_n \circ \Pi_\delta)$ is measurable.

Now, still holding δ fixed,

$$\begin{aligned}
 \sup_{h \in BL_1} |E_\xi h(\mathbb{G}'_n \circ \Pi_\delta) - E_\xi h(\mathbb{G}'_n)| &\leq \sup_{h \in BL_1} E_\xi |h(\mathbb{G}'_n \circ \Pi_\delta) - h(\mathbb{G}'_n)| \\
 &\leq E_\xi \|\mathbb{G}'_n \circ \Pi_\delta - \mathbb{G}'_n\|_{\mathcal{F}}^* \\
 &\leq E_\xi \|\mathbb{G}'_n\|_{\mathcal{F}_\delta}^*,
 \end{aligned}$$

where $\mathcal{F}_\delta \equiv \{f - g : \rho_P(f, g) < \delta, f, g \in \mathcal{F}\}$. Thus the outer expectation of the left-hand-side is bounded above by $E^* \|\mathbb{G}'_n\|_{\mathcal{F}_\delta}$. As we demonstrated in the proof of theorem 10.1, $E^* \|\mathbb{G}'_n\|_{\mathcal{F}_{\delta_n}} \rightarrow 0$, for any sequence $\delta_n \downarrow 0$. Now, we choose the sequence δ_n so that it goes to zero slowly enough to ensure that (10.3) still holds with δ replaced by δ_n . Combining this with (10.2), the desired result follows.

(ii) \Rightarrow (i): Let $h(\mathbb{G}'_n)^*$ and $h(\mathbb{G}'_n)_*$ denote measurable majorants and minorants with respect to $(\xi_1, \dots, \xi_n, X_1, \dots, X_n)$ jointly. We now have, by the triangle inequality and Fubini's theorem (lemma 6.14),

$$\begin{aligned}
 |E^* h(\mathbb{G}'_n) - Eh(\mathbb{G})| &\leq |E_X E_\xi h(\mathbb{G}'_n)^* - E_X^* E_\xi h(\mathbb{G}'_n)| \\
 &\quad + E_X^* |E_\xi h(\mathbb{G}'_n) - Eh(\mathbb{G})|,
 \end{aligned}$$

where E_X denotes taking the expectation over X_1, \dots, X_n . By (ii) and the dominated convergence theorem, the second term on the right side converges to zero for all $h \in BL_1$. Since the first term on the right is bounded above by $E_X E_\xi h(\mathbb{G}'_n)^* - E_X E_\xi h(\mathbb{G}'_n)_*$, it also converges to zero since \mathbb{G}'_n is asymptotically measurable. It is easy to see that the same result holds true if BL_1 is replaced by the class of all bounded, Lipschitz continuous nonnegative functions $h : \ell^\infty(\mathcal{F}) \mapsto \mathbb{R}$, and thus $\mathbb{G}'_n \rightsquigarrow \mathbb{G}$ unconditionally by the Portmanteau theorem (theorem 7.6). Hence \mathcal{F} is Donsker by the converse part of theorem 10.1.

(ii) \Rightarrow (iii): Since we can assume $\|P\|_{\mathcal{F}} = 0$, we have

$$(10.4) \quad |h(\mathbb{G}'_n) - h(\mathbb{G}''_n)| \leq |\bar{\xi} \mathbb{G}_n|.$$

Moreover, since (ii) also implies (i), we have that $E^* \|\bar{\xi} \mathbb{G}_n\|_{\mathcal{F}} \rightarrow 0$ by lemma 8.17. Thus $\sup_{h \in BL_1} |E_\xi h(\mathbb{G}'_n) - E_\xi h(\mathbb{G}''_n)| \rightarrow 0$ in outer probability. Since (10.4) also implies that \mathbb{G}''_n is asymptotically measurable, (iii) follows.

(iii) \Rightarrow (i): Arguing as we did in the proof that (ii) \Rightarrow (i), it is not hard to show that $\mathbb{G}''_n \rightsquigarrow \mathbb{G}$ unconditionally. Now theorem 10.1 yields that \mathcal{F} is Donsker. \square

We now present the outer-almost-sure conditional multiplier central limit theorem:

THEOREM 10.6 *Assume the conditions of theorem 10.4. Then the following are equivalent:*

$$(i) \quad \mathcal{F} \text{ is Donsker and } P^* \|f - Pf\|_{\mathcal{F}}^2 < \infty;$$

$$(ii) \quad \mathbb{G}'_n \overset{\text{as}^*}{\underset{\xi}{\rightsquigarrow}} \mathbb{G} \text{ in } \ell^\infty(\mathcal{F}).$$

(iii) $\mathbb{G}_n'' \xrightarrow[\xi]{\text{as*}} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$.

Proof. The equivalence of (i) and (ii) is given in theorem 2.9.7 of VW, and we omit its proof.

(ii) \Rightarrow (iii): As in the proof of theorem 10.4, we assume that $\|P\|_{\mathcal{F}} = 0$ without loss of generality. Since

$$(10.5) \quad |h(\mathbb{G}'_n) - h(\mathbb{G}''_n)| \leq |\sqrt{n\bar{\xi}}| \times \|\mathbb{P}_n\|_{\mathcal{F}},$$

for any $h \in BL_1$, we have

$$\sup_{h \in BL_1} |E_\xi h(\mathbb{G}'_n) - E_\xi h(\mathbb{G}''_n)| \leq E_\xi |\sqrt{n\bar{\xi}}| \times \|\mathbb{P}_n\|_{\mathcal{F}} \leq \|\mathbb{P}_n\|_{\mathcal{F}} \xrightarrow{\text{as*}} 0,$$

since the equivalence of (i) and (ii) implies that \mathcal{F} is both Donsker and Glivenko-Cantelli. Hence

$$\sup_{h \in BL_1} |E_\xi h(\mathbb{G}''_n) - Eh(\mathbb{G})| \xrightarrow{\text{as*}} 0.$$

The relation (10.5) also yields that $E_\xi h(\mathbb{G}''_n)^* - E_\xi h(\mathbb{G}''_n)_* \xrightarrow{\text{as*}} 0$, and thus (iii) follows.

(iii) \Rightarrow (ii): Let $h \in BL_1$. Since $E_\xi h(\mathbb{G}''_n)^* - Eh(\mathbb{G}) \xrightarrow{\text{as*}} 0$, we have $E^* h(\mathbb{G}''_n) \rightarrow Eh(\mathbb{G})$. Since this holds for all $h \in BL_1$, we now have that $\mathbb{G}''_n \xrightarrow{\text{as*}} \mathbb{G}$ unconditionally by the Portmanteau theorem (theorem 7.6). Now we can invoke theorem 10.4 to conclude that \mathcal{F} is both Donsker and Glivenko-Cantelli. Now (10.5) implies (ii) by using an argument almost identical to the one used in the previous paragraph. \square

10.1.3 Bootstrap Central Limit Theorems

Theorems 10.4 and 10.6 will now be used to prove theorems 2.6 and 2.7 from section 2.2.3. Recall that the multinomial bootstrap is obtained by resampling from the data X_1, \dots, X_n , with replacement, n times to obtain a bootstrapped sample X_1^*, \dots, X_n^* . The empirical measure $\hat{\mathbb{P}}_n^*$ of the bootstrapped sample has the same distribution—given the data—as the measure $\hat{\mathbb{P}}_n \equiv n^{-1} \sum_{i=1}^n W_{ni} \delta_{X_i}$, where $W_n \equiv (W_{n1}, \dots, W_{nn})$ is a multinomial(n, n^{-1}, \dots, n^{-1}) deviate independent of the data. As in section 2.2.3, let $\hat{\mathbb{P}}_n \equiv n^{-1} \sum_{i=1}^n W_{ni} \delta_{X_i}$ and $\hat{\mathbb{G}}_n \equiv \sqrt{n}(\hat{\mathbb{P}}_n - \mathbb{P}_n)$. Also recall the definitions $\tilde{\mathbb{P}}_n \equiv n^{-1} \sum_{i=1}^n (\xi/\bar{\xi}) \delta_{X_i}$ and $\tilde{\mathbb{G}}_n \equiv \sqrt{n}(\mu/\tau)(\tilde{\mathbb{P}}_n - \mathbb{P}_n)$, where the weights ξ_1, \dots, ξ_n are i.i.d. nonnegative, independent of X_1, \dots, X_n , with mean $0 < \mu < \infty$ and variance $0 < \tau^2 < \infty$, and with $\|\xi\|_{2,1} < \infty$. When $\bar{\xi} = 0$, we define $\tilde{\mathbb{P}}_n$ to be zero. Note that the weights ξ_1, \dots, ξ_n in this section must have μ subtracted from them and then divided by τ before they satisfy the criteria of the multiplier weights in the previous section.

Proof of theorem 2.6. The equivalence of (i) and (ii) follows from theorem 3.6.1 of VW, which proof we omit. We note, however, that a key component of this proof is a clever approximation of the multinomial weights with i.i.d. Poisson mean 1 weights. We will use this approximation in our proof of theorem 10.15 below.

We now prove the equivalence of (i) and (iii). Let $\xi_i^\circ \equiv \tau^{-1}(\xi_i - \mu)$, $i = 1, \dots, n$, and define $\mathbb{G}_n^\circ \equiv n^{-1/2} \sum_{i=1}^n (\xi_i^\circ - \bar{\xi}^\circ) \delta_{X_i}$, where $\bar{\xi}^\circ \equiv n^{-1} \sum_{i=1}^n \xi_i^\circ$. The basic idea is to show the asymptotic equivalence of $\tilde{\mathbb{G}}_n$ and \mathbb{G}_n° . Then theorem 10.4 can be used to establish the desired result. Accordingly,

$$(10.6) \quad \mathbb{G}_n^\circ - \tilde{\mathbb{G}}_n = \left(1 - \frac{\mu}{\xi}\right) \mathbb{G}_n^\circ = \left(\frac{\bar{\xi}}{\mu} - 1\right) \tilde{\mathbb{G}}_n.$$

First, assume that \mathcal{F} is Donsker. Since the weights $\xi_1^\circ, \dots, \xi_n^\circ$ satisfy the conditions of the unconditional multiplier central limit theorem, we have that $\mathbb{G}_n^\circ \rightsquigarrow \mathbb{G}$. Theorem 10.4 also implies that $\mathbb{G}_n^\circ \xrightarrow[\xi]{\mathbb{P}} \mathbb{G}$. Now (10.6) can be applied to verify that $\|\tilde{\mathbb{G}}_n - \mathbb{G}_n^\circ\|_{\mathcal{F}} \xrightarrow{\mathbb{P}} 0$, and thus $\tilde{\mathbb{G}}_n$ is asymptotically measurable and

$$\sup_{h \in BL_1} \left| \mathbb{E}_\xi h(\mathbb{G}_n^\circ) - \mathbb{E}_\xi h(\tilde{\mathbb{G}}_n) \right| \xrightarrow{\mathbb{P}} 0.$$

Thus (i) \Rightarrow (iii).

Second, assume that $\tilde{\mathbb{G}}_n \xrightarrow[\xi]{\mathbb{P}} \mathbb{G}$. It is not hard to show, arguing as we did in the proof of theorem 10.4 for the implication (ii) \Rightarrow (i), that $\tilde{\mathbb{G}}_n \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ unconditionally. By applying (10.6) again, we now have that $\|\mathbb{G}_n^\circ - \tilde{\mathbb{G}}_n\|_{\mathcal{F}} \xrightarrow{\mathbb{P}} 0$, and thus $\mathbb{G}_n^\circ \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ unconditionally. The unconditional multiplier central limit theorem now verifies that \mathcal{F} is Donsker, and thus (iii) \Rightarrow (i). \square

Proof of theorem 2.7. The equivalence of (i) and (ii) follows from theorem 3.6.2 of VW, which proof we again omit. We now prove the equivalence of (i) and (iii).

First, assume (i). Then $\mathbb{G}_n^\circ \xrightarrow[\xi]{\text{as}^*} \mathbb{G}$ by theorem 10.6. Fix $\rho > 0$, and note that by using the first equality in (10.6), we have for any $h \in BL_1$ that

$$(10.7) \quad \left| h(\tilde{\mathbb{G}}_n) - h(\mathbb{G}_n^\circ) \right| \leq 2 \times 1 \left\{ \left| 1 - \frac{\mu}{\xi} \right| > \rho \right\} + (\rho \|\mathbb{G}_n^\circ\|_{\mathcal{F}}) \wedge 1.$$

The first term on the right $\xrightarrow{\text{as}^*} 0$. Since the map $\|\cdot\|_{\mathcal{F}} \wedge 1 : \ell^\infty(\mathcal{F}) \mapsto \mathbb{R}$ is in BL_1 , we have by theorem 10.6 that $\mathbb{E}_\xi [(\rho \|\mathbb{G}_n^\circ\|_{\mathcal{F}}) \wedge 1] \xrightarrow{\text{as}^*} \mathbb{E} [\|\rho \mathbb{G}\|_{\mathcal{F}} \wedge 1]$. Let the sequence $0 < \rho_n \downarrow 0$ converge slowly enough so that the first term on the right in (10.7) $\xrightarrow{\text{as}^*} 0$ after replacing ρ with ρ_n . Since $\mathbb{E} [\|\rho_n \mathbb{G}\|_{\mathcal{F}} \wedge 1] \rightarrow 0$, we can apply \mathbb{E}_ξ to both sides of (10.7)—after replacing ρ with ρ_n —to obtain

$$\sup_{h \in BL_1} \left| h(\tilde{\mathbb{G}}_n) - h(\mathbb{G}_n^\circ) \right| \xrightarrow{\text{as}^*} 0.$$

Combining the fact that $h(\mathbb{G}_n^\circ)^* - h(\mathbb{G}_n^\circ)_* \xrightarrow{\text{as}^*} 0$ with additional applications of (10.7) yields $h(\tilde{\mathbb{G}}_n)^* - h(\tilde{\mathbb{G}}_n)_* \xrightarrow{\text{as}^*} 0$. Since h was arbitrary, we have established that $\tilde{\mathbb{G}}_n \xrightarrow[\xi]{\text{as}^*} \mathbb{G}$, and thus (iii) follows.

Second, assume (iii). Fix $\rho > 0$, and note that by using the second equality in (10.6), we have for any $h \in BL_1$ that

$$\left| h(\mathbb{G}_n^\circ) - h(\tilde{\mathbb{G}}_n) \right| \leq 2 \times 1 \left\{ \left| \frac{\bar{\xi}}{\mu} - 1 \right| > \rho \right\} + \left(\rho \|\tilde{\mathbb{G}}_n\|_{\mathcal{F}} \right) \wedge 1.$$

Since the first term on the right $\xrightarrow{\text{as}^*} 0$, we can use virtually identical arguments to those used in the previous paragraph—but with the roles of \mathbb{G}_n° and $\tilde{\mathbb{G}}_n$ reversed—to obtain that $\mathbb{G}_n^\circ \xrightarrow[\xi]{\text{as}^*} \mathbb{G}$. Now theorem 10.6 yields that \mathcal{F} is Donsker, and thus (i) follows. \square

10.1.4 Continuous Mapping Results

We now assume a more general set-up, where \hat{X}_n is a bootstrapped process in a Banach space $(\mathbb{D}, \|\cdot\|)$ and is composed of the sample data $\mathcal{X}_n \equiv (X_1, \dots, X_n)$ and a random weight vector $M_n \in \mathbb{R}^n$ independent of \mathcal{X}_n . We do not require that X_1, \dots, X_n be i.i.d. In this section, we obtain two continuous mapping results. The first result, proposition 10.7, is a simple continuous mapping results for the very special case of Lipschitz continuous maps. It is applicable to both the in-probability or outer-almost-sure versions of bootstrap consistency. An interesting special case is the map $g(x) = \|x\|$. In this case, the proposition validates the use of the bootstrap to construct asymptotically uniformly valid confidence bands for $\{Pf : f \in \mathcal{F}\}$ whenever Pf is estimated by $\mathbb{P}_n f$ and \mathcal{F} is P -Donsker.

Now assume that $\hat{X}_n \xrightarrow[M]{P} X$ and that the distribution of $\|X\|$ is continuous. Lemma 10.11 towards the end of this section reveals that $P(\|\hat{X}_n\| \leq t | \mathcal{X}_n)$ converges uniformly to $P(\|X\| \leq t)$, in probability. A parallel outer almost sure result holds when $\hat{X}_n \xrightarrow[M]{\text{as}^*} X$.

The second result, theorem 10.8, is a considerably deeper result for general continuous maps applied to bootstraps which are consistent in probability. Because of this generality, we must require certain measurability conditions on the map $M_n \mapsto \hat{X}_n$. Fortunately, based on the discussion in the paragraph following theorem 10.4 above, these measurability conditions are easily satisfied when either $\hat{X}_n = \hat{\mathbb{G}}_n$ or $\hat{X}_n = \tilde{\mathbb{G}}_n$. It appears that other continuous mapping results for bootstrapped empirical processes hold, such as for bootstraps which are outer almost surely consistent, but such results seem to be very challenging to verify.

PROPOSITION 10.7 *Let \mathbb{D} and \mathbb{E} be Banach spaces, X a tight random variable on \mathbb{D} , and $g : \mathbb{D} \mapsto \mathbb{E}$ Lipschitz continuous. We have the following:*

(i) *If $\hat{X}_n \overset{P}{\rightsquigarrow}_M X$, then $g(\hat{X}_n) \overset{P}{\rightsquigarrow}_M g(X)$.*

(ii) *If $\hat{X}_n \overset{as*}{\rightsquigarrow}_M X$, then $g(\hat{X}_n) \overset{as*}{\rightsquigarrow}_M g(X)$.*

Proof. Let $c_0 < \infty$ be the Lipschitz constant for g , and, without loss of generality, assume $c_0 \geq 1$. Note that for any $h \in BL_1(\mathbb{E})$, the map $x \mapsto h(g(x))$ is an element of $c_0 BL_1(\mathbb{D})$. Thus

$$\begin{aligned} \sup_{h \in BL_1(\mathbb{E})} \left| \mathbb{E}_M h(g(\hat{X}_n)) - \mathbb{E} h(g(X)) \right| &\leq \sup_{h \in c_0 BL_1(\mathbb{D})} \left| \mathbb{E}_M h(\hat{X}_n) - \mathbb{E} h(X) \right| \\ &= c_0 \sup_{h \in BL_1(\mathbb{D})} \left| \mathbb{E}_M h(\hat{X}_n) - \mathbb{E} h(X) \right|, \end{aligned}$$

and the desired result follows by the respective definitions of $\overset{P}{\rightsquigarrow}_M$ and $\overset{as*}{\rightsquigarrow}_M$. \square

THEOREM 10.8 *Let $g : \mathbb{D} \mapsto \mathbb{E}$ be continuous at all points in $\mathbb{D}_0 \subset \mathbb{D}$, where \mathbb{D} and \mathbb{E} are Banach spaces and \mathbb{D}_0 is closed. Assume that $M_n \mapsto h(\hat{X}_n)$ is measurable for every $h \in C_b(\mathbb{D})$ outer almost surely. Then if $\hat{X}_n \overset{P}{\rightsquigarrow}_M X$ in \mathbb{D} , where X is tight and $P^*(X \in \mathbb{D}_0) = 1$, $g(\hat{X}_n) \overset{P}{\rightsquigarrow}_M g(X)$.*

Proof. As in the proof of the implication (ii) \Rightarrow (i) of theorem 10.4, we can argue that $\hat{X}_n \rightsquigarrow X$ unconditionally, and thus $g(\hat{X}_n) \rightsquigarrow g(X)$ unconditionally by the standard continuous mapping theorem. Moreover, we can replace \mathbb{E} with its closed linear span so that the restriction of g to \mathbb{D}_0 has an extension $\tilde{g} : \mathbb{D} \mapsto \mathbb{E}$ which is continuous on all of \mathbb{D} by Dugundji's extension theorem (theorem 10.9 below). Thus $(g(\hat{X}_n), \tilde{g}(\hat{X}_n)) \rightsquigarrow (g(X), \tilde{g}(X))$, and hence $g(\hat{X}_n) - \tilde{g}(\hat{X}_n) \overset{P}{\rightarrow} 0$. Therefore we can assume without loss of generality that g is continuous on all of \mathbb{D} . We can also assume without loss of generality that \mathbb{D}_0 is a separable Banach space since X is tight. Hence $\mathbb{E}_0 \equiv g(\mathbb{D}_0)$ is also a separable Banach space.

Fix $\epsilon > 0$. There now exists a compact $K \subset \mathbb{E}_0$ such that $P(X \notin K) < \epsilon$. By theorem 10.9 below, the proof of which is given in section 10.4, we know there exists an integer $k < \infty$, elements $z_1, \dots, z_k \in C[0, 1]$, continuous functions $f_1, \dots, f_k : \mathbb{E} \mapsto \mathbb{R}$, and a Lipschitz continuous function $J : \overline{\text{lin}}(z_1, \dots, z_k) \mapsto \mathbb{E}$, such that the map $x \mapsto T_\epsilon(x) \equiv J\left(\sum_{j=1}^k z_j f_j(x)\right)$ has domain \mathbb{E} and range $\subset \mathbb{E}$ and satisfies $\sup_{x \in K} \|T_\epsilon(x) - x\| < \epsilon$. Let $BL_1 \equiv BL_1(\mathbb{E})$. We now have

$$\begin{aligned}
& \sup_{h \in BL_1} \left| \mathbb{E}_M h(g(\hat{X}_n)) - \mathbb{E} h(g(X)) \right| \\
& \leq \sup_{h \in BL_1} \left| \mathbb{E}_M h(T_\epsilon g(\hat{X}_n)) - \mathbb{E} h(T_\epsilon g(X)) \right| \\
& \quad + \mathbb{E}_M \left\{ \left\| T_\epsilon g(\hat{X}_n) - g(\hat{X}_n) \right\| \wedge 2 \right\} + \mathbb{E} \left\{ \left\| T_\epsilon g(X) - g(X) \right\| \wedge 2 \right\}.
\end{aligned}$$

However, the outer expectation of the second term on the right converges to the third term, as $n \rightarrow \infty$, by the usual continuous mapping theorem. Thus, provided

$$(10.8) \quad \sup_{h \in BL_1} \left| \mathbb{E}_M h(T_\epsilon g(\hat{X}_n)) - \mathbb{E} h(T_\epsilon g(X)) \right| \xrightarrow{P} 0,$$

we have that

$$\begin{aligned}
(10.9) \quad & \limsup_{n \rightarrow \infty} \mathbb{E}^* \left\{ \sup_{h \in BL_1} \left| \mathbb{E}_M h(g(\hat{X}_n)) - \mathbb{E} h(g(X)) \right| \right\} \\
& \leq 2\mathbb{E} \left\{ \left\| T_\epsilon g(X) - g(X) \right\| \vee 2 \right\} \\
& \leq 2\mathbb{E} \left\{ \left\| T_\epsilon g(X) - g(X) \right\| 1_{\{g(X) \in K\}} \right\} + 4\mathbb{P}(g(X) \notin K) \\
& < 6\epsilon.
\end{aligned}$$

Now note that for each $h \in BL_1$, $h \left(J \left(\sum_{j=1}^k z_j a_j \right) \right) = \tilde{h}(a_1, \dots, a_k)$ for all $(a_1, \dots, a_k) \in \mathbb{R}^k$ and some $\tilde{h} \in c_0 BL_1(\mathbb{R}^k)$, where $1 \leq c_0 < \infty$ (this follows since J is Lipschitz continuous and $\left\| \sum_{j=1}^k z_j a_j \right\| \leq \max_{1 \leq j \leq k} |a_j| \times \sum_{j=1}^k \|z_j\|$). Hence

$$\begin{aligned}
(10.10) \quad & \sup_{h \in BL_1} \left| \mathbb{E}_M h(T_\epsilon g(\hat{X}_n)) - \mathbb{E} h(T_\epsilon g(X)) \right| \\
& \leq \sup_{h \in c_0 BL_1(\mathbb{R}^k)} \left| \mathbb{E}_M h(u(\hat{X}_n)) - \mathbb{E} h(u(X)) \right| \\
& = c_0 \sup_{h \in BL_1(\mathbb{R}^k)} \left| \mathbb{E}_M h(u(\hat{X}_n)) - \mathbb{E} h(u(X)) \right|,
\end{aligned}$$

where $x \mapsto u(x) \equiv (f_1(g(x)), \dots, f_k(g(x)))$. Fix any $v : \mathbb{R}^k \mapsto [0, 1]$ which is Lipschitz continuous (the Lipschitz constant may be > 1). Then, since $\hat{X}_n \rightsquigarrow X$ unconditionally, $\mathbb{E}^* \left\{ \mathbb{E}_M v(u(\hat{X}_n))^* - \mathbb{E}_M v(u(\hat{X}_n))_* \right\} \leq \mathbb{E}^* \left\{ v(u(\hat{X}_n))^* - v(u(\hat{X}_n))_* \right\} \rightarrow 0$, where sub- and super- script $*$ denote measurable majorants and minorants, respectively, with respect to the joint probability space of (\mathcal{X}_n, M_n) . Thus

$$(10.11) \quad \left| \mathbb{E}_M v(u(\hat{X}_n)) - \mathbb{E}_M v(u(\hat{X}_n))^* \right| \xrightarrow{P} 0.$$

Note that we are using at this point the outer almost sure measurability of $M_n \mapsto v(u(\hat{X}_n))$ to ensure that $\mathbb{E}_M v(u(\hat{X}_n))$ is well defined, even if the resulting random expectation is not itself measurable.

Now, for every subsequence n' , there exists a further subsequence n'' such that $\hat{X}_{n''} \xrightarrow[M]{\text{as}^*} X$. This means that for this subsequence, the set B of data subsequences $\{\mathcal{X}_{n''} : n \geq 1\}$ for which $\mathbb{E}_M v(u(\hat{X}_{n''})) - \mathbb{E}v(u(X)) \rightarrow 0$ has inner probability 1. Combining this with (10.11) and proposition 7.22, we obtain that $\mathbb{E}_M v(u(\hat{X}_n)) - \mathbb{E}v(u(X)) \xrightarrow{P} 0$. Since v was an arbitrary real, Lipschitz continuous function on \mathbb{R}^k , we now have by part (i) of lemma 10.11 below followed by lemma 10.12 below, that

$$\sup_{h \in BL_1(\mathbb{R}^k)} \left| \mathbb{E}_M h(u(\hat{X}_n)) - \mathbb{E}h(u(X)) \right| \xrightarrow{P} 0.$$

Combining this with (10.10), we obtain that (10.8) is satisfied. The desired result now follows from (10.9), since $\epsilon > 0$ was arbitrary. \square

THEOREM 10.9 (*Dugundji's extension theorem*) *Let X be an arbitrary metric space, A a closed subset of X , L a locally convex linear space (which includes Banach vector spaces), and $f : A \mapsto L$ a continuous map. Then there exists a continuous extension of f , $F : X \mapsto L$. Moreover, $F(X)$ lies in the closed linear span of the convex hull of $f(A)$.*

Proof. This is theorem 4.1 of Dugundji (1951), and the proof can be found therein. \square

THEOREM 10.10 *Let $\mathbb{E}_0 \subset \mathbb{E}$ be Banach spaces with \mathbb{E}_0 separable and $\overline{\text{lin}} \mathbb{E}_0 \subset \mathbb{E}$. Then for every $\epsilon > 0$ and every compact $K \subset \mathbb{E}_0$, there exists an integer $k < \infty$, elements $z_1, \dots, z_k \in C[0, 1]$, continuous functions $f_1, \dots, f_k : \mathbb{E} \mapsto \mathbb{R}$, and a Lipschitz continuous function $J : \overline{\text{lin}}(z_1, \dots, z_k) \mapsto \mathbb{E}$, such that the map $x \mapsto T_\epsilon(x) \equiv J\left(\sum_{j=1}^k z_j f_j(x)\right)$ has domain \mathbb{E} and range $\subset \mathbb{E}$, is continuous, and satisfies $\sup_{x \in K} \|T_\epsilon(x) - x\| < \epsilon$.*

The proof of this theorem is given in section 10.4. For the next two lemmas, we use the usual partial ordering on \mathbb{R}^k to define relations between points, e.g., for any $s, t \in \mathbb{R}^k$, $s \leq t$ is equivalent to $s_1 \leq t_1, \dots, s_k \leq t_k$.

LEMMA 10.11 *Let X_n and X be random variables in \mathbb{R}^k for all $n \geq 1$. Define $\mathcal{S} \subset [\mathbb{R} \cup \{-\infty, \infty\}]^k$ to be the set of all continuity points of $t \mapsto F(t) \equiv \mathbb{P}(X \leq t)$ and H to be the set of all Lipschitz continuous functions $h : \mathbb{R}^k \mapsto [0, 1]$ (the Lipschitz constants may be > 1). Then, provided the expectations are well defined, we have:*

- (i) *If $\mathbb{E}[h(X_n)|\mathcal{Y}_n] \xrightarrow{P} \mathbb{E}h(X)$ for all $h \in H$, then $\sup_{t \in A} |\mathbb{P}(X_n \leq t|\mathcal{Y}_n) - F(t)| \xrightarrow{P} 0$ for all closed $A \subset \mathcal{S}$;*
- (ii) *If $\mathbb{E}[h(X_n)|\mathcal{Y}_n] \xrightarrow{\text{as}^*} \mathbb{E}h(X)$ for all $h \in H$, then $\sup_{t \in A} |\mathbb{P}(X_n \leq t|\mathcal{Y}_n) - F(t)| \xrightarrow{\text{as}^*} 0$ for all closed $A \subset \mathcal{S}$.*

Proof. Let $t_0 \in \mathcal{S}$. For every $\delta > 0$, there exists $h_1, h_2 \in H$, such that $h_1(u) \leq 1\{u \leq t_0\} \leq h_2(u)$ for all $u \in \mathbb{R}^k$ and $E[h_2(X) - h_1(X)] < \delta$. Under the condition in (i), we therefore have that $P(X_n \leq t_0 | \mathcal{Y}_n) \xrightarrow{P} F(t_0)$, since δ was arbitrary. The conclusion of (i) follows since this convergence holds for all $t_0 \in \mathcal{S}$, since both $P(X_n \leq t | \mathcal{Y}_n)$ and $F(t)$ are monotone in t with range $\subset [0, 1]$, and since $[0, 1]$ is compact. The proof for part (ii) follows similarly. \square

LEMMA 10.12 *Let $\{F_n\}$ and F be distribution functions on \mathbb{R}^k , and let $\mathcal{S} \subset [\mathbb{R} \cup \{-\infty, \infty\}]^k$ be the set of all continuity points of F . Then the following are equivalent:*

$$(i) \sup_{t \in A} |F_n(t) - F(t)| \rightarrow 0 \text{ for all closed } A \subset \mathcal{S}.$$

$$(ii) \sup_{h \in BL_1(\mathbb{R}^k)} \left| \int_{\mathbb{R}^k} h(dF_n - dF) \right| \rightarrow 0.$$

The relatively straightforward proof is saved as exercise 10.5.3.

10.2 The Bootstrap for Glivenko-Cantelli Classes

We now present several results for the bootstrap applied to Glivenko-Cantelli classes. The primary use of these results is to assist verification of consistency of bootstrapped estimators. The first theorem (theorem 10.13) consists of various multiplier bootstrap results, and it is followed by a corollary (corollary 10.14) which applies to certain weighted bootstrap results. The final theorem of this section (theorem 10.15) gives gives corresponding results for the multinomial bootstrap. On a first reading through this section, it might be best to skip the proofs and focus on the results and discussion between the proofs.

THEOREM 10.13 *Let \mathcal{F} be a class of measurable functions, and let ξ_1, \dots, ξ_n be i.i.d. nonconstant random variables with $0 < E|\xi| < \infty$ and independent of the sample data X_1, \dots, X_n . Let $\mathbb{W}_n \equiv n^{-1} \sum_{i=1}^n \xi_i (\delta_{X_i} - P)$ and $\tilde{\mathbb{W}}_n \equiv n^{-1} \sum_{i=1}^n (\xi_i - \bar{\xi}) \delta_{X_i}$, where $\bar{\xi} \equiv n^{-1} \sum_{i=1}^n \xi_i$. Then the following are equivalent:*

$$(i) \mathcal{F} \text{ is strong Glivenko-Cantelli};$$

$$(ii) \|\mathbb{W}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0;$$

$$(iii) E_{\xi} \|\mathbb{W}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0 \text{ and } P^* \|f - Pf\|_{\mathcal{F}} < \infty;$$

$$(iv) \text{For every } \eta > 0, P(\|\mathbb{W}_n\|_{\mathcal{F}} > \eta | \mathcal{X}_n) \xrightarrow{\text{as}^*} 0 \text{ and } P^* \|f - Pf\|_{\mathcal{F}} < \infty, \\ \text{where } \mathcal{X}_n \equiv (X_1, \dots, X_n);$$

- (v) For every $\eta > 0$, $P(\|\mathbb{W}_n\|_{\mathcal{F}}^* > \eta \mid \mathcal{X}_n) \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$, for some version of $\|\mathbb{W}_n\|_{\mathcal{F}}^*$, where the superscript $*$ denotes a measurable majorant with respect to $(\xi_1, \dots, \xi_n, X_1, \dots, X_n)$ jointly;
- (vi) $\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$;
- (vii) $E_{\xi}\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$;
- (viii) For every $\eta > 0$, $P(\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}} > \eta \mid \mathcal{X}_n) \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$;
- (ix) For every $\eta > 0$, $P(\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}}^* > \eta \mid \mathcal{X}_n) \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$, for some version of $\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}}^*$.

The lengthy proof is given in section 4 below. As is shown in the proof, the conditional expectations and conditional probabilities in (iii), (iv), (vii) and (viii) are well defined. This is because the quantities inside of the expectations in parts (iii) and (vii) (and in the conditional probabilities of (iv) and (viii)) are measurable as functions of ξ_1, \dots, ξ_n conditional on the data. The distinctions between (iv) and (v) and between (viii) and (ix) are not as trivial as they appear. This is because the measurable majorants involved are computed with respect to $(\xi_1, \dots, \xi_n, X_1, \dots, X_n)$ jointly, and thus the differences between $\|\mathbb{W}_n\|_{\mathcal{F}}$ and $\|\mathbb{W}_n\|_{\mathcal{F}}^*$ or between $\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}}$ and $\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}}^*$ may be nontrivial.

The following corollary applies to a class of weighted bootstraps which includes the Bayesian bootstrap mentioned earlier:

COROLLARY 10.14 *Let \mathcal{F} be a class of measurable functions, and let ξ_1, \dots, ξ_n be i.i.d. nonconstant, nonnegative random variables with $0 < E\xi < \infty$ and independent of X_1, \dots, X_n . Let $\tilde{\mathbb{P}}_n \equiv n^{-1} \sum_{i=1}^n (\xi_i/\bar{\xi})\delta_{X_i}$, where we set $\tilde{\mathbb{P}}_n = 0$ when $\bar{\xi} = 0$. Then the following are equivalent:*

- (i) \mathcal{F} is strong Glivenko-Cantelli.
- (ii) $\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$.
- (iii) $E_{\xi}\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$.
- (iv) For every $\eta > 0$, $P(\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}} > \eta \mid \mathcal{X}_n) \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$;
- (v) For every $\eta > 0$, $P(\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \mid \mathcal{X}_n) \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$, for some version of $\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^*$.

If in addition $P(\xi = 0) = 0$, then the requirement that $P^*\|f - Pf\|_{\mathcal{F}} < \infty$ in (ii) may be dropped.

Proof. Since the processes $\mathbb{P}_n - P$ and $\tilde{\mathbb{P}}_n - \mathbb{P}_n$ do not change when the class \mathcal{F} is replaced with $\dot{\mathcal{F}} \equiv \{f - Pf : f \in \mathcal{F}\}$, we can assume $\|P\|_{\mathcal{F}} = 0$ without loss of generality. Let the envelope of $\dot{\mathcal{F}}$ be denoted $\dot{F} \equiv \|f\|_{\dot{\mathcal{F}}}^*$. Since multiplying the ξ_i by a constant does not change $\xi_i/\bar{\xi}$, we can also assume $E\xi = 1$ without loss of generality. The fact that the conditional expressions in (iii) and (iv) are well defined can be argued as in the proof of theorem 10.13 (given below in section 4), and we do not repeat the details here.

(i) \Rightarrow (ii): Since

$$(10.12) \quad \tilde{\mathbb{P}}_n - \mathbb{P}_n - \tilde{\mathbb{W}}_n = \left(\frac{1}{\bar{\xi}} - 1\right) 1\{\bar{\xi} > 0\} \tilde{\mathbb{W}}_n - 1\{\bar{\xi} = 0\} \mathbb{P}_n,$$

(ii) follows by theorem 10.13 and the fact that $\bar{\xi} \xrightarrow{\text{as}^*} 1$.

(ii) \Rightarrow (i): Note that

$$(10.13) \quad \tilde{\mathbb{P}}_n - \mathbb{P}_n - \tilde{\mathbb{W}}_n = -(\bar{\xi} - 1) 1\{\bar{\xi} > 0\} (\tilde{\mathbb{P}}_n - \mathbb{P}_n) - 1\{\bar{\xi} = 0\} \mathbb{P}_n.$$

The first term on the right $\xrightarrow{\text{as}^*} 0$ by (ii), while the second term on the right is bounded in absolute value by $1\{\bar{\xi} = 0\} \|\mathbb{P}_n\|_{\dot{\mathcal{F}}} \leq 2 \times 1\{\bar{\xi} = 0\} \mathbb{P}_n \dot{F} \xrightarrow{\text{as}^*} 0$, by the moment condition.

(ii) \Rightarrow (iii): The method of proof will be to use the expansion (10.12) to show that $E_{\xi} \|\tilde{\mathbb{P}}_n - \mathbb{P}_n - \tilde{\mathbb{W}}_n\|_{\dot{\mathcal{F}}} \xrightarrow{\text{as}^*} 0$. Then (iii) will follow by theorem 10.13 and the established equivalence between (ii) and (i). Along this vein, we have by symmetry followed by an application of theorem 9.28 that

$$\begin{aligned} E_{\xi} \left\{ \left| \frac{1}{\bar{\xi}} - 1 \right| 1\{\bar{\xi} > 0\} \|\tilde{\mathbb{W}}_n\|_{\dot{\mathcal{F}}} \right\} &\leq \frac{1}{n} \sum_{i=1}^n \dot{F}(X_i) E_{\xi} \left\{ \left| \frac{1}{\bar{\xi}} - 1 \right| 1\{\bar{\xi} > 0\} \right\} \\ &= \mathbb{P}_n \dot{F} E_{\xi} \{1 - \bar{\xi} 1\{\bar{\xi} > 0\}\} \\ &\xrightarrow{\text{as}^*} 0. \end{aligned}$$

Since also $E_{\xi} [1\{\bar{\xi} = 0\}] \|\mathbb{P}_n\|_{\dot{\mathcal{F}}} \xrightarrow{\text{as}^*} 0$, the desired conclusion follows.

(iii) \Rightarrow (iv): This is obvious.

(iv) \Rightarrow (i): Consider again expansion (10.13). The moment conditions easily give us, conditional on X_1, X_2, \dots , that $1\{\bar{\xi} = 0\} \|\mathbb{P}_n\|_{\dot{\mathcal{F}}} \leq 1\{\bar{\xi} = 0\} \mathbb{P}_n \dot{F} \xrightarrow{P} 0$ for almost all sequences X_1, X_2, \dots . By (iv), we also obtain that $|\bar{\xi} - 1| 1\{\bar{\xi} > 0\} \|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\dot{\mathcal{F}}} \xrightarrow{P} 0$ for almost all sequences X_1, X_2, \dots . Thus assertion (viii) of theorem 10.13 follows, and we obtain (i).

If $P(\xi = 0) = 0$, then $1\{\bar{\xi} = 0\} \mathbb{P}_n = 0$ almost surely, and we no longer need the moment condition $P\dot{F} < \infty$ in the proofs of (ii) \Rightarrow (i) and (ii) \Rightarrow (iii), and thus the moment condition in (ii) can be dropped in this setting.

(ii) \Rightarrow (v): Assertion (ii) implies that there exists a measurable set B of infinite sequences $(\xi_1, X_1), (\xi_2, X_2), \dots$ with $P(B) = 1$ such that $\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\dot{\mathcal{F}}}^* \rightarrow 0$ on B for some version of $\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\dot{\mathcal{F}}}^*$. Thus by the bounded

convergence theorem, we have for any $\eta > 0$ and almost all sequences X_1, X_2, \dots ,

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbb{P} \left(\|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right) &= \limsup_{n \rightarrow \infty} \mathbb{E}_{\xi, \infty} 1 \left\{ \|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right\} 1\{B\} \\ &= \mathbb{E}_{\xi, \infty} \limsup_{n \rightarrow \infty} 1 \left\{ \|\tilde{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right\} 1\{B\} \\ &= 0. \end{aligned}$$

Thus (v) follows.

(v) \Rightarrow (iv): This is obvious. \square

The following theorem verifies consistency of the multinomial bootstrapped empirical measure defined in section 10.1.3, which we denote $\hat{\mathbb{P}}_n$, when \mathcal{F} is strong G-C. The proof is given in section 4 below.

THEOREM 10.15 *Let \mathcal{F} be a class of measurable functions, and let the multinomial vectors W_n in $\hat{\mathbb{P}}_n$ be independent of the data. Then the following are equivalent:*

(i) \mathcal{F} is strong Glivenko-Cantelli;

(ii) $\|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$;

(iii) $\mathbb{E}_W \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$;

(iv) For every $\eta > 0$, $\mathbb{P} \left(\|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}} > \eta \mid \mathcal{X}_n \right) \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$;

(v) For every $\eta > 0$, $\mathbb{P} \left(\|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \mid \mathcal{X}_n \right) \xrightarrow{\text{as}^*} 0$ and $P^*\|f - Pf\|_{\mathcal{F}} < \infty$, for some version of $\|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^*$.

10.3 A Simple Z-Estimator Master Theorem

Consider Z-estimation based on the estimating equation $\theta \mapsto \Psi_n(\theta) \equiv \mathbb{P}_n \psi_\theta$, where $\theta \in \Theta \subset \mathbb{R}^p$ and $x \mapsto \psi_\theta(x)$ is a measurable p -vector valued function for each θ . This is a special case of the more general Z-estimation approach discussed in section 2.2.5. Define the map $\theta \mapsto \Psi(\theta) \equiv P\psi_\theta$, and assume $\theta_0 \in \Theta$ satisfies $\Psi(\theta_0) = 0$. Let $\hat{\theta}_n$ be an approximate zero of Ψ_n , and let $\hat{\theta}_n^\circ$ be an approximate zero of the bootstrapped estimating equation $\theta \mapsto \Psi_n^\circ(\theta) \equiv \mathbb{P}_n^\circ \psi_\theta$, where \mathbb{P}_n° is either $\tilde{\mathbb{P}}_n$ of corollary 10.14—with ξ_1, \dots, ξ_n satisfying the conditions specified in the first paragraph of section 10.1.3 (the multiplier bootstrap)—or $\hat{\mathbb{P}}_n$ of theorem 10.15 (the multinomial bootstrap).

The goal of this section is to determine reasonably general conditions under which $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow Z$, where Z is mean zero normally distributed,

and $\sqrt{n}(\hat{\theta}_n^\circ - \hat{\theta}_n) \overset{P}{\rightsquigarrow} k_0 Z$. Here, we use $\overset{P}{\rightsquigarrow}$ to denote either $\overset{P}{\rightsquigarrow}_\xi$ or $\overset{P}{\rightsquigarrow}_W$ depending on which bootstrap is being used, and $k_0 = \tau/\mu$ for the multiplier bootstrap while $k_0 = 1$ for the multinomial bootstrap. One could also estimate the limiting variance rather than use the bootstrap, but there are many settings, such as least absolute deviation regression, where variance estimation may be more awkward than the bootstrap. For theoretical validation of the bootstrap approach, we have the following theorem, which is related to theorem 2.11 and which utilizes some of the bootstrap results of this chapter:

THEOREM 10.16 *Let $\Theta \subset \mathbb{R}^p$ be open, and assume $\theta_0 \in \Theta$ satisfies $\Psi(\theta_0) = 0$. Also assume the following:*

- (A) *For any sequence $\{\theta_n\} \in \Theta$, $\Psi(\theta_n) \rightarrow 0$ implies $\|\theta_n - \theta_0\| \rightarrow 0$;*
- (B) *The class $\{\psi_\theta : \theta \in \Theta\}$ is strong Glivenko-Cantelli;*
- (C) *For some $\eta > 0$, the class $\mathcal{F} \equiv \{\psi_\theta : \theta \in \Theta, \|\theta - \theta_0\| \leq \eta\}$ is Donsker and $P(\psi_\theta - \psi_{\theta_0})^2 \rightarrow 0$ as $\|\theta - \theta_0\| \rightarrow 0$;*
- (D) *$P\|\psi_{\theta_0}\|^2 < \infty$ and $\Psi(\theta)$ is differentiable at θ_0 with nonsingular derivative matrix V_{θ_0} ;*
- (E) *$\Psi_n(\hat{\theta}_n) = o_P(n^{-1/2})$ and $\Psi_n^\circ(\hat{\theta}_n^\circ) = o_P(n^{-1/2})$.*

Then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow Z \sim N(0, V_{\theta_0}^{-1} P[\psi_{\theta_0} \psi_{\theta_0}^T] (V_{\theta_0}^{-1})^T)$$

and $\sqrt{n}(\hat{\theta}_n^\circ - \hat{\theta}_n) \overset{P}{\rightsquigarrow} k_0 Z$.

Before giving the proof, we make a few comments about the conditions (A)–(E) of the theorem. Condition (A) is one of several possible identifiability conditions. Condition (B) is a sufficient condition, when combined with (A), to yield consistency of a zero of Ψ_n . This condition is generally reasonable to verify in practice. Condition (C) is needed for asymptotic normality of $\sqrt{n}(\hat{\theta}_n - \theta_0)$ and is also not hard to verify in practice. Condition (D) enables application of the delta method at the appropriate juncture in the proof below, and (E) is a specification of the level of approximation permitted in obtaining the zeros of the estimating equations. See exercise 10.5.5 below for a specific example of an estimation setting that satisfies these conditions.

Proof of theorem 10.16. By (B) and (E),

$$\|\Psi(\hat{\theta}_n)\| \leq \|\Psi_n(\hat{\theta}_n)\| + \sup_{\theta \in \Theta} \|\Psi_n(\theta) - \Psi(\theta)\| \leq o_P(1).$$

Thus $\hat{\theta}_n \xrightarrow{P} \theta_0$ by the identifiability condition (A). By assertion (ii) of either corollary 10.14 or theorem 10.15 (depending on which bootstrap is used),

condition (B) implies $\sup_{\theta \in \Theta} \|\Psi_n^\circ(\theta) - \Psi(\theta)\| \xrightarrow{\text{as}^*} 0$. Thus reapplication of conditions (A) and (E) yield $\hat{\theta}_n^\circ \xrightarrow{P} \theta_0$. Note that for the first part of the proof we will be using unconditional bootstrap results, while the associated conditional bootstrap results will be used only at the end.

By (C) and the consistency of $\hat{\theta}_0$, we have $\mathbb{G}_n \psi_{\hat{\theta}_n} - \mathbb{G}_n \psi_{\theta_0} \xrightarrow{P} 0$. Since (E) now implies that $\mathbb{G}_n \psi_{\hat{\theta}_n} = \sqrt{n}P(\psi_{\theta_0} - \psi_{\hat{\theta}_n}) + o_P(1)$, we can use the parametric (Euclidean) delta method plus differentiability of Ψ to obtain

$$(10.14) \quad \sqrt{n}V_{\theta_0}(\theta_0 - \hat{\theta}_n) + \sqrt{n}o_P(\|\hat{\theta}_n - \theta_0\|) = \mathbb{G}_n \psi_{\theta_0} + o_P(1).$$

Since V_{θ_0} is nonsingular, this yields that $\sqrt{n}\|\hat{\theta}_n - \theta_0\|(1 + o_P(1)) = O_P(1)$, and thus $\sqrt{n}(\hat{\theta}_n - \theta_0) = O_P(1)$. Combining this with (10.14), we obtain

$$(10.15) \quad \sqrt{n}(\hat{\theta}_n - \theta_0) = -V_{\theta_0}^{-1} \sqrt{n} \mathbb{P}_n \psi_{\theta_0} + o_P(1),$$

and thus $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow Z$ with the specified covariance.

The first part of condition (C) and theorem 2.6 imply that $\mathbb{G}_n^\circ \equiv k_0^{-1} \sqrt{n}(\mathbb{P}_n^\circ - \mathbb{P}_n) \rightsquigarrow \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ unconditionally, by arguments similar to those used in the (ii) \Rightarrow (i) part of the proof of theorem 10.4. Combining this with the second part of condition (C), we obtain $k_0 \mathbb{G}_n^\circ(\psi_{\hat{\theta}_n^\circ}) + \mathbb{G}_n(\psi_{\hat{\theta}_n^\circ}) - k_0 \mathbb{G}_n^\circ(\psi_{\theta_0}) - \mathbb{G}_n(\psi_{\theta_0}) \xrightarrow{P} 0$. Condition (E) now implies $\sqrt{n}P(\psi_{\theta_0} - \psi_{\hat{\theta}_n^\circ}) = \sqrt{n} \mathbb{P}_n^\circ \psi_{\theta_0} + o_P(1)$. Using similar arguments to those used in the previous paragraph, we obtain

$$\sqrt{n}(\hat{\theta}_n^\circ - \theta_0) = -V_{\theta_0}^{-1} \sqrt{n} \mathbb{P}_n^\circ \psi_{\theta_0} + o_P(1).$$

Combining with (10.15), we have

$$\sqrt{n}(\hat{\theta}_n^\circ - \hat{\theta}_n) = -V_{\theta_0}^{-1} \sqrt{n}(\mathbb{P}_n^\circ - \mathbb{P}_n) \psi_{\theta_0} + o_P(1).$$

The desired conditional bootstrap convergence now follows from theorem 2.6, part (ii) or part (iii) (depending on which bootstrap is used). \square

10.4 Proofs

Proof of theorem 10.10. Fix $\epsilon > 0$ and a compact $K \subset \mathbb{E}_0$. The proof stems from certain properties of separable Banach spaces which can be found in Megginson (1998). Specifically, the fact that every separable Banach space is isometrically isomorphic to a subspace of $C[0, 1]$, implies the existence of an isometric isomorphism $J_* : \mathbb{E}_0 \mapsto \mathbb{A}_0$, where \mathbb{A}_0 is a subspace of $C[0, 1]$. Since $C[0, 1]$ has a basis, we know by theorem 4.1.33 of Megginson (1998) that it also has the ‘‘approximation property.’’ This means by theorem 3.4.32 of Megginson (1998) that since $J_*(K)$ is compact, there exists a finite rank, bounded linear operator $T_* : \mathbb{A}_0 \mapsto \mathbb{A}_0$

such that $\sup_{y \in J_*(K)} \|T_*(y) - y\| < \epsilon$. Because T_* is finite rank, this means there exists elements $z_1, \dots, z_k \in \mathbb{A}_0 \subset C[0, 1]$ and bounded linear functionals $f_1^*, \dots, f_k^* : \mathbb{A}_0 \mapsto \mathbb{R}$ such that $T_*(y) = \sum_{j=1}^k z_j f_j^*(y)$. Note that since both J_* and $J_*^{-1} : \mathbb{A}_0 \mapsto \mathbb{E}_0$ are isometric isomorphisms, they are also both Lipschitz continuous with Lipschitz constant 1. This means that the map $x \mapsto \tilde{T}_\epsilon(x) \equiv J_*^{-1}(T_*(J_*(x)))$, with domain and range \mathbb{E}_0 , satisfies $\sup_{x \in K} \|\tilde{T}_\epsilon(x) - x\| < \infty$.

We now need to verify the existence of several important extensions. By Dugundji's extension theorem (theorem 10.9 above), there exists a continuous extension of J_* , $\tilde{J}_* : \mathbb{E} \mapsto \overline{\text{lin}} \mathbb{A}_0$. Also, by the Hahn-Banach extension theorem, there exist bounded linear extensions of f_1^*, \dots, f_k^* , $\tilde{f}_1^*, \dots, \tilde{f}_k^* : \overline{\text{lin}} \mathbb{A}_0 \mapsto \mathbb{R}$. Now let \tilde{J} denote the restriction of J_*^{-1} to the domain $\left\{ \sum_{j=1}^k z_j f_j^*(y) : y \in J_*(\mathbb{E}_0) \right\}$. Since \tilde{J} is Lipschitz continuous, as noted previously, we now have by theorem 10.17 below that there exists a Lipschitz continuous extension of \tilde{J} , $J : \overline{\text{lin}}(z_1, \dots, z_k) \mapsto \mathbb{E}$, with Lipschitz constant possibly larger than 1. Now define $x \mapsto f_j(x) \equiv \tilde{f}_j^*(\tilde{J}_*(x))$, for $j = 1, \dots, k$, and $x \mapsto T_\epsilon(x) \equiv J\left(\sum_{i=1}^k z_i f_i(x)\right)$; and note that T_ϵ is a continuous extension of \tilde{T}_ϵ . Now the quantities k , z_1, \dots, z_k , f_1, \dots, f_k , J and T_ϵ all satisfy the given requirements. \square

Proof of theorem 10.13. Since the processes $\mathbb{P}_n - P$, \mathbb{W}_n and $\tilde{\mathbb{W}}_n$ do not change when \mathcal{F} is replaced by $\dot{\mathcal{F}} \equiv \{f - Pf : f \in \mathcal{F}\}$, we can use for the index set either \mathcal{F} or $\dot{\mathcal{F}} \equiv \{f - Pf : f \in \mathcal{F}\}$ without changing the processes. Define also $\dot{F} \equiv \|f - Pf\|_{\mathcal{F}}^*$. We first need to show that the expectations and probabilities in (iii), (iv), (vii) and (viii) are well defined. Note that for fixed x_1, \dots, x_n and $a \equiv (a_1, \dots, a_n) \in \mathbb{R}^n$, the map

$$(a_1, \dots, a_n) \mapsto \left\| n^{-1} \sum_{i=1}^n a_i f(x_i) \right\|_{\dot{\mathcal{F}}} = \sup_{u \in B} |a^T u|,$$

where $B \equiv \{(f(x_1), \dots, f(x_n)) : f \in \dot{\mathcal{F}}\} \subset \mathbb{R}^n$. By the continuity of the map $(a, u) \mapsto |a^T u|$ and the separability of \mathbb{R} , this map is a measurable function even if the set B is not a Borel set. Thus the conditional expectations and conditional probabilities are indeed well defined.

(i) \Rightarrow (ii): Note that \mathcal{F} being P -G-C implies that $P\dot{F} < \infty$ by lemma 8.13. Because $\dot{\mathcal{F}}$ is G-C and ξ is trivially G-C, the desired result follows from corollary 9.26 (of section 9.3) and the fact that $\|\xi(f - Pf)\|_{\mathcal{F}}^* \leq |\xi| \times \dot{F}$ is integrable.

(ii) \Rightarrow (i): Since both $\text{sign}(\xi)$ and $\xi \cdot \dot{\mathcal{F}}$ are P -G-C, corollary 9.26 can be applied to verify that $\text{sign}(\xi) \cdot \xi \cdot \dot{\mathcal{F}} = |\xi| \cdot \dot{\mathcal{F}}$ is also P -G-C. We also have by lemma 8.13 that $P^*\dot{F} < \infty$ since $P|\xi| > 0$. Now we have for fixed X_1, \dots, X_n ,

$$(E\xi)\|\mathbb{P}_n\|_{\dot{\mathcal{F}}} = \|n^{-1} \sum_{i=1}^n (E\xi_i) \delta_{X_i}\|_{\dot{\mathcal{F}}} \leq E\xi\|\mathbb{W}_n\|_{\mathcal{F}},$$

and thus $(E\xi)E^*\|\mathbb{P}_n - P\|_{\mathcal{F}} \leq E^*\|\mathbb{W}_n\|_{\mathcal{F}}$. By applying theorem 9.28 twice, we obtain the desired result.

(ii) \Rightarrow (iii): Note first that (ii) immediately implies $P|\xi|\dot{F}(X) < \infty$ by lemma 8.13. Thus $P\dot{F} < \infty$ since $E|\xi| > 0$. Define $R_n \equiv n^{-1} \sum_{i=1}^n |\xi_i|\dot{F}(X_i)$, let B be the set of all infinite sequences $(\xi_1, X_1), (\xi_2, X_2), \dots$ such that $\|\mathbb{W}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ and $R_n \rightarrow E[|\xi|\dot{F}]$, and let $E_{\xi, \infty}$ be the expectation taken over the infinite sequence ξ_1, ξ_2, \dots holding the infinite sequence X_1, X_2, \dots fixed. Note that the set B has probability 1. Moreover, by the bounded convergence theorem,

$$\begin{aligned} \limsup E_{\xi} \|\mathbb{W}_n\|_{\mathcal{F}} 1\{R_n \leq K\} &= \limsup E_{\xi, \infty} \|\mathbb{W}_n\|_{\mathcal{F}} 1\{R_n \leq K\} 1\{B\} \\ &\leq E_{\xi, \infty} \limsup \|\mathbb{W}_n\|_{\mathcal{F}}^* 1\{R_n \leq K\} 1\{B\} \\ &= 0, \end{aligned}$$

outer almost surely, for any $K < \infty$. In addition, if we let $S_n = n^{-1} \sum_{i=1}^n |\xi_i|$, we have for any $0 < N < \infty$ that

$$\begin{aligned} E_{\xi} \|\mathbb{W}_n\|_{\mathcal{F}} 1\{R_n > K\} &\leq E_{\xi} R_n 1\{R_n > K\} \\ &\leq E_{\xi} [S_n 1\{R_n > K\}] \mathbb{P}_n \dot{F} \\ &\leq \frac{N(E|\xi|)[\mathbb{P}_n \dot{F}]^2}{K} + E_{\xi} [S_n 1\{S_n > N\}] \mathbb{P}_n \dot{F}, \end{aligned}$$

where the second-to-last inequality follows by symmetry. By exercise 10.5.4, the last line of the display has a $\limsup \leq N(P\dot{F})^2/K + (E|\xi|)^2/N$ outer almost surely. Thus, if we let $N = \sqrt{K}$ and allow $K \uparrow \infty$ slowly enough, we ensure that $\limsup_{n \rightarrow \infty} E_{\xi} \|\mathbb{W}_n\|_{\mathcal{F}} \rightarrow 0$, outer almost surely. Hence (iii) follows.

(iii) \Rightarrow (iv): This is obvious.

(iv) \Rightarrow (ii): (iv) clearly implies that $\|\mathbb{W}_n\|_{\mathcal{F}} \xrightarrow{P} 0$. Now lemma 8.16 implies that since the class $|\xi| \times \mathcal{F}$ has an integrable envelope, a version of $\|\mathbb{W}_n\|_{\mathcal{F}}^*$ must converge outer almost surely to a constant. Thus (ii) follows.

(ii) \Rightarrow (v): Assertion (ii) implies that there exists a measurable set B of infinite sequences $(\xi_1, X_1), (\xi_2, X_2), \dots$ with $P(B) = 1$ such that $\|\mathbb{W}_n\|_{\mathcal{F}}^* \rightarrow 0$ on B for some version of $\|\mathbb{W}_n\|_{\mathcal{F}}^*$. Thus by the bounded convergence theorem, we have for any $\eta > 0$ and almost all sequences X_1, X_2, \dots ,

$$\begin{aligned} \limsup_{n \rightarrow \infty} P(\|\mathbb{W}_n\|_{\mathcal{F}}^* > \eta) &= \limsup_{n \rightarrow \infty} E_{\xi, \infty} 1\{\|\mathbb{W}_n\|_{\mathcal{F}}^* > \eta\} 1\{B\} \\ &= E_{\xi, \infty} \limsup_{n \rightarrow \infty} 1\{\|\mathbb{W}_n\|_{\mathcal{F}}^* > \eta\} 1\{B\} \\ &= 0. \end{aligned}$$

Thus (v) follows.

(v) \Rightarrow (iv): This is obvious.

(ii) \Rightarrow (vi): Note that

$$(10.16) \quad \|\tilde{\mathbb{W}}_n - \mathbb{W}_n\|_{\mathcal{F}} \leq |\bar{\xi} - E\xi| \times |n^{-1} \sum_{i=1}^n \dot{F}(X_i)| + (E\xi) \|\mathbb{P}_n\|_{\mathcal{F}}.$$

Since (ii) \Rightarrow (i), $P^* \|f - Pf\|_{\mathcal{F}} < \infty$. Thus, since the centered weights $\xi_i - E\xi$ satisfy the conditions of the theorem as well as the original weights, the right side converges to zero outer almost surely. Hence $\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$, and (vi) follows.

(vi) \Rightarrow (ii): Since ξ_i can be replaced with $\xi_i - E\xi$ without changing $\tilde{\mathbb{W}}_n$, we will assume without loss of generality that $E\xi = 0$ (for this paragraph only). Accordingly, (10.16) implies

$$(10.17) \quad \|\tilde{\mathbb{W}}_n - \mathbb{W}_n\|_{\mathcal{F}} \leq |\bar{\xi} - E\xi| \times |n^{-1} \sum_{i=1}^n \dot{F}(X_i)|.$$

Thus (ii) will follow by the strong law of large numbers if we can show that $E\dot{F} < \infty$. Now let Y_1, \dots, Y_n be i.i.d. P independent of X_1, \dots, X_n , and let $\tilde{\xi}_1, \dots, \tilde{\xi}_n$ be i.i.d. copies of ξ_1, \dots, ξ_n independent of $X_1, \dots, X_n, Y_1, \dots, Y_n$. Define $\tilde{\mathbb{W}}_{2n} \equiv (2n)^{-1} \sum_{i=1}^n [(\xi_i - \bar{\xi})f(X_i) + (\tilde{\xi}_i - \bar{\xi})f(Y_i)]$ and $\tilde{\mathbb{W}}'_{2n} \equiv (2n)^{-1} \sum_{i=1}^n [(\tilde{\xi}_i - \bar{\xi})f(X_i) + (\xi_i - \bar{\xi})f(Y_i)]$, where $\bar{\xi} \equiv (2n)^{-1} \sum_{i=1}^n (\xi_i + \tilde{\xi}_i)$. Since both $\|\tilde{\mathbb{W}}_{2n}\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ and $\|\tilde{\mathbb{W}}'_{2n}\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$, we have that $\|\tilde{\mathbb{W}}_n - \tilde{\mathbb{W}}'_{2n}\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$. However,

$$\tilde{\mathbb{W}}_n - \tilde{\mathbb{W}}'_{2n} = \frac{1}{n} \sum_{i=1}^n \frac{\xi_i - \tilde{\xi}_i}{2} [f(X_i) - f(Y_i)],$$

and thus $\|n^{-1} \sum_{i=1}^n [f(X_i) - f(Y_i)]\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$ by the previously established equivalence between (i) and (ii) and the fact that the new weights $(\xi_i - \tilde{\xi}_i)/2$ satisfy the requisite conditions. Thus

$$E^* \dot{F} = E^* \|f(X) - Ef(Y)\|_{\mathcal{F}} \leq E^* \|f(X) - f(Y)\|_{\mathcal{F}} < \infty,$$

where the last inequality holds by lemma 8.13, and (ii) now follows.

(iii) \Rightarrow (vii): Since

$$E\xi \|\tilde{\mathbb{W}}_n - \mathbb{W}_n\|_{\mathcal{F}} \leq (E|\xi|) \|\mathbb{P}_n\|_{\mathcal{F}},$$

we have that $E\xi \|\tilde{\mathbb{W}}_n\|_{\mathcal{F}} \xrightarrow{\text{as}^*} 0$, because (iii) also implies (i).

(vii) \Rightarrow (viii): This is obvious.

(viii) \Rightarrow (vi): Since $\tilde{\mathbb{W}}_n$ does not change if the ξ_i s are replaced by $\xi_i - E\xi$, we will assume—as we did in the proof that (vi) \Rightarrow (ii)—that $E\xi = 0$ without loss of generality. By reapplication of (10.17) and the strong law of large numbers, we obtain that $\|\tilde{\mathbb{W}}_n\|_{\mathcal{F}} \xrightarrow{P} 0$. Since the class $|\xi| \times \mathcal{F}$ has an integrable envelope, reapplication of lemma 8.16 yields the desired result.

(vi) \Rightarrow (ix): The proof here is identical to the proof that (ii) \Rightarrow (v), after exchanging \mathbb{W}_n with $\tilde{\mathbb{W}}_n$.

(ix) \Rightarrow (viii): This is obvious. \square

Proof of theorem 10.15. The fact that the conditional expressions in assertions (iii) and (iv) are well defined can be argued as in the proof of theorem 10.13 above, and we omit the details.

(i) \Rightarrow (v): This follows from lemma 10.18 below since the vectors W_n/n are exchangeable and satisfy the other required conditions.

(v) \Rightarrow (iv): This is obvious.

(iv) \Rightarrow (i): For each integer $n \geq 1$, generate an infinite sequence of independent random row n -vectors $m_n^{(1)}, m_n^{(2)}, \dots$ as follows. Set $m_1^{(k)} = 1$ for all integers $k \geq 1$, and for each $n > 1$, generate an infinite sequence of i.i.d. Bernoullies $B_n^{(1)}, B_n^{(2)}, \dots$ with probability of success $1/n$, and set $m_n^{(k)} = [1 - B_n^{(k)}](m_{n-1}^{(k)}, 0) + B_n^{(k)}(0, \dots, 0, 1)$. Note that for each fixed n , $m_n^{(1)}, m_n^{(2)}, \dots$ are i.i.d. multinomial($1, n^{-1}, \dots, n^{-1}$) vectors. Independent of these random quantities, generate an infinite sequence U_1, U_2, \dots of i.i.d. Poisson random variables with mean 1, and set $N_n = \sum_{i=1}^n U_i$. Also make sure that all of these random quantities are independent of X_1, X_2, \dots . Without loss of generality assume $W_n = \sum_{i=1}^{N_n} m_n^{(i)}$ and define $\xi^{(n)} \equiv \sum_{i=1}^{N_n} m_n^{(i)}$. It is easy to verify that the W_n are indeed multinomial(n, n^{-1}, \dots, n^{-1}) vectors as claimed, and that $\xi_i^{(n)}, \dots, \xi_i^{(n)}$, where $(\xi_1^{(n)}, \dots, \xi_n^{(n)}) \equiv \xi^{(n)}$, are i.i.d. Poisson mean 1 random variables. Note also that these random weights are independent of X_1, X_2, \dots .

Let $\mathbb{W}_n \equiv n^{-1} \sum_{i=1}^n (\xi_i^{(n)} - 1)(\delta_{X_i} - P)$, and note that

$$\hat{\mathbb{P}}_n - \mathbb{P}_n - \mathbb{W}_n = n^{-1} \sum_{i=1}^n (W_{ni} - \xi_i^{(n)})(\delta_{X_i} - P).$$

Since the nonzero elements of $(W_{ni} - \xi_i^{(n)})$ all have the same sign by construction, we have that

$$\begin{aligned} E_{W, \xi} \|\hat{\mathbb{P}}_n - \mathbb{P}_n - \mathbb{W}_n\|_{\mathcal{F}} &\leq E_{W, \xi} \|n^{-1} \sum_{i=1}^n |W_{ni} - \xi_i^{(n)}| (\delta_{X_i} - P)\|_{\mathcal{F}} \\ &\leq \left(E \left| \frac{N_n - n}{n} \right| \right) [\mathbb{P}_n \dot{F} + P \dot{F}] \\ &\stackrel{\text{as}^*}{\rightarrow} 0, \end{aligned}$$

where the last inequality follows from the exchangeability result $E_{W, \xi} |W_{ni} - \xi_i^{(n)}| = E[|N_n - n|/n]$, $1 \leq i \leq n$, and the outer almost sure convergence to zero follows from the fact that $E[|N_n - n|/n] \leq n^{-1/2}$ combined with the moment conditions. In the forgoing, we have used $E_{W, \xi}$ to denote taking expectations over W_n and $\xi^{(n)}$ conditional on X_1, X_2, \dots . We have just

established that assertion (iv) holds in theorem 10.13 with weights $(\xi_1^{(n)} - 1, \dots, \xi_n^{(n)} - 1)$ that satisfy the necessary conditions. Thus \mathcal{F} is Glivenko-Cantelli.

(v) \Rightarrow (ii): Let $E_{W,\xi,\infty}$ be the expectation over the infinite sequences of the weights conditional on X_1, X_2, \dots . For fixed $\eta > 0$ and X_1, X_2, \dots , we have by the bounded convergence theorem

$$E_{W,\xi,\infty} \limsup_{n \rightarrow \infty} 1 \left\{ \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right\} = \limsup_{n \rightarrow \infty} E_{W,\xi,\infty} 1 \left\{ \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right\}.$$

But the right side $\rightarrow 0$ for almost all X_1, X_2, \dots by (v). This implies $1 \left\{ \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right\} \rightarrow 0$, almost surely. Now (ii) follows since η was arbitrary.

(ii) \Rightarrow (iii): Let B be the set of all sequences of weights and data for which $\|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* \rightarrow 0$. From (ii), we know that B is measurable, $P(B) = 1$ and, by the bounded convergence theorem, we have for every $\eta > 0$ and all X_1, X_2, \dots

$$\begin{aligned} 0 &= E_{W,\xi,\infty} \limsup_{n \rightarrow \infty} \left[1 \left\{ \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right\} 1\{B\} \right] \\ &= \limsup_{n \rightarrow \infty} E_{W,\xi,\infty} \left[1 \left\{ \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\mathcal{F}}^* > \eta \right\} 1\{B\} \right]. \end{aligned}$$

Since $P(B) = 1$, this last line implies (v), since η was arbitrary, and hence assertions (i) and (iv) also hold by the previously established equivalences. Fix $0 < K < \infty$, and note that the class $\dot{\mathcal{F}} \cdot 1\{\dot{F} \leq K\}$ is strong G-C by corollary 9.26. Now

$$\begin{aligned} E_W \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\dot{\mathcal{F}}} &\leq E_W \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\dot{\mathcal{F}} \cdot 1\{\dot{F} \leq K\}} \\ &\quad + E_W \left[(\hat{\mathbb{P}}_n + \mathbb{P}_n) \dot{F} 1\{\dot{F} > K\} \right] \\ &\leq E_W \|\hat{\mathbb{P}}_n - \mathbb{P}_n\|_{\dot{\mathcal{F}} \cdot 1\{\dot{F} \leq K\}} + 2\mathbb{P}_n[\dot{F} 1\{\dot{F} > K\}] \\ &\stackrel{\text{as}^*}{\rightarrow} 2P[\dot{F} 1\{\dot{F} > K\}], \end{aligned}$$

by assertion (iv). Since this last term can be made arbitrarily small by choosing K large enough, assertion (iii) follows.

(iii) \Rightarrow (iv): This is obvious. \square

THEOREM 10.17 *Let X, Z be metric spaces, with the dimension of X being finite, and let $Y \subset X$. For any Lipschitz continuous map $f : Y \mapsto Z$, there exists a Lipschitz continuous extension $F : X \mapsto Z$.*

Proof. This is a simplification of theorem 2 of Johnson, Lindenstrauss and Schechtman (1986), and the proof can be found therein. \square

LEMMA 10.18 *Let \mathcal{F} be a strong Glivenko-Cantelli class of measurable functions. For each n , let (M_{n1}, \dots, M_{nn}) be an exchangeable nonnegative*

random vector independent of X_1, X_2, \dots such that $\sum_{i=1}^n M_{ni} = 1$ and $\max_{1 \leq i \leq n} |M_{ni}| \xrightarrow{P} 0$. Then, for every $\eta > 0$,

$$P \left(\left\| \sum_{i=1}^n M_{ni} (\delta_{X_i} - P) \right\|_{\mathcal{F}}^* > \eta \mid \mathcal{X}_n \right) \xrightarrow{\text{as}^*} 0.$$

This is lemma 3.6.16 of VW, and we omit the proof.

10.5 Exercises

10.5.1. Show the following:

- (a) $\|\cdot\|_{2,1}$ is a norm on the space of real, square-integrable random variables.
- (b) For any real random variable ξ , and any $r > 2$, $(1/2)\|\xi\|_2 \leq \|\xi\|_{2,1} \leq (r/(r-2))\|\xi\|_r$. Hints: For the first inequality, show that $E|\xi|^2 \leq 2 \int_0^\infty P(|\xi| > u) u du \leq 2\|\xi\|_{2,1} \times \|\xi\|_2$. For the second inequality, show first that

$$\|\xi\|_{2,1} \leq a + \int_a^\infty \left(\frac{\|\xi\|_r^r}{x^r} \right)^{1/2} dx$$

for any $a > 0$.

10.5.2. Show that for any $p > 1$, and any real i.i.d. X_1, \dots, X_n with $E|X|^p < \infty$, we have

$$E \max_{1 \leq i \leq n} \frac{|X_i|}{n^{1/p}} \rightarrow 0,$$

as $n \rightarrow \infty$. Hint: Show first that for any $x > 0$,

$$\limsup_{n \rightarrow \infty} P \left(\max_{1 \leq i \leq n} \frac{|X_i|}{n^{1/p}} > x \right) \leq 1 - \exp \left(-\frac{E|X|^p}{x^p} \right).$$

10.5.3. Prove lemma 10.12.

10.5.4. Let ξ_1, \dots, ξ_n be i.i.d. nonnegative random variables with $E\xi < \infty$, and denote $S_n = n^{-1} \sum_{i=1}^n \xi_i$. Show that

$$\lim_{m \rightarrow \infty} \limsup_{n \rightarrow \infty} E[S_n 1\{S_n > m\}] = 0.$$

Hint: Theorem 9.28 may be useful here.

10.5.5. Assume that, given the covariate $Z \in \mathbb{R}^p$, Y is Bernoulli with probability of success $e^{\theta^T Z} / (1 + e^{\theta^T Z})$, where $\theta \in \Theta = \mathbb{R}^p$ and $E[Z Z^T]$ is positive definite. Assume that we observe an i.i.d. sample $(Y_1, Z_1), \dots, (Y_n, Z_n)$.

Z_n) generated from this model with true parameter $\theta_0 \in \mathbb{R}$. Show that the conditions of theorem 10.16 are satisfied for Z -estimators based on

$$\psi_\theta(y, z) = Z \left(Y - \frac{e^{\theta^T Z}}{1 + e^{\theta^T Z}} \right).$$

Note that one of the challenges here is the noncompactness of Θ .

10.6 Notes

Much of the material in section 10.1 is inspired by chapters 2.9 and 3.6 of VW, although the results for the weights $(\xi_1 - \bar{\xi}, \dots, \xi_n - \bar{\xi})$ and $(\xi_1/\bar{\xi} - 1, \dots, \xi_n/\bar{\xi} - 1)$ and the continuous mapping results are essentially new. The equivalence of assertions (i) and (ii) of theorem 10.4 is theorem 2.9.2 of VW, while the equivalence of (i) and (ii) of theorem 10.6 is theorem 2.9.6 of VW. Lemma 10.5 is lemma 2.9.5 of VW. Theorem 10.16 is an expansion of theorem 5.21 of van der Vaart (1998), and part (b) of exercise 10.5.1 is exercise 2.9.1 of VW.

11

Additional Empirical Process Results

In this chapter, we study several additional empirical process results that are useful but don't fall neatly into the framework of the other chapters. Because the contents of this chapter are somewhat specialized, some readers may want to skip it the first time they read the book. Although some of the results given herein will be used in later chapters, the results of this chapter are not really necessary for a philosophical understanding of the remainder of the book. On the other hand, this chapter does contain results and references that are useful for readers interested in the deeper potential of empirical process methods in genuinely hard statistical problems.

We first discuss bounding tail probabilities and moments of $\|\mathbb{G}_n\|_{\mathcal{F}}$. These results will be useful in chapter 14 for determining rates of convergence of M-estimators. We then discuss Donsker results for classes composed of sequences of functions and present several related statistical applications. After this, we discuss contiguous alternative probability measures P_n that get progressively closer to a fixed "null" probability measure P as n gets larger. These results will be useful in part III of the book, especially in chapter 18, where we discuss optimality of tests.

We then discuss weak convergence of sums of independent but not identically distributed stochastic processes which arise, for example, in clinical trials with non-independent randomization schemes such as biased coin designs (see, for example, Wei, 1978). We develop this topic in some depth, discussing both a central limit theorem and validity of a certain weighted bootstrap procedure. We also specialize these results to empirical processes based on i.i.d. data but with functions classes \mathcal{F}_n that change with n .

The final topic we cover is Donsker results for dependent observations. Here, our brief treatment is primarily meant to introduce the subject and point the interested reader toward the key results and references.

11.1 Bounding Moments and Tail Probabilities

We first consider bounding moments of $\|\mathbb{G}_n\|_{\mathcal{F}}^*$ under assumptions similar to the Donsker theorems of chapter 8. While these do not provide sharp bounds, the bounds are still useful for certain problems. We need to introduce some slightly modified entropy integrals. The first is based on a modified uniform entropy integral:

$$J^*(\delta, \mathcal{F}) \equiv \sup_Q \int_0^\delta \sqrt{1 + \log N(\epsilon \|F\|_{Q,2}, \mathcal{F}, L_2(Q))} d\epsilon,$$

where the supremum is over all finitely discrete probability measures Q with $\|F\|_{Q,2} > 0$. The only difference between this and the previously defined uniform entropy integral $J(\delta, \mathcal{F}, L_2)$, is the presence of the 1 under the radical. The following theorem, which we give without proof, is a subset of theorem 2.14.1 of VW:

THEOREM 11.1 *Let \mathcal{F} be a P -measurable class of measurable functions, with measurable envelope F . Then, for each $p \geq 1$,*

$$\|\|\mathbb{G}_n\|_{\mathcal{F}}^*\|_{P,p} \leq c_p J^*(1, \mathcal{F}) \|F\|_{P,2 \wedge p},$$

where the constant $c_p < \infty$ depends only on p .

We next provide an analogue of theorem 11.1 for bracketing entropy, based on the modified bracketing integral:

$$J_{[]}^*(\delta, \mathcal{F}) \equiv \int_0^\delta \sqrt{1 + \log N_{[]}(\epsilon \|F\|_{P,2}, \mathcal{F}, L_2(P))} d\epsilon.$$

The difference between this definition and the previously defined $J_{[]}(\delta, \mathcal{F}, L_2(P))$ is twofold: the presence of the 1 under the radical and a rescaling of ϵ by the factor $\|F\|_{P,2}$. The following theorem, a subset of theorem 2.14.2 of VW, is given without proof:

THEOREM 11.2 *Let \mathcal{F} be a class of measurable functions with measurable envelope F . Then*

$$\|\|\mathbb{G}_n\|_{\mathcal{F}}^*\|_{P,1} \leq c J_{[]}^*(1, \mathcal{F}) \|F\|_{P,2},$$

for some universal constant $c < \infty$.

For some problems, the previous moment bounds are sufficient. In other settings, more refined tail probability bounds are needed. To accomplish this, stronger assumptions are needed for the involved function classes. Recall the definition of pointwise measurable (PM) classes from section 8.2. The following tail probability results, the proofs of which are given in section 11.7 below, apply only to bounded and PM classes:

THEOREM 11.3 *Let \mathcal{F} be a pointwise measurable class of functions $f : \mathcal{X} \mapsto [-M, M]$, for some $M < \infty$, such that*

$$(11.1) \quad \sup_Q \log N(\epsilon, \mathcal{F}, L_2(Q)) \leq K \left(\frac{1}{\epsilon} \right)^W,$$

for all $0 < \epsilon \leq M$ and some constants $0 < W < 2$ and $K < \infty$, where the supremum is taken over all finitely discrete probability measures. Then

$$\| \|\mathbb{G}_n\|_{\mathcal{F}}^* \|_{\psi_2} \leq c,$$

for all $n \geq 1$, where $c < \infty$ depends only on K , W and M .

Examples of interesting function classes that satisfy the conditions of theorem 11.3 are bounded VC classes that are also PM, and the set of all non-decreasing distribution functions on \mathbb{R} . This follows from theorem 9.3 and lemma 9.11, since the class of distribution functions can be shown to be PM. To see this last claim, for each integer $m \geq 1$, let \mathcal{G}_m be the class of empirical distribution functions based on a sample of size m from the rationals union $\{-\infty, \infty\}$. It is not difficult to show that $\mathcal{G} \equiv \cup_{m \geq 1} \mathcal{G}_m$ is countable and that for each distribution function f , there exists a sequence $\{g_m\} \in \mathcal{G}$ such that $g_m(x) \rightarrow f(x)$, as $m \rightarrow \infty$, for each $x \in \mathbb{R}$.

THEOREM 11.4 *Let \mathcal{F} be a pointwise measurable class of functions $f : \mathcal{X} \mapsto [-M, M]$, for some $M < \infty$, such that*

$$(11.2) \quad N_{[]}(\epsilon, \mathcal{F}, L_2(P)) \leq K \left(\frac{1}{\epsilon} \right)^W,$$

for all $0 < \epsilon \leq M$ and some constants $K, W < \infty$. Then

$$\| \|\mathbb{G}_n\|_{\mathcal{F}}^* \|_{\psi_2} \leq c,$$

for all $n \geq 1$, where $c < \infty$ depends only on K , W and M .

An example of a function class that satisfies the conditions of theorem 11.4, are the Lipschitz classes of theorem 9.22 which satisfy condition 9.4, provided T is separable and $N(\epsilon, T, d) \leq K(1/\epsilon)^W$ for some constants $K, W < \infty$. This will certainly be true if (T, d) is a Euclidean space.

By lemma 8.1, if the real random variable X satisfies $\|X\|_{\psi_2} < \infty$, then the tail probabilities of X are “subgaussian” in the sense that $P(|X| >$

$x) \leq Ke^{-Cx^2}$ for some constants $K < \infty$ and $C > 0$. These results can be significantly refined under stronger conditions to yield more precise bounds on the constants. Some results along this line can be found in chapter 2.14 of VW. A very strong result applies to the empirical distribution function \mathbb{F}_n , where \mathcal{F} consists of left half-lines in \mathbb{R} :

THEOREM 11.5 *For any i.i.d. sample X_1, \dots, X_n with distribution F ,*

$$\mathbb{P} \left(\sup_{t \in \mathbb{R}} \sqrt{n} |\mathbb{F}_n(t) - F(t)| > x \right) \leq 2e^{-2x^2},$$

for all $x > 0$.

Proof. This is the celebrated result of Dvoretzky, Kiefer and Wolfowitz (1956), given in their lemma 2, as refined by Massart (1990) in his corollary 1. We omit the proof of their result but note that their result applies to the special case where F is continuous. We now show that it also applies when F may be discontinuous. Without loss of generality, assume that F has discontinuities, and let t_1, \dots, t_m be the locations of the discontinuities of F , where m may be infinity. Note that the number of discontinuities can be at most countable. Let p_1, \dots, p_m be the jump sizes of F at t_1, \dots, t_m . Now let U_1, \dots, U_n be i.i.d. uniform random variables independent of the X_1, \dots, X_n , and define new random variables

$$Y_i = X_i + \sum_{j=1}^m p_j [1\{X_i > t_j\} + U_i 1\{X_i = t_j\}],$$

$1 \leq i \leq n$. Define also the transformation $t \mapsto T(t) = t + \sum_{j=1}^m p_j 1\{t \geq t_j\}$; let \mathbb{F}_n^* be the empirical distribution of Y_1, \dots, Y_n ; and let F^* be the distribution of Y_1 . It is not hard to verify that

$$\begin{aligned} \sup_{t \in \mathbb{R}} |\mathbb{F}_n(t) - F(t)| &= \sup_{t \in \mathbb{R}} |\mathbb{F}_n^*(T(t)) - F^*(T(t))| \\ &\leq \sup_{s \in \mathbb{R}} |\mathbb{F}_n^*(s) - F^*(s)|, \end{aligned}$$

and the desired result now follows since F^* is continuous. \square

11.2 Sequences of Functions

Whether a countable class of functions \mathcal{F} is P -Donsker can be verified using the methods of chapters 9 and 10, but sometimes the special structure of certain countable classes simplifies the evaluation. This is true for certain classes composed of sequences of functions. The following is our first result in this direction:

THEOREM 11.6 *Let $\{f_i, i \geq 1\}$ be any sequence of measurable functions satisfying $\sum_{i=1}^{\infty} P(f_i - Pf_i)^2 < \infty$. Then the class*

$$\left\{ \sum_{i=1}^{\infty} c_i f_i : \sum_{i=1}^{\infty} |c_i| \leq 1 \text{ and the series converges pointwise} \right\}$$

is P -Donsker.

Proof. Since the class given in the conclusion of the theorem is the pointwise closure of the symmetric convex hull (see the comments given in section 9.1.1 just before theorem 9.4) of the class $\{f_i\}$, it is enough to verify that $\{f_i\}$ is Donsker, by theorem 9.29. To this end, fix $\epsilon > 0$ and define for each positive integer m , a partition $\{f_i\} = \cup_{i=1}^{m+1} \mathcal{F}_i$ as follows. For each $i = 1, \dots, m$, let \mathcal{F}_i consist of the single point f_i , and let $\mathcal{F}_{m+1} = \{f_{m_1}, f_{m+2}, \dots\}$. Since $\sup_{f, g \in \mathcal{F}_i} |\mathbb{G}_n(f-g)| = 0$ (trivially) for $i = 1, \dots, m$, we have, by Chebyshev's inequality,

$$\begin{aligned} P \left(\sup_i \sup_{f, g \in \mathcal{F}_i} |\mathbb{G}_n(f-g)| > \epsilon \right) &\leq P \left(\sup_{f \in \mathcal{F}_{m+1}} |\mathbb{G}_n f| > \frac{\epsilon}{2} \right) \\ &\leq \frac{4}{\epsilon^2} \sum_{i=m+1}^{\infty} P(f_i - Pf_i)^2. \end{aligned}$$

Since this last term can be made arbitrarily small by choosing m large enough, and since ϵ was arbitrary, the desired result now follows by theorem 2.1 via lemma 7.20. \square

When the sequence $\{f_i\}$ satisfies $Pf_i f_j = 0$ whenever $i \neq j$, theorem 11.6 can be strengthened as follows:

THEOREM 11.7 *Let $\{f_i, i \geq 1\}$ be any sequence of measurable functions satisfying $Pf_i f_j = 0$ for all $i \neq j$ and $\sum_{i=1}^{\infty} Pf_i^2 < \infty$. Then the class*

$$\left\{ \sum_{i=1}^{\infty} c_i f_i : \sum_{i=1}^{\infty} c_i^2 \leq 1 \text{ and the series converges pointwise} \right\}$$

is P -Donsker.

Proof. Since the conditions on $c \equiv (c_1, c_2, \dots)$ ensure $\sum_{i=1}^{\infty} c_i f_i \leq \sqrt{\sum_{i=1}^{\infty} f_i^2}$, we have by the dominated convergence theorem that pointwise converging sums also converge in $L_2(P)$. Now we argue that the class \mathcal{F} of all of these sequences is totally bounded in $L_2(P)$. This follows because \mathcal{F} can be arbitrarily closely approximated by a finite-dimensional set, since

$$P \left(\sum_{i>m} c_i f_i \right)^2 = \sum_{i>m} c_i^2 P f_i^2 \leq \sum_{i>m} P f_i^2 \rightarrow 0,$$

as $m \rightarrow \infty$. Thus the theorem is proved if we can show that the sequence \mathbb{G}_n , as a process indexed by \mathcal{F} , is asymptotically equicontinuous with respect to the $L_2(P)$ -seminorm. Accordingly, note that for any $f = \sum_{i=1}^{\infty} c_i f_i$, $g = \sum_{i=1}^{\infty} d_i f_i$, and integer $k \geq 1$,

$$\begin{aligned} |\mathbb{G}_n(f) - \mathbb{G}_n(g)|^2 &= \left| \sum_{i=1}^{\infty} (c_i - d_i) \mathbb{G}_n(f_i) \right|^2 \\ &\leq 2 \sum_{i=1}^k (c_i - d_i)^2 P f_i^2 \sum_{i=1}^k \frac{\mathbb{G}_n^2(f_i)}{P f_i^2} + 2 \sum_{i=k+1}^{\infty} (c_i - d_i)^2 \sum_{i=k+1}^{\infty} \mathbb{G}_n^2(f_i). \end{aligned}$$

Since, by assumption, $\|c - d\| \leq \|c\| + \|d\| \leq 2$ (here, $\|\cdot\|$ is the infinite-dimensional Euclidean norm), the above expression is bounded by

$$2\|f - g\|_{P,2}^2 \sum_{i=1}^k \frac{\mathbb{G}_n^2(f_i)}{P f_i^2} + 8 \sum_{i=k+1}^{\infty} \mathbb{G}_n^2(f_i).$$

Now take the supremum over all pairs of series f and g with $\|f - g\|_{P,2} < \delta$. Since also $E\mathbb{G}_n^2(f_i) \leq P f_i^2$, the expectation is bounded above by $2\delta^2 k + 8 \sum_{i=k+1}^{\infty} P f_i^2$. This quantity can now be made arbitrarily small by first choosing k large and then choosing δ small enough. \square

If the functions $\{f_i\}$ involved in the preceding theorem are an orthonormal sequence $\{\psi_i\}$ in $L_2(P)$, then the result can be reexpressed in terms of an *elliptical class* for a fixed sequence of constants $\{b_i\}$:

$$\mathcal{F} \equiv \left\{ \sum_{i=1}^{\infty} c_i \psi_i : \sum_{i=1}^{\infty} \frac{c_i^2}{b_i^2} \leq 1 \text{ and the series converges pointwise} \right\}.$$

More precisely, theorem 11.7 implies that \mathcal{F} is P -Donsker if $\sum_{i=1}^{\infty} b_i^2 < \infty$. Note that a sufficient condition for the stated pointwise convergence to hold at the point x for all $\{c_i\}$ satisfying $\sum_{i=1}^{\infty} c_i^2/b_i^2 \leq 1$ is for $\sum_{i=1}^{\infty} b_i^2 \psi_i^2(x) < \infty$. A very important property of an empirical process indexed by an elliptical class \mathcal{F} is the following:

$$(11.3) \quad \|\mathbb{G}_n\|_{\mathcal{F}}^2 = \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{\infty} c_i \mathbb{G}_n(\psi_i) \right|^2 = \sum_{i=1}^{\infty} b_i^2 \mathbb{G}_n^2(\psi_i).$$

In the central quantity, each function $f \in \mathcal{F}$ is represented by its series representation $\{c_i\}$. For the second equality, it is easy to see that the last term is an upper bound for the second term by the Cauchy-Schwartz inequality combined with the fact that $\sum_{i=1}^{\infty} c_i^2/b_i^2 \leq 1$. The next thing to note is that this maximum can be achieved by setting $c_i = b_i^2 \mathbb{G}_n(\psi_i) / \sqrt{\sum_{i=1}^{\infty} b_i^2 \mathbb{G}_n^2(\psi_i)}$.

An important use for elliptical classes is to characterize the limiting distribution of one- and two- sample Cramér-von Mises, Anderson-Darling,

and Watson statistics. We will now demonstrate this for both the Cramér-von Mises and Anderson-Darling statistics. For a study of the one-sample Watson statistic, see example 2.13.4 of VW. We now have the following key result, the proof of which we give in section 11.6. Note that the result (11.3) plays an important role in the proof.

THEOREM 11.8 *Let \mathbb{P}_n be the empirical distribution of an i.i.d. sample of uniform $[0, 1]$ random variables, let $\mathbb{G}_n = \sqrt{n}(\mathbb{P}_n - P)$ be the classical empirical process indexed by $t \in [0, 1]$, and let Z_1, Z_2, \dots be an i.i.d. sequence of standard normal deviates independent of \mathbb{P}_n . Also define the function classes*

$$\mathcal{F}_1 \equiv \left\{ \sum_{j=1}^{\infty} c_j \sqrt{2} \cos \pi j t : \sum_{j=1}^{\infty} c_j^2 \pi^2 j^2 \leq 1 \right\}$$

and

$$\mathcal{F}_2 \equiv \left\{ \sum_{j=1}^{\infty} c_j \sqrt{2} p_j(2t-1) : \sum_{j=1}^{\infty} c_j^2 j(j+1) \leq 1 \text{ and pointwise convergence} \right\},$$

where the functions $p_0(u) \equiv (1/2)\sqrt{2}$, $p_1(u) \equiv (1/2)\sqrt{6}u$, $p_2(u) \equiv (1/4)\sqrt{10} \times (3u^2 - 1)$, $p_3(u) \equiv (1/4)\sqrt{14}(5u^3 - 3u)$, and so on, are the orthonormalized Legendre polynomials in $L_2[-1, 1]$. Then the following are true:

(i) *The one-sample Cramér-von Mises statistic for uniform data satisfies*

$$\int_0^1 \mathbb{G}_n^2(t) dt = \|\mathbb{G}_n\|_{\mathcal{F}_1}^2 \rightsquigarrow \frac{1}{\pi^2} \sum_{j=1}^{\infty} \frac{Z_j^2}{j^2} \equiv T_1.$$

(ii) *The one-sample Anderson-Darling statistic for uniform data satisfies*

$$\int_0^1 \frac{\mathbb{G}_n^2(t)}{t(1-t)} dt = \|\mathbb{G}_n\|_{\mathcal{F}_2}^2 \rightsquigarrow \sum_{j=1}^{\infty} \frac{Z_j^2}{j(j+1)} \equiv T_2.$$

Theorem 11.8 applies to testing whether i.i.d. real data X_1, \dots, X_n comes from an arbitrary continuous distribution F . This is realized by replacing t with $F(x)$ throughout the theorem. A more interesting result can be obtained by applying the theorem to testing whether two samples have the same distribution. Let $\hat{F}_{n,j}$ be the empirical distribution of sample j of n_j i.i.d. real random variables, $j = 1, 2$, where the two samples are independent, and where $n = n_1 + n_2$. Let $\hat{F}_{n,0} \equiv (n_1 \hat{F}_{n,1} + n_2 \hat{F}_{n,2})/n$ be the pooled empirical distribution. The two-sample Cramér-von Mises statistic is

$$\hat{T}_1 \equiv \frac{n_1 n_2}{n_1 + n_2} \int_{-\infty}^{\infty} \left(\hat{F}_{n,1}(s) - \hat{F}_{n,2}(s) \right)^2 d\hat{F}_{n,0}(s),$$

while the two-sample Anderson-Darling statistics is

$$\hat{T}_2 \equiv \frac{n_1 n_2}{n_1 + n_2} \int_{-\infty}^{\infty} \frac{\left(\hat{F}_{n,1}(s) - \hat{F}_{n,2}(s) \right)^2}{\hat{F}_{n,0}(s) [1 - \hat{F}_{n,0}(s)]} d\hat{F}_{n,0}(s).$$

The proof of the following corollary is given in section 11.6:

COROLLARY 11.9 *Under the null hypothesis that the two samples come from the same continuous distribution F_0 , $\hat{T}_j \rightsquigarrow T_j$, as $n_1 \wedge n_2 \rightsquigarrow \infty$, for $j = 1, 2$.*

Since the limiting distributions do not depend on F_0 , critical values can be easily calculated by Monte Carlo simulation. Our own calculations resulted in critical values at the 0.05 level of 0.46 for T_1 and 2.50 for T_2 .

11.3 Contiguous Alternatives

For each $n \geq 1$, let X_{n1}, \dots, X_{nn} be i.i.d. random elements in a measurable space $(\mathcal{X}, \mathcal{A})$. Let P denote the common probability distribution under the “null hypothesis,” and let P_n be a “contiguous alternative hypothesis” distribution satisfying

$$(11.4) \quad \int \left[\sqrt{n}(dP_n^{1/2} - dP^{1/2}) - \frac{1}{2}h dP^{1/2} \right]^{1/2} \rightarrow 0,$$

as $n \rightarrow \infty$, for some measurable function $h : \mathcal{X} \mapsto \mathbb{R}$. The following lemma, which is part of lemma 3.10.11 of VW and which we give without proof, provides some properties for h :

LEMMA 11.10 *If the sequence of probability measures P_n satisfy (11.4), then necessarily $Ph = 0$ and $Ph^2 < \infty$.*

The following theorem gives very general weak convergence properties of the empirical process under the contiguous alternative P_n . Such weak convergence will be useful for studying efficiency of tests in chapter 18. This is theorem 3.10.12 of VW which we give without proof:

THEOREM 11.11 *Let \mathcal{F} be a P -Donsker class of measurable functions with $\|P\|_{\mathcal{F}} < \infty$, and assume the sequence of probability measures P_n satisfies (11.4). Then $\sqrt{n}(\mathbb{P}_n - P)$ converges under P_n in distribution in $\ell^\infty(\mathcal{F})$ to the process $f \mapsto \mathbb{G}(f) + Pfh$, where \mathbb{G} is a tight Brownian bridge. If, moreover, $\|P_n f^2\|_{\mathcal{F}} = O(1)$, then $\|\sqrt{n}(P_n - P)f - Pfh\|_{\mathcal{F}} \rightarrow 0$ and $\sqrt{n}(\mathbb{P}_n - P_n)$ converges under P_n in distribution to \mathbb{G} .*

We now present a bootstrap result for contiguous alternatives. For i.i.d. nonnegative random weights ξ_1, \dots, ξ_n with mean $0 < \mu < \infty$ and variance $0 < \tau^2 < \infty$, recall the bootstrapped empirical measures $\tilde{\mathbb{P}}_n f = n^{-1} \sum_{i=1}^n (\xi_i/\bar{\xi}) f(X_i)$ and $\tilde{\mathbb{G}}_n = \sqrt{n}(\mu/\tau)(\tilde{\mathbb{P}}_n - \mathbb{P}_n)$ from sections 2.2.3 and 10.1. We define several new symbols: $\xrightarrow{P_n}$ and $\xrightarrow[\sim]{P_n}$ denote convergence in probability and weak convergence, respectively, under the distribution P_n . In addition, $\hat{X}_n \xrightarrow[\sim]{P_n} X$ in a metric space \mathbb{D} denotes conditional bootstrap convergence in probability under P_n , i.e., $\sup_{g \in BL_1} |\mathbb{E}_M g(\hat{X}_n) - \mathbb{E} g(X)| \xrightarrow{P_n} 0$ and $\mathbb{E}_M g(\hat{X}_n)^* - \mathbb{E}_M g(\hat{X}_n)_* \xrightarrow{P_n} 0$, for all $g \in BL_1$, where BL_1 is the same bounded Lipschitz function space defined in section 2.2.3. Note that we require the weights ξ_1, \dots, ξ_n to have the same distribution and independence from X_{n1}, \dots, X_{nn} under both P_n and P .

THEOREM 11.12 *Let \mathcal{F} be a P -Donsker class of measurable functions, let P_n satisfy (11.4), and assume*

$$(11.5) \quad \lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} P_n(f - Pf)^2 1\{|f - Pf| > M\} = 0$$

for all $f \in \mathcal{F}$. Also let ξ_1, \dots, ξ_n be i.i.d. nonnegative random variables, independent of X_{n1}, \dots, X_{nn} , with mean $0 < \mu < \infty$, variance $0 < \tau^2 < \infty$, and with $\|\xi_1\|_{2,1} < \infty$. Then $\tilde{\mathbb{G}}_n \xrightarrow[\xi]{P_n} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$ and $\tilde{\mathbb{G}}_n$ is asymptotically measurable.

Proof. Let $\eta_i \equiv \tau^{-1}(\xi_i - \mu)$, $i = 1, \dots, n$, and note that

$$(11.6) \quad \begin{aligned} \tilde{\mathbb{G}}_n &= n^{-1/2}(\mu/\tau) \sum_{i=1}^n (\xi_i/\bar{\xi} - 1) \delta_{X_i} \\ &= n^{-1/2}(\mu/\tau) \sum_{i=1}^n (\xi_i/\bar{\xi} - 1) (\delta_{X_i} - P) \\ &= n^{-1/2} \sum_{i=1}^n \eta_i (\delta_{X_i} - P) \\ &\quad + \left(\frac{\mu}{\bar{\xi}} - 1\right) n^{-1/2} \sum_{i=1}^n \eta_i (\delta_{X_i} - P) \\ &\quad + \left(\frac{\mu}{\tau}\right) \left(\frac{\mu}{\bar{\xi}} - 1\right) n^{-1/2} \sum_{i=1}^n (\delta_{X_i} - P). \end{aligned}$$

Since \mathcal{F} is P -Donsker, we also have that $\dot{\mathcal{F}} \equiv \{f - Pf : f \in \mathcal{F}\}$ is P -Donsker. Thus by the unconditional multiplier central limit theorem (theorem 10.13), we have that $\eta \cdot \mathcal{F}$ is also P -Donsker. Now, by the fact that $\|P(f - Pf)\|_{\mathcal{F}} = 0$ (trivially) combined with theorem 11.11, both

$n^{-1/2} \sum_{i=1}^n \eta_i (\delta_{X_i} - P) \overset{P_n}{\rightsquigarrow} \mathbb{G}(f)$ and $n^{-1/2} \sum_{i=1}^n (\delta_{X_i} - P) \overset{P_n}{\rightsquigarrow} \mathbb{G}(f) + P(f - Pf)h$ in $\ell^\infty(\mathcal{F})$. The reason the first limiting process has mean zero is because η is independent of X and thus $P\eta(f - Pf)h = 0$ for all $f \in \mathcal{F}$. Thus the last two terms in (11.6) $\overset{P_n}{\rightsquigarrow} 0$ and $\tilde{\mathbb{G}}_n \overset{P_n}{\rightsquigarrow} \mathbb{G}$ in $\ell^\infty(\mathcal{F})$. This now implies the unconditional asymptotic tightness and desired asymptotic measurability of $\tilde{\mathbb{G}}_n$.

By the same kinds of arguments we used in the proof of theorem 10.4, all we need to verify now is that all finite dimensional collections $f_1, \dots, f_m \in \mathcal{F}$ converge under P_n in distribution, conditional on the data, to the appropriate limiting Gaussian process. Accordingly, let $Z_i = (f_1(X_i) - Pf_1, \dots, f_m(X_i) - Pf_m)'$, $i = 1, \dots, n$. What we need to show is that $n^{-1/2} \sum_{i=1}^n \eta_i Z_i$ converges weakly under P_n , conditional on the Z_1, \dots, Z_n , to a mean zero Gaussian process with variance $\Sigma \equiv PZ_1Z_1'$. By lemma 10.5, we are done if we can verify that $n^{-1} \sum_{i=1}^n Z_iZ_i' \overset{P_n}{\rightsquigarrow} \Sigma$ and $\max_{1 \leq i \leq n} \|Z_i\|/\sqrt{n} \overset{P_n}{\rightsquigarrow} 0$.

By the assumptions of the theorem,

$$\limsup_{n \rightarrow \infty} \{E_n(M) \equiv P_n[Z_1Z_1'1\{\|Z_1\| > M\}]\} \rightarrow 0$$

as $M \rightarrow \infty$. Note also that (11.4) can be shown to imply that $P_n h \rightarrow Ph$, as $n \rightarrow \infty$, for any bounded h (this verification is saved as an exercise). Thus, for any $M < \infty$, $P_n Z_1Z_1'1\{\|Z_1\| \leq M\}$ converges to $PZ_1Z_1'1\{\|Z_1\| \leq M\}$. Since M is arbitrary, this convergence continues to hold if M is replaced by a sequence M_n going to infinity slowly enough. Accordingly,

$$\begin{aligned} n^{-1} \sum_{i=1}^n Z_iZ_i' &= n^{-1} \sum_{i=1}^n Z_iZ_i'1\{\|Z_i\| > M_n\} \\ &\quad + n^{-1} \sum_{i=1}^n Z_iZ_i'1\{\|Z_i\| \leq M_n\} \\ &\overset{P_n}{\rightsquigarrow} PZ_1Z_1', \end{aligned}$$

as $n \rightarrow \infty$. Now we also have

$$\begin{aligned} \max_{1 \leq i \leq n} \frac{\|Z_i\|}{\sqrt{n}} &= \sqrt{\max_{1 \leq i \leq n} \frac{\|Z_i\|^2}{n}} \\ &\leq \sqrt{\max_{1 \leq i \leq n} \frac{\|Z_i\|^2}{n} 1\{\|Z_i\| \leq M\} + \frac{1}{n} \sum_{i=1}^n \|Z_i\|^2 1\{\|Z_i\| > M\}}. \end{aligned}$$

The first term under the last square root sign $\overset{P_n}{\rightsquigarrow} 0$ trivially, while the expectation under P_n of the second term, $E_n(M)$, goes to zero as $n \rightarrow \infty$ and $M \rightarrow \infty$ sufficiently slowly with n , as argued previously. Thus $\max_{1 \leq i \leq n} \|Z_i\|/\sqrt{n} \overset{P_n}{\rightsquigarrow} 0$, and the proof is complete. \square

11.4 Sums of Independent but not Identically Distributed Stochastic Processes

In this section, we are interested in deriving the limiting distribution of sums of the form $\sum_{i=1}^{m_n} f_{ni}(\omega, t)$, where the real-valued stochastic processes $\{f_{ni}(\omega, t), t \in T, 1 \leq i \leq m_n\}$, for all integers $n \geq 1$, are independent within rows on the probability space (Ω, \mathcal{A}, P) for some index set T . In addition to a central limit theorem, we will present a multiplier bootstrap result to aid in inference. An example using these techniques will be presented in the upcoming case studies II chapter. Throughout this section, function arguments or subscripts will sometimes be suppressed for notational clarity.

11.4.1 Central Limit Theorems

The notation and set-up are similar to that found in Pollard (1990) and Kosorok (2003). A slightly more general approach to the same question can be found in chapter 2.11 of VW. Part of the generality in VW is the ability to utilize bracketing entropy in addition to uniform entropy for establishing tightness. An advantage of Pollard's approach, on the other hand, is that total boundedness of the index set T is a conclusion rather than a condition. Both approaches have their merits and appear to be roughly equally useful in practice.

We need to introduce a few measurability conditions which are different from but related to conditions introduced in previous chapters. The first condition is *almost measurable Suslin*: Call a triangular array $\{f_{ni}(\omega, t), t \in T\}$ almost measurable Suslin (AMS) if for all integers $n \geq 1$, there exists a Suslin topological space $T_n \subset T$ with Borel sets \mathcal{B}_n such that

$$(i) \quad \mathbb{P}^* \left(\sup_{t \in T} \inf_{s \in T_n} \sum_{i=1}^{m_n} (f_{ni}(\omega, s) - f_{ni}(\omega, t))^2 > 0 \right) = 0,$$

(ii) For $i = 1 \dots m_n$, $f_{ni} : \Omega \times T_n \mapsto \mathbb{R}$ is $\mathcal{A} \times \mathcal{B}_n$ -measurable.

The second condition is stronger yet seems to be more easily verified in applications: Call a triangular array of processes $\{f_{ni}(\omega, t), t \in T\}$ *separable* if for every integer $n \geq 1$, there exists a countable subset $T_n \subset T$ such that

$$\mathbb{P}^* \left(\sup_{t \in T} \inf_{s \in T_n} \sum_{i=1}^{m_n} (f_{ni}(\omega, s) - f_{ni}(\omega, t))^2 > 0 \right) = 0.$$

The following lemma shows that separability implies AMS:

LEMMA 11.13 *If the triangular array of stochastic processes $\{f_{ni}(\omega, t), t \in T\}$ is separable, then it is AMS.*

Proof. The discrete topology applied to T_n makes it into a Suslin topology by countability, with resulting Borel sets \mathcal{B}_n . For $i = 1 \dots m_n$, $f_{ni} : \Omega \times T_n \mapsto \mathbb{R}$ is $\mathcal{A} \times \mathcal{B}_n$ -measurable since, for every $\alpha \in \mathbb{R}$,

$$\{(\omega, t) \in \Omega \times T_n : f_{ni}(\omega, t) > \alpha\} = \bigcup_{s \in T_n} \{(\omega, s) : f_{ni}(\omega, s) > \alpha\},$$

and the right-hand-side is a countable union of $\mathcal{A} \times \mathcal{B}_n$ -measurable sets. \square

The forgoing measurable Suslin condition is closely related to the definition given in Example 2.3.5 of VW, while the definition of separable arrays is similar in spirit to the definition of separable stochastic processes given in the discussion preceding lemma 7.2 in section 7.1 above. The modifications of these definitions presented in this section have been made to accommodate nonidentically distributed arrays for a broad scope of statistical applications. However, finding the best possible measurability conditions was not the primary goal.

We need the following definition of *manageability* (Definition 7.9 of Pollard, 1990, with minor modification). First, for any set $A \in \mathbb{R}^m$, let $D_m(x, A)$ be the packing number for the set A at Euclidean distance x , i.e., the largest k such that there exist k points in A with the smallest Euclidean distance between any two distinct points being greater than x . Also let $\mathcal{F}_{n\omega} \equiv \{[f_{n1}(\omega, t), \dots, f_{nm_n}(\omega, t)] \in \mathbb{R}^{m_n} : t \in T\}$; and for any vectors $u, v \in \mathbb{R}^m$, $u \odot v \in \mathbb{R}^m$ is the pointwise product and $\|\cdot\|$ denotes Euclidean distance. A triangular array of processes $\{f_{ni}(\omega, t)\}$ is manageable, with respect to the envelopes $F_n(\omega) \equiv [F_{n1}(\omega), \dots, F_{nm_n}(\omega)] \in \mathbb{R}^{m_n}$, if there exists a deterministic function λ (the *capacity bound*) for which

- (i) $\int_0^1 \sqrt{\log \lambda(x)} dx < \infty$,
- (ii) there exists $N \subset \Omega$ such that $P^*(N) = 0$ and for each $\omega \notin N$,

$$D_{m_n}(x \|\alpha \odot F_n(\omega)\|, \alpha \odot \mathcal{F}_{n\omega}) \leq \lambda(x),$$

for $0 < x \leq 1$, all vectors $\alpha \in \mathbb{R}^{m_n}$ of nonnegative weights, all $n \geq 1$, and where λ does not depend on ω or n .

We now state a minor modification of Pollard’s Functional Central Limit Theorem for the stochastic process sum

$$X_n(\omega, t) \equiv \sum_{i=1}^{m_n} [f_{ni}(\omega, t) - \mathbb{E}f_{ni}(\cdot, t)].$$

The modification is the inclusion of a sufficient measurability requirement which was omitted in Pollard’s (1990) version of the theorem.

THEOREM 11.14 *Suppose the triangular array $\{f_{ni}(\omega, t), t \in T\}$ consists of independent processes within rows, is AMS, and satisfies:*

- (A) the $\{f_{ni}\}$ are manageable, with envelopes $\{F_{ni}\}$ which are also independent within rows;
- (B) $H(s, t) = \lim_{n \rightarrow \infty} \mathbb{E} X_n(s) X_n(t)$ exists for every $s, t \in T$;
- (C) $\limsup_{n \rightarrow \infty} \sum_{i=1}^{m_n} \mathbb{E}^* F_{ni}^2 < \infty$;
- (D) $\lim_{n \rightarrow \infty} \sum_{i=1}^{m_n} \mathbb{E}^* F_{ni}^2 \mathbb{1}\{F_{ni} > \epsilon\} = 0$, for each $\epsilon > 0$;
- (E) $\rho(s, t) = \lim_{n \rightarrow \infty} \rho_n(s, t)$, where

$$\rho_n(s, t) \equiv \left(\sum_{i=1}^{m_n} \mathbb{E} |f_{ni}(\cdot, s) - f_{ni}(\cdot, t)|^2 \right)^{1/2},$$

exists for every $s, t \in T$, and for all deterministic sequences $\{s_n\}$ and $\{t_n\}$ in T , if $\rho(s_n, t_n) \rightarrow 0$ then $\rho_n(s_n, t_n) \rightarrow 0$.

Then

- (i) T is totally bounded under the ρ pseudometric;
- (ii) X_n converges weakly on $\ell^\infty(T)$ to a tight mean zero Gaussian process X concentrated on $UC(T, \rho)$, with covariance $H(s, t)$.

The proof is given in Kosorok (2003) who relies on chapter 10 of Pollard (1990) for some of the steps. We omit the details. We will use this theorem when discussing function classes changing with n , later in this chapter, as well as in one of the case studies of chapter 15. One can think of this theorem as a Lindeberg central limit theorem for stochastic processes, where condition (D) is a modified Lindeberg condition.

Note that the manageability condition (A) is an entropy condition quite similar to the BUEI condition of chapter 9. Pollard (1990) discussed several methods and preservation results for establishing manageability, including *bounded pseudodimension classes* which are very close in spirit to VC-classes of functions (see chapter 4 of Pollard): The set $\mathcal{F}_n \subset \mathbb{R}^n$ has pseudodimension of at most V if, for every point $t \in \mathbb{R}^{V+1}$, no *proper coordinate projection* of \mathcal{F}_n can surround t . A proper coordinate projection \mathcal{F}_n^k is obtained by choosing a subset i_1, \dots, i_k of indices $1, \dots, m_n$, where $k \leq m_n$, and then letting $\mathcal{F}_n^k = \{f_{ni_1}(t_1), \dots, f_{ni_k}(t_k) : (t_1, \dots, t_k) \in \mathbb{R}^k\}$.

It is not hard to verify that if $f_{n1}(t), \dots, f_{nm_n}(t)$ are always monotone increasing functions in t , then the resulting \mathcal{F}_n has pseudodimension 1. This happens because all two-dimensional projections always form monotone increasing trajectories, and thus can never surround any point in \mathbb{R}^2 . The proof is almost identical to the proof of lemma 9.10. By theorem 4.8 of Pollard (1990), every triangular array of stochastic processes for which \mathcal{F}_n has pseudodimension bounded above by $V < \infty$, for all $n \geq 1$, is manageable. As verified in the following theorem, complicated manageable classes can be built up from simpler manageable classes:

THEOREM 11.15 *Let f_{n1}, \dots, f_{nm_n} and g_{n1}, \dots, g_{nm_n} be manageable arrays with respective index sets T and U and with respective envelopes F_{n1}, \dots, F_{nm_n} and G_{n1}, \dots, G_{nm_n} . Then the following are true:*

- (i) $\{f_{n1}(t) + g_{n1}(u), \dots, f_{nm_n}(t) + g_{nm_n}(u) : (t, u) \in T \times U\}$, is manageable with envelopes $F_{n1} + G_{n1}, \dots, F_{nm_n} + G_{nm_n}$;
- (ii) $\{f_{n1}(t) \wedge g_{n1}(u), \dots, f_{nm_n}(t) \wedge g_{nm_n}(u) : (t, u) \in T \times U\}$ is manageable with envelopes $F_{n1} + G_{n1}, \dots, F_{nm_n} + G_{nm_n}$;
- (iii) $\{f_{n1}(t) \vee g_{n1}(u), \dots, f_{nm_n}(t) \vee g_{nm_n}(u) : (t, u) \in T \times U\}$ is manageable with envelopes $F_{n1} + G_{n1}, \dots, F_{nm_n} + G_{nm_n}$;
- (iv) $\{f_{n1}(t)g_{n1}(u), \dots, f_{nm_n}(t)g_{nm_n}(u) : (t, u) \in T \times U\}$ is manageable with envelopes $F_{n1}G_{n1}, \dots, F_{nm_n}G_{nm_n}$.

Proof. For any vectors $x_1, x_2, y_1, y_2 \in \mathbb{R}^{m_n}$, it is easy to verify that $\|x_1 \square y_1 - x_2 \square y_2\| \leq \|x_1 - x_2\| + \|y_1 - y_2\|$, where \square is any one of the operations \wedge, \vee , or $+$. Thus

$$N_{m_n}(\epsilon \|\alpha \odot F_n + G_n\|, \alpha \odot \mathcal{F}_n \square \mathcal{G}_n) \leq N_{m_n}(\epsilon \|\alpha \odot F_n\|, \alpha \odot \mathcal{F}_n) \\ \times N_{m_n}(\epsilon \|\alpha \odot G_n\|, \alpha \odot \mathcal{G}_n),$$

for any $0 < \epsilon \leq 1$ and all vectors $\alpha \in \mathbb{R}^{m_n}$ of nonnegative weights, where N_{m_n} denotes the covering number version of D_{m_n} , $\mathcal{G}_n \equiv \{g_{n1}(u), \dots, g_{nm_n}(u) : u \in U\}$, and $\mathcal{F}_n \square \mathcal{G}_n$ has the obvious interpretation. Thus parts (i)–(iii) of the theorem follow by the relationship between packing and covering numbers discussed in section 8.1.2 in the paragraphs preceding theorem 8.4.

Proving part (iv) requires a slightly different approach. Let x_1, x_2, y_1, y_2 be any vectors in \mathbb{R}^{m_n} , and choose the vectors $\tilde{x}, \tilde{y} \in \mathbb{R}^{m_n}$ with nonnegative components so that both $\tilde{x} - [x_1 \vee x_2 \vee (-x_1) \vee (-x_2)]$ and $\tilde{y} - [y_1 \vee y_2 \vee (-y_1) \vee (-y_2)]$ have only nonnegative components. It is not hard to verify that $\|x_1 \odot y_1 - x_2 \odot y_2\| \leq \tilde{y} \|x_1 - x_2\| + \tilde{x} \|y_1 - y_2\|$. From this, we can deduce

$$N_{m_n}(2\epsilon \|\alpha \odot F_n \odot G_n\|, \alpha \odot \mathcal{F}_n \odot \mathcal{G}_n) \\ \leq N_{m_n}(\epsilon \|\alpha \odot F_n \odot G_n\|, \alpha \odot G_n \odot \mathcal{F}_n) \\ \times N_{m_n}(\epsilon \|\alpha \odot F_n \odot G_n\|, \alpha \odot F_n \odot \mathcal{G}_n) \\ = N_{m_n}(\epsilon \|\alpha' \odot F_n\|, \alpha' \odot \mathcal{F}_n) \times N_{m_n}(\epsilon \|\alpha'' \odot G_n\|, \alpha'' \odot \mathcal{G}_n),$$

for any $0 < \epsilon \leq 1$ and all vectors $\alpha \in \mathbb{R}^{m_n}$ of nonnegative weights, where $\alpha' \equiv \alpha \odot G_n$ and $\alpha'' \equiv \alpha \odot F_n$. Since capacity bounds do not depend on the nonnegative weight vector (either α, α' or α''), part (iv) now follows, and the theorem is proved. \square

11.4.2 Bootstrap Results

We now present a weighted bootstrap for inference about the limiting process X of theorem 11.14. The basic idea shares some similarities with the wild bootstrap (Praestgaard and Wellner, 1993). Let $Z \equiv \{z_i, i \geq 1\}$ be a sequence of random variables satisfying

- (F) The $\{z_i\}$ are independent and identically distributed, on the probability space $\{\Omega_z, \mathcal{A}_z, \Pi_z\}$, with mean zero and variance 1.

Denote $\mu_{ni}(t) \equiv E f_{ni}(\cdot, t)$, and let $\hat{\mu}_{ni}(t)$ be estimators of $\mu_{ni}(t)$. The weighted bootstrapped process we propose for inference is

$$\hat{X}_{n\omega}(t) \equiv \sum_{i=1}^{m_n} z_i [f_{ni}(\omega, t) - \hat{\mu}_{ni}(\omega, t)],$$

which is defined on the product probability space $\{\Omega, \mathcal{A}, \Pi\} \times \{\Omega_z, \mathcal{A}_z, \Pi_z\}$, similar to what was done in chapter 9 for the weighted bootstrap in the i.i.d. case. What is unusual about this bootstrap is the need to estimate the μ_{ni} terms: this need is a consequence of the terms being non-identically distributed.

The proposed method of inference is to resample \hat{X}_n , using many realizations of z_1, \dots, z_{m_n} , to approximate the distribution of X_n . The following theorem gives us conditions under which this procedure is asymptotically valid:

THEOREM 11.16 *Suppose the triangular array $\{f_{ni}\}$ satisfies the conditions of theorem 11.14 and the sequence $\{z_i, i \geq 1\}$ satisfies condition (F) above. Suppose also that the array of estimators $\{\hat{\mu}_{ni}(\omega, t), t \in T, 1 \leq i \leq m_n, n \geq 1\}$ is AMS and satisfies the following:*

(G) $\sup_{t \in T} \sum_{i=1}^{m_n} [\hat{\mu}_{ni}(\omega, t) - \mu_{ni}(t)]^2 = o_P(1)$;

(H) the stochastic processes $\{\hat{\mu}_{ni}(\omega, t)\}$ are manageable with envelopes $\{\hat{F}_{ni}(\omega)\}$;

(I) $k \vee \sum_{i=1}^{m_n} [\hat{F}_{ni}(\omega)]^2$ converges to k in outer probability as $n \rightarrow \infty$, for some $k < \infty$.

Then the conclusions of Theorem 11.14 obtain, \hat{X}_n is asymptotically measurable, and $\hat{X}_n \xrightarrow[Z]{P} X$.

The main idea of the proof is to first study the conditional limiting distribution of $\hat{X}_{n\omega}(t) \equiv \sum_{i=1}^{m_n} z_i [f_{ni}(\omega, t) - \hat{\mu}_{ni}(t)]$, and then show that the limiting result is unchanged after replacing μ_{ni} with $\hat{\mu}_{ni}$. The first step is summarized in the following theorem:

THEOREM 11.17 *Suppose the triangular array $\{f_{ni}\}$ satisfies the conditions of theorem 11.14 and the sequence $\{z_i, i \geq 1\}$ satisfies condition (F) above. Then the conclusions of Theorem 11.14 obtain, \tilde{X}_n is asymptotically measurable, and $\tilde{X}_n \xrightarrow[Z]{P} X$.*

An interesting step in the proof of this theorem is verifying that manageability of the triangular array $z_1 f_{n1}, \dots, z_{m_n} f_{nm_n}$ follows directly from manageability of f_{n1}, \dots, f_{nm_n} . We now demonstrate this. For vectors $u \in \mathbb{R}^{m_n}$, let $|u|$ denote pointwise absolute value and $\text{sign}(u)$ denote pointwise sign. Now, for any nonnegative $\alpha \in \mathbb{R}^{m_n}$,

$$\begin{aligned} D_{m_n}(x \|\alpha \odot |z_n| \odot F_n(\omega)\|, \alpha \odot z_n \odot \mathcal{F}_{n\omega}) \\ = D_{m_n}(x \|\tilde{\alpha} \odot F_n(\omega)\|, \tilde{\alpha} \odot \text{sign}(z_n) \odot \mathcal{F}_{n\omega}) \\ = D_{m_n}(x \|\tilde{\alpha} \odot F_n(\omega)\|, \tilde{\alpha} \odot \mathcal{F}_{n\omega}), \end{aligned}$$

where $z_n \equiv \{z_1, \dots, z_{m_n}\}^T$, since the absolute value of the $\{z_i\}$ can be absorbed into the α to make $\tilde{\alpha}$ and since any coordinate change of sign does not effect the geometry of $\mathcal{F}_{n\omega}$. Thus the foregoing triangular array is manageable with envelopes $\{|z_i|F_{ni}(\omega)\}$. The remaining details of the proofs of both theorems 11.16 and 11.17, which we omit here, can be found in Kosorok (2003).

11.5 Function Classes Changing with n

We now return to the i.i.d. empirical process setting where the i.i.d. observations X_1, X_2, \dots are drawn from a measurable space $\{\mathcal{X}, \mathcal{A}\}$, with probability measure P . What is new, however, is that we allow the function class to depend on n . Specifically, we assume the function class has the form $\mathcal{F}_n \equiv \{f_{n,t} : t \in T\}$, where the functions $x \mapsto f_{n,t}(x)$ are indexed by a fixed T but are allowed to change with sample size n . Note that this trivially includes the standard empirical process set-up with an arbitrary but fixed function class \mathcal{F} by setting $T = \mathcal{F}$ and $f_{n,t} = t$ for all $n \geq 1$ and $t \in \mathcal{F}$. The approach we take is to specialize the results of section 11.4 after replacing manageability with a bounded uniform entropy integral condition. An alternative approach which can utilize either uniform or bracketing entropy is given in section 2.11.3 of VW, but we do not pursue this second approach here.

Let $X_n(t) \equiv n^{-1/2} \sum_{i=1}^n (f_{n,t}(X_i) - P f_{n,t})$, for all $t \in T$, and let F_n be an envelope for \mathcal{F}_n . We say that the sequence \mathcal{F}_n is AMS if for all $n \geq 1$, there exists a Suslin topological space $T_n \subset T$ with Borel sets \mathcal{B}_n such that

$$(11.7) \quad P^* \left(\sup_{t \in T} \sup_{s \in T_n} |f_{n,s}(X_1) - f_{n,t}(X_1)| > 0 \right) = 0$$

and $f_{n,\cdot} : \mathcal{X} \times T_n \mapsto \mathbb{R}$ is $\mathcal{A} \times \mathcal{B}_{n-}$ -measurable. Moreover, the sequence \mathcal{F}_n is said to be separable if, for all $n \geq 1$, there exists a countable subset $T_n \subset T$ such that (11.7) holds. The arguments in the proof of lemma 11.13 verify that separability implies AMS as is true for the more general setting of section 11.4.

We also require the following bounded uniform entropy integral condition:

$$(11.8) \quad \limsup_{n \rightarrow \infty} \sup_Q \int_0^1 \sqrt{\log N(\epsilon \|F_n\|_{Q,2}, \mathcal{F}_n, L_2(Q))} d\epsilon < \infty,$$

where, for each $n \geq 1$, \sup_Q is the supremum taken over all finitely discrete probability measures Q with $\|F_n\|_{Q,2} > 0$. We are now ready to present the following functional central limit theorem:

THEOREM 11.18 *Suppose \mathcal{F}_n is AMS and the following hold:*

- (A) \mathcal{F}_n satisfies (11.8) with envelop F_n ;
- (B) $H(s, t) = \lim_{n \rightarrow \infty} \mathbb{E} X_n(s) X_n(t)$ for every $s, t \in T$;
- (C) $\limsup_{n \rightarrow \infty} \mathbb{E}^* F_n^2 < \infty$;
- (D) $\lim_{n \rightarrow \infty} \mathbb{E}^* F_n^2 1\{F_n > \epsilon \sqrt{n}\} = 0$, for each $\epsilon > 0$;
- (E) $\rho(s, t) = \lim_{n \rightarrow \infty} \rho_n(s, t)$, where $\rho_n(s, t) \equiv \sqrt{\mathbb{E}[f_{n,s}(X_1) - f_{n,t}(X_2)]^2}$, exists for every $s, t \in T$, and for all deterministic sequences $\{s_n\}$ and $\{t_n\}$ in T , if $\rho(s_n, t_n) \rightarrow 0$ then $\rho_n(s_n, t_n) \rightarrow 0$.

Then

- (i) T is totally bounded under the ρ pseudometric;
- (ii) X_n converges weakly in $\ell^\infty(T)$ to a tight, mean zero Gaussian process X concentrated on $UC(T, \rho)$, with covariance $H(s, t)$.

Proof. The proof consists in showing that the current setting is just a special case of theorem 11.14. Specifically, we let $f_{ni}(t) = f_{n,t}(X_i)$ and $m_n = n$, and we study the array $\{f_{ni}(t), t \in T\}$. First, it is clear that \mathcal{F}_n being AMS implies that $\{f_{ni}(t), t \in T\}$ is AMS. Now let

$$\tilde{\mathcal{F}}_n \equiv \{[f_{n1}(t), \dots, f_{nn}(t)] \in \mathbb{R}^n : t \in T\}$$

and $\tilde{F}_n \equiv [\tilde{F}_{n1}, \dots, \tilde{F}_{nn}]$, where $\tilde{F}_{ni} \equiv F_n(X_i)/\sqrt{n}$; and note that for any $\alpha \in \mathbb{R}^n$,

$$D_n(\epsilon \|\alpha \odot \tilde{F}_n\|, \alpha \odot \tilde{\mathcal{F}}_n) \leq D(\epsilon \|F_n\|_{Q_{\alpha,2}}, \mathcal{F}_n, L_2(Q_\alpha)),$$

where $Q_\alpha \equiv (n\|\alpha\|)^{-1} \sum_{i=1}^n \alpha_i^2 \delta_{X_i}$ is a finitely discrete probability measure. Thus, by the relationship between packing and covering numbers given in section 8.1.2, we have that if we let

$$\lambda(x) = \limsup_{n \rightarrow \infty} \sup_Q N(x\|F_n\|_{Q,2}/2, \mathcal{F}_n, L_2(Q)),$$

where \sup_Q is taken over all finitely discrete probability measures, then condition 11.8 implies that

$$D_n(\epsilon\|\alpha \odot \tilde{F}_n\|, \alpha \odot \tilde{\mathcal{F}}_n) \leq \lambda(\epsilon),$$

for all $0 < \epsilon \leq 1$, all vectors $\alpha \in \mathbb{R}^n$ of nonnegative weights, and all $n \geq 1$; and that $\int_0^1 \sqrt{\log \lambda(\epsilon)} d\epsilon < \infty$. Note that without loss of generality, we can set $D_n(\cdot, \cdot) = 1$ whenever $\|\alpha \odot \tilde{F}_n\| = 0$ and let $\lambda(1) = 1$, and thus the foregoing arguments yield that condition (11.8) implies manageability of the triangular array $\{f_{n1}(t), \dots, f_{nn}(t), t \in T\}$.

Now the remaining conditions of the theorem can easily be shown to imply conditions (B) through (E) of theorem 11.14 for the new triangular array and envelope vector \tilde{F}_n . Hence the desired results follow from theorem 11.14. \square

The following lemma gives us an important example of a sequence of classes \mathcal{F}_n that satisfies condition (11.8):

LEMMA 11.19 *For fixed index set T , let $\mathcal{F}_n = \{f_{n,t} : t \in T\}$ be a VC class of measurable functions with VC-index V_n and integrable envelope F_n , for all $n \geq 1$, and assume $\sup_{n \geq 1} V_n = V < \infty$. Then the sequence \mathcal{F}_n satisfies condition (11.8).*

Proof. By theorem 9.3, there exists a universal constant K depending only on V such that

$$N(\epsilon\|F_n\|_{Q,2}, \mathcal{F}_n, L_2(Q)) \leq K \left(\frac{1}{\epsilon}\right)^{r(V-1)},$$

for all $0 < \epsilon \leq 1$. Note that we have extended the range of ϵ to include 1, but this presents no difficulty since $N(\|F_n\|_{Q,2}, \mathcal{F}_n, L_2(Q)) = 1$ always holds by the definition of an envelope function. The desired result now follows since $V_n \leq V$ for all $n \geq 1$. \square

As a simple example, let $T \subset \mathbb{R}$ and assume that $f_{n,t}(x)$ is always monotone increasing in t . Then \mathcal{F}_n always has VC-index 2 by lemma 9.10, and hence lemma 11.19 applies. Thus (11.8) holds. This particular situation will apply later in section 14.5.2 when we study the weak convergence of a certain monotone density estimator. This condition (11.8) is quite similar to the BUEI condition for fixed function classes, and most of the preservation results of section 9.1.2 will also apply. The following proposition, the proof of which is saved as an exercise, is one such preservation result:

PROPOSITION 11.20 Let \mathcal{G}_n and \mathcal{H}_n be sequences of classes of measurable functions with respective envelope sequences G_n and H_n , where condition (11.8) is satisfied for the sequences $(\mathcal{F}_n, F_n) = (\mathcal{G}_n, G_n)$ and $(\mathcal{F}_n, F_n) = (\mathcal{H}_n, H_n)$. Then condition (11.8) is also satisfied for the sequence of classes $\mathcal{F}_n = \mathcal{G}_n \cdot \mathcal{H}_n$ (consisting of all pairwise products) and the envelope sequence $F_n = G_n H_n$.

We now present a weighted bootstrap result for this setting. Let $Z \equiv \{z_i, i \geq 1\}$ be a sequence of random variables satisfying

- (F) The $\{z_i\}$ are positive, i.i.d. random variables which are independent of the data X_1, X_2, \dots and which have mean $0 < \mu < \infty$ and variance $0 < \tau^2 < \infty$.

The weighted bootstrapped process we propose for use here is

$$\hat{X}_n(t) \equiv \frac{\mu}{\tau} n^{-1/2} \sum_{i=1}^n \left(\frac{z_i}{\bar{z}_n} - 1 \right) f_{n,t}(X_i),$$

where $\bar{z}_n \equiv n^{-1} \sum_{i=1}^n z_i$. The following theorem tells us that this is a valid bootstrap procedure:

THEOREM 11.21 Suppose the class of functions \mathcal{F}_n , with envelope F_n , satisfies the conditions of theorem 11.18 and the sequence $\{z_i, i \geq 1\}$ satisfies condition (F) above. Then the conclusions of theorem 11.18 obtain, \hat{X}_n is asymptotically measurable, and $\hat{X}_n \xrightarrow[Z]{P} X$.

Proof. Note that

$$\begin{aligned} \frac{\tau}{\mu} \hat{X}_n(t) &= n^{-1/2} \sum_{i=1}^n \left(\frac{z_i}{\bar{z}_n} - 1 \right) (f_{n,t}(X_i) - P f_{n,t}) \\ &= n^{-1/2} \sum_{i=1}^n \left(\frac{z_i}{\mu} - 1 \right) (f_{n,t}(X_i) - P f_{n,t}) \\ &\quad + n^{-1/2} \left(\frac{\mu}{\bar{z}_n} - 1 \right) \sum_{i=1}^n \left(\frac{z_i}{\mu} - 1 \right) (f_{n,t}(X_i) - P f_{n,t}) \\ &\quad + n^{-1/2} \left(\frac{\mu}{\bar{z}_n} - 1 \right) \sum_{i=1}^n (f_{n,t}(X_i) - P f_{n,t}) \\ &= A_n(t) + B_n(t) + C_n(t). \end{aligned}$$

By theorem 11.18, $\mathbb{G}_n f_{n,t} = O_P^{*T}(1)$, where $O_P^{*T}(1)$ denotes a term bounded in outer probability uniformly over T . Since theorem 11.18 is really a special case of theorem 11.14, we have by theorem 11.17 that $n^{-1/2} \sum_{i=1}^n \frac{\mu}{\tau} \left(\frac{z_i}{\mu} - 1 \right) \times (f_{n,t}(X_i) - P f_{n,t}) \xrightarrow[Z]{P} X$, since $\frac{\mu}{\tau} \left(\frac{z_i}{\mu} - 1 \right)$ has mean zero and variance 1.

Hence $(\mu/\tau)A_n$ is asymptotically measurable and $(\mu/\tau)A_n \overset{P}{\rightsquigarrow} X$. Moreover, since $\bar{z}_n \xrightarrow{\text{as*}} \mu$, we now have that both $B_n = o_P^{*T}(1)$ and $C_n = o_P^{*T}(1)$, where $o_P^{*T}(1)$ denotes a term going to zero outer almost surely uniformly over T . The desired results now follow. \square

11.6 Dependent Observations

In the section, we will review a number of empirical process results for dependent observations. A survey of recent results on this subject is *Empirical Process Techniques for Dependent Data*, edited by Dehling, Mikosch and Sørensen (2002); and a helpful general reference on theory for dependent observations is *Dependence in Probability and Statistics: A Survey of Recent Results*, edited by Eberlein and Taqqu (1986). Our focus here will be on strongly mixing stationary sequences (see Bradley, 1986). For the interested reader, a few results for non-stationary dependent sequences can be found in Andrews (1991), while several results for long range dependent sequences can be found in Dehling and Taqqu (1989), Yu (1994) and Wu (2003), among other references.

Let X_1, X_2, \dots be a *stationary* sequence of possibly dependent random variables on a probability space (Ω, \mathcal{D}, Q) , and let \mathcal{M}_a^b be the σ -field generated by X_a, \dots, X_b . By stationary, we mean that for any set of positive integers m_1, \dots, m_k , the joint distribution of $X_{m_1+j}, X_{m_2+j}, \dots, X_{m_k+j}$ is unchanging for all integers $j \geq -m_1 + 1$. The sequence $\{X_i, i \geq 1\}$ is *strongly mixing* (also α -mixing) if

$$\alpha(k) \equiv \sup_{m \geq 1} \{|P(A \cap B) - P(A)P(B)| : A \in \mathcal{M}_1^m, B \in \mathcal{M}_{m+k}^\infty\} \rightarrow 0,$$

as $k \rightarrow \infty$, and it is *absolutely regular* (also β -mixing) if

$$\beta(k) \equiv E \sup_{m \geq 1} \{|P(B | \mathcal{M}_1^m) - P(B)| : B \in \mathcal{M}_{m+k}^\infty\} \rightarrow 0,$$

as $k \rightarrow \infty$. Other forms of mixing include ρ -mixing, ϕ -mixing, ψ -mixing and $*$ -mixing (see Definition 3.1 of Dehling and Philipp, 2002). Note that the stronger notion of m -dependence, where observations more than m lags apart are independent, implies that $\beta(k) = 0$ for all $k > m$ and therefore also implies absolute regularity. It is also known that absolute regularity implies strong mixing (see section 3.1 of Dehling and Philipp, 2002). Hereafter, we will restrict our attention to β -mixing sequences since these will be the most useful for our purposes.

We now present several empirical process Donsker and bootstrap results for absolutely regular stationary sequences. Let the values of X_1 lie in a Polish space \mathcal{X} with distribution P , and let \mathbb{G}_n be the empirical measure for the first n observations of the sequence, i.e., $\mathbb{G}_n f = n^{1/2} \sum_{i=1}^n (f(X_i) - Pf)$,

for any measurable $f : \mathcal{X} \mapsto \mathbb{R}$. We now present the following bracketing central limit theorem:

THEOREM 11.22 *Let X_1, X_2, \dots be a stationary sequence in a Polish space with marginal distribution P , and let \mathcal{F} be a class of functions in $L_2(P)$. Suppose there exists a $2 < p < \infty$ such that*

$$(a) \sum_{k=1}^{\infty} k^{2/(p-2)} \beta(k) < \infty, \text{ and}$$

$$(b) J_{[]}(\infty, \mathcal{F}, L_p(P)) < \infty.$$

Then $\mathbb{G}_n \rightsquigarrow \mathbb{H}$ in $\ell^\infty(\mathcal{F})$, where \mathbb{H} is a tight, mean zero Gaussian process with covariance

$$(11.9) \quad \Gamma(f, g) \equiv \lim_{k \rightarrow \infty} \sum_{i=1}^{\infty} \text{cov}(f(X_k), g(X_i)), \text{ for all } f, g \in \mathcal{F}.$$

Proof. The result follows through condition 2 of theorem 5.2 of Dedecker and Louhichi (2002), after noting that their condition 2 can be shown to be implied by conditions (a) and (b) above, via arguments contained in section 4.3 of Dedecker and Louhichi (2002). We omit the details. \square

We next present a result for VC classes \mathcal{F} . In this case, we need to address the issue of measurability with some care. For what follows, let \mathcal{B} be the σ -field of the measurable sets on \mathcal{X} . The class of functions \mathcal{F} is *permissible* if it can be indexed by some set T , i.e., $\mathcal{F} = \{f(\cdot, t) : t \in T\}$ (T could potentially be \mathcal{F} for this purpose), in such a way that the following holds:

$$(a) \quad T \text{ is a Suslin metric space with Borel } \sigma\text{-field } \mathcal{B}(T),$$

$$(b) \quad f(\cdot, \cdot) \text{ is a } \mathcal{B} \times \mathcal{B}(T)\text{-measurable function from } \mathbb{R} \times T \text{ to } \mathbb{R}.$$

Note that this definition is similar to the almost measurable Suslin condition of section 11.5. We now have the following theorem:

THEOREM 11.23 *Let X_1, X_2, \dots be a stationary sequence in a Polish space with marginal distribution P , and let \mathcal{F} be a class of functions in $L_2(P)$. Suppose there exists a $2 < p < \infty$ such that*

$$(a) \quad \lim_{k \rightarrow \infty} k^{2/(p-2)} (\log k)^{2(p-1)/(p-2)} \beta(k) = 0, \text{ and}$$

$$(b) \quad \mathcal{F} \text{ is permissible, VC, and has envelope } F \text{ satisfying } P^* F^p < \infty.$$

Then $\mathbb{G}_n \rightsquigarrow \mathbb{H}$ in $\ell^\infty(\mathcal{F})$, where \mathbb{H} is as defined in theorem 11.22 above.

Proof. This theorem is essentially theorem 2.1 of Arcones and Yu (1994), and the proof, under slightly different measurability conditions, can be found therein. The measurability issues, including the sufficiency of the permissibility condition, are addressed in the appendix of Yu (1994). We omit the details. \square

Now we consider the bootstrap. A problem with the usual nonparametric bootstrap is that samples of X_1, X_2, \dots randomly drawn with replacement will be independent and will lose the dependency structure of the stationary sequence. Hence the usual bootstrap will generally not work. A modified bootstrap, the *moving blocks bootstrap* (MBB), was independently introduced by Künsch (1989) and Liu and Singh (1992) to address this problem. The method works as follows for a stationary sample X_1, \dots, X_n : For a chosen block length $b \leq n$, extend the sample by defining $X_{n+i}, \dots, X_{n+b-1} = X_1, \dots, X_b$ and let k be the smallest integer such that $kb \geq n$. Now define blocks (as row vectors) $B_i = (X_i, X_{i+1}, \dots, X_{i+b-1})$, for $i = 1, \dots, n$, and sample from the B_i s with replacement to obtain k blocks $B_1^*, B_2^*, \dots, B_k^*$. The bootstrapped sample X_1^*, \dots, X_n^* consists of the first n observations from the row vector (B_1^*, \dots, B_k^*) . The bootstrapped empirical measure indexed by the class \mathcal{F} is then defined as

$$\mathbb{G}_n^* f \equiv n^{-1/2} \sum_{i=1}^n (f(X_i^*) - \mathbb{P}_n f),$$

for all $f \in \mathcal{F}$, where $\mathbb{P}_n f \equiv n^{-1} \sum_{i=1}^n f(X_i)$ is the usual empirical probability measure (except that the data are now potentially dependent).

For now, we will assume that X_1, X_2, \dots are real-valued, although the results probably could be extended to general Polish-valued random variables. MBB bootstrap consistency has been established for bracketing entropy in Bühlmann (1995), although the entropy requirements are much stronger than those of theorem 11.22 above, and also for VC classes in Radulović (1996). Other interesting, related references are Naik-Nimbalkar and Rajarshi (1994) and Peligrad (1998), among others. We conclude this section by presenting Radulović's (1996) theorem 1 (slightly modified to address measurability) without proof:

THEOREM 11.24 *Let X_1, X_2, \dots be a stationary sequence of real random variables with marginal distribution P , and let \mathcal{F} be a class of functions in $L_2(P)$. Also assume that X_1^*, \dots, X_n^* are generated by the MBB procedure with block size $b(n) \rightarrow \infty$, as $n \rightarrow \infty$, and that there exists $2 < p < \infty$, $q > p/(p-2)$, and $0 < \rho < (p-2)/[2(p-1)]$ such that*

$$(a) \limsup_{k \rightarrow \infty} k^q \beta(k) < \infty,$$

$$(b) \mathcal{F} \text{ is permissible, VC, and has envelope } F \text{ satisfying } P^* F^p < \infty, \text{ and}$$

$$(c) \limsup_{n \rightarrow \infty} n^{-\rho} b(n) < \infty.$$

Then $\mathbb{G}_n^* \xrightarrow{*} \mathbb{H}$ in $\ell^\infty(\mathcal{F})$, where \mathbb{H} is as defined in theorem 11.22 above.

11.7 Proofs

Proofs of theorems 11.3 and 11.4. Consider the class $\mathcal{F}' \equiv 1/2 + \mathcal{F}/(2M)$, and note that all functions $f \in \mathcal{F}'$ satisfy $0 \leq f \leq 1$. Moreover, $\|\mathbb{G}_n\|_{\mathcal{F}} = 2M\|\mathbb{G}_n\|_{\mathcal{F}'}$. Thus, if we prove the results for \mathcal{F}' , we are done.

For theorem 11.3, the condition (11.1) is also satisfied if we replace \mathcal{F} with \mathcal{F}' . This can be done without changing W , although we may need to change K to some $K' < \infty$. Theorem 2.14.10 of VW now yields that

$$(11.10) \quad \mathbb{P}^*(\|\mathbb{G}_n\|_{\mathcal{F}'} > t) \leq C e^{Dt^{U+\delta}} e^{-2t^2},$$

for every $t > 0$ and $\delta > 0$, where $U = W(6 - W)/(2 + W)$ and C and D depend only on K' , W , and δ . Since $0 < W < 2$, it can be shown that $0 < U < 2$. Accordingly, choose a $\delta > 0$ so that $U + \delta < 2$. Now, it can be shown that there exists a $C^* < \infty$ and $K^* > 0$ so that (11.10) $\leq C^* e^{-K^* t^2}$, for every $t > 0$. Theorem 11.3 now follows by lemma 8.1.

For theorem 11.4, the condition (11.2) implies the existence of a K' so that $N_{\square}(\epsilon, \mathcal{F}', L_2(P)) \leq (K'/\epsilon)^V$, for all $0 < \epsilon < K'$, where V is the one in (11.2). Now theorem 2.14.9 of VW yields that

$$(11.11) \quad \mathbb{P}^*(\|\mathbb{G}_n\|_{\mathcal{F}'} > t) \leq C t^V e^{-2t^2},$$

for every $t > 0$, where the constant C depends only on K' and V . Thus there exists a $C^* < \infty$ and $K^* > 0$ such that (11.11) $\leq C^* e^{-K^* t^2}$, for every $t > 0$. The desired result now follows. \square

Proof of theorem 11.8. For the proof of result (i), note that the cosines are bounded, and thus the series defining \mathcal{F}_1 is automatically pointwise convergent by the discussion prior to theorem 11.8. Now, the Cramér-von Mises statistic is the square of the $L_2[0, 1]$ norm of the function $t \mapsto \mathbb{G}_n(t) \equiv \mathbb{G}_n \mathbf{1}\{X \leq t\}$. Since the functions $\{g_i \equiv \sqrt{2} \sin \pi j t : j = 1, 2, \dots\}$ form an orthonormal basis for $L_2[0, 1]$, Parseval's formula tells us that the integral of the square of any function in $L_2[0, 1]$ can be replaced by the sum of the squares of the Fourier coefficients. This yields:

$$(11.12) \quad \int_0^1 \mathbb{G}_n^2(t) dt = \sum_{i=1}^{\infty} \left[\int_0^1 \mathbb{G}_n(t) g_i(t) dt \right]^2 = \sum_{i=1}^{\infty} \left[\mathbb{G}_n \int_0^X g_i(t) dt \right]^2.$$

But since $\int_0^x g_i(t) dt = -(\pi i)^{-1} \sqrt{2} \cos \pi i x$, the last term in (11.12) becomes $\sum_{i=1}^{\infty} \mathbb{G}_n^2(\sqrt{2} \cos \pi i X) / (\pi i)^2$. Standard methods can now be used to establish that this converges weakly to the appropriate limiting distribution. An alternative proof can be obtained via the relation (11.3) and the fact that \mathcal{F}_1 is an elliptical class and hence Donsker. Now (i) follows.

For result (ii), note that the orthonormalized Legendre polynomials can be obtained by applying the Gram-Schmidt procedure to the functions $1, u, u^2, \dots$. By problem 2.13.1 of VW, the orthonormalized Legendre

polynomials satisfy the differential equations $(1 - u^2)p_j''(u) - 2up_j'(u) = -j(j+1)p_j(u)$, for all $u \in [-1, 1]$ and integers $j \geq 1$. By change of variables, followed by partial integration and use of this differential identity, we obtain

$$\begin{aligned} & 2 \int_0^1 p_i'(2t-1)p_j'(2t-1)t(1-t)dt \\ &= \frac{1}{4} \int_{-1}^1 p_i'(u)p_j'(u)(1-u^2)du \\ &= -\frac{1}{4} \int_{-1}^1 p_i(u) [p_j''(u)(1-u^2) - 2up_j'(u)] du \\ &= \frac{1}{4}j(j+1)1\{i=j\}. \end{aligned}$$

Thus the functions $2\sqrt{2}p_j'(2t-1)\sqrt{t(1-t)}/\sqrt{j(j+1)}$, with j ranging over the positive integers, form an orthonormal basis for $L_2[0, 1]$. By Parseval's formula, we obtain

$$\begin{aligned} \int_0^1 \frac{\mathbb{G}_n^2(t)}{t(1-t)} dt &= \sum_{j=1}^{\infty} \left[\int_0^1 \mathbb{G}_n(t)p_j'(2t-1)dt \right]^2 \frac{8}{j(j+1)} \\ &= \sum_{j=1}^{\infty} \frac{2}{j(j+1)} \mathbb{G}_n^2(p_j(2t-1)). \end{aligned}$$

By arguments similar to those used to establish (i), we can now verify that (ii) follows. \square

Proof of corollary 11.9. By transforming the data using the transform $x \mapsto F_0(x)$, we can, without loss of generality, assume that the data are all i.i.d. uniform and reduce the interval of integration to $[0, 1]$. Now proposition 7.27 yields that $\hat{T}_1 \rightsquigarrow \int_0^1 \mathbb{G}^2(t)dt$, where $\mathbb{G}(t)$ is a standard Brownian bridge. Now Parseval's formula and arguments used in the proof of theorem 11.8 yield that

$$\int_0^1 \mathbb{G}^2(t)dt = \sum_{i=1}^{\infty} \mathbb{G}^2(\sqrt{2} \cos \pi x)/(\pi i)^2 \equiv T_1^*,$$

where now \mathbb{G} is a mean zero Gaussian Brownian bridge random measure, where the covariance between $\mathbb{G}(f)$ and $\mathbb{G}(g)$, where $f, g : [0, 1] \mapsto \mathbb{R}$, is $\int_0^1 f(s)g(s)ds - \int_0^1 f(s)ds \int_0^1 g(s)ds$. The fact that T_1^* is tight combined with the covariance structure of \mathbb{G} yields that T_1^* has the same distribution as T_1 , and the desired weak convergence result for \hat{T}_1 follows.

For \hat{T}_2 , apply the same transformation as above so that, without loss of generality, we can again assume that the data are i.i.d. uniform and that the interval of integration is $[0, 1]$. Let

$$\tilde{\mathbb{G}}_n \equiv \sqrt{\frac{n_1 n_2}{n_1 + n_2}} (\hat{F}_{n,1} - \hat{F}_{n,2}),$$

and fix $\epsilon \in (0, 1/2)$. We can now apply proposition 7.27 to verify that

$$\int_{\epsilon}^{1-\epsilon} \frac{\tilde{\mathbb{G}}_n^2(s)}{\hat{F}_{n,0}(s)[1-\hat{F}_{n,0}(s)]} d\hat{F}_{n,0}(s) \rightsquigarrow \int_{\epsilon}^{1-\epsilon} \frac{\mathbb{G}^2(s)}{s(1-s)} ds.$$

Note also that Fubini's theorem yields that both $E \left\{ \int_0^{\epsilon} \mathbb{G}^2(s)/[s(1-s)] ds \right\} = \epsilon$ and $E \left\{ \int_{1-\epsilon}^1 \mathbb{G}^2(s)/[s(1-s)] ds \right\} = \epsilon$.

We will now work towards bounding $\int_0^{\epsilon} \left(\tilde{\mathbb{G}}_n(s) / \left\{ \hat{F}_{n,0}(s)[1-\hat{F}_{n,0}(s)] \right\} \right) ds$. Fix $s \in (0, \epsilon)$ and note that, under the null hypothesis, the conditional distribution of $\tilde{\mathbb{G}}_n(s)$ given $\hat{F}_{n,0}(s) = m$ has the form

$$\sqrt{\frac{n_1 n_2}{n_1 + n_2}} \left(\frac{A}{n_1} - \frac{m - A}{n_2} \right),$$

where A is hypergeometric with density

$$P(A = a) = \binom{n_1}{a} \binom{n_1}{m-a} / \binom{n}{m},$$

where a is any integer between $(m - n_2) \vee 0$ and $m \wedge n_1$. Hence

$$\begin{aligned} E \left[\tilde{\mathbb{G}}_n^2(s) \mid \hat{F}_{n,0}(s) = m \right] &= \frac{n_1 + n_2}{n_1 n_2} \text{var}(A) \\ &= \frac{n}{n_1 n_2} \left(m \frac{n_1 n_2 (n - m)}{n^2 (n - 1)} \right) \\ &= \frac{n}{n - 1} \left[\frac{m(n - m)}{n^2} \right]. \end{aligned}$$

Thus

$$E \left[\int_0^{\epsilon} \frac{\tilde{\mathbb{G}}_n^2(s)}{\hat{F}_{n,0}(s)[1-\hat{F}_{n,0}(s)]} ds \right] = \frac{n}{n-1} E \hat{F}_{n,0}(\epsilon) \leq 2\epsilon,$$

for all $n \geq 2$. Similar arguments verify that

$$E \left[\int_{1-\epsilon}^1 \frac{\tilde{\mathbb{G}}_n^2(s)}{\hat{F}_{n,0}(s)[1-\hat{F}_{n,0}(s)]} ds \right] \leq 2\epsilon,$$

for all $n \geq 2$. Since ϵ was arbitrary, we now have that

$$\hat{T}_2 \rightsquigarrow \int_0^1 \frac{\mathbb{G}^2(s)}{s(1-s)} ds \equiv T_2^*,$$

where \mathbb{G} is the same Brownian bridge process used in defining T_1^* . Now we can again use arguments from the proof of theorem 11.8 to obtain that T_2^* has the same distribution as T_2 . \square

11.8 Exercises

11.8.1. Verify that F^* in the proof of theorem 11.5 is continuous.

11.8.2. Show that when P_n and P satisfy (11.4), we have that $P_n h \rightarrow Ph$, as $n \rightarrow \infty$, for all bounded and measurable h .

11.8.3. Prove proposition 11.20. Hint: Consider the arguments used in the proof of theorem 9.15.

11.9 Notes

Theorems 11.3 and 11.4 are inspired by theorem 2.14.9 of VW, and theorem 11.6 is derived from theorem 2.13.1 of VW. Theorem 11.7 is theorem 2.13.2 of VW, while theorem 11.8 is derived from examples 2.13.3 and 2.13.5 of VW. Much of the structure of section 11.4 comes from Kosorok (2003), although theorem 11.15 was derived from material in section 5 of Pollard (1990). Lemma 11.13 and theorems 11.14, 11.16 and 11.17, are lemma 2 and theorems 1 (with a minor modification), 3 and 2, respectively, of Kosorok (2003).

12

The Functional Delta Method

In this chapter, we build on the presentation of the functional delta method given in section 2.2.4. Recall the concept of Hadamard differentiability introduced in this section and also defined more precisely in section 6.3. The key result of section 2.2.4 is that the delta method and its bootstrap counterpart work provided the map ϕ is Hadamard differentiable tangentially to a suitable set \mathbb{D}_0 . We first present in section 12.1 clarifications and proofs of the two main theorems given in section 2.2.4, the functional delta method for Hadamard differentiable maps (theorem 2.8) and the conditional analog for the bootstrap (theorem 2.9). We then give in section 12.2 several important examples of Hadamard differentiable maps of use in statistics, along with specific illustrations of how those maps are utilized.

12.1 Main Results and Proofs

In this section, we first prove the functional delta method theorem (theorem 2.8) and then restate and prove theorem 2.9. Before proceeding, recall that X_n in the statement of theorem 2.8 is a random quantity that takes its values in a normed space \mathbb{D} .

Proof of theorem 2.8. Consider the map $h \mapsto r_n(\phi(\theta + r_n^{-1}h) - \phi(\theta)) \equiv g_n(h)$, and note that it is defined on the domain $\mathbb{D}_n \equiv \{h : \theta + r_n^{-1}h \in \mathbb{D}_\phi\}$ and satisfies $g_n(h_n) \rightarrow \phi'_\theta(h)$ for every $h_n \rightarrow h \in \mathbb{D}_0$ with $h_n \in \mathbb{D}_n$. Thus the conditions of the extended continuous mapping theorem (theorem 7.24)

are satisfied by $g(\cdot) = \phi'_\theta(\cdot)$. Hence conclusion (i) of that theorem implies $g_n(r_n(X_n - \theta)) \rightsquigarrow \phi'_\theta(X)$. \square

We now restate and prove theorem 2.9. The restatement clarifies the measurability condition. Before proceeding, recall the definitions of \mathbb{X}_n and $\hat{\mathbb{X}}_n$ in the statement of theorem 2.9. Specifically, $\mathbb{X}_n(X_n)$ is a sequence of random elements in a normed space \mathbb{D} based on the data sequence $\{X_n, n \geq 1\}$, while $\hat{\mathbb{X}}_n(X_n, W_n)$ is a bootstrapped version of \mathbb{X}_n based on both the data sequence and a sequence of weights $W = \{W_n, n \geq 1\}$. Note that the proof of this theorem utilizes the bootstrap continuous mapping theorem (theorem 10.8). Here is the restated version of theorem 2.9:

THEOREM 12.1 *For normed spaces \mathbb{D} and \mathbb{E} , let $\phi : \mathbb{D}_\phi \subset \mathbb{D} \mapsto \mathbb{E}$ be Hadamard-differentiable at μ tangentially to $\mathbb{D}_0 \subset \mathbb{D}$, with derivative ϕ'_μ . Let \mathbb{X}_n and $\hat{\mathbb{X}}_n$ have values in \mathbb{D}_ϕ , with $r_n(\mathbb{X}_n - \mu) \rightsquigarrow \mathbb{X}$, where \mathbb{X} is tight and takes its values in \mathbb{D}_0 for some sequence of constants $0 < r_n \rightarrow \infty$, the maps $W_n \mapsto h(\hat{\mathbb{X}}_n)$ are measurable for every $h \in C_b(\mathbb{D})$ outer almost surely, and where $r_n c(\hat{\mathbb{X}}_n - \mathbb{X}_n) \xrightarrow[W]{P} \mathbb{X}$, for a constant $0 < c < \infty$. Then $r_n c(\phi(\hat{\mathbb{X}}_n) - \phi(\mathbb{X}_n)) \xrightarrow[W]{P} \phi'_\mu(\mathbb{X})$.*

Proof. We can, without loss of generality, assume that \mathbb{D}_0 is complete and separable (since \mathbb{X} is tight), that \mathbb{D} and \mathbb{E} are both complete, and that $\phi'_\mu : \mathbb{D} \mapsto \mathbb{E}$ is continuous on all of \mathbb{D} , although it is permitted to not be bounded or linear off of \mathbb{D}_0 . To accomplish this, one can apply the Dugundji extension theorem (theorem 10.9) which extends any continuous operator defined on a closed subset to the entire space. It may be necessary to replace \mathbb{E} with its closed linear span to accomplish this.

We can now use arguments nearly identical to those used in the proof given in section 10.1 of theorem 10.4 to verify that, unconditionally, both $\hat{\mathbb{U}}_n \equiv r_n(\hat{\mathbb{X}}_n - \mathbb{X}_n) \rightsquigarrow c^{-1}\mathbb{X}$ and $r_n(\hat{\mathbb{X}}_n - \mu) \rightsquigarrow Z$, where Z is a tight random element. Fix some $h \in BL_1(\mathbb{D})$, define $\mathbb{U}_n \equiv r_n(\mathbb{X}_n - \mu)$, and let $\tilde{\mathbb{X}}_1$ and $\tilde{\mathbb{X}}_2$ be two independent copies of \mathbb{X} . We now have that

$$\begin{aligned} & \left| \mathbb{E}^* h(\hat{\mathbb{U}}_n + \mathbb{U}_n) - \mathbb{E} h(c^{-1}\tilde{\mathbb{X}}_1 + \tilde{\mathbb{X}}_2) \right| \\ & \leq \left| \mathbb{E}_{X_n} \mathbb{E}_{W_n} h(\hat{\mathbb{U}}_n + \mathbb{U}_n)^* - \mathbb{E}^* \mathbb{E}_{W_n} h(\hat{\mathbb{U}}_n + \mathbb{U}_n) \right| \\ & \quad + \mathbb{E}^* \left| \mathbb{E}_{W_n} h(\hat{\mathbb{U}}_n + \mathbb{U}_n) - \mathbb{E}_{\tilde{\mathbb{X}}_1} h(c^{-1}\tilde{\mathbb{X}}_1 + \mathbb{U}_n) \right| \\ & \quad + \left| \mathbb{E}^* \mathbb{E}_{\tilde{\mathbb{X}}_1} h(c^{-1}\tilde{\mathbb{X}}_1 + \mathbb{U}_n) - \mathbb{E}_{\tilde{\mathbb{X}}_2} \mathbb{E}_{\tilde{\mathbb{X}}_1} h(c^{-1}\tilde{\mathbb{X}}_1 + \tilde{\mathbb{X}}_2) \right|, \end{aligned}$$

where \mathbb{E}_{W_n} , $\mathbb{E}_{\tilde{\mathbb{X}}_1}$ and $\mathbb{E}_{\tilde{\mathbb{X}}_2}$ are expectations taken over the bootstrap weights, $\tilde{\mathbb{X}}_1$ and $\tilde{\mathbb{X}}_2$, respectively. The first term on the right in the above expression goes to zero by the asymptotic measurability of $\hat{\mathbb{U}}_n + \mathbb{U}_n = r_n(\hat{\mathbb{X}}_n - \mu)$. The second term goes to zero by the fact that $\hat{\mathbb{U}}_n \xrightarrow[W]{P} c^{-1}\mathbb{X}$ combined with

the fact that the map $x \mapsto h(x + \mathbb{U}_n)$ is Lipschitz continuous with Lipschitz constant 1 outer almost surely. The second term goes to zero since $\mathbb{U}_n \rightsquigarrow \mathbb{X}$ and the map $x \mapsto E_{\tilde{\mathbb{X}}_1} h(\tilde{\mathbb{X}}_1 + x)$ is also Lipschitz continuous with Lipschitz constant 1. Since h was arbitrary, we have by the Portmanteau theorem that, unconditionally,

$$r_n \begin{pmatrix} \hat{\mathbb{X}}_n - \mu \\ \mathbb{X}_n - \mu \end{pmatrix} \rightsquigarrow \begin{pmatrix} c^{-1}\tilde{\mathbb{X}}_1 + \tilde{\mathbb{X}}_2 \\ \tilde{\mathbb{X}}_2 \end{pmatrix}.$$

Now the functional delta method (theorem 2.8) yields

$$r_n \begin{pmatrix} \phi(\hat{\mathbb{X}}_n) - \phi(\mu) \\ \phi(\mathbb{X}_n) - \phi(\mu) \\ \hat{\mathbb{X}}_n - \mu \\ \mathbb{X}_n - \mu \end{pmatrix} \rightsquigarrow \begin{pmatrix} \phi'_\mu(c^{-1}\tilde{\mathbb{X}}_1 + \tilde{\mathbb{X}}_2) \\ \phi'_\mu(\tilde{\mathbb{X}}_2) \\ c^{-1}\tilde{\mathbb{X}}_1 + \tilde{\mathbb{X}}_2 \\ \tilde{\mathbb{X}}_2 \end{pmatrix},$$

since the map $(x, y) \mapsto (\phi(x), \phi(y), x, y)$ is Hadamard differentiable at (μ, μ) tangentially to \mathbb{D}_0 . This implies two things. First,

$$r_n c \begin{pmatrix} \phi(\hat{\mathbb{X}}_n) - \phi(\mathbb{X}_n) \\ \hat{\mathbb{X}}_n - \mathbb{X}_n \end{pmatrix} \rightsquigarrow \begin{pmatrix} \phi'_\mu(\mathbb{X}) \\ \mathbb{X} \end{pmatrix},$$

since ϕ'_μ is linear on \mathbb{D}_0 . Second, the usual continuous mapping theorem now yields that, unconditionally,

$$(12.1) \quad r_n c(\phi(\hat{\mathbb{X}}_n) - \phi(\mathbb{X}_n)) - \phi'_\mu(r_n c(\hat{\mathbb{X}}_n - \mathbb{X}_n)) \xrightarrow{P} 0,$$

since the map $(x, y) \mapsto x - \phi'_\mu(y)$ is continuous on all of $\mathbb{E} \times \mathbb{D}$.

Now for any map $h \in C_b(\mathbb{D})$, the map $x \mapsto h(r_n c(x - \mathbb{X}_n))$ is continuous and bounded for all $x \in \mathbb{D}$ outer almost surely. Thus the maps $W_n \mapsto h(r_n c(\hat{\mathbb{X}}_n - \mathbb{X}_n))$ are measurable for every $h \in C_b(\mathbb{D})$ outer almost surely. Hence the bootstrap continuous mapping theorem, theorem 10.8, yields that $\phi'_\mu(r_n c(\hat{\mathbb{X}}_n - \mathbb{X}_n)) \xrightarrow[W]{P} \phi'_\mu(\mathbb{X})$. The desired result now follows from (12.1). \square

12.2 Examples

We now give several important examples of Hadamard differentiable maps, along with illustrations of how these maps are utilized in statistical applications.

12.2.1 Composition

Recall from section 2.2.4 the map $\phi : \mathbb{D}_\phi \mapsto D[0, 1]$, where $\phi(f) = 1/f$ and $\mathbb{D}_\phi = \{f \in D[0, 1] : |f| > 0\}$. In that section, we established that ϕ was

Hadamard differentiable, tangentially to $D[0, 1]$, with derivative at $\theta \in \mathbb{D}_\phi$ equal to $h \mapsto -h/\theta^2$. This is a simple example of the following general composition result:

LEMMA 12.2 *Let $g : B \subset \bar{\mathbb{R}} \equiv [-\infty, \infty] \mapsto \mathbb{R}$ be differentiable with derivative continuous on all closed subsets of B , and let $\mathbb{D}_\phi = \{A \in \ell^\infty(\mathcal{X}) : \{R(A)\}^\delta \subset B \text{ for some } \delta > 0\}$, where \mathcal{X} is a set, $R(C)$ denotes the range of the function $C \in \ell^\infty(\mathcal{X})$, and D^δ is the δ -enlargement of the set D . Then $A \mapsto g \circ A$ is Hadamard-differentiable as a map from $\mathbb{D}_\phi \subset \ell^\infty(\mathcal{X})$ to $\ell^\infty(\mathcal{X})$, at every $A \in \mathbb{D}_\phi$. The derivative is given by $\phi'_A(\alpha) = g'(A)\alpha$, where g' is the derivative of g .*

Before giving the proof, we briefly return to our simple example of the reciprocal map $A \mapsto 1/A$. The differentiability of this map easily generalizes from $D[0, 1]$ to $\ell^\infty(\mathcal{X})$, for arbitrary \mathcal{X} , provided we restrict the domain of the reciprocal map to $\mathbb{D}_\phi = \{A \in \ell^\infty(\mathcal{X}) : \inf_{x \in \mathcal{X}} |A(x)| > 0\}$. This follows after applying lemma 12.2 to the set $B = [-\infty, 0) \cup (0, \infty]$.

Proof of lemma 12.2. Note that $\mathbb{D} = \ell^\infty(\mathcal{X})$ in this case, and that the tangent set for the derivative is all of \mathbb{D} . Let t_n be any real sequence with $t_n \rightarrow 0$, let $\{h_n\} \in \ell^\infty(\mathcal{X})$ be any sequence converging to $\alpha \in \ell^\infty(\mathcal{X})$ uniformly, and define $A_n = A + t_n h_n$. Then, by the conditions of the theorem, there exists a closed $B_1 \subset B$ such that $\{R(A) \cup R(A_n)\}^\delta \subset B_1$ for some $\delta > 0$ and all n large enough. Hence

$$\sup_{x \in \mathcal{X}} \left| \frac{g(A(x) + t_n h_n(x)) - g(A(x))}{t_n} - g'(A(x))\alpha(x) \right| \rightarrow 0,$$

as $n \rightarrow \infty$, since continuous functions on closed sets are bounded and thus g' is uniformly continuous on B_1 . \square

12.2.2 Integration

For an $M < \infty$ and an interval $[a, b] \in \bar{\mathbb{R}}$, let $BV_M[a, b]$ be the set of all functions $A \in D[a, b]$ with total variation $\int_{(a,b]} |dA(s)| \leq M$. In this section, we consider, for given functions $A \in D[a, b]$ and $B \in BV_M[a, b]$ and domain $\mathbb{D}_M \equiv D[a, b] \times BV_M[a, b]$, the maps $\phi : \mathbb{D}_M \mapsto \mathbb{R}$ and $\psi : \mathbb{D}_M \mapsto D[a, b]$ defined by

$$(12.2) \quad \phi(A, B) = \int_{(a,b]} A(s)dB(s) \quad \text{and} \quad \psi(A, B)(t) = \int_{(a,t]} A(s)dB(s).$$

The following lemma verifies that these two maps are Hadamard differentiable:

LEMMA 12.3 *For each fixed $M < \infty$, the maps $\phi : \mathbb{D}_M \mapsto \mathbb{R}$ and $\psi : \mathbb{D}_M \mapsto D[a, b]$ defined in (12.2) are Hadamard differentiable at each $(A, B) \in \mathbb{D}_M$ with $\int_{(a,b]} |dA| < \infty$. The derivatives are given by*

$$\begin{aligned}\phi'_{A,B}(\alpha, \beta) &= \int_{(a,b]} Ad\beta + \int_{(a,b]} \alpha dB, \quad \text{and} \\ \psi'_{A,B}(\alpha, \beta)(t) &= \int_{(a,t]} Ad\beta + \int_{(a,t]} \alpha dB.\end{aligned}$$

Note that in the above lemma we define $\int_{(a,t]} Ad\beta = A(t)\beta(t) - A(a)\beta(a) - \int_{(a,t]} \beta(s-)dA(s)$ so that the integral is well defined even when β does not have bounded variation. We will present the proof of this lemma at the end of this section.

We now look at two statistical applications of lemma 12.3, the two-sample Wilcoxon rank sum statistic, and the Nelson-Aalen integrated hazard estimator. Consider first the Wilcoxon statistic. Let X_1, \dots, X_m and Y_1, \dots, Y_n be independent samples from distributions F and G on the reals. If \mathbb{F}_m and \mathbb{G}_n are the respective empirical distribution functions, the Wilcoxon rank sum statistic for comparing F and G has the form

$$T_1 = m \int_{\mathbb{R}} (m\mathbb{F}_m(x) + n\mathbb{G}_n(x))d\mathbb{F}_m(x).$$

If we temporarily assume that F and G are continuous, then

$$\begin{aligned}T_1 &= mn \int_{\mathbb{R}} \mathbb{G}_n(x)d\mathbb{F}_m(x) + m^2 \int_{\mathbb{R}} \mathbb{F}_m(x)d\mathbb{F}_m(x) \\ &= mn \int_{\mathbb{R}} \mathbb{G}_n(x)d\mathbb{F}_m(x) + \frac{m^2 + m}{2} \\ &\equiv mnT_2 + \frac{m^2 + m}{2},\end{aligned}$$

where T_2 is the Mann-Whitney statistic. When F or G have atoms, the relationship between the Wilcoxon and Mann-Whitney statistics is more complex. We will now study the asymptotic properties of the Mann-Whitney version of the rank sum statistic, T_2 .

For arbitrary F and G , $T_2 = \phi(\mathbb{G}_n, \mathbb{F}_m)$, where ϕ is as defined in lemma 12.3. Note that F , G , \mathbb{F}_m and \mathbb{G}_n all have total variation ≤ 1 . Thus lemma 12.3 applies, and we obtain that the Hadamard derivative of ϕ at $(A, B) = (G, F)$ is the map $\phi'_{G,F}(\alpha, \beta) = \int_{\mathbb{R}} Gd\beta + \int_{\mathbb{R}} \alpha dF$, which is continuous and linear over $\alpha, \beta \in D[-\infty, \infty]$. If we assume that $m/(m+n) \rightarrow \lambda \in [0, 1]$, as $m \wedge n \rightarrow \infty$, then

$$\sqrt{\frac{mn}{m+n}} \begin{pmatrix} \mathbb{G}_n - G \\ \mathbb{F}_m - F \end{pmatrix} \rightsquigarrow \begin{pmatrix} \sqrt{\lambda} \mathbb{B}_1(G) \\ \sqrt{1-\lambda} \mathbb{B}_2(F) \end{pmatrix},$$

where \mathbb{B}_1 and \mathbb{B}_2 are independent standard Brownian bridges. Hence $\mathbb{G}_G(\cdot) \equiv \mathbb{B}_1(G(\cdot))$ and $\mathbb{G}_F(\cdot) \equiv \mathbb{B}_2(F(\cdot))$ both live in $D[-\infty, \infty]$. Now theorem 2.8 yields

$$T_2 \rightsquigarrow \sqrt{\lambda} \int_{\mathbb{R}} G d\mathbb{G}_F + \sqrt{1-\lambda} \int_{\mathbb{R}} \mathbb{G}_F dG,$$

as $m \wedge n \rightarrow \infty$. When $F = G$ and F is continuous, this limiting distribution is mean zero normal with variance $1/12$. The delta method bootstrap, theorem 12.1, is also applicable and can be used to obtain an estimate of the limiting distribution under more general hypotheses on F and G .

We now shift our attention to the Nelson-Aalen estimator under right censoring. In the right censored survival data setting, an observation consists of the pair (X, δ) , where $X = T \wedge C$ is the minimum of a failure time T and censoring time C , and $\delta = 1\{T \leq C\}$. T and C are assumed to be independent. Let F be the distribution function for T , and define the integrated baseline hazard for F to be $\Lambda(t) = \int_0^t dF(s)/S(s-)$, where $S \equiv 1 - F$ is the survival function. The Nelson-Aalen estimator for Λ , based on the i.i.d. sample $(X_1, \delta_1), \dots, (X_n, \delta_n)$, is

$$\hat{\Lambda}_n(t) \equiv \int_{[0,t]} \frac{d\hat{N}_n(s)}{\hat{Y}_n(s)},$$

where $\hat{N}_n(t) \equiv n^{-1} \sum_{i=1}^n \delta_i 1\{X_i \leq t\}$ and $\hat{Y}_n(t) \equiv n^{-1} \sum_{i=1}^n 1\{X_i \geq t\}$. It is easy to verify that the classes $\{\delta 1\{X \leq t\}, t \geq 0\}$ and $\{1\{X \geq t\} : t \geq 0\}$ are both Donsker and hence that

$$(12.3) \quad \sqrt{n} \begin{pmatrix} \hat{N}_n - N_0 \\ \hat{Y}_n - Y_0 \end{pmatrix} \rightsquigarrow \begin{pmatrix} \mathbb{G}_1 \\ \mathbb{G}_2 \end{pmatrix},$$

where $N_0(t) \equiv P(T \leq t, C \geq T)$, $Y_0(t) \equiv P(X \geq t)$, and \mathbb{G}_1 and \mathbb{G}_2 are tight Gaussian processes with respective covariances $N_0(s \wedge t) - N_0(s)N_0(t)$ and $Y_0(s \vee t) - Y_0(s)Y_0(t)$ and with cross-covariance $1\{s \geq t\} [N_0(s) - N_0(t-)] - N_0(s)Y_0(t)$. Note that while we have already seen this survival set-up several times (eg., sections 2.2.5 and 4.2.2), we are choosing to use slightly different notation than previously used to emphasize certain features of the underlying empirical processes.

The Nelson-Aalen estimator depends on the pair (\hat{N}_n, \hat{Y}_n) through the two maps

$$(A, B) \mapsto \left(A, \frac{1}{B} \right) \mapsto \int_{[0,t]} \frac{1}{B} dA.$$

From section 12.1.1, lemma 12.3, and the chain rule (lemma 6.19), it is easy to see that this composition map is Hadamard differentiable on a domain of the type $\{(A, B) : \int_{[0,\tau]} |dA(t)| \leq M, \inf_{t \in [0,\tau]} |B(t)| \geq \epsilon\}$ for a given $M < \infty$ and $\epsilon > 0$, at every point (A, B) such that $1/B$ has bounded variation. Note that the interval of integration we are using, $[0, \tau]$, is left-closed rather than left-open as in the definition of ψ given in (12.2). However, if we pick some $\eta > 0$, then in fact integrals over $[0, t]$, for any $t > 0$, of functions which have zero variation over $(-\infty, 0)$ are unchanged if

we replace the interval of integration with $(-\eta, t]$. Thus we will still be able to utilize lemma 12.3 in our current set-up. In this case, the point (A, B) of interest is $A = N_0$ and $B = Y_0$. Thus if t is restricted to the interval $[0, \tau]$, where τ satisfied $Y_0(\tau) > 0$, then it is easy to see that the pair (\hat{N}_n, \hat{Y}_n) is contained in the given domain with probability tending to 1 as $n \rightarrow \infty$. The derivative of the composition map is given by

$$(\alpha, \beta) \mapsto \left(\alpha, \frac{-\beta}{Y_0^2} \right) \mapsto \int_{[0,t]} \frac{d\alpha}{Y_0} - \int_{[0,t]} \frac{\beta dN_0}{Y_0^2}.$$

Thus from (12.3), we obtain via theorem 2.8 that

$$(12.4) \quad \sqrt{n}(\hat{\Lambda}_n - \Lambda) \rightsquigarrow \int_{[0,(\cdot)]} \frac{d\mathbb{G}_1}{Y_0} - \int_{[0,(\cdot)]} \frac{\mathbb{G}_2 dN_0}{Y_0^2}.$$

The Gaussian process on the right side of (12.4) is equal to $\int_{[0,(\cdot)]} d\mathbb{M}/Y_0$, where $\mathbb{M}(t) \equiv \mathbb{G}_1(t) - \int_{[0,t]} \mathbb{G}_2 d\Lambda$ can be shown to be a Gaussian martingale with independent increments and covariance $\int_{[0,s \wedge t]} (1 - \Delta\Lambda) d\Lambda$, where $\Delta A(t) \equiv A(t) - A(t-)$ is the mass at t of a signed-measure A . This means that the Gaussian process on the right side of (12.4) is also a Gaussian martingale with independent increments but with covariance $\int_{[0,s \wedge t]} (1 - \Delta\Lambda) d\Lambda / Y_0$. A useful discussion of continuous time martingales arising in right censored survival data can be found in Fleming and Harrington (1991).

The delta method bootstrap, theorem 12.1, is also applicable here and can be used to obtain an estimate of the limiting distribution. However, when Λ is continuous over $[0, \tau]$, the independent increments structure implies that the limiting distribution is time-transformed Brownian motion. More precisely, the limiting process can be expressed as $\mathbb{W}(v(t))$, where \mathbb{W} is standard Brownian motion on $[0, \infty)$ and $v(t) \equiv \int_{(0,t]} d\Lambda / Y_0$. As discussed in chapter 7 of Fleming and Harrington (1991), this fact can be used to compute asymptotically exact simultaneous confidence bands for Λ .

Proof of lemma 12.3. For sequences $t_n \rightarrow 0$, $\alpha_n \rightarrow \alpha$, and $\beta_n \rightarrow \beta$, define $A_n \equiv A + t_n \alpha_n$ and $B_n \equiv B + t_n \beta_n$. Since we require that $(A_n, B_n) \in \mathbb{D}_M$, we know that the total variation of B_n is bounded by M . Consider first the derivative of ψ , and note that

$$(12.5) \quad \frac{\int_{(0,t]} A_n dB_n - \int_{(0,t]} A dB}{t_n} - \psi'_{A,B}(\alpha_n, \beta_n) = \int_{(0,t]} \alpha d(B_n - B) + \int_{(0,t]} (\alpha_n - \alpha) d(B_n - B).$$

Since it is easy to verify that the map $(\alpha, \beta) \mapsto \psi'_{A,B}(\alpha, \beta)$ is continuous and linear, the desired Hadamard differentiability of ψ will follow provided the right side of (12.5) goes to zero. To begin with, the second term on the

right side goes to zero uniformly over $t \in (a, b]$, since both B_n and B have total variation bounded by M .

Now, for the first term on the right side of (12.5), fix $\epsilon > 0$. Since α is cadlag, there exists a partition $a = t_0 < t_1 < \dots < t_m = b$ such that α varies less than ϵ over each interval $[t_{i-1}, t_i)$, $1 \leq i \leq m$, and $m < \infty$. Now define the function $\tilde{\alpha}$ to be equal to $\alpha(t_{i-1})$ over the interval $[t_{i-1}, t_i)$, $1 \leq i \leq m$, with $\tilde{\alpha}(b) = \alpha(b)$. Thus

$$\begin{aligned} & \left\| \int_{(0,t]} \alpha d(B_n - B) \right\|_{\infty} \\ & \leq \left\| \int_{(0,t]} (\alpha - \tilde{\alpha}) d(B_n - B) \right\|_{\infty} + \left\| \int_{(0,t]} \tilde{\alpha} d(B_n - B) \right\|_{\infty} \\ & \leq \|\alpha - \tilde{\alpha}\|_{\infty} 2M + \sum_{i=1}^m |\alpha(t_{i-1})| \times |(B_n - B)(t_i) - (B_n - B)(t_{i-1})| \\ & \quad + |\alpha(b)| \times |(B_n - B)(b)| \\ & \leq \epsilon 2M + (2m + 1) \|B_n - B\|_{\infty} \|\alpha\|_{\infty} \\ & \rightarrow \epsilon 2M, \end{aligned}$$

as $n \rightarrow \infty$. Since ϵ was arbitrary, we have that the first term on the right side of (12.5) goes to zero, as $n \rightarrow \infty$, and the desired Hadamard differentiability of ψ follows.

Now the desired Hadamard differentiability of ϕ follows from the trivial but useful lemma 12.4 below, by taking the extraction map $f : D[a, b] \mapsto \mathbb{R}$ defined by $f(x) = x(b)$, noting that $\phi = f(\psi)$, and then applying the chain rule for Hadamard derivatives (lemma 6.19). \square

LEMMA 12.4 *Let T be a set and fix $T_0 \subset T$. Define the extraction map $f : \ell^{\infty}(T) \mapsto \ell^{\infty}(T_0)$ as $f(x) = \{x(t) : t \in T_0\}$. Then f is Hadamard differentiable at all $x \in \ell^{\infty}(T)$ with derivative $f'_x(h) = \{h(t) : t \in T_0\}$.*

Proof. Let t_n be any real sequence with $t_n \rightarrow 0$, and let $\{h_n\} \in \ell^{\infty}(T)$ be any sequence converging to $h \in \ell^{\infty}(T)$. The desired conclusion follows after noting that $t_n^{-1}[f(x + t_n h_n) - f(x)] = \{h_n(t) : t \in T_0\} \rightarrow \{h(t) : t \in T_0\}$, as $n \rightarrow \infty$. \square

12.2.3 Product Integration

For a function $A \in D(0, b]$, let $A^c(t) \equiv A(t) - \sum_{0 < s \leq t} \Delta A(s)$, where ΔA is as defined in the previous section, be the continuous part of A . We define the product integral to be the map $A \mapsto \phi(A)$, where

$$\phi(A)(t) \equiv \prod_{0 < s \leq t} (1 + dA(s)) = \exp(A^c(t)) \prod_{0 < s \leq t} (1 - \Delta A(s)).$$

The first product is merely notation, but it is motivated by the mathematical definition of the product integral:

$$\phi(A)(t) = \lim_{\max_i |t_i - t_{i-1}| \rightarrow 0} \prod_i \{1 + [A(t_i) - A(t_{i-1})]\},$$

where the limit is over partitions $0 = t_0 < t_1 < \cdots < t_m = t$ with maximum separation decreasing to zero. We will also use the notation

$$\phi(A)(s, t] = \prod_{s < u \leq t} (1 + dA(u)) \equiv \frac{\phi(A)(t)}{\phi(A)(s)},$$

for all $0 \leq s < t$. The two terms on the left are defined by the far right term. Three alternative definitions of the product integral, as solutions of two different Volterra integral equations and as a “Peano series,” are given in exercise 12.3.2.

The following lemma verifies that product integration is Hadamard differentiable:

LEMMA 12.5 *For fixed constants $0 < b, M < \infty$, the product integral map $\phi : BV_M[0, b] \subset D[0, b] \mapsto D[0, b]$ is Hadamard differentiable with derivative*

$$\phi'_A(\alpha)(t) = \int_{(0, t]} \phi(A)(0, u) \phi(A)(u, t] d\alpha(u).$$

When $\alpha \in D[0, b]$ has unbounded variation, the above quantity is well-defined by integration by parts.

We give the proof later on in this section, after first discussing an important statistical application. From the discussion of the Nelson-Aalen estimator $\hat{\Lambda}_n$ in section 12.2.2, it is not hard to verify that in the right-censored survival analysis setting $S(t) = \phi(-\Lambda)(t)$, where ϕ is the product integration map. Moreover, it is easily verified that the Kaplan-Meier estimator \hat{S}_n discussed in sections 2.2.5 and 4.3 satisfies $\hat{S}_n(t) = \phi(-\hat{\Lambda}_n)(t)$.

We can now use lemma 12.5 to derive the asymptotic limiting distribution of $\sqrt{n}(\hat{S}_n - S)$. As in section 12.2.2, we will restrict our time domain to $[0, \tau]$, where $P(X > \tau) > 0$. Under these circumstances, there exists an $M < \infty$, such that $\Lambda(\tau) < M$ and $\hat{\Lambda}_n(\tau) < M$ with probability tending to 1 as $n \rightarrow \infty$. Now lemma 12.5, combined with (12.4) and the discussion immediately following, yields

$$\begin{aligned} \sqrt{n}(\hat{S}_n - S) &\rightsquigarrow - \int_{(0, (\cdot)]} \phi(-\Lambda)(0, u) \phi(-\Lambda)(u, t] \frac{d\mathbb{M}}{Y_0} \\ &= -S(t) \int_{(0, (\cdot)]} \frac{d\mathbb{M}}{(1 - \Delta\Lambda)Y_0}, \end{aligned}$$

where \mathbb{M} is a Gaussian martingale with independent increments and covariance $\int_{(0, s \wedge t]} (1 - \Delta\Lambda) d\Lambda / Y_0$. Thus $\sqrt{n}(\hat{S}_n - S)/S$ is asymptotically time-transformed Brownian motion $\mathbb{W}(w(t))$, where \mathbb{W} is standard Brownian

motion on $[0, \infty)$ and where $w(t) \equiv \int_{(0,t]} [(1 - \Delta\Lambda)Y_0]^{-1} d\Lambda$. Along the lines discussed in the Nelson-Aalen example of section 12.2.2, the form of the limiting distribution can be used to obtain asymptotically exact simultaneous confidence bands for S . The delta method bootstrap, theorem 12.1, can also be used for inference on S .

Before giving the proof of lemma 12.5, we present the following lemma which we will need and which includes the important *Duhamel equation* for the difference between two product integrals:

LEMMA 12.6 For $A, B \in D(0, b]$, we have for all $0 \leq s < t \leq b$ the following, where M is the sum of the total variation of A and B :

(i) (the Duhamel equation)

$$(\phi(B) - \phi(A))(s, t] = \int_{(s,t]} \phi(A)(0, u)\phi(B)(u, t](B - A)(du).$$

(ii) $\|\phi(A) - \phi(B)\|_{(s,t]} \leq e^M(1 + M)^2\|A - B\|_{(s,t]}.$

Proof. For any $u \in (s, t]$, consider the function $C_u \in D(s, t]$ defined as

$$C_u(x) = \begin{cases} A(x) - A(s), & \text{for } s \leq x < u, \\ A(u-) - A(s), & \text{for } x = u, \\ A(u-) - A(s) + B(x) - B(u), & \text{for } u < x \leq t. \end{cases}$$

Using the Peano series expansion of exercise 12.3.2, part (b), we obtain:

$$\begin{aligned} \phi(A)(s, u)\phi(B)(u, t] &= \phi(C_u)(s, t] = 1 \\ &+ \sum_{m,n \geq 0: m+n \geq 1} \int_{s < t_1 < \dots < t_m < u < t_{m+1} < \dots < t_{m+n} \leq t} A(dt_1) \cdots A(dt_m) \\ &\times B(dt_{m+1}) \cdots B(dt_{m+n}). \end{aligned}$$

Thus

$$\begin{aligned} &\int_{(s,t]} \phi(A)(s, u)\phi(B)(u, t](B - A)(du) \\ &= \sum_{n \geq 1} \int_{s < x_1 < \dots < x_n \leq t} \left[1 + \sum_{m \geq 1} \int_{s < t_1 < \dots < t_m < x_1} A(dt_1) \cdots A(dt_m) \right] \\ &\quad \times B(dx_1) \cdots B(dx_n) \\ &\quad - \sum_{n \geq 1} \int_{s < t_1 < \dots < t_n \leq t} \left[1 + \sum_{m \geq 1} \int_{t_n < x_1 < \dots < x_m \leq t} B(dx_1) \cdots B(dx_m) \right] \\ &\quad \times A(dt_1) \cdots A(dt_n) \end{aligned}$$

$$\begin{aligned}
&= \sum_{n \geq 1} \int_{s < x_1 < \dots < x_n \leq t} B(dx_1) \cdots B(dx_n) \\
&\quad - \sum_{n \geq 1} \int_{s < t_1 < \dots < t_n \leq t} A(dt_1) \cdots A(dt_n) \\
&= \phi(B)(s, t] - \phi(A)(s, t].
\end{aligned}$$

This proves part (i).

For part (ii), we need to derive an integration by parts formula for the Duhamel equation. Define $G = B - A$ and $H(u) \equiv \int_0^u \phi(B)(v, t] G(dv)$. Now integration by parts gives us

$$\begin{aligned}
(12.6) \quad &\int_{(s, t]} \phi(A)(0, u) \phi(B)(u, t] G(du) \\
&= \phi(A)(t) H(t) - \phi(A)(s) H(s) - \int_{(s, t]} H(u) \phi(A)(du).
\end{aligned}$$

From the backwards integral equation (part (c) of exercise 12.3.2), we know that $\phi(B)(dv, t] = -\phi(B)(v, t] B(dv)$, and thus, by integration by parts, we obtain

$$H(u) = G(u) \phi(B)(u, t] + \int_{(0, u]} G(v-) \phi(B)(v, t] B(dv).$$

Combining this with (12.6) and the fact that $\phi(A)(du) = \phi(A)(u-) A(du)$, we get

$$\begin{aligned}
(12.7) \quad &\int_{(s, t]} \phi(A)(0, u) \phi(B)(u, t] G(du) \\
&= \phi(A)(t) \int_{(0, t]} G(u-) \phi(B)(u, t] B(du) \\
&\quad - \phi(A)(s) \phi(B)(s, t] G(s) \\
&\quad - \phi(A)(s) \int_{(0, s]} G(u-) \phi(B)(u, t] B(du) \\
&\quad - \int_{(s, t]} G(u) \phi(B)(u, t] \phi(A)(u-) A(du) \\
&\quad - \int_{(s, t]} \int_{(0, u]} G(v-) \phi(B)(v, t] B(dv) \phi(A)(u-) A(du).
\end{aligned}$$

From exercise 12.3.3, we know that $\phi(A)$ and $\phi(B)$ are bounded by the exponentiation of the respective total variations of A and B . Now the desired result follows. \square

Proof of lemma 12.5. Set $A_n = A + t_n \alpha_n$ for a sequence $\alpha_n \rightarrow \alpha$ with the total variation of both A and A_n bounded by M . In view of the Duhamel equation (part (i) of lemma 12.6 above), it suffices to show that

$$\int_{(0,t]} \phi(A)(0, u)\phi(A_n)(u, t]d\alpha_n(u) \rightarrow \int_{(0,t]} \phi(A)(0, u)\phi(A)(u, t]d\alpha(u),$$

uniformly in $0 \leq t \leq b$. Fix $\epsilon > 0$. Since $\alpha \in D[0, b]$, there exists a function $\tilde{\alpha}$ with total variation $V < \infty$ such that $\|\alpha - \tilde{\alpha}\|_\infty \leq \epsilon$.

Now recall that the derivation of the integration by parts formula (12.7) from the proof of lemma 12.6 did not depend on the definition of G , other than the necessity of G being right-continuous. If we replace G with $\alpha - \tilde{\alpha}$, we obtain from (12.7) that

$$\begin{aligned} \left\| \int_{(0,(\cdot)]} \phi(A)(0, u)\phi(A_n)(u, t]d(\alpha_n - \tilde{\alpha})(u) \right\|_\infty &\leq e^{2M}(1 + 2M)^2\|\alpha_n - \tilde{\alpha}\|_\infty \\ &\rightarrow e^{2M}(1 + 2M)^2\epsilon, \end{aligned}$$

as $n \rightarrow \infty$, since $\|\alpha_n - \alpha\|_\infty \rightarrow 0$. Moreover,

$$\begin{aligned} \left\| \int_{(0,(\cdot)]} \phi(A)(0, u) [\phi(A_n) - \phi(A)](u, t]d\tilde{\alpha}(u) \right\|_\infty &\leq \|\phi(A_n) - \phi(A)\|_\infty \|\phi(A)\|_\infty V \\ &\rightarrow 0, \end{aligned}$$

as $n \rightarrow \infty$. Now using again the integration by parts formula (12.7), but with $G = \tilde{\alpha} - \alpha$, we obtain

$$\left\| \int_{(0,(\cdot)]} \phi(A)(0, u)\phi(A)(u, t]d(\tilde{\alpha} - \alpha)(u) \right\|_\infty \leq e^{2M}(1 + 2M)^2\epsilon.$$

Thus the desired result follows since ϵ was arbitrary. \square

12.2.4 Inversion

Recall the derivation given in the paragraphs following theorem 2.8 of the Hadamard derivative of the inverse of a distribution function F . Note that this derivation did not depend on F being a distribution function per se. In fact, the derivation will carry through unchanged if we replace the distribution function F with any nondecreasing, cadlag function A satisfying mild regularity conditions. For a non-decreasing function $B \in D(-\infty, \infty)$, define the left-continuous inverse $r \mapsto B^{-1}(r) \equiv \inf\{x : B(x) \geq r\}$. We will hereafter use the notation $\tilde{D}[u, v]$ to denote all left-continuous functions with right-hand limits (caglad) on $[u, v]$ and $D_1[u, v]$ to denote the restriction of all non-decreasing functions in $D(-\infty, \infty)$ to the interval $[u, v]$. Here is a precise statement of the general Hadamard differentiation result for non-decreasing functions:

LEMMA 12.7 Let $-\infty < p \leq q < \infty$, and let the non-decreasing function $A \in D(-\infty, \infty)$ be continuously differentiable on the interval $[u, v] \equiv [A^{-1}(p) - \epsilon, A^{-1}(q) + \epsilon]$, for some $\epsilon > 0$, with derivative A' being strictly positive and bounded over $[u, v]$. Then the inverse map $B \mapsto B^{-1}$ as a map $D_1[u, v] \subset D[u, v] \mapsto \tilde{D}[p, q]$ is Hadamard differentiable at A tangentially to $C[u, v]$, with derivative $\alpha \mapsto -(\alpha/A') \circ A^{-1}$.

We will give the proof of lemma 12.7 at the end of this section. We now restrict ourselves to the setting where A is a distribution function which we will now denote by F . The following lemma provides two results similar to lemma 12.7 but which utilize knowledge about the support of the distribution function F . Let $D_2[u, v]$ be the subset of distribution functions in $D_1[u, v]$ with support only on $[u, \infty)$, and let $D_3[u, v]$ be the subset of distribution functions in $D_2[u, v]$ which have support only on $[u, v]$.

LEMMA 12.8 Let F be a distribution function. We have the following:

- (i) Let $F \in D_2[u, \infty)$, for finite u , and let $q \in (0, 1)$. Assume F is continuously differentiable on the interval $[u, v] = [u, F^{-1}(q) + \epsilon]$, for some $\epsilon > 0$, with derivative f being strictly positive and bounded over $[u, v]$. Then the inverse map $G \mapsto G^{-1}$ as a map $D_2[u, v] \subset D[u, v] \mapsto \tilde{D}(0, q)$ is Hadamard differentiable at F tangentially to $C[u, v]$.
- (ii) Let $F \in D_3[u, v]$, for $[u, v]$ compact, and assume that F is continuously differentiable on $[u, v]$ with derivative f strictly positive and bounded over $[u, v]$. Then the inverse map $G \mapsto G^{-1}$ as a map $D_3[u, v] \subset D[u, v] \mapsto \tilde{D}(0, 1)$ is Hadamard differentiable at F tangentially to $C[u, v]$.

In either case, the derivative is the map $\alpha \mapsto -(\alpha/f) \circ F^{-1}$.

Before giving the proof of the above two lemmas, we will discuss some important statistical applications. As discussed in section 2.2.4, an important application of these results is to estimation and inference for the quantile function $p \mapsto F^{-1}(p)$ based on the usual empirical distribution function for i.i.d. data. Lemma 12.8 is useful when some information is available on the support of F , since it allows the range of p to extend as far as possible. These results are applicable to other estimators of the distribution function F besides the usual empirical distribution, provided the standardized estimators converge to a tight limiting process over the necessary intervals. Several examples of such estimators include the Kaplan-Meier estimator, the self-consistent estimator of Chang (1990) for doubly-censored data, and certain estimators from dependent data as mentioned in Kosorok (1999).

We now apply lemma 12.8 to the construction of quantile processes based on the Kaplan-Meier estimator discussed in section 12.2.3 above. Since it is known that the support of a survival function is on $[0, \infty)$, we can utilize part (i) of this lemma. Define the Kaplan-Meier quantile process

$\{\hat{\xi}(p) \equiv \hat{F}_n^{-1}(p), 0 < p \leq q\}$, where $\hat{F}_n = 1 - \hat{S}_n$, \hat{S}_n is the Kaplan-Meier estimator, and where $0 < q < F(\tau)$ for τ as defined in the previous section. Assume that F is continuously differentiable on $[0, \tau]$ with density f bounded below by zero and finite. Combining the results of the previous section with part (i) of lemma 12.8 and theorem 2.8, we obtain

$$\sqrt{n}(\hat{\xi} - \xi)(\cdot) \rightsquigarrow \frac{S(\xi(\cdot))}{f(\xi(\cdot))} \int_{(0, \xi(\cdot)]} \frac{d\mathbb{M}}{(1 - \Delta\Lambda)Y_0},$$

in $\tilde{D}(0, q]$, where $\xi(p) \equiv \xi_p$ and \mathbb{M} is the Gaussian martingale described in the previous section. Thus $\sqrt{n}(\hat{\xi} - \xi)f(\xi)/S(\xi)$ is asymptotically time-transformed Brownian motion with time-transform $w(\xi)$, where w is as defined in the previous section, over the interval $(0, q]$. As described in Kosorok (1999), one can construct kernel estimators for f —which can be shown to be uniformly consistent—to facilitate inference. An alternative approach is the bootstrap which can be shown to be valid in this setting based on theorem 12.1.

Proof of lemma 12.7. The arguments are essentially identical to those used in the paragraphs following theorem 2.8, except that the distribution function F is replaced by a more general, non-decreasing function A . \square

Proof of lemma 12.8. To prove part (i), let $\alpha_n \rightarrow \alpha$ uniformly in $D[u, v]$ and $t_n \rightarrow 0$, where α is continuous and $F + t_n\alpha_n$ is contained in $D_2[u, v]$ for all $n \geq 1$. Abbreviate $F^{-1}(p)$ and $(F + t_n\alpha_n)^{-1}(p)$ to ξ_p and ξ_{pn} , respectively. Since F and $F + t_n\alpha_n$ have domains (u, ∞) (the lower bound by assumption), we have that $\xi_p, \xi_{pn} > u$ for all $0 < p \leq q$. Moreover, $\xi_p, \xi_{pn} \leq v$ for all n large enough. Thus the numbers $\epsilon_{pn} \equiv t_n^2 \wedge (\xi_{pn} - u)$ are positive for all $0 < p \leq q$, for all n large enough. Hence, by definition, we have for all $p \in (0, q]$ that

$$(12.8) \quad (F + t_n\alpha_n)(\xi_{pn} - \epsilon_{pn}) \leq p \leq (F + t_n\alpha_n)(\xi_{pn}),$$

for all sufficiently large n .

By the smoothness of F , we have $F(\xi_p) = p$ and $F(\xi_{pn} - \epsilon_{pn}) = F(\xi_{pn}) + O(\epsilon_{pn})$, uniformly over $p \in (0, q]$. Thus from (12.8) we obtain

$$(12.9) \quad -t_n\alpha(\xi_{pn}) + o(t_n) \leq F(\xi_{pn}) - F(\xi_p) \leq -t_n\alpha(\xi_{pn} - \epsilon_{pn}) + o(t_n),$$

where the $o(t_n)$ terms are uniform over $0 < p \leq q$. Both the far left and far right sides are $O(t_n)$, while the middle term is bounded above and below by constants times $|\xi_{pn} - \xi_p|$, for all $0 < p \leq q$. Hence $|\xi_{pn} - \xi_p| = O(t_n)$, uniformly over $0 < p \leq q$. The part (i) result now follows from (12.9), since $F(\xi_{pn}) - F(\xi_p) = -f(\xi_p)(\xi_{pn} - \xi_p) + E_n$, where $E_n = o(\sup_{0 < p \leq q} |\xi_{pn} - \xi_p|)$ by the uniform differentiability of F over $(u, v]$.

Note that part (ii) of this lemma is precisely part (ii) of lemma 3.9.23 of VW, and the details of the proof (which are quite similar to the proof of part (i)) can be found therein. \square

12.2.5 Other Mappings

We now mention briefly a few additional interesting examples. The first example is the *copula map*. For a bivariate distribution function H , with marginals $F_H(x) \equiv H(x, \infty)$ and $G_H(y) \equiv H(\infty, y)$, the copula map is the map ϕ from bivariate distributions on \mathbb{R}^2 to bivariate distributions on $[0, 1]^2$ defined as follows:

$$H \mapsto \phi(H)(u, v) = H(F_H^{-1}(u), G_H^{-1}(v)), \quad (u, v) \in [0, 1]^2,$$

where the inverse functions are the left-continuous quantile functions defined in the previous section. Section 3.9.4.4 of VW verifies that this map is Hadamard differentiable in a manner which permits developing inferential procedures for the copula function based on i.i.d. bivariate data.

The second example is multivariate trimming. Let P be a given probability distribution on \mathbb{R}^d , fix $\alpha \in (0, 1/2]$, and define \mathcal{H} to be the collection of all closed half-spaces in \mathbb{R}^d . The set $K_P \equiv \cap\{H \in \mathcal{H} : P(H) \geq 1 - \alpha\}$ can be easily shown to be compact and convex (see exercise 12.3.4). The α -trimmed mean is the quantity

$$T(P) \equiv \frac{1}{\lambda(K_P)} \int_{K_P} x d\lambda(x),$$

where λ is the Lebesgue measure on \mathbb{R}^d . Using non-trivial arguments, section 3.9.4.6 of VW shows how $P \mapsto T(P)$ can be formulated as a Hadamard differentiable functional of P and how this formulation can be applied to develop inference for $T(P)$ based on i.i.d. data from P .

There are many other important examples in statistics, some of which we will explore later on in this book, including a delta method formulation of Z-estimator theory which we will describe in the next chapter (chapter 13) and several other examples in the case studies of chapter 15.

12.3 Exercises

12.3.1. In the Wilcoxon statistic example of section 12.2.2, verify explicitly that every hypothesis of theorem 2.8 is satisfied.

12.3.2. Show that the product integral of A , $\phi(A)(s, t]$, is equivalent to the following:

- (a) The unique solution B of the following Volterra integral equation:

$$B(s, t] = 1 + \int_{(s, t]} B(s, u) A(du).$$

(b) The following **Peano series** representation:

$$\phi(A)(s, t] = 1 + \sum_{m=1}^{\infty} \int_{s < t_1 < \dots < t_m \leq t} A(dt_1) \cdots A(dt_m),$$

where the signed-measure interpretation of A is being used. Hint: Use the uniqueness from part (a).

(c) The unique solution B of the “backward” Volterra integral equation:

$$B(s, t] = 1 + \int_{(s, t]} B(u, t] A(du).$$

Hint: Start at t and go backwards in time to s .

12.3.3. Let ϕ be the product integral map of section 12.2.3. Show that if the total variation of A over the interval $(s, t]$ is M , then $|\phi(A)(s, t]| \leq e^M$. Hint: Recall that $\log(1 + x) \leq x$ for all $x > 0$.

12.3.4. Show that the set $K_P \subset \mathbb{R}^d$ defined in section 12.2.5 is compact and convex.

12.4 Notes

Much of the material of this chapter is inspired by chapter 3.9 of VW, although there is some new material and the method of presentation is different. Section 12.1 contains results from sections 3.9.1 and 3.9.3 of VW, although our results are specialized to Banach spaces (rather than the more general topological vector spaces). The examples of sections 12.2.1 through 12.2.4 are modified versions of the examples of sections 3.9.4.3, 3.9.4.1, 3.9.4.5 and 3.9.4.2, respectively, of VW. The order has been changed to emphasize a natural progression leading up to quantile inference based on the Kaplan-Meier estimator. Lemma 12.2 is a generalization of lemma 3.9.25 of VW, while lemmas 12.3 and 12.5 correspond to lemmas 3.9.17 and 3.9.30 of VW. Lemma 12.7 is a generalization of part (i) of lemma 3.9.23 of VW, while part (ii) of lemma 12.8 corresponds to part (ii) of lemma 3.9.23 of VW. Exercise 12.3.2 is based on exercises 3.9.5 and 3.9.6 of VW.

13

Z-Estimators

Recall from section 2.2.5 that Z-estimators are approximate zeros of data-dependent functions. These data-dependent functions, denoted Ψ_n , are maps between a possibly infinite dimensional normed parameter space Θ and a normed space \mathbb{L} , where the respective norms are $\|\cdot\|$ and $\|\cdot\|_{\mathbb{L}}$. The Ψ_n are frequently called estimating equations. A quantity $\hat{\theta}_n \in \Theta$ is a Z-estimator if $\|\Psi_n(\hat{\theta}_n)\|_{\mathbb{L}} \xrightarrow{P} 0$. In this chapter, we extend and prove the results of section 2.2.5. As part of this, we extend the Z-estimator master theorem, theorem 10.16, to the infinite dimensional case, although we divide the result into two parts, one for consistency and one for weak convergence.

We first discuss consistency and present a Z-estimator master theorem for consistency. We then discuss weak convergence and examine closely the special case of Z-estimators which are empirical measures of Donsker classes. We then use this structure to develop a Z-estimator master theorem for weak convergence. Both master theorems, the one for consistency and the one for weak convergence, will include results for the bootstrap. Finally, we demonstrate how Z-estimators can be viewed as Hadamard differentiable functionals of the involved estimating equations and how this structure enables use of a modified delta method to obtain very general results for Z-estimators. Recall from section 2.2.5 that the Kaplan-Meier estimator is an important and instructive example of a Z-estimator. A more sophisticated example, which will be presented later in the case studies of chapter 15, is the nonparametric maximum likelihood estimator for the proportional odds survival model.

13.1 Consistency

The main consistency result we have already presented in theorem 2.10 of section 2.2.5, and the proof of this theorem was given as an exercise (exercise 2.4.2). We will now extend this result to the bootstrapped Z-estimator. First, we restate the identifiability condition of theorem 2.10: The map $\Psi : \Theta \mapsto \mathbb{L}$ is identifiable at $\theta_0 \in \Theta$ if

$$(13.1) \quad \|\Psi(\theta_n)\|_{\mathbb{L}} \rightarrow 0 \text{ implies } \|\theta_n - \theta_0\| \text{ for any } \{\theta_n\} \in \Theta.$$

Note that there are alternative identifiability conditions that will also work, including the stronger condition that both $\Psi(\theta_0) = 0$ and $\Psi : \Theta \mapsto \mathbb{L}$ be one-to-one. Nevertheless, condition (13.1) seems to be the most efficient for our purposes.

In what follows, we will use the bootstrap-weighted empirical process \mathbb{P}_n° to denote either the nonparametric bootstrapped empirical process (with multinomial weights) or the multiplier bootstrapped empirical process defined by $f \mapsto \mathbb{P}_n^\circ f = n^{-1} \sum_{i=1}^n (\xi_i/\bar{\xi}) f(X_i)$, where ξ_1, \dots, ξ_n are i.i.d. positive weights with $0 < \mu = \mathbb{E}\xi_1 < \infty$ and $\bar{\xi} = n^{-1} \sum_{i=1}^n \xi_i$. Note that this is a special case of the weighted bootstrap introduced in theorem 10.13 but with the addition of $\bar{\xi}$ in the denominator. We leave it as an exercise (exercise 13.4.1) to verify that the conclusions of theorem 10.13 are not affected by this addition. Let $\mathcal{X}_n \equiv \{X_1, \dots, X_n\}$ as given in theorem 10.13. The following is the main result of this section:

THEOREM 13.1 (*Master Z-estimator theorem for consistency*) *Let $\theta \mapsto \Psi(\theta) = P\psi_\theta$, $\theta \mapsto \Psi_n(\theta) = \mathbb{P}_n\psi_\theta$ and $\theta \mapsto \Psi_n^\circ(\theta) = \mathbb{P}_n^\circ\psi_\theta$, where Ψ satisfies (13.1) and the class $\{\psi_\theta : \theta \in \Theta\}$ is P-Glivenko-Cantelli. Then, provided $\|\Psi_n(\hat{\theta}_n)\|_{\mathbb{L}} = o_P(1)$ and*

$$(13.2) \quad \mathbb{P} \left(\|\Psi_n^\circ(\hat{\theta}_n)\|_{\mathbb{L}} > \eta \mid \mathcal{X}_n \right) = o_P(1) \text{ for every } \eta > 0,$$

we have both $\|\hat{\theta}_n - \theta_0\| = o_P(1)$ and $\mathbb{P} \left(\|\hat{\theta}_n^\circ - \theta_0\| > \eta \mid \mathcal{X}_n \right) = o_P(1)$ for every $\eta > 0$.

Note in (13.2) the absence of an outer probability on the left side. This is because, as argued in section 2.2.3, a Lipschitz continuous map of either of these bootstrapped empirical processes is measurable with respect to the random weights conditional on the data.

Proof of theorem 13.1. The result that $\|\hat{\theta}_n - \theta_0\| = o_P(1)$ is a conclusion from theorem 2.10. For the conditional bootstrap result, (13.2) implies that for some sequence $\eta_n \downarrow 0$, $P \left(\|\Psi(\hat{\theta}_n^\circ)\|_{\mathbb{L}} > \eta_n \mid \mathcal{X}_n \right) = o_P(1)$, since

$$P \left(\sup_{\theta \in \Theta} \|\Psi_n^\circ(\theta) - \Psi(\theta)\| > \eta \mid \mathcal{X}_n \right) = o_P(1)$$

for all $\eta > 0$, by theorems 10.13 and 10.15. Thus, for any $\epsilon > 0$,

$$\begin{aligned} \mathbb{P}\left(\|\hat{\theta}_n^\circ - \theta_0\| > \epsilon \mid \mathcal{X}_n\right) &\leq \mathbb{P}\left(\|\hat{\theta}_n^\circ - \theta_0\| > \epsilon, \|\Psi(\hat{\theta}_n^\circ)\|_{\mathbb{L}} \leq \eta_n \mid \mathcal{X}_n\right) \\ &\quad + \mathbb{P}\left(\|\Psi(\hat{\theta}_n^\circ)\|_{\mathbb{L}} > \eta_n \mid \mathcal{X}_n\right) \\ &\xrightarrow{\mathbb{P}} 0, \end{aligned}$$

since the identifiability condition (13.1) implies that for all $\delta > 0$ there exists an $\eta > 0$ such that $\|\Psi(\theta)\|_{\mathbb{L}} < \eta$ implies $\|\theta - \theta_0\| < \delta$. Hence it is impossible for there to exist any $\theta \in \Theta$ such that both $\|\theta - \theta_0\| > \epsilon$ and $\|\Psi(\theta)\|_{\mathbb{L}} < \eta_n$ for all $n \geq 1$. The conclusion now follows since ϵ was arbitrary. \square

Note that we might have worked toward obtaining outer almost sure results since we are making a strong Glivenko-Cantelli assumption for the class of functions involved. However, we only need convergence in probability for statistical applications. Notice also that we only assumed $\|\Psi_n^\circ(\hat{\theta}_n^\circ)\|$ goes to zero conditionally rather than unconditionally as done in theorem 10.16. However, it seems to be easier to check the conditional version in practice. Moreover, the unconditional version is actually stronger than the conditional version, since

$$\mathbb{E}^* \mathbb{P}\left(\|\Psi_n^\circ(\hat{\theta}_n^\circ)\|_{\mathbb{L}} > \eta \mid \mathcal{X}_n\right) \leq \mathbb{P}^*\left(\|\Psi_n^\circ(\hat{\theta}_n^\circ)\|_{\mathbb{L}} > \eta\right)$$

by the version of Fubini's theorem given as theorem 6.14. It is unclear how to extend this argument to the outer almost sure setting. This is another reason for restricting our attention to the convergence in probability results. Nevertheless, we still need the strong Glivenko-Cantelli assumption since this enables the use of theorems 10.13 and 10.15.

13.2 Weak Convergence

In this section, we first provide general results for Z-estimators which may not be based on i.i.d. data. We then present sufficient conditions for the i.i.d. case when the estimating equation is an empirical measure ranging over a Donsker class. Finally, we give a master theorem for Z-estimators based on i.i.d. data which includes bootstrap validity.

13.2.1 The General Setting

We now prove theorem 2.11 and give a method of weakening the differentiability requirement for Ψ . An important thing to note is that no assumptions about the data being i.i.d. are required. The proof follows closely the proof of theorem 3.3.1 given in VW.

Proof of theorem 2.11. By the definitions of $\hat{\theta}_n$ and θ_0 ,

$$\begin{aligned} (13.3) \quad \sqrt{n} \left(\Psi(\hat{\theta}_n) - \Psi(\theta_0) \right) &= -\sqrt{n} \left(\Psi_n(\hat{\theta}_n) - \Psi(\hat{\theta}_n) \right) + o_P(1) \\ &= -\sqrt{n}(\Psi_n - \Psi)(\theta_0) + o_P(1 + \sqrt{n}\|\hat{\theta}_n - \theta_0\|), \end{aligned}$$

by assumption (2.12). Note the error terms throughout this theorem are with respect to the norms of the spaces, e.g. Θ or \mathbb{L} , involved. Since $\dot{\Psi}_{\theta_0}$ is continuously invertible, we have by part (i) of lemma 6.16 that there exists a constant $c > 0$ such that $\|\dot{\Psi}_{\theta_0}(\theta - \theta_0)\| \geq c\|\theta - \theta_0\|$ for all θ and θ_0 in $\overline{\text{lin}} \Theta$. Combining this with the differentiability of Ψ yields $\|\Psi(\theta) - \Psi(\theta_0)\| \geq c\|\theta - \theta_0\| + o(\|\theta - \theta_0\|)$. Combining this with (13.3), we obtain

$$\sqrt{n}\|\hat{\theta}_n - \theta_0\|(c + o_P(1)) \leq O_P(1) + o_P(1 + \sqrt{n}\|\hat{\theta}_n - \theta_0\|).$$

We now have that $\hat{\theta}_n$ is \sqrt{n} -consistent for θ_0 with respect to $\|\cdot\|$. By the differentiability of Ψ , the left side of (13.3) can be replaced by $\sqrt{n}\dot{\Psi}_{\theta_0}(\hat{\theta}_n - \theta_0) + o_P(1 + \sqrt{n}\|\hat{\theta}_n - \theta_0\|)$. This last error term is now $o_P(1)$ as also is the error term on the right side of (13.3). Now the result (2.13) follows. Next the continuity of $\dot{\Psi}_{\theta_0}^{-1}$ and the continuous mapping theorem yield $\sqrt{n}\hat{\theta}_n - \theta_0 \rightsquigarrow -\dot{\Psi}_{\theta_0}^{-1}(Z)$ as desired. \square

The following lemma allows us to weaken the Fréchet differentiability requirement to Hadamard differentiability when it is also known that $\sqrt{n}(\hat{\theta}_n - \theta_0)$ is asymptotically tight:

LEMMA 13.2 *Assume the conditions of theorem 2.11 except that consistency of $\hat{\theta}_n$ is strengthened to asymptotic tightness of $\sqrt{n}(\hat{\theta}_n - \theta_0)$ and the Fréchet differentiability of Ψ is weakened to Hadamard differentiability at θ_0 . Then the results of theorem 2.11 still hold.*

Proof. The asymptotic tightness of $\sqrt{n}(\hat{\theta}_n - \theta_0)$ enables expression (13.3) to imply $\sqrt{n} \left(\Psi(\hat{\theta}_n) - \Psi(\theta_0) \right) = -\sqrt{n}(\Psi_n - \Psi)(\theta_0) + o_P(1)$. The Hadamard differentiability of Ψ yields $\sqrt{n} \left(\Psi(\hat{\theta}_n) - \Psi(\theta_0) \right) = \sqrt{n}\dot{\Psi}_{\theta_0}(\hat{\theta}_n - \theta_0) + o_P(1)$. Combining, we now have $\sqrt{n}\dot{\Psi}_{\theta_0}(\hat{\theta}_n - \theta_0) = -\sqrt{n}(\Psi_n - \Psi)(\theta_0) + o_P(1)$, and all of the results of the theorem follow. \square

13.2.2 Using Donsker Classes

We now consider the special case where the data involved are i.i.d., i.e., $\Psi_n(\theta)(h) = \mathbb{P}_n \psi_{\theta,h}$ and $\Psi(\theta)(h) = P \psi_{\theta,h}$, for measurable functions $\psi_{\theta,h}$, where h ranges over an index set \mathcal{H} . The following lemma gives us reasonably verifiable sufficient conditions for (2.12) to hold:

LEMMA 13.3 *Suppose the class of functions*

$$(13.4) \quad \{\psi_{\theta,h} - \psi_{\theta_0,h} : \|\theta - \theta_0\| < \delta, h \in \mathcal{H}\}$$

is P -Donsker for some $\delta > 0$ and

$$(13.5) \quad \sup_{h \in \mathcal{H}} P(\psi_{\theta, h} - \psi_{\theta_0, h})^2 \rightarrow 0, \text{ as } \theta \rightarrow \theta_0.$$

Then if $\Psi_n(\hat{\theta}_n) = o_P(n^{-1/2})$ and $\hat{\theta}_n \xrightarrow{P} \theta_0$, $\sqrt{n}(\Psi_n - \Psi)(\hat{\theta}_n) - \sqrt{n}(\Psi_n - \Psi)(\theta_0) = o_P(1)$.

Before giving the proof of this lemma, we make the somewhat trivial observation that the conclusion of this lemma implies (2.12).

Proof of lemma 13.3. Let $\Theta_\delta \equiv \{\theta : \|\theta - \theta_0\| < \delta\}$ and define the extraction function $f : \ell^\infty(\Theta_\delta \times \mathcal{H}) \times \Theta_\delta \mapsto \ell^\infty(\mathcal{H})$ as $f(z, \theta)(h) \equiv z(\theta, h)$, where $z \in \ell^\infty(\Theta_\delta \times \mathcal{H})$. Note that f is continuous at every point (z, θ_1) such that $\sup_{h \in \mathcal{H}} |z(\theta, h) - z(\theta_1, h)| \rightarrow 0$ as $\theta \rightarrow \theta_1$. Define the stochastic process $Z_n(\theta, h) \equiv \mathbb{G}_n(\psi_{\theta, h} - \psi_{\theta_0, h})$ indexed by $\Theta_\delta \times \mathcal{H}$. As assumed, the process Z_n converges weakly in $\ell^\infty(\Theta_\delta \times \mathcal{H})$ to a tight Gaussian process Z_0 with continuous sample paths with respect to the metric ρ defined by $\rho^2((\theta_1, h_1), (\theta_2, h_2)) = P(\psi_{\theta_1, h_1} - \psi_{\theta_0, h_1} - \psi_{\theta_2, h_2} + \psi_{\theta_0, h_2})^2$. Since, $\sup_{h \in \mathcal{H}} \rho((\theta, h), (\theta_0, h)) \rightarrow 0$ by assumption, we have that f is continuous at almost all sample paths of Z_0 . By Slutsky's theorem (theorem 7.15), $(Z_n, \hat{\theta}_n) \rightsquigarrow (Z_0, \theta_0)$. The continuous mapping theorem (theorem 7.7) now implies that $Z_n(\hat{\theta}_n) = f(Z_n, \hat{\theta}_n) \rightsquigarrow f(Z_0, \theta_0) = 0$. \square

If, in addition to the assumptions of lemma 13.3, we are willing to assume

$$(13.6) \quad \{\psi_{\theta_0, h} : h \in \mathcal{H}\}$$

is P -Donsker, then $\sqrt{n}(\Psi_n - \Psi)(\theta_0) \rightsquigarrow Z$, and all of the weak convergence assumptions of theorem 2.11 are satisfied. Alternatively, we could just assume that

$$(13.7) \quad \mathcal{F}_\delta \equiv \{\psi_{\theta, h} : \|\theta - \theta_0\| < \delta, h \in \mathcal{H}\}$$

is P -Donsker for some $\delta > 0$, then both (13.4) and (13.6) are P -Donsker for some $\delta > 0$. We are now well poised for a Z -estimator master theorem for weak convergence.

13.2.3 A Master Theorem and the Bootstrap

In this section, we augment the results of the previous section to achieve a general Z -estimator master theorem that includes both weak convergence and validity of the bootstrap. Here we consider the two bootstrapped Z -estimators described in section 13.1, except that for the multiplier bootstrap we make the additional requirements that $0 < \tau^2 = \text{var}(\xi_1) < \infty$ and $\|\xi_1\|_{2,1} < \infty$. We use $\overset{P}{\rightsquigarrow}$ to denote either $\overset{P}{\rightsquigarrow}_\xi$ or $\overset{P}{\rightsquigarrow}_W$ depending on which bootstrap is being used, and we let the constant $k_0 = \tau/\mu$ for the multiplier bootstrap and $k_0 = 1$ for the multinomial bootstrap. Here is the main result:

THEOREM 13.4 Assume $\Psi(\theta_0) = 0$ and the following hold:

- (A) $\theta \mapsto \Psi(\theta)$ satisfies (13.1);
- (B) The class $\{\psi_{\theta,h}; \theta \in \Theta, h \in \mathcal{H}\}$ is P -Glivenko-Cantelli;
- (C) The class \mathcal{F}_δ in (13.7) is P -Donsker for some $\delta > 0$;
- (D) Condition (13.5) holds;
- (E) $\|\Psi_n(\hat{\theta}_n)\|_{\mathbb{L}} = o_P(n^{-1/2})$ and $P\left(\sqrt{n}\|\Psi_n^\circ(\hat{\theta}_n^\circ)\|_{\mathbb{L}} > \eta \mid \mathcal{X}_n\right) = o_P(1)$ for every $\eta > 0$;
- (F) $\theta \mapsto \Psi(\theta)$ is Fréchet-differentiable at θ_0 with continuously invertible derivative $\dot{\Psi}_{\theta_0}$.

Then $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow -\dot{\Psi}_{\theta_0}^{-1}Z$, where $Z \in \ell^\infty(\mathcal{H})$ is the tight, mean zero Gaussian limiting distribution of $\sqrt{n}(\Psi_n - \Psi)(\theta_0)$, and $\sqrt{n}(\hat{\theta}_n^\circ - \hat{\theta}_n) \overset{P}{\rightsquigarrow} k_0 Z$.

Condition (A) is identifiability. Conditions (B) and (C) are consistency and asymptotic normality conditions for the estimating equation. Condition (D) is an asymptotic equicontinuity condition for the estimating equation at θ_0 . Condition (E) simply states that the estimators are approximate zeros of the estimating equation, while condition (F) specifies the smoothness and invertibility requirements of the derivative of Ψ . Except for the last half of condition (E), all of the conditions are requirements for asymptotic normality of $\sqrt{n}(\hat{\theta}_n - \theta_0)$. What is perhaps surprising is how little additional assumptions are needed to obtain bootstrap validity. Only an assurance that the bootstrapped estimator is an approximate zero of the bootstrapped estimating equation is required. Thus bootstrap validity is almost an automatic consequence of asymptotic normality.

Proof of theorem 13.4. The consistency of $\hat{\theta}_n$ and weak convergence of $\sqrt{n}(\hat{\theta}_n - \theta_0)$ follow from theorems 13.1 and 13.4 and lemma 13.3. Theorem 13.1 also yields that there exists a decreasing sequence $0 < \eta_n \downarrow 0$ such that

$$P\left(\|\hat{\theta}_n^\circ - \theta_0\| > \eta_n \mid \mathcal{X}_n\right) = o_P(1).$$

Now we can use the same arguments used in the proof of lemma 13.3, in combination with theorem 2.6, to obtain that $\sqrt{n}(\Psi_n^\circ - \Psi)(\hat{\theta}_n^\circ) - \sqrt{n}(\Psi_n^\circ - \Psi)(\theta_0) = E_n$, where $P(E_n > \eta \mid \mathcal{X}_n) = o_P(1)$ for all $\eta > 0$. Combining this with arguments used in the proof of theorem 2.11, we can deduce that $\sqrt{n}(\hat{\theta}_n^\circ - \theta_0) = -\dot{\Psi}_{\theta_0}^{-1}\sqrt{n}(\Psi_n^\circ - \Psi)(\theta_0) + E'_n$, where $P(E'_n > \eta \mid \mathcal{X}_n) = o_P(1)$ for all $\eta > 0$. Combining this with the conclusion of theorem 2.11, we obtain $\sqrt{n}(\hat{\theta}_n^\circ - \hat{\theta}_n) = -\dot{\Psi}_{\theta_0}^{-1}\sqrt{n}(\Psi_n^\circ - \Psi_n)(\theta_0) + E''_n$, where $P(E''_n > \eta \mid \mathcal{X}_n) = o_P(1)$ for all $\eta > 0$. The final conclusion now follows from reapplication of theorem 2.6. \square

13.3 Using the Delta Method

There is an alternative approach to Z-estimators which may be more effective for more general data settings, including non-i.i.d. and dependent data settings. The idea is to view the extraction of the zero from the estimating equation as a continuous mapping. Our approach is closely related to the approach in section 3.9.4.7 of VW but with some important modifications which simplify the required assumptions. We require Θ to be the subset of a Banach space and \mathbb{L} to be a Banach space. Let $\ell^\infty(\Theta, \mathbb{L})$ be the Banach space of all uniformly norm-bounded functions $z : \Theta \mapsto \mathbb{L}$. Let $Z(\Theta, \mathbb{L})$ be the subset consisting of all maps with at least one zero, and let $\Phi(\Theta, \mathbb{L})$ be the collection of all maps (or algorithms) $\phi : Z(\Theta, \mathbb{L}) \mapsto \Theta$ that for each element $z \in Z(\Theta, \mathbb{L})$ extract one of its zeros $\phi(z)$. This structure allows for multiple zeros.

The following lemma gives us a kind of uniform Hadamard differentiability of members of $\Phi(\Theta, \mathbb{L})$ which we will be able to use to obtain a delta method result for Z-estimators $\hat{\theta}_n$ that satisfy $\Psi_n(\hat{\theta}_n) = o_P(r_n^{-1})$ for some sequence $r_n \rightarrow \infty$ for which $X_n(\theta) \equiv r_n(\Psi_n - \Psi)(\theta)$ converges weakly to a tight, random element in $X \in \ell^\infty(\Theta_0, \mathbb{L})$, where $\Theta_0 \subset \Theta$ is an open neighborhood of θ_0 and $\|X(\theta) - X(\theta_0)\|_{\mathbb{L}} \rightarrow 0$ as $\theta \rightarrow \theta_0$ almost surely, i.e., X has continuous sample paths in θ . Define $\ell_0^\infty(\Theta, \mathbb{L})$ to be the elements $x \in \ell^\infty(\Theta, \mathbb{L})$ for which $\|x(\theta) - x(\theta_0)\|_{\mathbb{L}} \rightarrow 0$ as $\theta \rightarrow \theta_0$.

THEOREM 13.5 *Assume $\Psi : \Theta \mapsto \mathbb{L}$ is uniformly norm-bounded over Θ , $\Psi(\theta_0) = 0$, and condition (13.1) holds. Let Ψ also be Fréchet differentiable at θ_0 with continuously invertible derivative $\dot{\Psi}_{\theta_0}$. Then the continuous linear operator $\phi'_\Psi : \ell_0^\infty(\Theta, \mathbb{L}) \mapsto \text{lin } \Theta$ defined by $z \mapsto \phi'_\Psi(z) \equiv -\dot{\Psi}_{\theta_0}^{-1}(z(\theta_0))$ satisfies:*

$$\sup_{\phi \in \Phi(\Theta, \mathbb{L})} \left\| \frac{\phi(\Psi + t_n z_n) - \phi(\Psi)}{t_n} - \phi'_\Psi(z(\theta_0)) \right\| \rightarrow 0,$$

as $n \rightarrow \infty$, for any sequences $(t_n, z_n) \in (0, \infty) \times \ell^\infty(\Theta, \mathbb{L})$ such that $t_n \downarrow 0$, $\Psi + t_n z_n \in Z(\Theta, \mathbb{L})$, and $z_n \rightarrow z \in \ell_0^\infty(\Theta, \mathbb{L})$.

Proof. Let $0 < t_n \downarrow 0$ and $z_n \rightarrow z \in \ell_0^\infty(\Theta, \mathbb{L})$ such that $\Psi + t_n z_n \in Z(\Theta, \mathbb{L})$. Choose any sequence $\{\phi_n\} \in \Phi(\Theta, \mathbb{L})$, and note that $\theta_n \equiv \phi_n(\Psi + t_n z_n)$ satisfies $\Psi(\theta_n) + t_n z_n = 0$ by construction. Hence $\Psi(\theta_n) = O(t_n)$. By condition (13.1), $\theta_n \rightarrow \theta_0$. By the Fréchet differentiability of Ψ ,

$$\liminf_{n \rightarrow \infty} \frac{\|\Psi(\theta_n) - \Psi(\theta_0)\|_{\mathbb{L}}}{\|\theta_n - \theta_0\|} \geq \liminf_{n \rightarrow \infty} \frac{\|\dot{\Psi}_{\theta_0}(\theta_n - \theta_0)\|_{\mathbb{L}}}{\|\theta_n - \theta_0\|} \geq \inf_{\|g\|=1} \|\dot{\Psi}_{\theta_0}(g)\|_{\mathbb{L}},$$

where g ranges over $\text{lin } \Theta$. Since the inverse of $\dot{\Psi}_{\theta_0}$ is continuous, the right side of the above is positive. Thus there exists a universal constant $c < \infty$ (depending only on $\dot{\Psi}_{\theta_0}$ and $\text{lin } \Theta$) for which $\|\theta_n - \theta_0\| < c\|\Psi(\theta_n) -$

$\Psi(\theta_0)\|_{\mathbb{L}} = c\|t_n z_n(\theta_n)\|_{\mathbb{L}}$ for all n large enough. Hence $\|\theta_n - \theta_0\| = O(t_n)$. By Fréchet differentiability, $\Psi(\theta_n) - \Psi(\theta_0) = \dot{\Psi}_{\theta_0}(\theta_n - \theta_0) + o(\|\theta_n - \theta_0\|)$, where $\dot{\Psi}_{\theta_0}$ is linear and continuous on $\text{lin } \Theta$. The remainder term is $o(t_n)$ by previous arguments. Combining this with the fact that $t_n^{-1}(\Psi(\theta_n) - \Psi(\theta_0)) = -z_n(\theta_n) \rightarrow z(\theta_0)$, we obtain

$$\frac{\theta_n - \theta_0}{t_n} = \dot{\Psi}_{\theta_0}^{-1} \left(\frac{\Psi(\theta_n) - \Psi(\theta_0)}{t_n} + o(1) \right) \rightarrow -\dot{\Psi}_{\theta_0}^{-1}(z(\theta_0)).$$

The conclusion now follows since the sequence ϕ_n was arbitrary. \square

The following simple corollary allows the delta method to be applied to Z-estimators. Let $\tilde{\phi} : \ell^\infty(\Theta, \mathbb{L}) \mapsto \Theta$ be a map such that for each $x \in \ell^\infty(\Theta, \mathbb{L})$, $\tilde{\phi}(x) = \theta_1 \neq \theta_0$ when $x \notin Z(\Theta, \mathbb{L})$ and $\tilde{\phi}(x) = \phi(x)$ for some $\phi \in \Phi(\Theta, \mathbb{L})$ otherwise.

COROLLARY 13.6 *Suppose Ψ satisfies the conditions of theorem 13.5, $\hat{\theta}_n = \tilde{\phi}(\Psi_n)$, and Ψ_n has at least one zero for all n large enough, outer almost surely. Suppose also that $r_n(\Psi_n - \Psi) \rightsquigarrow X$ in $\ell^\infty(\Theta, \mathbb{L})$, with X tight and $\|X(\theta)\|_{\mathbb{L}} \rightarrow 0$ as $\theta \rightarrow \theta_0$ almost surely. Then $r_n(\hat{\theta}_n - \theta_0) \rightsquigarrow -\dot{\Psi}_{\theta_0}^{-1}X(\theta_0)$.*

Proof. Since Ψ_n has a zero for all n large enough, outer almost surely, we can, without loss of generality, assume that Ψ_n has a zero for all n . Thus we can assume that $\tilde{\phi} \in \Phi(\Theta, \mathbb{L})$. By theorem 13.5, we know that $\tilde{\phi}$ is Hadamard differentiable tangentially to $\ell_0^\infty(\Theta, \mathbb{L})$, which, by assumption, contains X with probability 1. Thus theorem 2.8 applies, and the desired result follows. \square

We leave it as an exercise (see exercise 13.4.2 below) to develop a corollary which utilizes $\tilde{\phi}$ to obtain a bootstrap result for Z-estimators. A drawback with this approach is that root finding algorithms in practice are seldom exact, and room needs to be allowed for computational error. The following corollary of theorem 13.5 yields a very general Z-estimator result based on a modified delta method. We make the fairly realistic assumption that the Z-estimator $\hat{\theta}_n$ is computed from Ψ_n using a deterministic algorithm (e.g., a computer program) that is allowed to depend on n and which is not required to yield an exact root of Ψ_n .

COROLLARY 13.7 *Suppose Ψ satisfies the conditions of theorem 13.5, and $\hat{\theta}_n = A_n(\Psi_n)$ for some sequence of deterministic algorithms $A_n : \ell^\infty(\Theta, \mathbb{L}) \mapsto \Theta$ and random sequence $\Psi_n : \Theta \mapsto \mathbb{L}$ of estimating equations such that $\Psi_n \xrightarrow{P} \Psi$ in $\ell^\infty(\Theta, \mathbb{L})$ and $\Psi_n(\hat{\theta}_n) = o_P(r_n^{-1})$, where $0 < r_n \rightarrow \infty$ is a sequence of constants for which $r_n(\Psi_n - \Psi) \rightsquigarrow X$ in $\ell^\infty(\Theta_0, \mathbb{L})$ for some closed $\Theta_0 \subset \Theta$ containing an open neighborhood of θ_0 , with X tight and $\|X(\theta)\|_{\mathbb{L}} \rightarrow 0$ as $\theta \rightarrow \theta_0$ almost surely. Then $r_n(\hat{\theta}_n - \theta_0) \rightsquigarrow -\dot{\Psi}_{\theta_0}^{-1}X(\theta_0)$.*

Proof. Let $X_n \equiv r_n(\Psi_n - \Psi)$. By theorem 7.26, there exists a new probability space $(\tilde{\Omega}, \tilde{\mathcal{A}}, \tilde{P})$ on which: $E^*f(\tilde{\Psi}_n) = E^*f(\Psi_n)$ for all bounded

$f : \ell^\infty(\Theta, \mathbb{L}) \mapsto \mathbb{R}$ and all $n \geq 1$; $\tilde{\Psi}_n \xrightarrow{\text{as}^*} \Psi$ in $\ell^\infty(\Theta, \mathbb{L})$; $r_n(\tilde{\Psi}_n - \Psi) \xrightarrow{\text{as}^*} \tilde{X}$ in $\ell^\infty(\Theta_0, \mathbb{L})$; \tilde{X} and X have the same distributions; and $r_n(A_n(\tilde{\Psi}_n) - \theta_0)$ and $r_n(A_n(\Psi_n) - \theta_0)$ have the same distributions. Note that for any bounded f , $\tilde{\Psi}_n \mapsto f(\tilde{\Psi}_n(A_n(\tilde{\Psi}_n))) = g(\tilde{\Psi}_n)$ for some bounded g . Thus $\tilde{\Psi}_n(\tilde{\theta}_n) = o_{\tilde{P}}(r_n^{-1})$ for $\tilde{\theta}_n \equiv A_n(\tilde{\Psi}_n)$.

Hence for any subsequence n' there exists a further subsequence n'' such that $\tilde{\Psi}_{n''} \xrightarrow{\text{as}^*} \Psi$ in $\ell^\infty(\Theta, \mathbb{L})$, $r_{n''}(\tilde{\Psi}_{n''} - \Psi) \xrightarrow{\text{as}^*} \tilde{X}$ in $\ell^\infty(\Theta_0, \mathbb{L})$, and $\tilde{\Psi}_{n''}(\tilde{\theta}_{n''}) \xrightarrow{\text{as}^*} 0$ in \mathbb{L} . Thus also $\Psi(\tilde{\theta}_{n''}) \xrightarrow{\text{as}^*} 0$, which implies $\tilde{\theta}_{n''} \xrightarrow{\text{as}^*} 0$. Note that $\tilde{\theta}_{n''}$ is a zero of $\tilde{\Psi}_{n''}(\theta) - \tilde{\Psi}(\tilde{\theta}_{n''})$ by definition and is contained in Θ_0 for all n large enough. Hence, for all n large enough, $r_{n''}(\tilde{\theta}_{n''} - \theta_0) = r_{n''}(\phi_{n''}(\tilde{\Psi}_{n''} - \tilde{\Psi}_{n''}(\tilde{\theta}_{n''})) - \phi_{n''}(\Psi))$ for some sequence $\phi_n \in \Phi(\Theta_0, \mathbb{L})$ possibly dependent on sample realization $\tilde{\omega} \in \tilde{\Omega}$. This implies that for all n large enough,

$$\begin{aligned}
 & \left\| r_{n''}(\tilde{\theta}_{n''} - \theta_0) - \phi'_\Psi(\tilde{X}) \right\| \\
 & \leq \sup_{\phi \in \Phi(\Theta_0, \mathbb{L})} \left| r_{n''}(\phi(\tilde{\Psi}_{n''} - \tilde{\Psi}_{n''}(\tilde{\theta}_{n''})) - \phi(\Psi)) - \phi'_\Psi(\tilde{X}) \right| \\
 & \xrightarrow{\text{as}^*} 0,
 \end{aligned}$$

by theorem 13.5 (with Θ_0 replacing Θ). This implies $\left\| r_{n''}(A_{n''}(\tilde{\Psi}_{n''}) - \theta_0) - \phi'_\Psi(\tilde{X}) \right\| \xrightarrow{\text{as}^*} 0$. Since this holds for every subsequence, we have $r_n(A_n(\tilde{\Psi}_n) - \theta_0) \rightsquigarrow \phi'_\Psi(\tilde{X})$. This of course implies $r_n(A_n(\Psi_n) - \theta_0) \rightsquigarrow \phi'_\Psi(X)$. \square

The following corollary extends the previous result to generalized bootstrapped processes. Let Ψ_n° be a bootstrapped version of Ψ_n based on both the data sequence X_n (the data used in Ψ_n) and a sequence of weights $W = \{W_n, n \geq 1\}$.

COROLLARY 13.8 *Assume the conditions of corollary 13.7 and, in addition, that $\hat{\theta}_n^\circ = A_n(\Psi_n^\circ)$ for a sequence of bootstrapped estimating equations $\Psi_n^\circ(X_n, W_n)$, with $\Psi_n^\circ - \Psi \xrightarrow[W]{P} 0$ and $r_n \Psi_n^\circ(\hat{\theta}_n^\circ) \xrightarrow[W]{P} 0$ in $\ell^\infty(\Theta, \mathbb{L})$, and with $r_n c(\Psi_n^\circ - \Psi_n) \xrightarrow[W]{P} X$ in $\ell^\infty(\Theta_0, \mathbb{L})$ for some $0 < c < \infty$, where the maps $W_n \mapsto h(\Psi_n^\circ)$ are measurable for every $h \in C_b(\ell^\infty(\Theta, \mathbb{L}))$ outer almost surely. Then $r_n c(\hat{\theta}_n^\circ - \hat{\theta}_n) \xrightarrow[W]{P} \phi'_\Psi(X)$.*

Proof. This proof shares many similarities with the proof of theorem 12.1. To begin with, by using the same arguments used in the beginning of that proof, we can obtain that, unconditionally,

$$r_n \begin{pmatrix} \Psi_n^\circ - \Psi \\ \Psi_n - \Psi \end{pmatrix} \rightsquigarrow \begin{pmatrix} c^{-1} X'_1 + X'_2 \\ X'_2 \end{pmatrix}$$

in $\ell^\infty(\Theta_0, \mathbb{L})$, where X'_1 and X'_2 are two independent copies of X . We can also obtain that $\Psi_n^\circ \xrightarrow{P} \Psi$ in $\ell^\infty(\Theta, \mathbb{L})$ unconditionally. Combining this with

a minor adaptation of the above corollary 13.7 (see exercise 13.4.3 below) we obtain unconditionally that

$$r_n \begin{pmatrix} \hat{\theta}_n^\circ - \theta_0 \\ \hat{\theta}_n - \theta_0 \\ \Psi_n^\circ(\theta_0) - \Psi(\theta_0) \\ \Psi_n(\theta_0) - \Psi(\theta_0) \end{pmatrix} \rightsquigarrow \begin{pmatrix} \phi'_\Psi(c^{-1}X'_1 + X'_2) \\ \phi'_\Psi(X'_2) \\ c^{-1}X'_1(\theta_0) + X'_2(\theta_0) \\ X'_2(\theta_0) \end{pmatrix}.$$

This implies two things. First,

$$r_n c \begin{pmatrix} \hat{\theta}_n^\circ - \hat{\theta}_n \\ (\Psi_n^\circ - \Psi_n)(\theta_0) \end{pmatrix} \rightsquigarrow \begin{pmatrix} \phi'_\Psi(X) \\ X(\theta_0) \end{pmatrix}$$

unconditionally, since ϕ'_Ψ is linear on \mathbb{L} . Second, the usual continuous mapping theorem now yields unconditionally that

$$(13.8) \quad r_n c(\hat{\theta}_n^\circ - \hat{\theta}_n) + \dot{\Psi}_{\theta_0}^{-1}(r_n c(\Psi_n^\circ - \Psi_n)(\theta_0)) \xrightarrow{P} 0,$$

since the map $(x, y) \mapsto x + \dot{\Psi}_{\theta_0}^{-1}(y)$ is continuous on all of $\text{lin } \Theta \times \mathbb{L}$ (recall that $x \mapsto \phi'_\Psi(x) = -\dot{\Psi}_{\theta_0}^{-1}(x(\theta_0))$).

Now for any map $h \in C_b(\mathbb{L})$, the map $x \mapsto h(r_n c(x - \Psi_n(\theta_0)))$ is continuous and bounded for all $x \in \mathbb{L}$ outer almost surely. Thus the maps $W_n \mapsto h(r_n c(\Psi_n^\circ - \Psi_n)(\theta_0))$ are measurable for every $h \in C_b(\mathbb{L})$ outer almost surely. Hence the bootstrap continuous mapping theorem, theorem 10.8, yields that $\dot{\Psi}_{\theta_0}^{-1}(r_n c(\hat{\Psi}_n^\circ - \Psi_n)(\theta_0)) \xrightarrow[W]{P} \dot{\Psi}_{\theta_0}^{-1}(X)$. The desired result now follows from (13.8). \square

An interesting application of the above results is to estimating equations for empirical processes from dependent data as discussed in section 11.6. Specifically, suppose $\Psi_n(\theta)(h) = \mathbb{P}_n \psi_{\theta,h}$, where the stationary sample data X_1, X_2, \dots and $\mathcal{F} = \{\psi_{\theta,h} : \theta \in \Theta, h \in \mathcal{H}\}$ satisfy the conditions of theorem 11.22 with marginal distribution P , and let $\Psi(\theta)(h) = P\psi_{\theta,h}$. Then the conclusion of theorem 11.22 is that $\sqrt{n}(\Psi_n - \Psi) \rightsquigarrow \mathbb{H}$ in $\ell^\infty(\Theta \times \mathcal{H})$, where \mathbb{H} is a tight, mean zero Gaussian process. Provided Ψ satisfies the conditions of theorem 13.1, and provided a few other conditions hold, corollary 13.7 will give us weak convergence of the standardized Z-estimators $\sqrt{n}(\hat{\theta}_n - \theta_0)$ based on Ψ_n . Under certain regularity conditions, a moving blocks bootstrapped estimation equation Ψ_n° can be shown by theorem 11.24 to satisfy the requirements of corollary 13.8. This enables valid bootstrap estimation of the limiting distribution of $\sqrt{n}(\hat{\theta}_n - \theta_0)$. These results can also be extended to stationary sequences with long range dependence, where the normalizing rate r_n may differ from \sqrt{n} .

13.4 Exercises

13.4.1. Show that the addition of $\bar{\xi}$ in the denominator of the weights in the weighted bootstrap introduced in theorem 10.13, as discussed in section 13.1, does not affect the conclusions of that theorem.

13.4.2. Develop a bootstrap central limit theorem for Z -estimators based theorem 12.1 which utilizes the Hadamard differentiability of the zero-extraction map $\tilde{\phi}$ used in corollary 13.6.

13.4.3. Verify the validity of the “minor adaptation” of corollary 13.7 used in the proof of corollary 13.8.

13.5 Notes

Theorem 2.11 and lemma 13.2 are essentially a decomposition of theorem 3.3.1 of VW into two parts. Lemma 13.3 is lemma 3.3.5 of VW.

14

M-Estimators

M-estimators, as introduced in section 2.2.6, are approximate maximizers of objective functions computed from data. Note the estimators based on minimizing objective functions are trivially also M-estimators after taking the negative of the objective function. In some ways, M-estimators are more basic than Z-estimators since Z-estimators can always be expressed as M-estimators. The reverse is not true, however, since there are M-estimators which cannot be effectively formulated as Z-estimators. Nevertheless, Z-estimator theory is usually much easier to use whenever it can be applied. The focus, then, of this chapter is on M-estimator settings for which it is not practical to directly use Z-estimator theory. The usual issues for M-estimation theory are to establish consistency, determine the correct rate of convergence, establish weak convergence, and, finally, to conduct inference.

We first present a key result central to M-estimation theory, the argmax theorem, which permits deriving weak limits of M-estimators as the the argmax of the limiting process. This is useful for both consistency, which we discuss next, and weak convergence. The section on consistency includes a proof of theorem 2.12. We then discuss how to determine the correct rate of convergence which is necessary for establishing weak convergence based on the argmax theorem. We then present general results for “regular estimators,” i.e., estimators whose rate of convergence is \sqrt{n} . We then give several examples which illustrate M-estimation theory for non-regular estimators which have rates distinct from \sqrt{n} . Much of the content of this chapter is inspired by chapter 3.2 of VW.

The M-estimators in both the regular and non-regular examples we present will be i.i.d. empirical processes of the form $M_n(\theta) = \mathbb{P}_n m_\theta$ for

a class of measurable, real valued functions $\mathcal{M} = \{m_\theta : \theta \in \Theta\}$, where the parameter space Θ is usually a subset of a semimetric space. The entropy of the class \mathcal{M} plays a crucial role in determining the proper rate of convergence. The aspect of the rate of convergence determining process is often the most difficult part technically in M-estimation theory and usually requires fairly precise bounds on moments of the empirical processes involved, such as those described in section 11.1. We note that our presentation involves only a small amount of the useful results in the area. Much more of these kinds of results can be found in chapter 3.4 of VW and in van de Geer (2000).

Inference for M-estimators is more challenging than it is for Z-estimators because the bootstrap is not in general automatically valid, especially when the convergence rate is not \sqrt{n} . On the other hand, subsampling m out of n observations (see Politis and Romano, 1994) can be shown to be universally valid, provided $m \rightarrow \infty$ and $m/n \rightarrow 0$. However, even this result is not entirely satisfactory because it requires n to be quite large since m must also be large yet small relative to n . Bootstrapping and other methods of inference for M-estimators is currently an area of active research, but we do not pursue it further in this chapter.

14.1 The Argmax Theorem

We now consider a sequence $\{M_n(h) : h \in H\}$ of stochastic processes indexed by a metric space H . Let \hat{h}_n denote an M-estimator obtained by nearly-maximizing M_n . The idea of the argmax theorem presented below is that under reasonable regularity conditions, when $M_n \rightsquigarrow M$, where M is another stochastic process in $\ell^\infty(H)$, that $\hat{h}_n \rightsquigarrow \hat{h}$, where \hat{h} is the argmax of M . If we know that the rate of convergence of an M-estimator $\hat{\theta}_n$ is r_n (a nondecreasing, positive sequence), where $\hat{\theta}_n$ is the argmax of $\theta \mapsto M_n(\theta)$, then $\hat{h}_n = r_n(\hat{\theta}_n - \theta_0)$ can be expressed as the argmax of $h \mapsto \tilde{M}_n(h) \equiv r_n [M_n(\theta_0 + h/r_n) - M_n(\theta_0)]$. Provided $\tilde{M}_n \rightsquigarrow M$, and the regularity conditions of the argmax theorem apply, $\hat{h}_n = r_n(\hat{\theta}_n - \theta_0) \rightsquigarrow \hat{h}$, where \hat{h} is the argmax of M . Consistency results will follow for the choice $r_n = 1$ for all $n \geq 1$.

Note that in the theorem, we require the sequence \hat{h}_n to be uniformly tight. This is stronger than asymptotic tightness, as pointed out in lemma 7.10, but is also quite easy to establish for Euclidean parameters which will be our main focus in this chapter. For finite Euclidean estimators that are measurable, uniform tightness follows automatically from asymptotic tightness (see exercise 14.6.1). This is a reasonable restriction, since, in practice, most infinite-dimensional estimators that converge weakly can usually be expressed as Z-estimators. Our consistency results that we present later on will not require uniform tightness and will thus be more readily applicable

to infinite dimensional estimators. Returning to the theorem at hand, most weak convergence results for non-regular estimators apply to finite dimensional parameters, and thus the theorem below will be applicable. We also note that it is not hard to modify these results for applicability to specific settings, including some infinite dimensional settings. For interesting examples in this direction, see Ma and Kosorok (2005) and Kosorok and Song (2006). We now present the argmax theorem, which utilizes upper semicontinuity as defined in section 6.1:

THEOREM 14.1 (Argmax theorem) *Let M_n, M be stochastic processes indexed by a metric space H such that $M_n \rightsquigarrow M$ in $\ell^\infty(K)$ for every compact $K \subset H$. Suppose also that almost all sample paths $h \mapsto M(h)$ are upper semicontinuous and possess a unique maximum at a (random) point \hat{h} , which as a random map in H is tight. If the sequence \hat{h}_n is uniformly tight and satisfies $M_n(\hat{h}_n) \geq \sup_{h \in H} M_n(h) - o_P(1)$, then $\hat{h}_n \rightsquigarrow \hat{h}$ in H .*

Proof. Fix $\epsilon > 0$. By uniform tightness of \hat{h}_n and tightness of \hat{h} , there exists a compact set $K \subset H$ such that $\limsup_{n \rightarrow \infty} P^*(\hat{h}_n \in K) \geq 1 - \epsilon$ and $P(\hat{h} \in K) \geq 1 - \epsilon$. Then almost surely

$$M(\hat{h}) > \sup_{h \notin G, h \in K} M(h),$$

for every open $G \ni \hat{h}$, by upper semicontinuity of M . To see this, suppose it were not true. Then by the compactness of K , there would exist a convergent sequence $h_m \in G^c \cap K$, for some open $G \ni \hat{h}$, with $h_m \rightarrow h$ and $M(h_m) \rightarrow M(\hat{h})$. The upper semicontinuity forces $M(h) \geq M(\hat{h})$ which contradicts the uniqueness of the maximum.

Now apply lemma 14.2 below for the sets $A = B = K$, to obtain

$$\begin{aligned} \limsup_{n \rightarrow \infty} P^*(\hat{h}_n \in F) &\leq \limsup_{n \rightarrow \infty} P^*(\hat{h}_n \in F \cap K) + \limsup_{n \rightarrow \infty} P^*(\hat{h}_n \notin K) \\ &\leq P(\hat{h} \in F \cup K^c) + \epsilon \\ &\leq P(\hat{h} \in F) + P(\hat{h} \in K^c) + \epsilon \\ &\leq P(\hat{h} \in F) + 2\epsilon. \end{aligned}$$

The desired result now follows from the Portmanteau theorem since ϵ was arbitrary. \square

LEMMA 14.2 *Let M_n, M be stochastic processes indexed by a metric space H , and let $A, B \subset H$ be arbitrary. Suppose there exists a random element \hat{h} such that almost surely*

$$(14.1) \quad M(\hat{h}) > \sup_{h \notin G, A \in K} M(h), \text{ for every open } G \ni \hat{h}.$$

Suppose the sequence \hat{h}_n satisfies $M_n(\hat{h}_n) \geq \sup_{h \in H} M_n(h) - o_P(1)$. Then, if $M_n \rightsquigarrow M$ in $\ell^\infty(A \cup B)$, we have for every closed set F ,

$$\limsup_{n \rightarrow \infty} \mathbb{P}^*(\hat{h}_n \in F \cap A) \leq \mathbb{P}(\hat{h} \in F \cup B^c).$$

Proof. By the continuous mapping theorem,

$$\sup_{h \in F \cap A} M_n(h) - \sup_{h \in B} M_n(h) \rightsquigarrow \sup_{h \in F \cap A} M(h) - \sup_{h \in B} M(h),$$

and thus

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbb{P}^*(\hat{h}_n \in F \cap A) &\leq \limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\sup_{h \in F \cap A} M_n(h) \geq \sup_{h \in B} M_n(h) - o_P(1) \right) \\ &\leq \limsup_{n \rightarrow \infty} \mathbb{P}^* \left(\sup_{h \in F \cap A} M_n(h) \geq \sup_{h \in B} M_n(h) - o_P(1) \right) \\ &\leq \mathbb{P} \left(\sup_{h \in F \cap A} M(h) \geq \sup_{h \in B} M(h) \right), \end{aligned}$$

by Slutsky's theorem (to get rid of the $o_P(1)$ part) followed by the Portmanteau theorem. Note that the event E in the last probability can't happen when $\hat{h} \in F^c \cap B$ because of assumption (14.1) and the fact that F^c is open. Thus E is contained in the set $\{\hat{h} \in F\} \cup \{\hat{h} \notin B\}$, and the conclusion of the lemma follows. \square

14.2 Consistency

We can obtain a consistency result by specializing the argmax theorem to the setting where M is fixed. This will not yield as general a result as theorem 2.12 because of the uniform tightness requirement. The primary goal of this section is to prove theorem 2.12. Before giving the proof, we want to present a result comparing a few different ways of establishing identifiability. We assume throughout this section that (Θ, d) is a metric space. In the following lemma, the condition given in (i) is the identifiability condition assumed in theorem 2.12, while the condition (ii) is often called the “well-separated maximum” condition:

LEMMA 14.3 *Let $M : \Theta \mapsto \mathbb{R}$ be a map and $\theta_0 \in \Theta$ a point. The following conditions are equivalent:*

- (i) *For any sequence $\{\theta_n\} \in \Theta$, $\liminf_{n \rightarrow \infty} M(\theta_n) \geq M(\theta_0)$ implies $d(\theta_n, \theta_0) \rightarrow 0$.*
- (ii) *For every open $G \ni \theta_0$, $M(\theta_0) > \sup_{\theta \notin G} M(\theta)$.*

The following condition implies both (i) and (ii):

(iii) M is upper semicontinuous with a unique maximum at θ_0 .

Proof. Suppose (i) is true but (ii) is not. Then there exists an open $G \ni \theta_0$ such that $\sup_{\theta \notin G} M(\theta) \geq M(\theta_0)$. This implies the existence of a sequence θ_n with both $\liminf_{n \rightarrow \infty} M(\theta_n) \geq M(\theta_0)$ and $d(\theta_n, \theta_0) \rightarrow \tau > 0$. Which is a contradiction. Thus (i) implies (ii). Now assume (ii) is true but (i) is not. Then there exists a sequence with $\liminf_{n \rightarrow \infty} M(\theta_n) \geq M(\theta_0)$ but with $\theta_n \notin G$ for all n large enough and some open $G \ni \theta_0$. Of course, this contradicts (ii), and thus (ii) implies (i). Now suppose M is upper semicontinuous with a unique maximum at θ_0 but (ii) does not hold. Then there exists an open $G \ni \theta_0$ for which $\sup_{\theta \notin G} M(\theta) \geq M(\theta_0)$. But this implies that the set $\{\theta : M(\theta) \geq M(\theta_0)\}$ contains at least one point in addition to θ_0 since G^c is closed. This contradiction completes the proof. \square

Any one of the three identifiability conditions given in the above lemma are sufficient for theorem 2.12. The most convenient condition in practice will depend on the setting. Here is the awaited proof:

Proof of theorem 2.12. Since $\liminf_{n \rightarrow \infty} M(\theta_n) \geq M(\theta_0)$ implies $d(\theta_n, \theta_0) \rightarrow 0$ for any sequence $\{\theta_n\} \in \Theta$, we know that there exists a non-decreasing cadlag function $f : [0, \infty] \mapsto [0, \infty]$ that satisfies both $f(0) = 0$ and $d(\theta, \theta_0) \leq f(|M(\theta) - M(\theta_0)|)$ for all $\theta \in \Theta$. The details for constructing such an f are left as an exercise (see exercise 14.6.2).

For part (i), note that $M(\theta_0) \geq M(\hat{\theta}_n) \geq M_n(\hat{\theta}_n) - \|M_n - M\|_{\Theta} \geq M_n(\theta_0) - o_P(1) \geq M(\theta_0) - o_P(1)$. By the previous paragraph, this implies $d(\hat{\theta}_n, \theta_0) \leq f(|M(\hat{\theta}_n) - M(\theta_0)|) \xrightarrow{P} 0$. An almost identical argument yields part (ii). \square

14.3 Rate of Convergence

In this section, we relax the requirement that (Θ, d) be a metric space to only requiring that it to be a semimetric space. If $\theta \mapsto M(\theta)$ is two times differentiable at a point of maximum θ_0 , then the first derivative of M at θ_0 must vanish while the second derivative should be negative definite. Thus it is not unreasonable to require that $M(\theta) - M(\theta_0) \leq -cd^2(\theta, \theta_0)$ for all θ in a neighborhood of θ_0 and some $c > 0$. The following theorem shows that an upper bound for the rate of convergence of a near-maximizer of a random objection function M_n can be obtained from the modulus of continuity of $M_n - M$ at θ_0 . In practice, one may need to try several rates that satisfy the conditions of this theorem before finding the right r_n for which the weak limit of $r_n(\hat{\theta}_n - \theta_0)$ is nontrivial.

THEOREM 14.4 (Rate of convergence) *Let M_n be a sequence of stochastic processes indexed by a semimetric space (Θ, d) and $M : \Theta \mapsto \mathbb{R}$ a deterministic function such that for every θ in a neighborhood of θ_0 , there exists a $c_1 > 0$ such that*

$$(14.2) \quad M(\theta) - M(\theta_0) \leq -c_1 \tilde{d}^2(\theta, \theta_0),$$

where $\tilde{d} : \Theta \times \Theta \mapsto [0, \infty)$ satisfies $\tilde{d}(\theta_n, \theta_0) \rightarrow 0$ whenever $d(\theta_n, \theta_0) \rightarrow 0$. Suppose that for all n large enough and sufficiently small δ , the centered process $M_n - M$ satisfies

$$(14.3) \quad \mathbb{E}^* \sup_{\tilde{d}(\theta, \theta_0) < \delta} \sqrt{n} |(M_n - M)(\theta) - (M_n - M)(\theta_0)| \leq c_2 \phi_n(\delta),$$

for $c_2 < \infty$ and functions ϕ_n such that $\delta \mapsto \phi_n(\delta)/\delta^\alpha$ is decreasing for some $\alpha < 2$ not depending on n . Let

$$(14.4) \quad r_n^2 \phi_n(r_n^{-1}) \leq c_3 \sqrt{n}, \text{ for every } n \text{ and some } c_3 < \infty.$$

If the sequence $\hat{\theta}_n$ satisfies $M_n(\hat{\theta}_n) \geq \sup_{\theta \in \Theta} M_n(\theta) - O_P(r_n^{-2})$ and converges in outer probability to θ_0 , then $r_n \tilde{d}(\hat{\theta}_n, \theta_0) = O_P(1)$.

Proof. We will use a modified “peeling device” (see, for example, section 5.3 of van de Geer, 2000) for the proof. For every $\eta > 0$, let $\eta' > 0$ be a number for which $\tilde{d}(\theta, \theta_0) \leq \eta$ whenever $\theta \in \Theta$ satisfies $d(\theta, \theta_0) \leq \eta'$ and also $\tilde{d}(\theta, \theta_0) \leq \eta/2$ whenever $\theta \in \Theta$ satisfies $d(\theta, \theta_0) \leq \eta'/2$. Such an η' always exists for each η by the assumed relationship between d and \tilde{d} . Note also that $M_n(\hat{\theta}_n) \geq \sup_{\theta \in \Theta} M_n(\theta) - O_P(r_n^{-2}) \geq M_n(\theta_0) - O_P(r_n^{-2})$. Now fix $\epsilon > 0$, and choose $K < \infty$ such that the probability that $M_n(\hat{\theta}_n) - M_n(\theta_0) < -Kr_n^{-2}$ is $\leq \epsilon$.

For each n , the parameter space minus the point θ_0 can be partitioned into “peels” $S_{j,n} = \{\theta : 2^{j-1} < r_n \tilde{d}(\theta, \theta_0) \leq 2^j\}$ with j ranging over the integers. Assume that $M_n(\hat{\theta}_n) - M_n(\theta_0) \geq -Kr_n^{-2}$, and note that if $r_n \tilde{d}(\hat{\theta}_n, \theta_0)$ is $> 2^M$ for a given integer M , then $\hat{\theta}_n$ is in one of the peels $S_{j,n}$, with $j > M$. In that situation, the supremum of the map $\theta \mapsto M_n(\theta) - M_n(\theta_0) + Kr_n^{-2}$ is nonnegative by the property of $\hat{\theta}_n$. Conclude that for every $\eta > 0$,

$$(14.5) \quad \begin{aligned} & \mathbb{P}^* \left(r_n \tilde{d}(\hat{\theta}_n, \theta_0) > 2^M \right) \\ & \leq \sum_{j \geq M, 2^j \leq \eta r_n} \mathbb{P}^* \left(\sup_{\theta \in S_{j,n}} [M_n(\theta) - M_n(\theta_0) + Kr_n^{-2}] \geq 0 \right) \\ & \quad + \mathbb{P}^* \left(2d(\hat{\theta}_n, \theta_0) \geq \eta' \right) + \mathbb{P}^* \left(M_n(\hat{\theta}_n) - M_n(\theta_0) < -Kr_n^{-2} \right). \end{aligned}$$

The $\limsup_{n \rightarrow \infty}$ of the sum of the two probabilities after the summation on the right side is $\leq \epsilon$ by the consistency of $\hat{\theta}_n$ and the choice of K . Now choose η small enough so that (14.2) holds for all $d(\theta, \theta_0) \leq \eta'$ and (14.3) holds for all $\delta \leq \eta$. Then for every j involved in the sum, we have for every $\theta \in S_{j,n}$, $M(\theta) - M(\theta_0) \leq -c_1 \tilde{d}(\theta, \theta_0) \leq -c_1 2^{2j-2} r_n^{-2}$. In terms of the centered process $W_n \equiv M_n - M$, the summation on the right side of (14.5) may be bounded by

$$\begin{aligned}
\sum_{j \geq M, 2^j \leq \eta r_n} \mathbb{P}^* \left(\|W_n(\theta) - W_n(\theta_0)\|_{S_{j,n}} \geq \frac{c_1 2^{2j-2} - K}{r_n^2} \right) \\
\leq \sum_{j \geq M} \frac{c_2 \phi_n(2^j/r_n) r_n^2}{\sqrt{n} (c_1 2^{2j-2} - K)} \\
\leq \sum_{j \geq M} \frac{c_2 c_3 2^{j\alpha}}{(c_1 2^{2j-2} - K)},
\end{aligned}$$

by Markov's inequality, the conditions on r_n , and the fact that $\phi_n(c\delta) \leq c^\alpha \phi_n(\delta)$ for every $c > 1$ as a consequence of the assumptions on ϕ_n . It is not difficult to show that this last sum goes to zero as $M \rightarrow \infty$ (verifying this is saved as an exercise). Thus we can choose an $M < \infty$ such that the $\limsup_{n \rightarrow \infty}$ of the left side of (14.5) is $\leq 2\epsilon$. The desired result now follows since ϵ was arbitrary. \square

Consider the i.i.d. setting with criterion functions of the form $M_n(\theta) = \mathbb{P}_n m_\theta$ and $M(\theta) = P m_\theta$. The scaled and centered process $\sqrt{n}(M_n - M) = \mathbb{G}_n m_\theta$ equals the empirical process at m_θ . The assertion (14.3) involves assessing the suprema of the empirical process by classes of functions $\mathcal{M}_\delta \equiv \{m_\theta - m_{\theta_0} : \tilde{d}(\theta, \theta_0) < \delta\}$. Taking this view, establishing (14.3) will require fairly precise—but not unreasonably precise—knowledge of the involved empirical process. The moment results in section 11.1 will be useful here, and we will illustrate this with several examples later on in this chapter. We note that the problems we address in this book represent only a small subset of the scope and capabilities of empirical process techniques for determining rates of M-estimators. We close this section with the following corollary which essentially specializes theorem 14.4 to the i.i.d. setting. Because the specialization is straightforward, we omit the somewhat trivial proof. Recall that the relation $a \lesssim b$ means that a is less than or equal b times a universal finite and positive constant.

COROLLARY 14.5 *In the i.i.d. setting, assume that for every θ in a neighborhood of θ_0 , $P(m_\theta - m_{\theta_0}) \lesssim -\tilde{d}^2(\theta, \theta_0)$, where \tilde{d} satisfies the conditions given in theorem 14.4. Assume moreover that there exists a function ϕ such that $\delta \mapsto \phi(\delta)/\delta^\alpha$ is decreasing for some $\alpha < 2$ and, for every n , $\mathbb{E}^* \|\mathbb{G}_n\|_{\mathcal{M}_\delta} \lesssim \phi(\delta)$. If the sequence $\hat{\theta}_n$ satisfies $\mathbb{P}_n m_{\hat{\theta}_n} \geq \sup_{\theta \in \Theta} \mathbb{P}_n m_\theta - O_P(r_n^{-2})$ and converges in outer probability to θ_0 , then $r_n \tilde{d}(\hat{\theta}_n, \theta_0) = O_P(1)$ for every sequence r_n for which $r_n^2 \phi(1/r_n) \lesssim \sqrt{n}$ for all $n \geq 1$.*

14.4 Regular Euclidean M-Estimators

A general result for Euclidean M-estimators based on i.i.d. data was given in theorem 2.13 of section 2.2.6. We now prove this theorem. In section 2.2.6, the theorem was used to establish asymptotic normality of a least-absolute-

deviation regression estimator. Establishing asymptotic normality in this situation is quite difficult without empirical process methods.

Proof of theorem 2.13. We first utilize corollary 14.5 to verify that \sqrt{n} is the correct rate of convergence. We will use Euclidean distance as both the discrepancy measure and distance through, i.e. $\tilde{d}(\theta_1, \theta_2) = d(\theta_1, \theta_2) = \|\theta_1 - \theta_2\|$. Condition (2.18) of the theorem indicates that \mathcal{M}_δ in this instance is a Lipschitz class, and thus theorem 9.22 implies that

$$(14.6) \quad N_{[]} (2\epsilon \|F_\delta\|_{P,2}, \mathcal{M}_\delta, L_2(P)) \lesssim \epsilon^{-p},$$

where $F_\delta \equiv \delta \dot{m}$ is an envelope for \mathcal{M}_δ . To see this, it may be helpful to rewrite condition (2.18) as

$$|m_{\theta_1}(x) - m_{\theta_2}(x)| \leq \frac{\|\theta_1 - \theta_2\|}{\delta} F_\delta(x).$$

Now (14.6) can be utilized in theorem 11.2 to obtain

$$E^* \|\mathbb{G}_n\|_{\mathcal{M}_\delta} \lesssim \|F_\delta\|_{P,2} \lesssim \delta.$$

Hence the modulus of continuity condition in corollary 14.5 is satisfied for $\phi(\delta) = \delta$. Combining condition (2.19), the maximality of θ_0 , and the fact that the second derivative matrix V is nonsingular and continuous, yields that $M(\theta) - M(\theta_0) \lesssim -\|\theta - \theta_0\|^2$. Since $\phi(\delta)/\delta^\alpha = \delta^{1-\alpha}$ is decreasing for any $\alpha \in (1, 2)$ and $n\phi(1/\sqrt{n}) = n/\sqrt{n} = \sqrt{n}$, the remaining conditions of corollary 14.5 are satisfied for $r_n = \sqrt{n}$, and thus $\sqrt{n}(\hat{\theta}_n - \theta_0) = O_P(1)$.

The next step is to apply the argmax theorem (theorem 14.1) to the process $h \mapsto U_n(h) \equiv n(M_n(\theta_0 + h/\sqrt{n}) - M_n(\theta_0))$. Now fix a compact $K \subset \mathbb{R}^p$, and note that

$$\begin{aligned} U_n(h) &= \mathbb{G}_n [\sqrt{n}(m_{\theta_0+h/\sqrt{n}} - m_{\theta_0}) - h^T \dot{m}_{\theta_0}] \\ &\quad + h^T \mathbb{G}_n \dot{m}_{\theta_0} + n(M(\theta_0 + h/\sqrt{n}) - M(\theta_0)) \\ &\equiv E_n(h) + h^T \mathbb{G}_n \dot{m}_{\theta_0} + \frac{1}{2} h^T V h + o(1), \end{aligned}$$

where $o(1)$ denotes a quantity going to zero uniformly over K . Note that $\hat{h}_n \equiv \sqrt{n}(\hat{\theta}_n - \theta_0)$ satisfies $U_n(\hat{h}_n) \geq \sup_{h \in \mathbb{R}^p} U_n(h) - o_P(1)$. Thus, provided we can establish that $\|E_n\|_K = o_P(1)$, the argmax theorem will yield that $\hat{h}_n \rightsquigarrow \hat{h}$, where \hat{h} is the argmax of $h \mapsto U(h) \equiv h^T Z + (1/2)h^T V h$, where Z is the Gaussian limiting distribution of $\mathbb{G}_n \dot{m}_{\theta_0}$. Hence $\hat{h} = -V^{-1}Z$ and the desired result will follow.

We now prove $\|E_n\|_K = o_P(1)$ for all compact $K \subset \mathbb{R}^p$. let $u_h^n(x) \equiv \sqrt{n}(m_{\theta_0+h/\sqrt{n}}(x) - m_{\theta_0}(x)) - h^T \dot{m}_{\theta_0}(x)$, and note that by (2.18),

$$|u_{h_1}^n(x) - u_{h_2}^n(x)| \leq (\dot{m}(x) + \|\dot{m}_{\theta_0}(x)\|) \|h_1 - h_2\|,$$

for all $h_1, h_2 \in \mathbb{R}^p$ and all $n \geq 1$. Fix a compact $K \subset \mathbb{R}^p$, and let $\mathcal{F}_n \equiv \{u_h^n : h \in K\}$. Applying theorem 9.22 once again, but with $\|\cdot\| = \|\cdot\|_{Q,2}$ (for any probability measure Q on \mathcal{X}) instead of $\|\cdot\|_{P,2}$, we obtain

$$N_{[]} (2\epsilon \|F_n\|_{Q,2}, \mathcal{F}_n, L_2(Q)) \leq k\epsilon^{-p},$$

where the envelope $F_n \equiv (\dot{m} + \|\dot{m}_{\theta_0}\|)\|h\|_K$ and $k < \infty$ does not depend on K or n . Lemma 9.18 in chapter 9 now yields that

$$\sup_{n \geq 1} \sup_Q N(\epsilon \|F_n\|_{Q,2}, \mathcal{F}_n, L_2(Q)) \leq k \left(\frac{2}{\epsilon}\right)^p,$$

where the second supremum is taken over all finitely discrete probability measures on \mathcal{X} . This implies that condition (A) of theorem 11.18 holds for \mathcal{F}_n and F_n . In addition, condition (2.19) implies that condition (B) of theorem 11.18 also holds with $H(s, t) = 0$ for all $s, t \in K$. It is not difficult to verify that all of the remaining conditions of theorem 11.18 also hold (we save this as an exercise), and thus $\mathbb{G}_n u_h^n \rightsquigarrow 0$ in $\ell^\infty(K)$. This, of course, is the desired result. \square

14.5 Non-Regular Examples

We now present two examples in detail that illustrate the techniques presented in this chapter for parameter estimation with non-regular rates of convergence. The first example considers a simple change-point model with three parameters wherein two of the parameter estimates converge at the regular rate while one of the parameter estimates converges at the n -rate, i.e., it converges faster than \sqrt{n} . The second example is monotone density estimation based on the Grenander estimator which is shown to yield convergence at the cube-root rate.

14.5.1 A Change-Point Model

For this model, we observe i.i.d. realizations of $X = (Y, Z)$, where $Y = \alpha 1\{Z \leq \zeta\} + \beta 1\{Z > \zeta\} + \epsilon$, Z and ϵ are independent with ϵ continuous, $E\epsilon = 0$ and $\sigma^2 \equiv E\epsilon^2 < \infty$, $\gamma \equiv (\alpha, \beta) \in \mathbb{R}^2$ and ζ is known to lie in a bounded interval $[a, b]$. The unknown parameters can be collected as $\theta = (\gamma, \zeta)$, and the subscript zero will be used to denote the true parameter values. We make the very important assumption that $\alpha_0 \neq \beta_0$ and also assume that Z has a strictly bounded and positive density f over $[a, b]$ with $P(Z < a) > 0$ and $P(Z > b) > 0$. Our goal is to estimate θ through least squares. This is the same as maximizing $M_n(\theta) = \mathbb{P}_n m_\theta$, where

$$m_\theta(x) \equiv -(y - \alpha 1\{z \leq \zeta\} - \beta 1\{z > \zeta\})^2.$$

Let $\hat{\theta}_n$ be maximizers of $M_n(\theta)$, where $\hat{\theta}_n \equiv (\hat{\gamma}_n, \hat{\zeta}_n)$ and $\hat{\gamma}_n \equiv (\hat{\alpha}_n, \hat{\beta}_n)$.

Since we are not assuming that γ is bounded, we first need to prove the existence of $\hat{\gamma}_n$, i.e., we need to prove that $\|\hat{\gamma}_n\| = O_P(1)$. We then

need to provide consistency of all parameters and then establish the rates of convergence for the parameters. Finally, we need to obtain the joint limiting distribution of the parameter estimates.

Existence.

Note that the covariate Z and parameter ζ can be partitioned into four mutually exclusive sets: $\{Z \leq \zeta \wedge \zeta_0\}$, $\{\zeta < Z \leq \zeta_0\}$, $\{\zeta_0 < Z \leq \zeta\}$ and $\{Z > \zeta \vee \zeta_0\}$. Since also $1\{Z < a\} \leq 1\{Z \leq \zeta \wedge \zeta_0\}$ and $1\{Z > b\} \leq 1\{Z > \zeta \vee \zeta_0\}$ by assumption, we obtain $-\mathbb{P}_n \epsilon^2 = M_n(\theta_0) \leq M_n(\hat{\theta}_n)$

$$\leq -\mathbb{P}_n \left[(\epsilon - \hat{\alpha}_n + \alpha_0)^2 1\{Z < a\} + (\epsilon - \hat{\beta}_n + \beta_0)^2 1\{Z > b\} \right].$$

By decomposing the squares, we now have

$$\begin{aligned} & (\hat{\alpha}_n - \alpha_0)^2 \mathbb{P}_n[\epsilon^2 1\{Z < a\}] + (\hat{\beta}_n - \beta_0)^2 \mathbb{P}_n[\epsilon^2 1\{Z > b\}] \\ & \leq \mathbb{P}_n[\epsilon^2 1\{a \leq z \leq b\}] \\ & \quad + 2|\hat{\alpha}_n - \alpha_0| \mathbb{P}_n[\epsilon 1\{Z < a\}] + 2|\hat{\beta}_n - \beta_0| \mathbb{P}_n[\epsilon 1\{Z > b\}] \\ & \leq O_p(1) + o_P(1) \|\hat{\gamma}_n - \gamma_0\|. \end{aligned}$$

Since $\mathbb{P}(Z < a) \wedge \mathbb{P}(Z > b) > 0$, the above now implies that $\|\hat{\gamma}_n - \gamma_0\|^2 = O_P(1 + \|\hat{\gamma}_n - \gamma_0\|)$ and hence that $\|\hat{\gamma}_n - \gamma_0\| = O_P(1)$. Thus all the parameters are bounded in probability and therefore exist.

Consistency.

Our approach to establishing consistency will be to utilize the argmax theorem (theorem 14.1). We first need to establish that $M_n \rightsquigarrow M$ in $\ell^\infty(K)$ for all compact $K \subset H \equiv \mathbb{R}^2 \times [a, b]$, where $M(\theta) \equiv Pm_\theta$. We then need to show that $\theta \mapsto M(\theta)$ is upper semicontinuous with a unique maximum at θ_0 . We already know from the previous paragraph that $\hat{\theta}_n$ is asymptotically tight (i.e., $\|\hat{\theta}_n\| = O_P(1)$). The argmax theorem will then yield that $\hat{\theta}_n \rightsquigarrow \theta_0$ as desired.

Fix a compact $K \subset H$. We now verify that $\mathcal{F}_K \equiv \{m_\theta : \theta \in K\}$ is Glivenko-Cantelli. Note that

$$\begin{aligned} m_\theta(X) &= -(\epsilon - \alpha + \alpha_0)^2 1\{Z \leq \zeta \wedge \zeta_0\} - (\epsilon - \beta + \alpha_0)^2 1\{\zeta < Z \leq \zeta_0\} \\ &\quad - (\epsilon - \alpha + \beta_0)^2 1\{\zeta_0 < Z \leq \zeta\} - (\epsilon - \beta + \beta_0)^2 1\{Z > \zeta \vee \zeta_0\}. \end{aligned}$$

It is not difficult to verify that $\{(\epsilon - \alpha + \alpha_0)^2 : \theta \in K\}$ and $1\{Z \leq \zeta \wedge \zeta_0 : \theta \in K\}$ are separately Glivenko-Cantelli classes. Thus the product of the two class is also Glivenko-Cantelli by corollary 9.26 since the product of the two envelopes is integrable. Similar arguments reveal that the remaining components of the sum are also Glivenko-Cantelli, and reapplication of corollary 9.26 yields that \mathcal{F}_K itself is Glivenko-Cantelli. Thus $M_n \rightsquigarrow M$ in $\ell^\infty(K)$ for all compact K .

We now establish upper semicontinuity of $\theta \mapsto M(\theta)$ and uniqueness of the maximum. Using the decomposition of the sets for (Z, ζ) used in the *Existence* paragraph above, we have

$$\begin{aligned} M(\theta) &= -P\epsilon^2 - (\alpha - \alpha_0)^2 P(Z \leq \zeta \wedge \zeta_0) - (\beta - \alpha_0)^2 P(\zeta < Z \leq \zeta_0) \\ &\quad - (\alpha - \beta_0)^2 P(\zeta_0 < Z \leq \zeta) - (\beta - \beta_0)^2 P(Z > \zeta \vee \zeta_0) \\ &\leq -P\epsilon^2 = M(\theta_0). \end{aligned}$$

Because Z has a bounded density on $[a, b]$, we obtain that M is continuous. It is also clear that M has a unique maximum at θ_0 because the density of Z is bounded below and $\alpha_0 \neq \beta_0$ (see exercise 14.6.5 below). Now the conditions of the argmax theorem are met, and the desired consistency follows.

Rate of convergence.

We will utilize corollary 14.5 to obtain the convergence rates via the discrepancy function $\tilde{d}(\theta, \theta_0) \equiv \|\gamma - \gamma_0\| + \sqrt{|\zeta - \zeta_0|}$. Note that this is not a norm since it does not satisfy the triangle inequality. Nevertheless, $\tilde{d}(\theta, \theta_0) \rightarrow 0$ if and only if $\|\theta - \theta_0\| \rightarrow 0$. Moreover, from the *Consistency* paragraph above, we have that

$$\begin{aligned} M(\theta) - M(\theta_0) &= -P\{Z \leq \zeta \wedge \zeta_0\}(\alpha - \alpha_0)^2 - P\{Z > \zeta \vee \zeta_0\}(\beta - \beta_0)^2 \\ &\quad - P\{\zeta < Z \leq \zeta_0\}(\beta - \alpha_0)^2 \\ &\quad - P\{\zeta_0 < Z \leq \zeta\}(\alpha - \beta_0)^2 \\ &\leq -P\{Z < a\}(\alpha - \alpha_0)^2 - P\{Z > b\}(\beta - \beta_0)^2 \\ &\quad - k_1(1 - o(1))|\zeta - \zeta_0| \\ &\leq -(k_1 \wedge \delta_1 - o(1))\tilde{d}^2(\theta, \theta_0), \end{aligned}$$

where the first inequality follows from the fact that the product of the density of Z and $(\alpha_0 - \beta_0)^2$ is bounded below by some $k_1 > 0$, and the second inequality follows from both $P(Z < a)$ and $P(Z > b)$ being bounded below by some $\delta_1 > 0$. Thus $M(\theta) - M(\theta_0) \lesssim -\tilde{d}^2(\theta, \theta_0)$ for all $\|\theta - \theta_0\|$ small enough, as desired.

Consider now the class of functions $\mathcal{M}_\delta \equiv \{m_\theta - m_{\theta_0} : \tilde{d}(\theta, \theta_0) < \delta\}$. Using previous calculations, we have

$$\begin{aligned} (14.7) \quad m_\theta - m_{\theta_0} &= 2(\alpha - \alpha_0)\epsilon 1\{Z \leq \zeta \wedge \zeta_0\} + 2(\beta - \beta_0)\epsilon 1\{Z > \zeta \vee \zeta_0\} \\ &\quad + 2(\beta - \alpha_0)\epsilon 1\{\zeta < Z \leq \zeta_0\} + 2(\alpha - \beta_0)\epsilon 1\{\zeta_0 < Z \leq \zeta\} \\ &\quad - (\alpha - \alpha_0)^2 1\{Z \leq \zeta \wedge \zeta_0\} - (\beta - \beta_0)^2 1\{Z > \zeta \vee \zeta_0\} \\ &\quad - (\beta - \alpha_0)^2 1\{\zeta < Z \leq \zeta_0\} - (\alpha - \beta_0)^2 1\{\zeta_0 < Z \leq \zeta\} \\ &\equiv A_1(\theta) + A_2(\theta) + B_1(\theta) + B_2(\theta) \\ &\quad - C_1(\theta) - C_2(\theta) - D_1(\theta) - D_2(\theta). \end{aligned}$$

Consider first A_1 . Since $\{1\{Z \leq t\} : t \in [a, b]\}$ is a VC class, it is easy to compute that

$$E^* \sup_{\tilde{d}(\theta, \theta_0) < \delta} |\mathbb{G}_n A_1(\theta)| \lesssim \delta,$$

as a consequence of lemma 8.17. Similar calculations apply to A_2 . Similar calculations also apply to C_1 and C_2 , except that the upper bounds will be $\lesssim \delta^2$ instead of $\lesssim \delta$. Now we consider B_1 . An envelope for the class $\mathcal{F} = \{B_1(\theta) : \tilde{d}(\theta, \theta_0) < \delta\}$ is $F = 2(|\beta_0 - \alpha_0| + \delta)|\epsilon|1\{\zeta_0 - \delta^2 < Z \leq \zeta_0\}$. It is not hard to verify that

$$(14.8) \quad \log N_{[]}(\eta \|F\|_{P,2}, \mathcal{F}, L_2(P)) \lesssim \log(1/\eta)$$

(see exercise 14.6.6). Now theorem 11.4 yields that

$$E^* \sup_{\tilde{d}(\theta, \theta_0) < \delta} |\mathbb{G}_n B_1(\theta)| = E^* \|\mathbb{G}_n\|_{\mathcal{F}} \lesssim \|F\|_{P,2} \lesssim \delta.$$

Similar calculations apply also to B_2 , D_1 and D_2 . Combining all of these results with the fact that $O(\delta^2) = O(\delta)$, we obtain $E^* \|\mathbb{G}_n\|_{\mathcal{M}_\delta} \lesssim \delta$.

Now when $\delta \mapsto \phi(\delta) = \delta$, $\phi(\delta)/\delta^\alpha$ is decreasing for any $\alpha \in (1, 2)$. Thus the conditions of corollary 14.5 are satisfied with $\phi(\delta) = \delta$. Since $r_n^2 \phi(1/r_n) = r_n$, we obtain that $\sqrt{n} \tilde{d}(\hat{\theta}_n, \theta_0) = O_P(1)$. By the form of \tilde{d} , this now implies that $\sqrt{n} \|\hat{\gamma}_n - \gamma_0\| = O_P(1)$ and $n|\hat{\zeta}_n - \zeta_0| = O_P(1)$.

Weak convergence.

We will utilize a minor modification of the argmax theorem and the rate result above to obtain the limiting distribution of $\hat{h}_n = (\sqrt{n}(\hat{\gamma}_n - \gamma_0), n(\hat{\zeta}_n - \zeta_0))$. From the rate result, we know that \hat{h}_n is uniformly tight and is the smallest argmax of $h \mapsto Q_n(h) \equiv n\mathbb{P}_n(m_{\theta_{n,h}} - m_{\theta_0})$, where $\theta_{n,h} \equiv \theta_0 + (h_1/\sqrt{n}, h_2/\sqrt{n}, h_3/n)$ and $h \equiv (h_1, h_2, h_3) \in \mathbb{R}^3 \equiv H$. Note that we have qualified \hat{h}_n as being the smallest argmax, which is interpreted componentwise since H is three dimensional. This is because if we hold (h_1, h_2) fixed, $M_n(\theta_{n,h})$ does not vary in h_3 over the interval $n[Z_{(j)} - \zeta_0, Z_{(j+1)} - \zeta_0]$, for $j = 1, \dots, n$, where $Z_{(1)}, \dots, Z_{(n)}$ are the order statistics for Z_1, \dots, Z_n , $Z_{(0)} \equiv -\infty$, and $Z_{(n+1)} \equiv \infty$. Because $P(Z < a) > 0$, we only need to consider h_3 at the values $n(Z_{(j)} - \zeta_0)$, $j = 1, \dots, n$, provided n is large enough.

Let \mathcal{D}_K be the space of functions $q : K \subset H \mapsto \mathbb{R}$, that are continuous in the first two arguments (h_1, h_2) and right-continuous and piecewise constant in the third argument h_3 . For each $q \in \mathcal{D}_K$, let $h_3 \mapsto J_q(h_3)$ be the cadlag counting process with $J_q(0-) = 0$, jumps of size positive 1 at each jump point in $q(\cdot, \cdot, h_3)$ for $h_3 \geq 0$, and with $h_3 \mapsto J_q(-h_3)$ also having jumps of size positive 1 (but left-continuous) at each jump point in $q(\cdot, \cdot, h_3)$ for $h_3 < 0$ (the left-continuity comes from the reversed time scale). Thus $J_q(h_3)$ is decreasing for $h_3 < 0$ and increasing for $h_3 \geq 0$. For $q_1, q_2 \in \mathcal{D}_K$, define the distance $d_K(q_1, q_2)$ to be the sum of the uniform

distance $\|q_1 - q_2\|_K$ and the Skorohod distance between J_{q_1} and J_{q_2} . Now it is not difficult to see that the smallest argmax function is continuous on \mathcal{D}_K with respect to d_K . We will argue that $Q_n(h) \rightsquigarrow Q(h)$ in (\mathcal{D}_K, d_K) , for some limiting process Q , and for each compact $K \subset H$. By the continuous mapping theorem, the smallest argmax of the restriction of Q_n to K will converge weakly to the smallest argmax of the restriction of Q to K . Since \hat{h}_n is uniformly tight, we obtain $\hat{h}_n \rightsquigarrow \hat{h}$, where \hat{h} is the smallest argmax of Q .

All that remains is to establish the specified weak convergence and to characterize Q . We first argue that $Q_n - \tilde{Q}_n = o_P^K(1)$ in (\mathcal{D}_K, d_K) for each compact $K \subset H$, where

$$\begin{aligned} \tilde{Q}_n(h) &\equiv 2h_1\sqrt{n}\mathbb{P}_n[\epsilon 1\{Z \leq \zeta_0\}] - h_1^2\mathbb{P}(Z \leq \zeta_0) \\ &\quad + 2h_2\sqrt{n}\mathbb{P}_n[\epsilon 1\{Z > \zeta_0\}] - h_2^2\mathbb{P}(Z > \zeta_0) \\ &\quad + n\mathbb{P}_n[-2(\alpha_0 - \beta_0)\epsilon - (\alpha_0 - \beta_0)^2]1\{\zeta_0 + h_3/n < Z \leq \zeta_0\} \\ &\quad + n\mathbb{P}_n[2(\alpha_0 - \beta_0)\epsilon - (\alpha_0 - \beta_0)^2]1\{\zeta_0 < Z \leq \zeta_0 + h_3/n\} \\ &\equiv \tilde{A}_n(h_1) + \tilde{B}_n(h_2) + \tilde{C}_n(h_3) + \tilde{D}_n(h_3). \end{aligned}$$

The superscript K in $o_P^K(1)$ indicates that the error is in terms of d_K . Fix a compact $K \subset H$. Note that by (14.7), $\tilde{A}_n(h_1) = n\mathbb{P}_n[A_1(\theta_{n,h}) - C_1(\theta_{n,h}) - A_1(\theta_0) + C_1(\theta_0)] + \tilde{E}_n(h)$, where

$$\begin{aligned} \tilde{E}_n(h) &= 2h_1\mathbb{G}_n[1\{Z \leq (\zeta_0 + h_3/n) \wedge \zeta_0\} - 1\{Z \leq \zeta_0\}] \\ &\quad - h_1^2[\mathbb{P}_n 1\{Z \leq (\zeta_0 + h_3/n) \wedge \zeta_0\} - \mathbb{P}(Z \leq \zeta_0)] \\ &\rightarrow 0 \end{aligned}$$

in probability, as $n \rightarrow \infty$, uniformly over $h \in K$. A similar analysis reveals the uniform equivalence of $\tilde{B}_n(h_2)$ and $n\mathbb{P}_n[A_2(\theta_{n,h}) - C_2(\theta_{n,h}) - A_2(\theta_0) + C_2(\theta_0)]$. It is fairly easy to see that $\tilde{C}_n(h_3)$ and $n\mathbb{P}_n[B_1(\theta_{n,h}) - D_1(\theta_{n,h}) - B_1(\theta_0) + D_1(\theta_0)]$ are asymptotically uniformly equivalent in probability as also $\tilde{D}_n(h_3)$ and $n\mathbb{P}_n[B_2(\theta_{n,h}) - D_2(\theta_{n,h}) - B_2(\theta_0) + C_2(\theta_0)]$. Thus $Q_n - \tilde{Q}_n$ goes to zero, in probability, uniformly over $h \in K$. Note that the potential jump points in h_3 for Q_n and \tilde{Q}_n remain the same, and thus $Q_n - \tilde{Q}_n = o_P^K(1)$ as desired.

Lemma 14.6 below shows that $\tilde{Q}_n \rightsquigarrow Q \equiv 2h_1Z_1 - h_1^2\mathbb{P}(Z \leq \zeta_0) + 2h_2Z_2 - h_2^2\mathbb{P}(Z > \zeta_0) + Q^+(h_3) + Q^-(h_3)1\{h_3 < 0\}$ in (\mathcal{D}_K, d_K) , where Z_1, Z_2, Q^+ and Q^- are all independent and Z_1 and Z_2 are mean zero Gaussian with respective variances $\sigma^2\mathbb{P}(Z \leq \zeta_0)$ and $\sigma^2\mathbb{P}(Z > \zeta_0)$. Let $s \mapsto \nu^+(s)$ be a right-continuous homogeneous Poisson process on $[0, \infty)$ with intensity parameter $f(\zeta_0)$ (recall that f is the density of ϵ), and let $s \mapsto \nu^-(s)$ be another Poisson process, independent of ν^+ , on $[-\infty, 0)$ which is left-continuous and goes backward in time with intensity $f(\zeta_0)$. Let $(V_k^+)_{k \geq 1}$ and $(V_k^-)_{k \geq 1}$ be independent sequences of i.i.d. random variables with V_1^+ being a realization of $2(\alpha_0 - \beta_0)\epsilon - (\alpha_0 - \beta_0)^2$ and V_1^- being

a realization of $-2(\alpha_0 - \beta_0)\epsilon - (\alpha_0 - \beta_0)^2$. Also define $V_0^+ = V_0^- = 0$ for convenience. Then $h_3 \mapsto Q^+(h_3) \equiv 1\{h_3 > 0\} \sum_{0 \leq k \leq \nu^+(h_3)} V_k^+$ and $h_3 \mapsto Q^-(h_3) \equiv 1\{h_3 < 0\} \sum_{0 \leq k \leq \nu^-(h_3)} V_k^-$.

Putting this all together, we conclude that $\hat{h} = (\hat{h}_1, \hat{h}_2, \hat{h}_3)$, where all three components are mutually independent, \hat{h}_1 and \hat{h}_2 are both mean zero Gaussian with respective variances $\sigma^2/P(Z \leq \zeta_0)$ and $\sigma^2/P(Z > \zeta_0)$, and where \hat{h}_3 is the smallest argmax of $h_3 \mapsto Q^+(h_3) + Q^-(h_3)$. Note that the expected value of both V_1^+ and V_1^- is $-(\alpha_0 - \beta_0)^2$. Thus $Q^+ + Q^-$ will be zero at $h_3 = 0$ and eventually always negative for all h_3 far enough away from zero. This means that the smallest argmax of $Q^+ + Q^-$ will be bounded in probability as desired.

LEMMA 14.6 For each compact $K \subset H$, $\tilde{Q}_n \rightsquigarrow Q$ in (\mathcal{D}_K, d_K) .

Proof. The details of the proof of weak convergence in the uniform norm follow very closely the proof of theorem 5 in Kosorok and Song (2006), and we omit the details. The required convergence of the jump locations follows from the assumed continuity of ϵ which ensures that the jump sizes will never be tied combined with the joint independence of ν^+ , ν^- , $(V_k^+)_{k \geq 1}$ and $(V_k^-)_{k \geq 1}$. \square

14.5.2 Monotone Density Estimation

This example is a special case of the cube root asymptotic results of Kim and Pollard (1990) which was analyzed in detail in section 3.2.14 of VW. Let X_1, \dots, X_n be a sample of size n from a Lebesgue density f on $[0, \infty)$ that is known to be decreasing. The maximum likelihood estimator \hat{f}_n of f is the non-increasing step function equal to the left derivative of the *least concave majorant* of the empirical distribution function \mathbb{F}_n . This \hat{f}_n is the celebrated Grenander estimator (Grenander, 1956). For a fixed value of $t > 0$, we will study the properties of $\hat{f}_n(t)$ under the assumption that f is differentiable at t with derivative $-\infty < f'(t) < 0$. Specifically, we will establish consistency of $\hat{f}_n(t)$, verify that the rate of convergence of \hat{f}_n is $n^{1/3}$, and derive weak convergence of $n^{1/3}(\hat{f}_n(t) - f(t))$. Existence of \hat{f}_n will be verified automatically as a consequence of consistency.

Consistency.

Let \hat{F}_n denote the least concave majorant of \mathbb{F}_n . In general, the least concave majorant of a function g is the smallest concave function h such that $h \geq g$. One can construct \hat{F}_n by imagining a a string tied at $(x, y) = (0, 0)$ which is pulled tight over the top of the function graph $(x, y = \mathbb{F}_n(x))$. The slope of each of the piecewise linear segments will be non-increasing, and the string (\hat{F}_n) will touch \mathbb{F}_n at two or points (x_j, y_j) , $j = 0, \dots, k$, where $k \geq 1$, $(x_0, y_0) \equiv (0, 0)$ and x_k is the last observation in the sample. For all

$x > x_k$, we set $\hat{F}_n(x) = 1$. Note also that \hat{F}_n is continuous. We leave it as an exercise to verify that this algorithm does indeed produce the least concave majorant of \mathbb{F}_n . The following lemma (Marshall's lemma) yields that \hat{F}_n is uniformly consistent for F . We save the proof as another exercise.

LEMMA 14.7 (*Marshall's lemma*) *Under the give conditions, $\sup_{t \geq 0} |\hat{F}_n(t) - F(t)| \leq \sup_{t \geq 0} |\mathbb{F}_n(t) - F(t)|$.*

Now fix $0 < \delta < t$, and note that by definition of \hat{f}_n ,

$$\frac{\hat{F}_n(t + \delta) - \hat{F}_n(t)}{\delta} \leq \hat{f}_n(t) \leq \frac{\hat{F}_n(t) - \hat{F}_n(t - \delta)}{\delta}.$$

By Marshall's lemma, the upper and lower bounds converge almost surely to $\delta^{-1}(F(t) - F(t - \delta))$ and $\delta^{-1}(F(t + \delta) - F(t))$, respectively. By the assumptions on F and the arbitrariness of δ , we obtain $\hat{f}_n(t) \xrightarrow{\text{as}^*} f(t)$.

Rate of convergence.

To determine the rate of convergence, we need to perform an interesting inverse transformation of the problem that will also be useful for obtaining the weak limiting distribution. Define the stochastic process $\{\hat{s}_n(a) : a > 0\}$ by $\hat{s}_n(a) = \operatorname{argmax}_{s \geq 0} \{\mathbb{F}_n(s) - as\}$, where the largest value is selected when multiple maximizers exist. The function \hat{s}_n is a sort of inverse of the function \hat{f}_n in the sense that $\hat{f}_n(t) \leq a$ if and only if $\hat{s}_n(a) \leq t$ for every $t \geq 0$ and $a > 0$. To see this, first assume that $\hat{f}_n(t) \leq a$. This means that the left derivative of \hat{F}_n is $\leq a$ at t . Hence a line of slope a which is moved down vertically from $+\infty$ will first touch \hat{F}_n at a point s_0 to the left of (or equal to) t . That point is also the point at which \hat{F}_n is furthest away from the line $s \mapsto as$ passing through the origin. Thus $s_0 = \operatorname{argmax}_{s \geq 0} \{\mathbb{F}_n(s) - as\}$, and hence $\hat{s}_n(a) \leq t$. Now suppose $\hat{s}_n(a) \leq t$. Then the argument can be taken in reverse to see that the slope of the line that touches \hat{F}_n at $\hat{s}_n(a)$ is less than or equal to the left derivative of \hat{F}_n at t , and thus $\hat{f}_n(t) \leq a$. Hence,

$$(14.9) \quad \mathbb{P}(n^{1/3}(\hat{f}_n(t) - f(t)) \leq x) = \mathbb{P}(\hat{s}_n(f(t) + xn^{-1/3}) \leq t),$$

and the desired rate and weak convergence result can be deduced from the argmax values of $x \mapsto \hat{s}_n(f(t) + xn^{-1/3})$. Applying the change of variable $s \mapsto t + g$ in the definition of \hat{s}_n , we obtain

$$\hat{s}_n(f(t) + xn^{-1/3}) - t = \operatorname{argmax}_{\{g > -t\}} \{\mathbb{F}_n(t + g) - (f(t) + xn^{-1/3})(t + g)\}.$$

In this manner, the probability on the left side of (14.9) is precisely $\mathbb{P}(\hat{g}_n \leq 0)$, where \hat{g}_n is the argmax above.

Now, by the previous argmax expression combined with the fact that the location of the maximum of a function does not change when the function is shifted vertically, we have $\hat{g}_n \equiv \operatorname{argmax}_{\{g > -t\}} \{M_n(g) \equiv \mathbb{F}_n(t + g) -$

$\mathbb{F}_n(t) - f(t)g - xgn^{-1/3}$. It is not hard to see that $\hat{g}_n = O_P(1)$ and that $M_n(g) \xrightarrow{P} M(g) \equiv F(t+g) - F(t) - f(t)g$ uniformly on compacts, and thus $\hat{g}_n = o_P(1)$. We now utilize theorem 14.4 to obtain the rate for \hat{g}_n , with the metric $d(\theta_1, \theta_2) = |\theta_1 - \theta_2|$, $\theta = g$, $\theta_0 = 0$ and $\bar{d} = d$. Note the fact that $M_n(0) = M(0) = 0$ will simplify the calculations. It is now easy to see that $M(g) \lesssim -g^2$, and by using theorem 11.2, that

$$\begin{aligned} \mathbb{E}^* \sup_{|g| < \delta} \sqrt{n} |M_n(g) - M(g)| &\leq \mathbb{E}^* \sup_{|g| < \delta} |\mathbb{G}_n(1\{X \leq t+g\} - 1\{X \leq t\})| \\ &\quad + O(\sqrt{n}\delta n^{-1/3}) \\ &\lesssim \phi_n(\delta) \equiv \delta^{1/2} + \sqrt{n}\delta n^{-1/3}. \end{aligned}$$

Clearly, $\phi_n(\delta)/\delta^\alpha$ is decreasing in δ for $\alpha = 3/2$. Since $n^{2/3}\phi_n(n^{-1/3}) = n^{1/2} + n^{1/6}n^{-1/3} = O(n^{1/2})$, theorem 14.4 yields $n^{1/3}\hat{g}_n = O_P(1)$. We show in the next section how this enables weak convergence of $n^{1/3}(\hat{f}(t) - f(t))$.

Weak convergence.

Let $\hat{h}_n = n^{1/3}\hat{g}_n$, and note that since the maximum of a function does not change when the function is multiplied by a constant, we have that \hat{h}_n is the argmax of the process

$$\begin{aligned} (14.10) \quad h &\mapsto n^{2/3}M_n(n^{-1/3}h) \\ &= n^{2/3}(\mathbb{P}_n - P) \left(1\{X \leq t + hn^{-1/3}\} - 1\{X \leq t\} \right) \\ &\quad + n^{2/3} \left[F(t + hn^{-1/3}) - F(t) - f(t)hn^{-1/3} \right] - xh. \end{aligned}$$

Fix $0 < K < \infty$, and apply theorem 11.18 to the sequence of classes $\mathcal{F}_n = \{n^{1/6} (1\{X \leq t + hn^{-1/3}\} - 1\{X \leq t\}) : -K \leq h \leq K\}$ with envelope sequence $F_n = n^{1/6}1\{t - Kn^{-1/3} \leq X \leq t + Kn^{-1/3}\}$, to obtain that the process on the right side of (14.10) converges in $\ell^\infty(-K, K)$ to

$$h \mapsto \mathbb{H}(h) \equiv \sqrt{f(t)}\mathbb{Z}(h) + \frac{1}{2}f'(t)h^2 - xh,$$

where \mathbb{Z} is a two-sided Brownian motion originating at zero (two independent Brownian motions starting at zero, one going to the right of zero and the other going to the left). From the previous paragraph, we know that $\hat{h}_n = O_P(1)$. Since it is not hard to verify that \mathbb{H} is continuous with a unique maximum, the argmax theorem now yields by the arbitrariness of K that $\hat{h}_n \rightsquigarrow \hat{h}$, where $\hat{h} = \operatorname{argmax} \mathbb{H}$. By exercise 14.6.9 below, we can simplify the form of \hat{h} to $|4f'(t)f(t)|^{1/3} \operatorname{argmax}_h \{\mathbb{Z}(h) - h^2\}$. Thus by (14.9), we obtain that

$$n^{1/3}(\hat{f}_n(t) - f(t)) \rightsquigarrow |f'(t)f(t)|^{1/3}\mathbb{C},$$

where the random variable $\mathbb{C} \equiv \operatorname{argmax}_h \{\mathbb{Z}(h) - h^2\}$ has Chernoff's distribution (see Groeneboom, 1989).

14.6 Exercises

14.6.1. Show that for a sequence X_n of measurable, Euclidean random variables which are finite almost surely, measurability plus asymptotic tightness implies uniform tightness.

14.6.2. For a metric space (\mathbb{D}, d) , let $H : \mathbb{D} \mapsto [0, \infty]$ be a function such that $H(x_0) = 0$ for a point $x_0 \in \mathbb{D}$ and $H(x_n) \rightarrow 0$ implies $d(x_n, x_0) \rightarrow 0$ for any sequence $\{x_n\} \in \mathbb{D}$. Show that there exists a non-decreasing cadlag function $f : [0, \infty] \mapsto [0, \infty]$ that satisfies both $f(0) = 0$ and $d(x, x_0) \leq f(|H(x)|)$ for all $x \in \mathbb{D}$. Hint: Use the fact that the given conditions on H imply the existence of a decreasing sequence $0 < \tau_n \downarrow 0$ such that $H(x) < \tau_n$ implies $d(x, x_0) < 1/n$, and note that it is permissible to have $f(u) = \infty$ for all $u \geq \tau_1$.

14.6.3. In the proof of theorem 14.4, verify that for fixed $c < \infty$ and $\alpha < 2$,

$$\sum_{j \geq M} \frac{2^{j\alpha}}{2^{2j} - c} \rightarrow 0,$$

as $M \rightarrow \infty$.

14.6.4. In the context of the last paragraph of the proof of theorem 2.13, given in section 14.4, complete the verification of the conditions of theorem 11.18.

14.6.5. Consider the function $\theta \mapsto M(\theta)$ defined in section 14.5.1. Show that it has a unique maximum over $\mathbb{R}^2 \times [a, b]$. Also show that the maximum is not unique if $\alpha_0 = \beta_0$.

14.6.6. Verify (14.8).

14.6.7. Verify that the algorithm described in the second paragraph of section 14.5.2 does indeed generate the least concave majorant of \mathbb{F}_n .

14.6.8. The goal of this exercise is to prove Marshall's lemma given in section 14.5.2. Denoting $A_n(t) \equiv \hat{F}_n(t) - F(t)$ and $B_n(t) \equiv \mathbb{F}_n(t) - F(t)$, the proof can be broken into the following steps:

(a) Show that $0 \geq \inf_{t \geq 0} A_n(t) \geq \inf_{t \geq 0} B_n(t)$.

(b) Show that

i. $\sup_{t \geq 0} A_n(t) \geq 0$ and $\sup_{t \geq 0} B_n(t) \geq 0$.

ii. If $\sup_{t \geq 0} B_n(t) = 0$, then $\sup_{t \geq 0} A_n(t) = 0$.

iii. If $\sup_{t \geq 0} B_n(t) > 0$, then $\sup_{t \geq 0} A_n(t) \leq \sup_{t \geq 0} B_n(t)$ (this last step is tricky).

Now verify that $0 \leq \sup_{t \geq 0} A_n(t) \leq \sup_{t \geq 0} B_n(t)$.

(c) Now complete the proof.

14.6.9. Let $\{\mathbb{Z}(h) : h \in \mathbb{R}\}$ be a standard two-sided Brownian motion with $\mathbb{Z}(0) = 0$. (The process is zero-mean Gaussian and the increment $\mathbb{Z}(g) - \mathbb{Z}(h)$ has variance $|g - h|$.) Then $\operatorname{argmax}_h \{a\mathbb{Z}(h) - bh^2 - ch\}$ is equal in distribution to $(a/b)^{2/3} \operatorname{argmax}_g \{\mathbb{Z}(g) - g^2\} - c/(2b)$, where $a, b, c > 0$. Hint: The process $h \mapsto \mathbb{Z}(\sigma h - \mu)$ is equal in distribution to the process $h \mapsto \sqrt{\sigma}\mathbb{Z}(g) - \mathbb{Z}(\mu)$, where $\sigma \geq 0$ and $\mu \in \mathbb{R}$. Apply the change of variable $h = (a/b)^{2/3}g - c/(2b)$ and note that the location of a maximum does not change by multiplication by a positive constant or a vertical shift.

14.7 Notes

Theorem 14.1 and lemma 14.2 are theorem 3.2.2 and lemma 3.2.1, respectively, of VW, while theorem 14.4 and corollary 14.5 are modified versions of theorem 3.2.5 and corollary 3.2.6 of VW. The monotone density estimation example in section 14.5.2 is a variation of example 3.2.14 of VW. The limiting behavior of the Grenander estimator of this example was obtained by Prakasa Rao (1969). Exercise 14.6.8 is an expanded version of exercise 24.5 of van der Vaart (1998) and exercise 14.6.9 is an expanded version of exercise 3.2.5 of VW.

References

- Andersen, P. K., Borgan, Ø., Gill, R. D., and Keiding, N. (1993). *Statistical Models Based on Counting Processes*. Springer-Verlag, New York.
- Andrews, D. W. K. (1991). An empirical process central limit theorem for dependent non-identically distributed random variables. *Journal of Multivariate Analysis*, 38:187–203.
- Arcones, M. A. and Yu, B. (1994). Central limit theorems for empirical and U -processes of stationary mixing sequences. *Journal of Theoretical Probability*, 7:47–71.
- Bassett, Jr., G. and Koenker, R. (1978). Asymptotic theory of least absolute error regression. *Journal of the American Statistical Association*, 73:618–622.
- Bickel, P. J., Klaassen, C. A., Ritov, Y., and Wellner, J. A. (1997). *Efficient and Adaptive Estimation for Semiparametric Models*. Springer-Verlag, New York.
- Bickel, P. J. and Rosenblatt, M. (1973). On some global measures of the deviations of density function estimates. *Annals of Statistics*, 1:1071–1095.
- Bilius, Y., Gu, M., and Ying, Z. (1997). Towards a general theory for Cox model with staggered entry. *Annals of Statistics*, 25:662–682.
- Billingsley, P. (1968). *Convergence of Probability Measures*. Wiley, New York.

- Billingsley, P. (1986). *Probability and Measure*. Wiley, New York, second edition.
- Bradley, R. C. (1986). Basic properties of strong mixing conditions. In Eberlein, E. and Taqqu, M. S., editors, *Dependence in Probability and Statistics: A Survey of Recent Results*, pages 165–192. Birkhäuser, Basel.
- Bühlmann, P. (1995). The blockwise bootstrap for general empirical processes of stationary sequences. *Stochastic Processes and Their Applications*, 58:247–265.
- Cantelli, F. P. (1933). Sulla determinazione empirica delle leggi di probabilità. *Giornale dell'Istituto Italiano degli Attuari*, 4:421–424.
- Chang, M. N. (1990). Weak convergence of a self-consistent estimator of the survival function with doubly censored data. *Annals of Statistics*, 18:391–404.
- Cox, D. R. (1975). Partial likelihood. *Biometrika*, 62:269–276.
- Dedecker, J. and Louhichi, S. (2002). Maximal inequalities and empirical central limit theorems. In Dehling, H., Mikosch, T., and Sørensen, M., editors, *Empirical Process Techniques for Dependent Data*, pages 137–159. Birkhäuser, Boston.
- Dehling, H., Mikosch, T., and Sørensen, M., editors (2002). *Empirical Process Techniques for Dependent Data*. Birkhäuser, Boston.
- Dehling, H. and Philipp, W. (2002). Empirical process techniques for dependent data. In Dehling, H., Mikosch, T., and Sørensen, M., editors, *Empirical Process Techniques for Dependent Data*, pages 3–113. Birkhäuser, Boston.
- Dehling, H. and Taqqu, M. (1989). The empirical process of some long-range dependent sequences with an application to u-statistics. *Annals of Statistics*, 17:1767–1783.
- Donsker, M. D. (1952). Justification and extension of doob's heuristic approach to the kolmogorov-smirnov theorems. *Annals of Mathematical Statistics*, 23:277–281.
- Dudley, R. M. and Philipp, W. (1983). Invariance principles for sums of banach space valued random elements and empirical processes. *Probability Theory and Related Fields*, 62:509–552.
- Dugundji, J. (1951). An extension of Tietze's theorem. *Pacific Journal of Mathematics*, 1:353–367.

- Dvoretzky, A., Kiefer, J., and Wolfowitz, J. (1956). Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. *Annals of Mathematical Statistics*, 27:642–669.
- Eberlein, E. and Taqqu, M. S., editors (1986). *Dependence in Probability and Statistics: A Survey of Recent Results*. Birkhäuser, Basel.
- Efron, B. (1967). The two sample problem with censored data. *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 4:831–855.
- Fine, J. P., Yan, J., and Kosorok, M. R. (2004). Temporal process regression. *Biometrika*, 91:683–703.
- Fleming, T. R. and Harrington, D. P. (1991). *Counting Processes and Survival Analysis*. Wiley, New York.
- Glivenko, V. (1933). Sulla determinazione empirica della leggi di probabilità. *Giornale dell'Istituto Italiano Degli Attuari*, 4:92–99.
- Godambe, V. P. (1960). An optimum property of regular maximum likelihood estimation. *Annals of Mathematical Statistics*, 31:1208–1211.
- Grenander, U. (1956). On the theory of mortality measurement, part ii. *Skandinavisk Aktuarietidskrift*, 39:125–153.
- Groeneboom, P. (1989). Brownian motion with a parabolic drift and airy functions. *Probability Theory and Related Fields*, 81:79–109.
- Hall, P. (1992). *The Bootstrap and Edgeworth Expansion*. Springer-Verlag, New York.
- Huber, P. J. (1981). *Robust Statistics*. Wiley, New York.
- Jameson, J. O. (1974). *Topology and Normed Spaces*. Chapman and Hall, London.
- Johnson, W. B., Lindenstrauss, J., and Schechtman, G. (1986). Extensions of lipschitz maps into banach spaces. *Israel Journal of Mathematics*, 54:129–138.
- Karatzas, I. and Shreve, S. E. (1991). *Brownian Motion and Stochastic Calculus*. Springer, New York, 2 edition.
- Kim, J. and Pollard, D. (1990). Cube root asymptotics. *Annals of Statistics*, 18:191–219.
- Komlós, J., Major, P., and Tusnády, G. (1975). An approximation of partial sums of independent rv's and the sample df. i. *Probability Theory and Related Fields*, 32:111–131.

- Komlós, J., Major, P., and Tusnády, G. (1976). An approximation of partial sums of independent rv's and the sample df. ii. *Probability Theory and Related Fields*, 34:33–58.
- Kosorok, M. R. (1999). Two-sample quantile tests under general conditions. *Biometrika*, 86:909–921.
- Kosorok, M. R. (2003). Bootstraps of sums of independent but not identically distributed stochastic processes. *Journal of Multivariate Analysis*, 84:299–318.
- Kosorok, M. R. and Song, R. Inference under right censoring for transformation models with a change-point based on a covariate threshold. *Annals of Statistics*. To appear.
- Kress, R. (1999). *Linear Integral Equations*. Springer, New York, second edition.
- Künsch, H. R. (1989). The jackknife and the bootstrap for general stationary observation. *Annals of Statistics*, 17:1217–1241.
- Liu, R. Y. and Singh, K. (1992). Moving blocks jackknife and bootstrap capture weak dependence. In LePage, R. and Billard, L., editors, *Exploring the Limits of Bootstrap*, pages 225–248. Wiley, New York.
- Mammen, E. and van de Geer, S. (1997). Penalized quasi-likelihood estimation in partial linear models. *Annals of Statistics*, 25:1014–1035.
- Massart, P. (1990). The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *Annals of Probability*, 18:1269–1283.
- Meggison, R. E. (1998). *An Introduction to Banach Space Theory*. Springer, New York.
- Murphy, S. A. and van der Vaart, A. W. (2000). On profile likelihood. with comments and a rejoinder by the authors. *Journal of the American Statistical Association*, 95:449–485.
- Naik-Nimbalkar, U. V. and Rajarshi, M. B. (1994). Validity of blockwise bootstrap for empirical processes with stationary observations. *Annals of Statistics*, 22:980–994.
- Peligrad, M. (1998). On the blockwise bootstrap for empirical processes for stationary sequences. *Annals of Probability*, 26:877–901.
- Politis, D. N. and Romano, J. P. (1994). Large sample confidence regions based on subsampling under minimal assumptions. *Annals of Statistics*, 22:2031–2050.

- Pollard, D. (1990). *Empirical Processes: Theory and Applications*. NSF-CBMS Regional Conference Series in Probability and Statistics 2. Institute of Mathematical Statistics and American Statistical Association, Hayward, California.
- Pollard, D. (1991). Asymptotics for least absolute deviation regression estimators. *Econometric Theory*, 7:186–199.
- Polonik, W. (1995). Measuring mass concentrations and estimating density contour clusters—An excess mass approach. *Annals of Statistics*, 23:855–881.
- Praestgaard, J. and Wellner, J. A. (1993). Exchangeably weighted bootstraps of the general empirical process. *Annals of Probability*, 21:2053–2086.
- Radulović, D. (1996). The bootstrap for empirical processes based on stationary observations. *Stochastic Processes and Their Applications*, 65:259–279.
- Rao, B. L. S. P. (1969). Estimation of a unimodal density. *Sankya, Series A*, 31:23–36.
- Rubin, D. B. (1981). The Bayesian bootstrap. *Annals of Statistics*, 9:130–134.
- Rudin, W. (1987). *Real and Complex Analysis*. McGraw-Hill, Inc., New York, third edition.
- Rudin, W. (1991). *Functional Analysis*. McGraw-Hill, Inc., New York, second edition.
- Shao, J. and Tu, D. (1995). *The Jackknife and Bootstrap*. Springer-Verlag, New York.
- Strassen, V. (1964). An invariance principle for the law of the iterated logarithm. *Probability Theory and Related Fields*, 3:211–226.
- van de Geer, S. A. (2000). *Empirical Processes in M-Estimation*. Cambridge University Press, New York.
- van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge University Press, New York.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Processes: With Applications in Statistics*. Springer-Verlag, New York.
- Wei, L. J. (1978). The adaptive biased coin design for sequential experiments. *Annals of Statistics*, 6:92–100.

- Wu, W. B. (2003). Empirical processes of long-memory sequences. *Bernoulli*, 9:809–831.
- Yu, B. (1994). Rates of convergence for empirical processes of stationary mixing sequences. *Annals of Probability*, 22:94–116.

List of Symbols

<p>B^*: adjunct of the operator B, 39</p> <p>BL_1: space of functions with Lipschitz norm bounded by 1, 19</p> <p>$\overset{\text{as*}}{\rightsquigarrow}_M$: conditional convergence of bootstrap outer almost surely, 20</p> <p>$\overset{\text{P}}{\rightsquigarrow}_M$: conditional convergence of bootstrap in probability, 19</p> <p>\mathbb{B}: standard Brownian bridge process, 11</p> <p>$C[a, b]$: continuous, real functions on $[a, b]$, 23</p> <p>$C_b(\mathbb{D})$: space of bounded, continuous maps $f : \mathbb{D} \mapsto \mathbb{R}$, 14</p> <p>$\overset{\text{as}}{\rightarrow}$: convergence almost surely, 10</p> <p>$\overset{\text{as*}}{\rightarrow}$: convergence outer almost surely, 14</p> <p>$\overset{\text{P}}{\rightarrow}$: convergence in probability, 14</p> <p>\rightsquigarrow: weak convergence, 11</p> <p>$\text{conv}\mathcal{F}$: convex hull of a class \mathcal{F}, 152</p>	<p>$\overline{\text{conv}}\mathcal{F}$: closed convex hull of a class \mathcal{F}, 152</p> <p>$\text{sconv}\mathcal{F}$: symmetric convex hull of a class \mathcal{F}, 152</p> <p>$\overline{\text{sconv}}\mathcal{F}$: symmetric closed convex hull of a class \mathcal{F}, 152</p> <p>$\text{cov}[X, Y]$: covariance of X and Y, 11</p> <p>$D[a, b]$: space of cadlag functions on $[a, b]$, 22</p> <p>δB: boundary of the set B, 104</p> <p>δ_x: point mass at x or Dirac measure, 10</p> <p>E_*: inner expectation, 14</p> <p>E^*: outer expectation, 14</p> <p>F: distribution function, 9</p> <p>\mathcal{F}, \mathcal{G}: collections of functions, 10, 19</p> <p>\mathbb{F}_n: empirical distribution function, 10</p> <p>\mathbb{G}: general Brownian bridge, 11</p> <p>\mathbb{G}_n: empirical process, 11</p>
---	---

- G_n : standardized empirical distribution function, 11
 $1\{A\}$: indicator of A , 4
 $\langle \cdot, \cdot \rangle$: inner product, 41
 $J_{[]}$: bracketing integral, 17
 $\ell^\infty(T)$: set of all uniformly bounded real functions on T , 11
 L_r : equivalence classes of r -integrable functions, 16
 $L_2^0(P)$: mean zero subspace of $L_2(P)$, 35
 $\text{lin } \mathbb{H}$: linear span of \mathbb{H} , 39
 $\overline{\text{lin}} \mathbb{H}$: closed linear span of \mathbb{H} , 39
 \mapsto : function specifier (“maps to”), 11
 T_* : maximal measurable minorant of T , 14
 $a \vee b$: maximum of a and b , 19
 $(\mathbb{D}, d), (\mathbb{E}, e)$: metric spaces, 13
 T^* : minimal measurable majorant of T , 14
 $a \wedge b$: minimum of a and b , 6
 $\|\cdot\|_{2,1}$: special variant of L_2 norm, 20
 $\|\cdot\|_{r,P}$: $L_r(P)$ norm, 16
 $\|\cdot\|_\infty$: uniform norm, 19
 $N_{[]}$: bracketing number, 16
 Ω, Ω_n : sample spaces, 14, 103
 $\|\cdot\|_\psi$: Orlicz ψ -norm, 124
 \mathbb{P}_n° : symmetrized empirical process, 134
 $\dot{\mathcal{P}}_P$: tangent set, 35
 \mathcal{P} : collection of probability measures, 33
 \mathbb{P}_n : empirical measure, 10
 P, Q : probability measures on underlying sample space \mathcal{X} , 10, 17
 P_* : inner probability, 14
 P^* : outer probability, 14
 ψ_p : the function $x \mapsto \psi_p(x) = \frac{e^{x^p} - 1}{e^{x^p} + 1}$, 125
 \mathbb{R} : real numbers, 3
 ρ : semimetric on a set T , 15
 T : index set for a stochastic process, 9
 $UC(T, \rho)$: set of uniformly continuous functions from (T, ρ) to \mathbb{R} , 15
 $V(\mathcal{C}), V(\mathcal{F})$: VC-index of a set \mathcal{C} or function class \mathcal{F} , 150, 151
 $\text{var}[X]$: variance of X , 15
 \mathcal{X} : sample space, 10
 (X, \mathcal{O}) : topological space, 78