

SAS Programming in Clinical Trials

Chapter 3. SAS STAT

Dongfeng Li

Autumn 2010

Chapter Contents

- ▶ **Some statistics background;**
- ▶ Descriptive statistics: concepts and programs;
- ▶ Comparing means and proportions;
- ▶ Analysis of variance.
- ▶ Students should master the basic concepts, descriptive statistics measures and graphs, basic hypothesis testing, basic analysis of variance.

Chapter Contents

- ▶ **Some statistics background;**
- ▶ **Descriptive statistics: concepts and programs;**
- ▶ Comparing means and proportions;
- ▶ Analysis of variance.
- ▶ Students should master the basic concepts, descriptive statistics measures and graphs, basic hypothesis testing, basic analysis of variance.

Chapter Contents

- ▶ Some statistics background;
- ▶ Descriptive statistics: concepts and programs;
- ▶ Comparing means and proportions;
- ▶ Analysis of variance.
- ▶ Students should master the basic concepts, descriptive statistics measures and graphs, basic hypothesis testing, basic analysis of variance.

Chapter Contents

- ▶ Some statistics background;
- ▶ Descriptive statistics: concepts and programs;
- ▶ Comparing means and proportions;
- ▶ Analysis of variance.
- ▶ Students should master the basic concepts, descriptive statistics measures and graphs, basic hypothesis testing, basic analysis of variance.

Chapter Contents

- ▶ Some statistics background;
- ▶ Descriptive statistics: concepts and programs;
- ▶ Comparing means and proportions;
- ▶ Analysis of variance.
- ▶ Students should master the basic concepts, descriptive statistics measures and graphs, basic hypothesis testing, basic analysis of variance.

Section Contents

- ▶ Review statistics concepts:
- ▶ Distribution, discrete distribution, continuous distribution, PDF, CDF, quantile. Normal distribution.
- ▶ Mean, median, variance, standard deviation, interquartile range, skewness, kurtosis.
- ▶ Population, parameter, sample, statistics, estimates, sampling distribution.
- ▶ MLE, standard error.
- ▶ Hypothesis tests, two types of errors, p-value.

Section Contents

- ▶ Review statistics concepts:
- ▶ Distribution, discrete distribution, continuous distribution, PDF, CDF, quantile. Normal distribution.
- ▶ Mean, median, variance, standard deviation, interquantile range, skewness, kurtosis.
- ▶ Population, parameter, sample, statistics, estimates, sampling distribution.
- ▶ MLE, standard error.
- ▶ Hypothesis tests, two types of errors, p-value.

Section Contents

- ▶ Review statistics concepts:
- ▶ Distribution, discrete distribution, continuous distribution, PDF, CDF, quantile. Normal distribution.
- ▶ Mean, median, variance, standard deviation, interquantile range, skewness, kurtosis.
- ▶ Population, parameter, sample, statistics, estimates, sampling distribution.
- ▶ MLE, standard error.
- ▶ Hypothesis tests, two types of errors, p-value.

Section Contents

- ▶ Review statistics concepts:
- ▶ Distribution, discrete distribution, continuous distribution, PDF, CDF, quantile. Normal distribution.
- ▶ Mean, median, variance, standard deviation, interquantile range, skewness, kurtosis.
- ▶ Population, parameter, sample, statistics, estimates, sampling distribution.
- ▶ MLE, standard error.
- ▶ Hypothesis tests, two types of errors, p-value.

Section Contents

- ▶ Review statistics concepts:
- ▶ Distribution, discrete distribution, continuous distribution, PDF, CDF, quantile. Normal distribution.
- ▶ Mean, median, variance, standard deviation, interquantile range, skewness, kurtosis.
- ▶ Population, parameter, sample, statistics, estimates, sampling distribution.
- ▶ MLE, standard error.
- ▶ Hypothesis tests, two types of errors, p-value.

Section Contents

- ▶ Review statistics concepts:
- ▶ Distribution, discrete distribution, continuous distribution, PDF, CDF, quantile. Normal distribution.
- ▶ Mean, median, variance, standard deviation, interquantile range, skewness, kurtosis.
- ▶ Population, parameter, sample, statistics, estimates, sampling distribution.
- ▶ MLE, standard error.
- ▶ Hypothesis tests, two types of errors, p-value.

- ▶ **Random Variable:**
 - ▶ **Discrete**, such as sex, patient/control, age group.
 - ▶ **Continuous**, such as weight, blood pressure.
- ▶ **Distribution:** used to describe the relative chance of taking some value.
 - ▶ For discrete variable X , use $P(X = x_i)$, where $\{x_i\}$ are the value set of X . Called **probability mass function(PMF)**.
 - ▶ For continuous variable X , use the **probability density function(PDF)** $f(x)$, where $P(X \in (x - \epsilon, x + \epsilon)) \propto f(x)(2\epsilon)$.

- ▶ **Random Variable:**
 - ▶ **Discrete**, such as sex, patient/control, age group.
 - ▶ **Continuous**, such as weight, blood pressure.
- ▶ **Distribution:** used to describe the relative chance of taking some value.
 - ▶ For discrete variable X , use $P(X = x_i)$, where $\{x_i\}$ are the value set of X . Called **probability mass function(PMF)**.
 - ▶ For continuous variable X , use the **probability density function(PDF)** $f(x)$, where $P(X \in (x - \epsilon, x + \epsilon)) \propto f(x)(2\epsilon)$.

- ▶ **Random Variable:**
 - ▶ **Discrete**, such as sex, patient/control, age group.
 - ▶ **Continuous**, such as weight, blood pressure.
- ▶ **Distribution:** used to describe the relative chance of taking some value.
 - ▶ For discrete variable X , use $P(X = x_i)$, where $\{x_i\}$ are the value set of X . Called **probability mass function(PMF)**.
 - ▶ For continuous variable X , use the **probability density function(PDF)** $f(x)$, where $P(X \in (x - \epsilon, x + \epsilon)) \propto f(x)(2\epsilon)$.

- ▶ **Random Variable:**
 - ▶ **Discrete**, such as sex, patient/control, age group.
 - ▶ **Continuous**, such as weight, blood pressure.
- ▶ **Distribution:** used to describe the relative chance of taking some value.
 - ▶ For discrete variable X , use $P(X = x_i)$, where $\{x_i\}$ are the value set of X . Called **probability mass function(PMF)**.
 - ▶ For continuous variable X , use the **probability density function(PDF)** $f(x)$, where $P(X \in (x - \epsilon, x + \epsilon)) \propto f(x)(2\epsilon)$.

- ▶ **Random Variable:**
 - ▶ **Discrete**, such as sex, patient/control, age group.
 - ▶ **Continuous**, such as weight, blood pressure.
- ▶ **Distribution:** used to describe the relative chance of taking some value.
 - ▶ For discrete variable X , use $P(X = x_i)$, where $\{x_i\}$ are the value set of X . Called **probability mass function(PMF)**.
 - ▶ For continuous variable X , use the **probability density function(PDF)** $f(x)$, where $P(X \in (x - \epsilon, x + \epsilon)) \propto f(x)(2\epsilon)$.

- ▶ **Random Variable:**
 - ▶ **Discrete**, such as sex, patient/control, age group.
 - ▶ **Continuous**, such as weight, blood pressure.
- ▶ **Distribution:** used to describe the relative chance of taking some value.
 - ▶ For discrete variable X , use $P(X = x_i)$, where $\{x_i\}$ are the value set of X . Called **probability mass function(PMF)**.
 - ▶ For continuous variable X , use the **probability density function(PDF)** $f(x)$, where $P(X \in (x - \epsilon, x + \epsilon)) \propto f(x)(2\epsilon)$.

- Cumulative distribution function(CDF) $F(x)$:

$$F(x) = P(X \leq x)$$

$$P(x \in (a, b]) = F(b) - F(a)$$

- For discrete distribution,

$$F(x) = \sum_{x_i \leq x} P(X = x_i)$$

- For continuous distribution,

$$F(x) = \int_{-\infty}^x f(t)dt$$

- Cumulative distribution function(CDF) $F(x)$:

$$F(x) = P(X \leq x)$$

$$P(x \in (a, b]) = F(b) - F(a)$$

- For discrete distribution,

$$F(x) = \sum_{x_i \leq x} P(X = x_i)$$

- For continuous distribution,

$$F(x) = \int_{-\infty}^x f(t)dt$$

- Cumulative distribution function(CDF) $F(x)$:

$$F(x) = P(X \leq x)$$

$$P(x \in (a, b]) = F(b) - F(a)$$

- For discrete distribution,

$$F(x) = \sum_{x_i \leq x} P(X = x_i)$$

- For continuous distribution,

$$F(x) = \int_{-\infty}^x f(t)dt$$

- ▶ **Quantile function**: the inverse of CDF

$$q(p) = F^{-1}(p), p \in (0, 1)$$

if $F(x)$ is 1-1 mapping.

- ▶ Generally, $q(p) = x_p$ where

$$P(X \leq x_p) \geq p, P(X \geq x_p) \geq 1 - p$$

(x_p can be non-unique.)

- ▶ **Quantile function**: the inverse of CDF

$$q(p) = F^{-1}(p), p \in (0, 1)$$

if $F(x)$ is 1-1 mapping.

- ▶ Generally, $q(p) = x_p$ where

$$P(X \leq x_p) \geq p, P(X \geq x_p) \geq 1 - p$$

(x_p can be non-unique.)

The normal distribution

- ▶ Standard normal distribution(N(0,1)), PDF

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

- ▶ Normal distribution $N(\mu, \sigma^2)$, PDF

$$f(x) = \varphi\left(\frac{x - \mu}{\sigma}\right) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}$$

The normal distribution

- ▶ Standard normal distribution(N(0,1)), PDF

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

- ▶ Normal distribution $N(\mu, \sigma^2)$, PDF

$$f(x) = \varphi\left(\frac{x - \mu}{\sigma}\right) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}$$

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **median**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **median**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **median**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **median**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **meadian**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **meadian**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **meadian**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ PDF is a curve with infinite number of points.
- ▶ We can use some numbers to describe the key part of a distribution.
- ▶ Firstly, the **location measurement**.
 - ▶ The **mean** EX .
 - ▶ The **meadian**, $x_{0.5} = q(0.5)$ where

$$P(x \leq x_{0.5}) \geq 0.5, P(x \geq x_{0.5}) \geq 0.5$$

(can be non-unique).

- ▶ **Variability measurement**.
 - ▶ The **standard deviation** $\sigma_X = \sqrt{E(X - EX)^2}$.
 - ▶ **Interquantile range** $q(0.75) - q(0.25)$.

- ▶ Shape characteristics. **Skewness**

$$E \left(\frac{X - EX}{\sigma_X} \right)^3$$

- ▶ Symmetric, left-skewed, right-skewed density.
- ▶ Shape characteristics. **Kurtosis**

$$E \left(\frac{X - EX}{\sigma_X} \right)^4 - 3$$

- ▶ Heavy tail problem.
- ▶ Multimodal distribution. E.g., the weight of 10 year old and 20 year old mixed together.

- ▶ Shape characteristics. **Skewness**

$$E \left(\frac{X - EX}{\sigma_X} \right)^3$$

- ▶ Symmetric, left-skewed, right-skewed density.

- ▶ Shape characteristics. **Kurtosis**

$$E \left(\frac{X - EX}{\sigma_X} \right)^4 - 3$$

- ▶ Heavy tail problem.
- ▶ Multimodal distribution. E.g., the weight of 10 year old and 20 year old mixed together.

- ▶ Shape characteristics. **Skewness**

$$E \left(\frac{X - EX}{\sigma_X} \right)^3$$

- ▶ Symmetric, left-skewed, right-skewed density.
- ▶ Shape characteristics. **Kurtosis**

$$E \left(\frac{X - EX}{\sigma_X} \right)^4 - 3$$

- ▶ Heavy tail problem.
- ▶ Multimodel distribution. E.g., the weight of 10 year old and 20 year old mixed together.

- ▶ Shape characteristics. **Skewness**

$$E \left(\frac{X - EX}{\sigma_X} \right)^3$$

- ▶ Symmetric, left-skewed, right-skewed density.
- ▶ Shape characteristics. **Kurtosis**

$$E \left(\frac{X - EX}{\sigma_X} \right)^4 - 3$$

- ▶ Heavy tail problem.
- ▶ Multimodel distribution. E.g., the weight of 10 year old and 20 year old mixed together.

- ▶ Shape characteristics. **Skewness**

$$E \left(\frac{X - EX}{\sigma_X} \right)^3$$

- ▶ Symmetric, left-skewed, right-skewed density.
- ▶ Shape characteristics. **Kurtosis**

$$E \left(\frac{X - EX}{\sigma_X} \right)^4 - 3$$

- ▶ Heavy tail problem.
- ▶ Multimodel distribution. E.g., the weight of 10 year old and 20 year old mixed together.

- ▶ A **population** is modeled by a random variable or a distribution.
- ▶ **Population parameters**: unknown numbers which could decide the distribution. E.g., for $N(\mu, \sigma^2)$ population, (μ, σ^2) are the two unknown population parameters.

- ▶ A **population** is modeled by a random variable or a distribution.
- ▶ **Population parameters**: unknown numbers which could decide the distribution. E.g., for $N(\mu, \sigma^2)$ population, (μ, σ^2) are the two unknown population parameters.

- ▶ A **sample**(typically) is n observations draw independently from the population, X_1, X_2, \dots, X_n .
- ▶ A sample can be regarded as n numbers, or n iid random variables.
- ▶ **Statistics**: calculate some values to estimate the distribution of the population. Can be regarded as a random variable. Such as sample mean, sample standard deviation. Can be used to estimate unknown population parameters, called **estimates**.
- ▶ **Sampling distribution**: The distribution of an estimator or other statistic. E.g., if X_1, \dots, X_n is a sample from $N(\mu, \sigma^2)$, then $\bar{X} \sim N(\mu, \sigma^2/n)$.

- ▶ A **sample** (typically) is n observations drawn independently from the population, X_1, X_2, \dots, X_n .
- ▶ A sample can be regarded as n numbers, or n iid random variables.
- ▶ **Statistics**: calculate some values to estimate the distribution of the population. Can be regarded as a random variable. Such as sample mean, sample standard deviation. Can be used to estimate unknown population parameters, called **estimates**.
- ▶ **Sampling distribution**: The distribution of an estimator or other statistic. E.g., if X_1, \dots, X_n is a sample from $N(\mu, \sigma^2)$, then $\bar{X} \sim N(\mu, \sigma^2/n)$.

- ▶ A **sample** (typically) is n observations drawn independently from the population, X_1, X_2, \dots, X_n .
- ▶ A sample can be regarded as n numbers, or n iid random variables.
- ▶ **Statistics**: calculate some values to estimate the distribution of the population. Can be regarded as a random variable. Such as sample mean, sample standard deviation. Can be used to estimate unknown population parameters, called **estimates**.
- ▶ **Sampling distribution**: The distribution of an estimator or other statistic. E.g., if X_1, \dots, X_n is a sample from $N(\mu, \sigma^2)$, then $\bar{X} \sim N(\mu, \sigma^2/n)$.

- ▶ A **sample** (typically) is n observations drawn independently from the population, X_1, X_2, \dots, X_n .
- ▶ A sample can be regarded as n numbers, or n iid random variables.
- ▶ **Statistics**: calculate some values to estimate the distribution of the population. Can be regarded as a random variable. Such as sample mean, sample standard deviation. Can be used to estimate unknown population parameters, called **estimates**.
- ▶ **Sampling distribution**: The distribution of an estimator or other statistic. E.g., if X_1, \dots, X_n is a sample from $N(\mu, \sigma^2)$, then $\bar{X} \sim N(\mu, \sigma^2/n)$.

MLE(Maximum Likelihood Estimate)

- MLE is a commonly used parameter estimation method. For random sample X_1, X_2, \dots, X_n from population X , if PDF of PMF of X is $f(x; \beta)$, then

$$\hat{\beta} = \arg \max_{\beta} L(\beta) = \arg \max_{\beta} \prod_{i=1}^n f(X_i; \beta)$$

is called the MLE of the unknow parameter β .

- Under some conditions, when the sample size $n \rightarrow \infty$, $\hat{\beta}$ has a limiting(approximate) distribution $N(\beta, \sigma_{\beta}^2)$.

MLE(Maximum Likelihood Estimate)

- MLE is a commonly used parameter estimation method. For random sample X_1, X_2, \dots, X_n from population X , if PDF of PMF of X is $f(x; \beta)$, then

$$\hat{\beta} = \arg \max_{\beta} L(\beta) = \arg \max_{\beta} \prod_{i=1}^n f(X_i; \beta)$$

is called the MLE of the unknow parameter β .

- Under some conditions, when the sample size $n \rightarrow \infty$, $\hat{\beta}$ has a limiting(approximate) distribution $N(\beta, \sigma_{\beta}^2)$.

Standard Error

- ▶ If $\hat{\beta} = \psi(X_1, \dots, X_n)$ is a estimate of an population parameter β , it has a sampling distribution $F_{\hat{\beta}}(x)$.
- ▶ The standard deviation of the sampling distribution is $\sigma_{\hat{\beta}}$.
- ▶ An estimate of $\sigma_{\hat{\beta}}$ is called the **standard error(SE)** of $\hat{\beta}$.
- ▶ If $\hat{\beta}$ has a limiting normal distribution, then $\hat{\beta}$ is approximately distributed $N(\beta, SE(\hat{\beta}))$.
- ▶ SE can measure the precision of estimation.
- ▶ SE can be used to construct (approximate) confidence intervals.

Standard Error

- ▶ If $\hat{\beta} = \psi(X_1, \dots, X_n)$ is a estimate of an population parameter β , it has a sampling distribution $F_{\hat{\beta}}(x)$.
- ▶ The standard deviation of the sampling distribution is $\sigma_{\hat{\beta}}$.
- ▶ An estimate of $\sigma_{\hat{\beta}}$ is called the **standard error(SE)** of $\hat{\beta}$.
- ▶ If $\hat{\beta}$ has a limiting normal distribution, then $\hat{\beta}$ is approximately distributed $N(\beta, SE(\hat{\beta}))$.
- ▶ SE can measure the precision of estimation.
- ▶ SE can be used to construct (approximate) confidence intervals.

Standard Error

- ▶ If $\hat{\beta} = \psi(X_1, \dots, X_n)$ is a estimate of an population parameter β , it has a sampling distribution $F_{\hat{\beta}}(x)$.
- ▶ The standard deviation of the sampling distribution is $\sigma_{\hat{\beta}}$.
- ▶ An estimate of $\sigma_{\hat{\beta}}$ is called the **standard error(SE)** of $\hat{\beta}$.
- ▶ If $\hat{\beta}$ has a limiting normal distribution, then $\hat{\beta}$ is approximately distributed $N(\beta, SE(\hat{\beta}))$.
- ▶ SE can measure the precision of estimation.
- ▶ SE can be used to construct (approximate) confidence intervals.

Standard Error

- ▶ If $\hat{\beta} = \psi(X_1, \dots, X_n)$ is a estimate of an population parameter β , it has a sampling distribution $F_{\hat{\beta}}(x)$.
- ▶ The standard deviation of the sampling distribution is $\sigma_{\hat{\beta}}$.
- ▶ An estimate of $\sigma_{\hat{\beta}}$ is called the **standard error(SE)** of $\hat{\beta}$.
- ▶ If $\hat{\beta}$ has a limiting normal distribution, then $\hat{\beta}$ is approximately distributed $N(\beta, SE(\hat{\beta}))$.
- ▶ SE can measure the precision of estimation.
- ▶ SE can be used to construct (approximate) confidence intervals.

Standard Error

- ▶ If $\hat{\beta} = \psi(X_1, \dots, X_n)$ is a estimate of an population parameter β , it has a sampling distribution $F_{\hat{\beta}}(x)$.
- ▶ The standard deviation of the sampling distribution is $\sigma_{\hat{\beta}}$.
- ▶ An estimate of $\sigma_{\hat{\beta}}$ is called the **standard error(SE)** of $\hat{\beta}$.
- ▶ If $\hat{\beta}$ has a limiting normal distribution, then $\hat{\beta}$ is approximately distributed $N(\beta, SE(\hat{\beta}))$.
- ▶ SE can measure the precision of estimation.
- ▶ SE can be used to construct (approximate) confidence intervals.

Standard Error

- ▶ If $\hat{\beta} = \psi(X_1, \dots, X_n)$ is a estimate of an population parameter β , it has a sampling distribution $F_{\hat{\beta}}(x)$.
- ▶ The standard deviation of the sampling distribution is $\sigma_{\hat{\beta}}$.
- ▶ An estimate of $\sigma_{\hat{\beta}}$ is called the **standard error(SE)** of $\hat{\beta}$.
- ▶ If $\hat{\beta}$ has a limiting normal distribution, then $\hat{\beta}$ is approximately distributed $N(\beta, SE(\hat{\beta}))$.
- ▶ SE can measure the precision of estimation.
- ▶ SE can be used to construct (approximate) confidence intervals.

Hypothesis Tests

- ▶ To test the **null hypothesis** H_0 against the **alternative hypothesis** H_a on the population;
- ▶ Given sample X_1, X_2, \dots, X_n from population $F(x; \theta)$, construct some statistic ξ , whose distribution does not depend on θ , but its value can indicate the possible choice of H_0 or H_a .
- ▶ Traditionally, we first choose an **significance level** α , then we find a **rejection area** W , where $\sup_{H_0} P(\xi \in W) \leq \alpha$. We reject H_0 when $\xi \in W$.

Hypothesis Tests

- ▶ To test the **null hypothesis** H_0 against the **alternative hypothesis** H_a on the population;
- ▶ Given sample X_1, X_2, \dots, X_n from population $F(x; \theta)$, construct some statistic ξ , whose distribution does not depend on θ , but its value can indicate the possible choice of H_0 or H_a .
- ▶ Traditionally, we first choose an **significance level** α , then we find a **rejection area** W , where $\sup_{H_0} P(\xi \in W) \leq \alpha$. We reject H_0 when $\xi \in W$.

Hypothesis Tests

- ▶ To test the **null hypothesis** H_0 against the **alternative hypothesis** H_a on the population;
- ▶ Given sample X_1, X_2, \dots, X_n from population $F(x; \theta)$, construct some statistic ξ , whose distribution does not depend on θ , but its value can indicate the possible choice of H_0 or H_a .
- ▶ Traditionally, we first choose an **significance level** α , then we find a **rejection area** W , where $\sup_{H_0} P(\xi \in W) \leq \alpha$. We reject H_0 when $\xi \in W$.

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don't choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don't choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don't choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don’t choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don’t choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don’t choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don’t choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Errors and P-Values

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don’t choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

- ▶ Two types of possible errors in a hypothesis test:
 - ▶ **Type I error**, H_0 true but rejected. Error rate is at most the significance level α .
 - ▶ **Type II error**, H_0 false but accepted. Error rate can be as large as $1 - \alpha$.
- ▶ To reduce type II error:
 - ▶ Construct theoretically “good” tests.
 - ▶ Don’t choose α too small.
 - ▶ Choose a big enough sample size n .
- ▶ **p-value**: the minimum α we can use if we want to reject H_0 , after the test statistic value is known.
- ▶ The smaller the p-value is, the more confidence we have when we reject H_0 . Reject H_0 if and only if the p-value is less than α .

Example hypothesis test

- ▶ $X \sim N(\mu, \sigma^2)$. Sample is X_1, \dots, X_n .

- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.

- ▶ Test statistic

$$T = \frac{\bar{X} - \mu_0}{SE(\bar{X})}$$

where $SE(\bar{X}) = S/\sqrt{n}$, S is the sample standard deviation.

- ▶ Rejection area: $W = \{\xi > \lambda\}$, $\lambda = F^{-1}(1 - \alpha, n - 1)$, $F^{-1}(p, n)$ is the quantile function of the $t(n - 1)$ distribution.
- ▶ P-value is $1 - F(T; n - 1)$, where $F(x; n)$ is the CDF of the $t(n)$ distribution.
- ▶ The larger T is, the smaller the p-value is, the more confidence we have when we reject H_0 .

Example hypothesis test

- ▶ $X \sim N(\mu, \sigma^2)$. Sample is X_1, \dots, X_n .

- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.

- ▶ Test statistic

$$T = \frac{\bar{X} - \mu_0}{SE(\bar{X})}$$

where $SE(\bar{X}) = S/\sqrt{n}$, S is the sample standard deviation.

- ▶ Rejection area: $W = \{\xi > \lambda\}$, $\lambda = F^{-1}(1 - \alpha, n - 1)$, $F^{-1}(p, n)$ is the quantile function of the $t(n - 1)$ distribution.
- ▶ P-value is $1 - F(T; n - 1)$, where $F(x; n)$ is the CDF of the $t(n)$ distribution.
- ▶ The larger T is, the smaller the p-value is, the more confidence we have when we reject H_0 .

Example hypothesis test

- ▶ $X \sim N(\mu, \sigma^2)$. Sample is X_1, \dots, X_n .
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.
- ▶ Test statistic

$$T = \frac{\bar{X} - \mu_0}{SE(\bar{X})}$$

where $SE(\bar{X}) = S/\sqrt{n}$, S is the sample standard deviation.

- ▶ Rejection area: $W = \{\xi > \lambda\}$, $\lambda = F^{-1}(1 - \alpha, n - 1)$, $F^{-1}(p, n)$ is the quantile function of the $t(n - 1)$ distribution.
- ▶ P-value is $1 - F(T; n - 1)$, where $F(x; n)$ is the CDF of the $t(n)$ distribution.
- ▶ The larger T is, the smaller the p-value is, the more confidence we have when we reject H_0 .

Example hypothesis test

- ▶ $X \sim N(\mu, \sigma^2)$. Sample is X_1, \dots, X_n .

- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.

- ▶ Test statistic

$$T = \frac{\bar{X} - \mu_0}{SE(\bar{X})}$$

where $SE(\bar{X}) = S/\sqrt{n}$, S is the sample standard deviation.

- ▶ Rejection area: $W = \{\xi > \lambda\}$, $\lambda = F^{-1}(1 - \alpha, n - 1)$, $F^{-1}(p, n)$ is the quantile function of the $t(n - 1)$ distribution.
- ▶ P-value is $1 - F(T; n - 1)$, where $F(x; n)$ is the CDF of the $t(n)$ distribution.
- ▶ The larger T is, the smaller the p-value is, the more confidence we have when we reject H_0 .

Example hypothesis test

- ▶ $X \sim N(\mu, \sigma^2)$. Sample is X_1, \dots, X_n .

- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.

- ▶ Test statistic

$$T = \frac{\bar{X} - \mu_0}{SE(\bar{X})}$$

where $SE(\bar{X}) = S/\sqrt{n}$, S is the sample standard deviation.

- ▶ Rejection area: $W = \{\xi > \lambda\}$, $\lambda = F^{-1}(1 - \alpha, n - 1)$, $F^{-1}(p, n)$ is the quantile function of the $t(n - 1)$ distribution.
- ▶ P-value is $1 - F(T; n - 1)$, where $F(x; n)$ is the CDF of the $t(n)$ distribution.
- ▶ The larger T is, the smaller the p-value is, the more confidence we have when we reject H_0 .

Example hypothesis test

- ▶ $X \sim N(\mu, \sigma^2)$. Sample is X_1, \dots, X_n .
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.
- ▶ Test statistic

$$T = \frac{\bar{X} - \mu_0}{SE(\bar{X})}$$

where $SE(\bar{X}) = S/\sqrt{n}$, S is the sample standard deviation.

- ▶ Rejection area: $W = \{\xi > \lambda\}$, $\lambda = F^{-1}(1 - \alpha, n - 1)$, $F^{-1}(p, n)$ is the quantile function of the $t(n - 1)$ distribution.
- ▶ P-value is $1 - F(T; n - 1)$, where $F(x; n)$ is the CDF of the $t(n)$ distribution.
- ▶ The larger T is, the smaller the p-value is, the more confidence we have when we reject H_0 .

Section Contents

- ▶ **Measurement level of variables.**
- ▶ Descriptive statistics for nominal variables.
- ▶ Descriptive statistics for interval variables.
- ▶ Histograms, boxplots, QQ plots, probability plots, stem-leaf plots.
- ▶ Using PROC FREQ, PROC MEANS, PROC UNIVARIATE.

Section Contents

- ▶ Measurement level of variables.
- ▶ Descriptive statistics for nominal variables.
- ▶ Descriptive statistics for interval variables.
- ▶ Histograms, boxplots, QQ plots, probability plots, stem-leaf plots.
- ▶ Using PROC FREQ, PROC MEANS, PROC UNIVARIATE.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics
Summary statistics

Comparing the
Mean

Analysis of
Variance

Section Contents

- ▶ Measurement level of variables.
- ▶ Descriptive statistics for nominal variables.
- ▶ Descriptive statistics for interval variables.
- ▶ Histograms, boxplots, QQ plots, probability plots, stem-leaf plots.
- ▶ Using PROC FREQ, PROC MEANS, PROC UNIVARIATE.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics
Summary statistics

Comparing the
Mean

Analysis of
Variance

Section Contents

- ▶ Measurement level of variables.
- ▶ Descriptive statistics for nominal variables.
- ▶ Descriptive statistics for interval variables.
- ▶ Histograms, boxplots, QQ plots, probability plots, stem-leaf plots.
- ▶ Using PROC FREQ, PROC MEANS, PROC UNIVARIATE.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics
Summary statistics

Comparing the
Mean

Analysis of
Variance

Section Contents

- ▶ Measurement level of variables.
- ▶ Descriptive statistics for nominal variables.
- ▶ Descriptive statistics for interval variables.
- ▶ Histograms, boxplots, QQ plots, probability plots, stem-leaf plots.
- ▶ Using PROC FREQ, PROC MEANS, PROC UNIVARIATE.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics
Summary statistics

Comparing the
Mean

Analysis of
Variance

Measurement levels

- ▶ **Nominal, such as sex, job type, race.**
- ▶ Ordinal, such as age group(child, youth, medium, old), dose level(none, low, high).
- ▶ Interval, such as weight, blood pressure, heart rate.

Measurement levels

- ▶ Nominal, such as sex, job type, race.
- ▶ Ordinal, such as age group(child, youth, medium, old), dose level(none, low, high).
- ▶ Interval, such as weight, blood pressure, heart rate.

Measurement levels

- ▶ Nominal, such as sex, job type, race.
- ▶ Ordinal, such as age group(child, youth, medium, old), dose level(none, low, high).
- ▶ Interval, such as weight, blood pressure, heart rate.

Statistics for nominal variables

- ▶ The value set.
- ▶ The count of each value(called frequency), and the percentage.
- ▶ Use a bar chart to show the distribution.

Statistics for nominal variables

- ▶ The value set.
- ▶ The count of each value(called frequency), and the percentage.
- ▶ Use a bar chart to show the distribution.

Statistics for nominal variables

- ▶ The value set.
- ▶ The count of each value(called frequency), and the percentage.
- ▶ Use a bar chart to show the distribution.

Example: frequency table

Use PROC FREQ to list the value set and the frequency counts, percents.

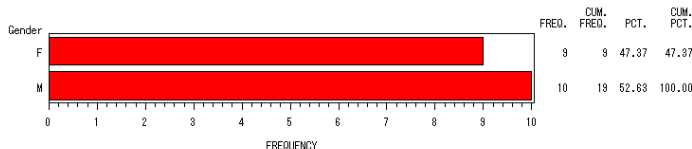
```
proc freq data=sasuser.class;  
  tables sex;  
run;
```

Gender				
sex	Frequency	Percent	Cumulative Frequency	Cumulative Percent
F	9	47.37	9	47.37
M	10	52.63	19	100.00

Example: bar chart

Use PROC GCHART to make a bar chart. `hbar` could be replaced with `vbar`, `hbar3d`, `vbar3d`.

```
proc gchart data=sasuser.class;  
  hbar sex;  
run;
```



Descriptive statistics for interval variables

► Location statistics: sample mean, sample median.

► Variability statistics: sample standard deviation, sample interquantile range, range, coefficient of variation.

► Sample skewness $\frac{n}{(n-1)(n-2)} \sum \left(\frac{y_i - \bar{y}}{s} \right)^3$.

► Sample kurtosis $\frac{n(n+1)}{(n-1)(n-2)(n-3)} \frac{\sum (y_i - \bar{y})^4}{s^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$

► Moments, quantiles.

Descriptive statistics for interval variables

- ▶ Location statistics: sample mean, sample median.
- ▶ Variability statistics: sample standard deviation, sample interquantile range, range, coefficient of variation.

- ▶ Sample skewness $\frac{n}{(n-1)(n-2)} \sum \left(\frac{y_i - \bar{y}}{s} \right)^3$.

- ▶ Sample kurtosis $\frac{n(n+1)}{(n-1)(n-2)(n-3)} \frac{\sum (y_i - \bar{y})^4}{s^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$

- ▶ Moments, quantiles.

Descriptive statistics for interval variables

- ▶ Location statistics: sample mean, sample median.
- ▶ Variability statistics: sample standard deviation, sample interquantile range, range, coefficient of variation.

- ▶ Sample skewness $\frac{n}{(n-1)(n-2)} \sum \left(\frac{y_i - \bar{y}}{s} \right)^3$.

- ▶ Sample kurtosis $\frac{n(n+1)}{(n-1)(n-2)(n-3)} \frac{\sum (y_i - \bar{y})^4}{s^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$

- ▶ Moments, quantiles.

Descriptive statistics for interval variables

- ▶ Location statistics: sample mean, sample median.
- ▶ Variability statistics: sample standard deviation, sample interquantile range, range, coefficient of variation.

- ▶ Sample skewness $\frac{n}{(n-1)(n-2)} \sum \left(\frac{y_i - \bar{y}}{s} \right)^3$.

- ▶ Sample kurtosis $\frac{n(n+1)}{(n-1)(n-2)(n-3)} \frac{\sum (y_i - \bar{y})^4}{s^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$

- ▶ Moments, quantiles.

Descriptive statistics for interval variables

- ▶ Location statistics: sample mean, sample median.
- ▶ Variability statistics: sample standard deviation, sample interquantile range, range, coefficient of variation.

- ▶ Sample skewness $\frac{n}{(n-1)(n-2)} \sum \left(\frac{y_i - \bar{y}}{s} \right)^3$.

- ▶ Sample kurtosis $\frac{n(n+1)}{(n-1)(n-2)(n-3)} \frac{\sum (y_i - \bar{y})^4}{s^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$

- ▶ Moments, quantiles.

PROC UNIVARIATE

- ▶ Use PROC UNIVARIATE to compute sample statistics for an interval variable.
- ▶ Example:

```
proc univariate data=sasuser.class;  
  var height;  
run;
```

- ▶ Use PROC UNIVARIATE to compute sample statistics for an interval variable.
- ▶ Example:

```
proc univariate data=sasuser.class;  
  var height;  
run;
```

Dongfeng Li

Moments			
N	19	Sum Weights	19
Mean	62.3368421	Sum Observations	1184.4
Std Deviation	5.12707525	Variance	26.2869006
Skewness	-0.2596696	Kurtosis	-0.1389692
Uncorrected SS	74304.92	Corrected SS	473.164211
Coeff Variation	8.22479143	Std Error Mean	1.17623173

Basic Statistical Measures			
Location		Variability	
Mean	62.33684	Std Deviation	5.12708
Median	62.80000	Variance	26.28690
Mode	62.50000	Range	20.70000
		Interquartile Range	9.00000

Note	The mode displayed is the smallest of 2 modes with a count of 2.
-------------	------------------------------------------------------------------

Descriptive statistics

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

Quantiles (Definition 5)	
Quantile	Estimate
100% Max	72.0
99%	72.0
95%	72.0
90%	69.0
75% Q3	66.5
50% Median	62.8
25% Q1	57.5
10%	56.3
5%	51.3
1%	51.3
0% Min	51.3

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
51.3	7	66.5	6
56.3	4	66.5	19
56.5	1	67.0	12
57.3	13	69.0	10
57.5	18	72.0	16

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics

Summary statistics

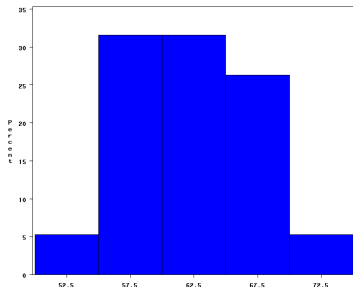
Comparing the
Mean

Analysis of
Variance

Histograms

- ▶ The **histogram** is a **nonparametric** estimate of the population density function.
- ▶ It divide the value set into intervals, and count the percent of observations in each interval.
- ▶ Example:

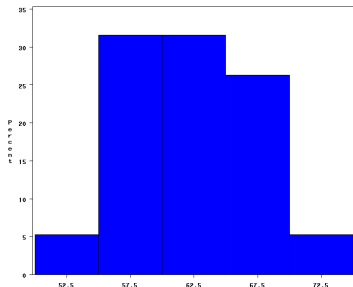
```
proc univariate data=sasuser.class  
    noprint;  
    var height;    histogram;  
run;
```



Histograms

- ▶ The **histogram** is a **nonparametric** estimate of the population density function.
- ▶ It divide the value set into intervals, and count the percent of observations in each interval.
- ▶ Example:

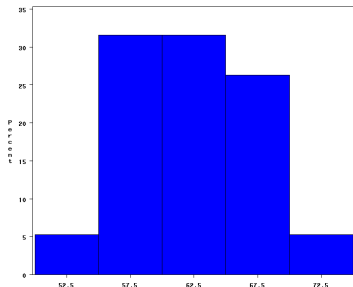
```
proc univariate data=sasuser.class  
    noprint;  
    var height;    histogram;  
run;
```



Histograms

- ▶ The **histogram** is a **nonparametric** estimate of the population density function.
- ▶ It divide the value set into intervals, and count the percent of observations in each interval.
- ▶ Example:

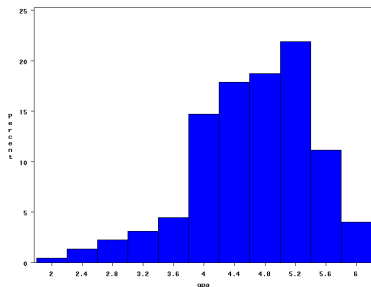
```
proc univariate data=sasuser.class  
    noprint;  
    var height;    histogram;  
run;
```



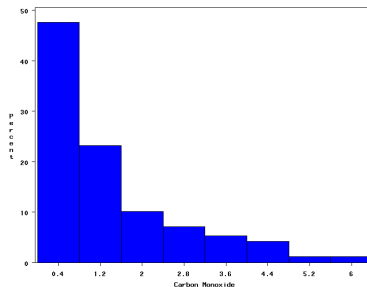
Examples of histograms

Skewness shown in histograms.

Left skewed histogram



Right skewed histogram



Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics

Summary statistics

Comparing the
Mean

Analysis of
Variance

Box plot

- ▶ Quantiles could be used to describe key distribution characteristics.
- ▶ The **box plot** shows the median, 1/4 and 3/4 quantile, and the minimum and maximum of a variable.
- ▶ The box holds the middle 50% range. Length is the IQR(inter quantile range).
- ▶ The whiskers extend to the minimum and the maximum; but the length of the whisker is not allowed to exceed 1.5 times the IQR.
- ▶ Points beyond whiskers are plotted as individual points. They are candidate for outliers.

Box plot

- ▶ Quantiles could be used to describe key distribution characteristics.
- ▶ The **box plot** shows the median, 1/4 and 3/4 quantile, and the minimum and maximum of a variable.
- ▶ The box holds the middle 50% range. Length is the IQR(inter quantile range).
- ▶ The whiskers extend to the minimum and the maximum; but the length of the whisker is not allowed to exceed 1.5 times the IQR.
- ▶ Points beyond whiskers are plotted as individual points. They are candidate for outliers.

Box plot

- ▶ Quantiles could be used to describe key distribution characteristics.
- ▶ The **box plot** shows the median, 1/4 and 3/4 quantile, and the minimum and maximum of a variable.
- ▶ The box holds the middle 50% range. Length is the IQR(inter quantile range).
- ▶ The whiskers extend to the minimum and the maximum; but the length of the whisker is not allowed to exceed 1.5 times the IQR.
- ▶ Points beyond whiskers are plotted as individual points. They are candidate for outliers.

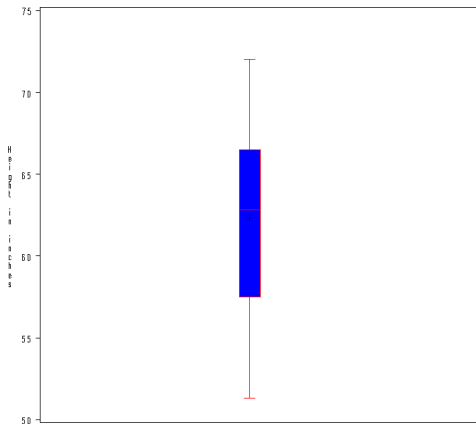
Box plot

- ▶ Quantiles could be used to describe key distribution characteristics.
- ▶ The **box plot** shows the median, 1/4 and 3/4 quantile, and the minimum and maximum of a variable.
- ▶ The box holds the middle 50% range. Length is the IQR(inter quantile range).
- ▶ The whiskers extend to the minimum and the maximum; but the length of the whisker is not allowed to exceed 1.5 times the IQR.
- ▶ Points beyond whiskers are plotted as individual points. They are candidate for outliers.

Box plot

- ▶ Quantiles could be used to describe key distribution characteristics.
- ▶ The **box plot** shows the median, 1/4 and 3/4 quantile, and the minimum and maximum of a variable.
- ▶ The box holds the middle 50% range. Length is the IQR(inter quantile range).
- ▶ The whiskers extend to the minimum and the maximum; but the length of the whisker is not allowed to exceed 1.5 times the IQR.
- ▶ Points beyond whiskers are plotted as individual points. They are candidate for outliers.

Example: boxplot of height



Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics

Summary statistics

Comparing the
Mean

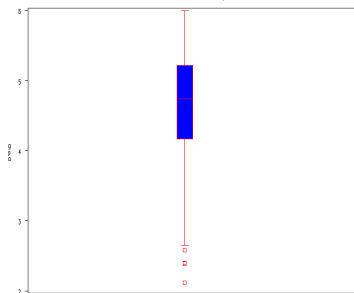
Analysis of
Variance


```
proc sql;  
  create view work._tmp_0 as  
    select height, 1 as _dummy_  
      from sasuser.class;  
title ;  
axis1 major=none value=none label=none;  
proc boxplot data=_tmp_0;  
  plot height*_dummy_ /  
    boxstyle=skematic  
    cboxfill=blue haxis=axis1;  
run;
```

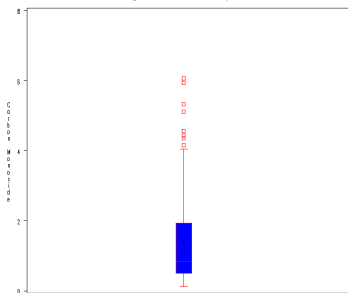
Example: boxplot of GPA and CO

Skewness and possible outliers shown in boxplot.

Left Skewed Boxplot



Right Skewed Boxplot



Some Statistics
Background

Descriptive
statistics: concepts
and programs

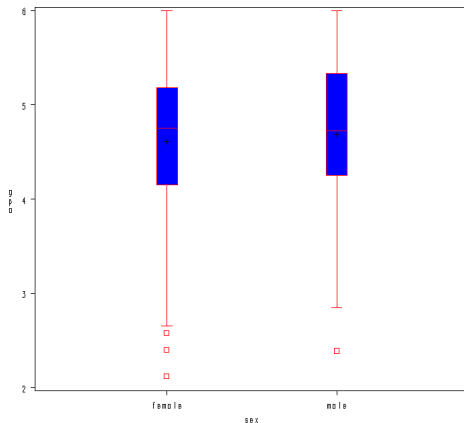
Descriptive statistics

Summary statistics

Comparing the
Mean

Analysis of
Variance

Grouped boxplot



Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics

Summary statistics

Comparing the
Mean

Analysis of
Variance

```
proc sort data=sasuser.gpa out=_tmp_1;  
  by sex;  
proc boxplot data=_tmp_1;  
  plot gpa*sex /  
    boxstyle=skematic  
    cboxfill=blue;  
run;
```

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics
Summary statistics

Comparing the
Mean

Analysis of
Variance

- ▶ The normal distribution is the most commonly used distribution.
- ▶ Graphs are designed to show discrepancy from the normal distribution. **QQ plot** is one of them.
- ▶ Suppose Y_1, Y_2, \dots, Y_n is from $N(\mu, \sigma^2)$ and sorted ascendingly.
- ▶ CDF $F(x) = \Phi(\frac{x-\mu}{\sigma})$, $F(Y_i) \approx \frac{i}{n}$, $\Phi(\frac{Y_i-\mu}{\sigma}) \approx \frac{i}{n}$, $Y_i \approx \mu + \sigma\Phi^{-1}(\frac{i}{n})$.
- ▶ Let $x_i = \Phi^{-1}(\frac{i}{n})$, plot $(x_i, Y_i), i = 1, \dots, n$. The points should lie around the line $y = \mu + \sigma x$.
- ▶ Continuity adjustment: $x_i = \Phi^{-1}(\frac{i-0.375}{n+0.25})$.

- ▶ The normal distribution is the most commonly used distribution.
- ▶ Graphs are designed to show discrepancy from the normal distribution. **QQ plot** is one of them.
- ▶ Suppose Y_1, Y_2, \dots, Y_n is from $N(\mu, \sigma^2)$ and sorted ascendingly.
- ▶ CDF $F(x) = \Phi(\frac{x-\mu}{\sigma})$, $F(Y_i) \approx \frac{i}{n}$, $\Phi(\frac{Y_i-\mu}{\sigma}) \approx \frac{i}{n}$,
 $Y_i \approx \mu + \sigma\Phi^{-1}(\frac{i}{n})$.
- ▶ Let $x_i = \Phi^{-1}(\frac{i}{n})$, plot $(x_i, Y_i), i = 1, \dots, n$. The points should lie around the line $y = \mu + \sigma x$.
- ▶ Continuity adjustment: $x_i = \Phi^{-1}(\frac{i-0.375}{n+0.25})$.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Descriptive statistics
Summary statistics

Comparing the
Mean

Analysis of
Variance

- ▶ The normal distribution is the most commonly used distribution.
- ▶ Graphs are designed to show discrepancy from the normal distribution. **QQ plot** is one of them.
- ▶ Suppose Y_1, Y_2, \dots, Y_n is from $N(\mu, \sigma^2)$ and sorted ascendingly.
- ▶ CDF $F(x) = \Phi(\frac{x-\mu}{\sigma})$, $F(Y_i) \approx \frac{i}{n}$, $\Phi(\frac{Y_i-\mu}{\sigma}) \approx \frac{i}{n}$,
 $Y_i \approx \mu + \sigma\Phi^{-1}(\frac{i}{n})$.
- ▶ Let $x_i = \Phi^{-1}(\frac{i}{n})$, plot $(x_i, Y_i), i = 1, \dots, n$. The points should lie around the line $y = \mu + \sigma x$.
- ▶ Continuity adjustment: $x_i = \Phi^{-1}(\frac{i-0.375}{n+0.25})$.

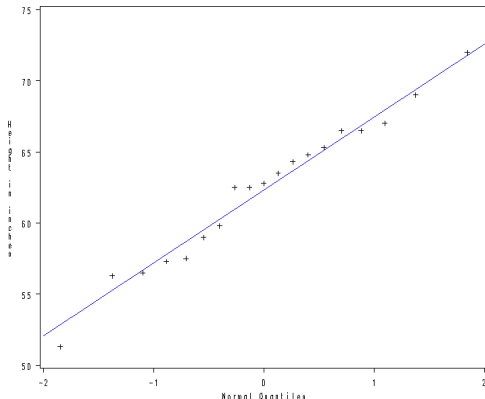
- ▶ The normal distribution is the most commonly used distribution.
- ▶ Graphs are designed to show discrepancy from the normal distribution. **QQ plot** is one of them.
- ▶ Suppose Y_1, Y_2, \dots, Y_n is from $N(\mu, \sigma^2)$ and sorted ascendingly.
- ▶ CDF $F(x) = \Phi(\frac{x-\mu}{\sigma})$, $F(Y_i) \approx \frac{i}{n}$, $\Phi(\frac{Y_i-\mu}{\sigma}) \approx \frac{i}{n}$,
 $Y_i \approx \mu + \sigma\Phi^{-1}(\frac{i}{n})$.
- ▶ Let $x_i = \Phi^{-1}(\frac{i}{n})$, plot $(x_i, Y_i), i = 1, \dots, n$. The points should lie around the line $y = \mu + \sigma x$.
- ▶ Continuity adjustment: $x_i = \Phi^{-1}(\frac{i-0.375}{n+0.25})$.

- ▶ The normal distribution is the most commonly used distribution.
- ▶ Graphs are designed to show discrepancy from the normal distribution. **QQ plot** is one of them.
- ▶ Suppose Y_1, Y_2, \dots, Y_n is from $N(\mu, \sigma^2)$ and sorted ascendingly.
- ▶ CDF $F(x) = \Phi(\frac{x-\mu}{\sigma})$, $F(Y_i) \approx \frac{i}{n}$, $\Phi(\frac{Y_i-\mu}{\sigma}) \approx \frac{i}{n}$,
 $Y_i \approx \mu + \sigma\Phi^{-1}(\frac{i}{n})$.
- ▶ Let $x_i = \Phi^{-1}(\frac{i}{n})$, plot $(x_i, Y_i), i = 1, \dots, n$. The points should lie around the line $y = \mu + \sigma x$.
- ▶ Continuity adjustment: $x_i = \Phi^{-1}(\frac{i-0.375}{n+0.25})$.

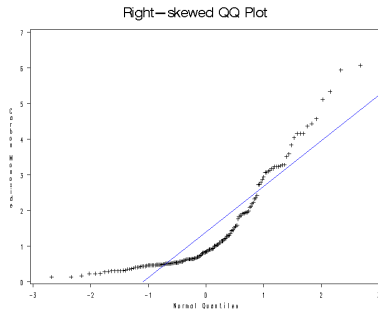
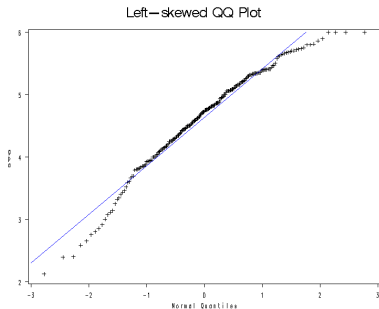
- ▶ The normal distribution is the most commonly used distribution.
- ▶ Graphs are designed to show discrepancy from the normal distribution. **QQ plot** is one of them.
- ▶ Suppose Y_1, Y_2, \dots, Y_n is from $N(\mu, \sigma^2)$ and sorted ascendingly.
- ▶ CDF $F(x) = \Phi(\frac{x-\mu}{\sigma})$, $F(Y_i) \approx \frac{i}{n}$, $\Phi(\frac{Y_i-\mu}{\sigma}) \approx \frac{i}{n}$,
 $Y_i \approx \mu + \sigma\Phi^{-1}(\frac{i}{n})$.
- ▶ Let $x_i = \Phi^{-1}(\frac{i}{n})$, plot $(x_i, Y_i), i = 1, \dots, n$. The points should lie around the line $y = \mu + \sigma x$.
- ▶ Continuity adjustment: $x_i = \Phi^{-1}(\frac{i-0.375}{n+0.25})$.

Example: QQ Plot of HEIGHT

```
proc univariate  
    data=sasuser.class noprint;  
    qqplot height /  
        normal(mu=est sigma=est);  
run;
```



Example: Skewness in QQ Plot



Probability Plot

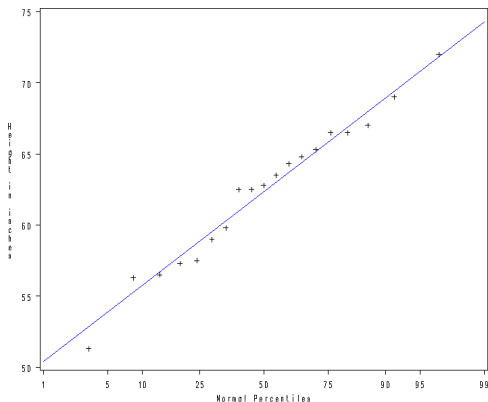
- ▶ Probability plot is the same plot as a QQ plot, except that the label of the x axis $\Phi(x_i)$ instead of x_i .
- ▶ Probability plots are preferable for graphical estimation of percentiles, whereas Q-Q plots are preferable for graphical estimation of distribution parameters.

Probability Plot

- ▶ Probability plot is the same plot as a QQ plot, except that the label of the x axis $\Phi(x_i)$ instead of x_i .
- ▶ Probability plots are preferable for graphical estimation of percentiles, whereas Q-Q plots are preferable for graphical estimation of distribution parameters.

Example: Probability Plot of HEIGHT

```
proc univariate  
    data=sasuser.class noprint;  
    probplot height /  
        normal(mu=est sigma=est);  
run;
```



The Stem-leaf Plot

- ▶ The **stem-leaf plot** is a text-based plot, which display information like the histogram, but with detail on each data value. Each “leaf” corresponds to one data value.
- ▶ PROC UNIVARIATE has an option PLOT, which generates text-based stem-leaf plot, boxplot and QQ plot.

The Stem-leaf Plot

- ▶ The **stem-leaf plot** is a text-based plot, which display information like the histogram, but with detail on each data value. Each “leaf” corresponds to one data value.
- ▶ PROC UNIVARIATE has an option PLOT, which generates text-based stem-leaf plot, boxplot and QQ plot.

```
proc univariate data=sasuser.class
    plot;
    var height;
run;
```

Stem	Leaf	#
7	2	1
6	556679	6
6	022344	6
5	66789	5
5	1	1

-----+-----+-----+-----+

Multiply Stem.Leaf by $10^{**}+1$

PROC MEANS and PROC SUMMARY

- ▶ **PROC MEANS and PROC SUMMARY are used to produce summary statistics.**
- ▶ They can display overall summary statistics and classified summary statistics. PROC MEANS displays result by default; PROC SUMMARY need the PRINT option to display result.
- ▶ They can output data sets with the summary statistics. PROC SUMMARY is designed to do this, although PROC MEANS can do the same.

PROC MEANS and PROC SUMMARY

- ▶ PROC MEANS and PROC SUMMARY are used to produce summary statistics.
- ▶ They can display overall summary statistics and classified summary statistics. PROC MEANS displays result by default; PROC SUMMARY need the PRINT option to display result.
- ▶ They can output data sets with the summary statistics. PROC SUMMARY is designed to do this, although PROC MEANS can do the same.

PROC MEANS and PROC SUMMARY

- ▶ PROC MEANS and PROC SUMMARY are used to produce summary statistics.
- ▶ They can display overall summary statistics and classified summary statistics. PROC MEANS displays result by default; PROC SUMMARY need the PRINT option to display result.
- ▶ They can output data sets with the summary statistics. PROC SUMMARY is designed to do this, although PROC MEANS can do the same.

Example of overall statistics

- ▶ Use the VAR statement to specify which variables to summarize.
- ▶ Example of overall statistics:

```
proc means data=sasuser.class;  
  var height weight;  
run;  
proc summary data=sasuser.class print;  
  var height weight;  
run;
```

Example of overall statistics

- ▶ Use the VAR statement to specify which variables to summarize.
- ▶ Example of overall statistics:

```
proc means data=sasuser.class;  
    var height weight;  
run;  
proc summary data=sasuser.class print;  
    var height weight;  
run;
```

Classified summary

- ▶ Use the CLASS statement to specify one or more class variables.
- ▶ Example:

```
proc means data=sasuser.class ;  
  var height weight;  
  class sex;  
run;
```


Classified summary

- ▶ Use the CLASS statement to specify one or more class variables.
- ▶ Example:

```
proc means data=sasuser.class ;  
  var height weight;  
  class sex;  
run;
```

Other statistical measures

- ▶ PROC MEANS calculate summaries N, Mean, Standard Deviation, Minimum, Maximum by default.
- ▶ Use PROC MEANS options to specify other statistical measures.
- ▶ Example:

```
proc means data=sasuser.class MEAN VAR CV;  
  var height weight;  
  class sex;  
run;
```

Other statistical measures

- ▶ PROC MEANS calculate summaries N, Mean, Standard Deviation, Minimum, Maximum by default.
- ▶ Use PROC MEANS options to specify other statistical measures.
- ▶ Example:

```
proc means data=sasuser.class MEAN VAR CV;  
  var height weight;  
  class sex;  
run;
```

Other statistical measures

- ▶ PROC MEANS calculate summaries N, Mean, Standard Deviation, Minimum, Maximum by default.
- ▶ Use PROC MEANS options to specify other statistical measures.
- ▶ Example:

```
proc means data=sasuser.class MEAN VAR CV;  
  var height weight;  
  class sex;  
run;
```

Producing output datasets

- ▶ Use the OUTPUT statement to save the summary statistics to an output dataset.

- ▶ Syntax:

```
OUTPUT OUT=output-dataset keyword=  
keyword= ... / AUTONAME;
```

- ▶ Where *keyword* is a name of some statistical measure, like MEAN, STD, N, MIN, MAX, etc.

- ▶ Example:

```
proc means data=sasuser.class MEAN VAR CV;  
  var height weight;  
  class sex;  
  output out=res N= CV= / AUTONAME;  
run;  
proc print data=res;run;
```


Producing output datasets

- ▶ Use the OUTPUT statement to save the summary statistics to an output dataset.

- ▶ Syntax:

OUTPUT OUT=*output-dataset* keyword=
keyword= ... / AUTONAME;

- ▶ Where *keyword* is a name of some statistical measure, like MEAN, STD, N, MIN, MAX, etc.

- ▶ Example:

```
proc means data=sasuser.class MEAN VAR CV;  
  var height weight;  
  class sex;  
  output out=res N= CV= / AUTONAME;  
run;  
proc print data=res;run;
```


Comparing the Mean

- ▶ One sample Z tests, t tests.
- ▶ Two sample t tests.
- ▶ The Wilcoxon rank sum test.
- ▶ Paired t tests.
- ▶ Comparing proportions: one sample and two sample.

Comparing the Mean

- ▶ One sample Z tests, t tests.
- ▶ Two sample t tests.
- ▶ The Wilcoxon rank sum test.
- ▶ Paired t tests.
- ▶ Comparing proportions: one sample and two sample.

Comparing the Mean

- ▶ One sample Z tests, t tests.
- ▶ Two sample t tests.
- ▶ The Wilcoxon rank sum test.
- ▶ Paired t tests.
- ▶ Comparing proportions: one sample and two sample.

Comparing the Mean

- ▶ One sample Z tests, t tests.
- ▶ Two sample t tests.
- ▶ The Wilcoxon rank sum test.
- ▶ Paired t tests.
- ▶ Comparing proportions: one sample and two sample.

Comparing the Mean

- ▶ One sample Z tests, t tests.
- ▶ Two sample t tests.
- ▶ The Wilcoxon rank sum test.
- ▶ Paired t tests.
- ▶ Comparing proportions: one sample and two sample.

One Sample Z Test—Two-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$.
- ▶ $W = \{|Z| > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha/2)$.
- ▶ p-value = $2(1 - \Phi(|Z|))$.
- ▶ **Wald test:** Based on the central limit theorem, to test $H_0 : \theta = \theta_0 \longleftrightarrow H_a : \theta \neq \theta_0$, use $W = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$ as the test statistic, if $W \rightarrow N(0, 1)$.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

One Sample Z Test—Two-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$.
- ▶ $W = \{|Z| > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha/2)$.
- ▶ p-value = $2(1 - \Phi(|Z|))$.
- ▶ **Wald test:** Based on the central limit theorem, to test $H_0 : \theta = \theta_0 \longleftrightarrow H_a : \theta \neq \theta_0$, use $W = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$ as the test statistic, if $W \rightarrow N(0, 1)$.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

One Sample Z Test—Two-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$.
- ▶ $W = \{ |Z| > \lambda \}$, $\lambda = \Phi^{-1}(1 - \alpha/2)$.
- ▶ p-value = $2(1 - \Phi(|Z|))$.
- ▶ **Wald test:** Based on the central limit theorem, to test $H_0 : \theta = \theta_0 \longleftrightarrow H_a : \theta \neq \theta_0$, use $W = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$ as the test statistic, if $W \rightarrow N(0, 1)$.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

One Sample Z Test—Two-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$.
- ▶ $W = \{|Z| > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha/2)$.
- ▶ p-value = $2(1 - \Phi(|Z|))$.
- ▶ **Wald test:** Based on the central limit theorem, to test $H_0 : \theta = \theta_0 \longleftrightarrow H_a : \theta \neq \theta_0$, use $W = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$ as the test statistic, if $W \rightarrow N(0, 1)$.

One Sample Z Test—Two-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$.
- ▶ $W = \{|Z| > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha/2)$.
- ▶ p-value = $2(1 - \Phi(|Z|))$.
- ▶ **Wald test:** Based on the central limit theorem, to test $H_0 : \theta = \theta_0 \longleftrightarrow H_a : \theta \neq \theta_0$, use $W = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$ as the test statistic, if $W \rightarrow N(0, 1)$.

One Sample Z Test—Two-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$.
- ▶ $W = \{|Z| > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha/2)$.
- ▶ p-value = $2(1 - \Phi(|Z|))$.
- ▶ **Wald test:** Based on the central limit theorem, to test $H_0 : \theta = \theta_0 \longleftrightarrow H_a : \theta \neq \theta_0$, use $W = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})}$ as the test statistic, if $W \rightarrow N(0, 1)$.

One Sample Z Test—One-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{\mu = \mu_0}{\sim} N(0, 1)$.
- ▶ $W = \{Z > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - \Phi(Z)$.

One Sample Z Test—One-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{\mu = \mu_0}{\sim} N(0, 1)$.
- ▶ $W = \{Z > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - \Phi(Z)$.

One Sample Z Test—One-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{\mu = \mu_0}{\sim} N(0, 1)$.
- ▶ $W = \{Z > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - \Phi(Z)$.

One Sample Z Test—One-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{\mu = \mu_0}{\sim} N(0, 1)$.
- ▶ $W = \{Z > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - \Phi(Z)$.

One Sample Z Test—One-sided

- ▶ $X \sim N(\mu, \sigma_0^2)$, σ_0 known.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$
- ▶ $Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}} \stackrel{\mu = \mu_0}{\sim} N(0, 1)$.
- ▶ $W = \{Z > \lambda\}$, $\lambda = \Phi^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - \Phi(Z)$.

Example: Z test

- For HEIGHT in SASUSER.CLASS, suppose $\sigma_0 = 5$, $\mu_0 = 65$. One sample Z two-sided Z test:

```
proc means data=sasuser.class  
    mean std n;  
    var height;  
    output out=_tmp_1  
        mean=mu std=sigma n=n;  
run;
```


Example: Wald test

- Use the Wald test for HEIGHT mean.

```
data _null_;  
  set _tmp_1;  
  file print;  
  mu0 = 65;  
  sigma0 = sigma;  
  z = (mu - mu0)/(sigma0 / sqrt(n));  
  pvalue = 2*(1 -  
    cdf('normal', abs(z)));  
  put 'Z: ' Z 12.4  
      '      Pr>|Z|: ' pvalue PVALUE.;  
run;  
Z:          -2.2641      Pr>|Z|: 0.0236
```

One-sample t Test for the Mean—Two-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{H_0}{\sim} t(n-1)$.
- ▶ $W = \{|T| > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha/2)$.
- ▶ $\text{p-value} = 2(1 - F_{t(n-1)}(|T|))$.

One-sample t Test for the Mean—Two-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{H_0}{\sim} t(n-1)$.
- ▶ $W = \{|T| > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha/2)$.
- ▶ $\text{p-value} = 2(1 - F_{t(n-1)}(|T|))$.

One-sample t Test for the Mean—Two-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{H_0}{\sim} t(n-1)$.
- ▶ $W = \{|T| > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha/2)$.
- ▶ $\text{p-value} = 2(1 - F_{t(n-1)}(|T|))$.

One-sample t Test for the Mean—Two-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{H_0}{\sim} t(n-1)$.
- ▶ $W = \{|T| > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha/2)$.
- ▶ $\text{p-value} = 2(1 - F_{t(n-1)}(|T|))$.

One-sample t Test for the Mean—Two-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{H_0}{\sim} t(n-1)$.
- ▶ $W = \{|T| > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha/2)$.
- ▶ $\text{p-value} = 2(1 - F_{t(n-1)}(|T|))$.

One-sample t Test for the Mean—One-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{\mu = \mu_0}{\sim} t(n-1)$.
- ▶ $W = \{T > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - F_{t(n-1)}(T)$.

One-sample t Test for the Mean—One-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{\mu = \mu_0}{\sim} t(n-1)$.
- ▶ $W = \{T > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - F_{t(n-1)}(T)$.

One-sample t Test for the Mean—One-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{\mu = \mu_0}{\sim} t(n-1)$.
- ▶ $W = \{T > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - F_{t(n-1)}(T)$.

One-sample t Test for the Mean—One-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{\mu = \mu_0}{\sim} t(n-1)$.
- ▶ $W = \{T > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - F_{t(n-1)}(T)$.

One-sample t Test for the Mean—One-sided

- ▶ $X \sim N(\mu, \sigma^2)$, μ and σ unknown.
- ▶ $H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$.
- ▶ $T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \stackrel{\mu = \mu_0}{\sim} t(n-1)$.
- ▶ $W = \{T > \lambda\}$, $\lambda = F_{t(n-1)}^{-1}(1 - \alpha)$.
- ▶ p-value = $1 - F_{t(n-1)}(T)$.

Example: Two-sided T Test

- ▶ For HEIGHT in SASUSER.CLASS, to test if its mean is 65.
- ▶ Use PROC TTEST, H0= to set μ_0 , VAR to define the variable to test.
- ▶ Result: the mean height is significantly different from 65 at significance level $\alpha = 0.05$.

Example: Two-sided T Test

- ▶ For HEIGHT in SASUSER.CLASS, to test if its mean is 65.
- ▶ Use PROC TTEST, H0= to set μ_0 , VAR to define the variable to test.
- ▶ Result: the mean height is significantly different from 65 at significance level $\alpha = 0.05$.

Example: Two-sided T Test

- ▶ For HEIGHT in SASUSER.CLASS, to test if its mean is 65.
- ▶ Use PROC TTEST, H0= to set μ_0 , VAR to define the variable to test.
- ▶ Result: the mean height is significantly different from 65 at significance level $\alpha = 0.05$.


```
proc ttest data=sasuser.class H0=65;
  var height;
run;
```

Variable	N	Mean	Std Dev
height	19	62.337	5.1271

Variable	DF	t Value	Pr > t
height	18	-2.26	0.0362

Example: One-sided T Test

- ▶ For the previous example, what if our question is “is the mean height greater than 65?”
- ▶ $H_a : \mu > 65$, which gives $H_0 : \mu \leq 65$.
- ▶ The program remains the same, but the p-value is only for the two-sided test.
 - ▶ If $\bar{X} \geq \mu_0$, one sided p-value is half the two-sided p-value.
 - ▶ If $\bar{X} < \mu_0$, which means H_a is not sensible, we just accept H_0 .
- ▶ Since mean height is $62.337 < 65$, H_a is not sensible, we cannot conclude that mean height is significantly larger than 65.

Example: One-sided T Test

- ▶ For the previous example, what if our question is “is the mean height greater than 65?”
- ▶ $H_a : \mu > 65$, which gives $H_0 : \mu \leq 65$.
- ▶ The program remains the same, but the p-value is only for the two-sided test.
 - ▶ If $\bar{X} \geq \mu_0$, one sided p-value is half the two-sided p-value.
 - ▶ If $\bar{X} < \mu_0$, which means H_a is not sensible, we just accept H_0 .
- ▶ Since mean height is $62.337 < 65$, H_a is not sensible, we cannot conclude that mean height is significantly larger than 65.

Example: One-sided T Test

- ▶ For the previous example, what if our question is “is the mean height greater than 65?”
- ▶ $H_a : \mu > 65$, which gives $H_0 : \mu \leq 65$.
- ▶ The program remains the same, but the p-value is only for the two-sided test.
 - ▶ If $\bar{X} \geq \mu_0$, one sided p-value is half the two-sided p-value.
 - ▶ If $\bar{X} < \mu_0$, which means H_a is not sensible, we just accept H_0 .
- ▶ Since mean height is $62.337 < 65$, H_a is not sensible, we cannot conclude that mean height is significantly larger than 65.

Example: One-sided T Test

- ▶ For the previous example, what if our question is “is the mean height greater than 65?”
- ▶ $H_a : \mu > 65$, which gives $H_0 : \mu \leq 65$.
- ▶ The program remains the same, but the p-value is only for the two-sided test.
 - ▶ If $\bar{X} \geq \mu_0$, one sided p-value is half the two-sided p-value.
 - ▶ If $\bar{X} < \mu_0$, which means H_a is not sensible, we just accept H_0 .
- ▶ Since mean height is $62.337 < 65$, H_a is not sensible, we cannot conclude that mean height is significantly larger than 65.

Example: One-sided T Test

- ▶ For the previous example, what if our question is “is the mean height greater than 65?”
- ▶ $H_a : \mu > 65$, which gives $H_0 : \mu \leq 65$.
- ▶ The program remains the same, but the p-value is only for the two-sided test.
 - ▶ If $\bar{X} \geq \mu_0$, one sided p-value is half the two-sided p-value.
 - ▶ If $\bar{X} < \mu_0$, which means H_a is not sensible, we just accept H_0 .
- ▶ Since mean height is $62.337 < 65$, H_a is not sensible, we cannot conclude that mean height is significantly larger than 65.

Example: One-sided T Test

- ▶ For the previous example, what if our question is “is the mean height greater than 65?”
- ▶ $H_a : \mu > 65$, which gives $H_0 : \mu \leq 65$.
- ▶ The program remains the same, but the p-value is only for the two-sided test.
 - ▶ If $\bar{X} \geq \mu_0$, one sided p-value is half the two-sided p-value.
 - ▶ If $\bar{X} < \mu_0$, which means H_a is not sensible, we just accept H_0 .
- ▶ Since mean height is $62.337 < 65$, H_a is not sensible, we cannot conclude that mean height is significantly larger than 65.

- ▶ X and Y are two independent populations,
 $X \sim N(\mu_1, \sigma^2)$, $Y \sim N(\mu_2, \sigma^2)$.
- ▶ Sample from X : X_1, \dots, X_{n_1} ; Sample from Y :
 Y_1, \dots, Y_{n_2} .
- ▶ Two sided test: $H_0 : \mu_1 = \mu_2 \longleftrightarrow H_a : \mu_1 \neq \mu_2$.
- ▶ Test statistic:

Two Sample T-Test

- ▶ X and Y are two independent populations,
 $X \sim N(\mu_1, \sigma^2)$, $Y \sim N(\mu_2, \sigma^2)$.
- ▶ Sample from X : X_1, \dots, X_{n_1} ; Sample from Y :
 Y_1, \dots, Y_{n_2} .
- ▶ Two sided test: $H_0 : \mu_1 = \mu_2 \longleftrightarrow H_a : \mu_1 \neq \mu_2$.
- ▶ Test statistic:

$$S_p^2 = \frac{1}{n_1 + n_2 - 2} \left(\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{j=1}^{n_2} (Y_j - \bar{Y})^2 \right)$$

$$T = \frac{\bar{X} - \bar{Y}}{S_p / \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \stackrel{H_0}{\sim} t(n_1 + n_2 - 2)$$

- ▶ Rejection field: $\{|T| > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha/2)$.
- ▶ p-value: $2[1 - F_{t(n_1+n_2-2)}(|T|)]$.

- ▶ X and Y are two independent populations,
 $X \sim N(\mu_1, \sigma^2)$, $Y \sim N(\mu_2, \sigma^2)$.
- ▶ Sample from X : X_1, \dots, X_{n_1} ; Sample from Y :
 Y_1, \dots, Y_{n_2} .
- ▶ Two sided test: $H_0 : \mu_1 = \mu_2 \longleftrightarrow H_a : \mu_1 \neq \mu_2$.
- ▶ Test statistic:

Two Sample T-Test

- ▶ X and Y are two independent populations,
 $X \sim N(\mu_1, \sigma^2)$, $Y \sim N(\mu_2, \sigma^2)$.
- ▶ Sample from X : X_1, \dots, X_{n_1} ; Sample from Y :
 Y_1, \dots, Y_{n_2} .
- ▶ Two sided test: $H_0 : \mu_1 = \mu_2 \longleftrightarrow H_a : \mu_1 \neq \mu_2$.
- ▶ Test statistic:

$$S_p^2 = \frac{1}{n_1 + n_2 - 2} \left(\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{j=1}^{n_2} (Y_j - \bar{Y})^2 \right)$$

$$T = \frac{\bar{X} - \bar{Y}}{S_p / \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \stackrel{H_0}{\sim} t(n_1 + n_2 - 2)$$

- ▶ Rejection field: $\{|T| > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha/2)$.
- ▶ p-value: $2[1 - F_{t(n_1+n_2-2)}(|T|)]$.

Two Sample T-Test

- ▶ X and Y are two independent populations,
 $X \sim N(\mu_1, \sigma^2)$, $Y \sim N(\mu_2, \sigma^2)$.
- ▶ Sample from X : X_1, \dots, X_{n_1} ; Sample from Y :
 Y_1, \dots, Y_{n_2} .
- ▶ Two sided test: $H_0 : \mu_1 = \mu_2 \longleftrightarrow H_a : \mu_1 \neq \mu_2$.
- ▶ Test statistic:

$$S_p^2 = \frac{1}{n_1 + n_2 - 2} \left(\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{j=1}^{n_2} (Y_j - \bar{Y})^2 \right)$$

$$T = \frac{\bar{X} - \bar{Y}}{S_p / \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \stackrel{H_0}{\sim} t(n_1 + n_2 - 2)$$

- ▶ Rejection field: $\{|T| > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha/2)$.
- ▶ p-value: $2[1 - F_{t(n_1+n_2-2)}(|T|)]$.

Two Sample T-Test

- ▶ X and Y are two independent populations,
 $X \sim N(\mu_1, \sigma^2)$, $Y \sim N(\mu_2, \sigma^2)$.
- ▶ Sample from X : X_1, \dots, X_{n_1} ; Sample from Y :
 Y_1, \dots, Y_{n_2} .
- ▶ Two sided test: $H_0 : \mu_1 = \mu_2 \longleftrightarrow H_a : \mu_1 \neq \mu_2$.
- ▶ Test statistic:

$$S_p^2 = \frac{1}{n_1 + n_2 - 2} \left(\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{j=1}^{n_2} (Y_j - \bar{Y})^2 \right)$$

$$T = \frac{\bar{X} - \bar{Y}}{S_p / \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \stackrel{H_0}{\sim} t(n_1 + n_2 - 2)$$

- ▶ Rejection field: $\{|T| > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha/2)$.
- ▶ p-value: $2[1 - F_{t(n_1+n_2-2)}(|T|)]$.

One-sided Two-sample T-Test

- ▶ $H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$.
- ▶ Rejection field: $\{T > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha)$.
- ▶ p-value: $1 - F_{t(n_1+n_2-2)}(T)$.
- ▶ Prerequisites for two-sample t-test:
 - ▶ Independence;
 - ▶ Normality.
 - ▶ Variance equality(homogeneity);

One-sided Two-sample T-Test

- ▶ $H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$.
- ▶ Rejection field: $\{T > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha)$.
- ▶ p-value: $1 - F_{t(n_1+n_2-2)}(T)$.
- ▶ Prerequisites for two-sample t-test:
 - ▶ Independence;
 - ▶ Normality.
 - ▶ Variance equality(homogeneity);

One-sided Two-sample T-Test

- ▶ $H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$.
- ▶ Rejection field: $\{T > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha)$.
- ▶ p-value: $1 - F_{t(n_1+n_2-2)}(T)$.
- ▶ Prerequisites for two-sample t-test:
 - ▶ Independence;
 - ▶ Normality.
 - ▶ Variance equality(homogeneity);

One-sided Two-sample T-Test

- ▶ $H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$.
- ▶ Rejection field: $\{T > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha)$.
- ▶ p-value: $1 - F_{t(n_1+n_2-2)}(T)$.
- ▶ Prerequisites for two-sample t-test:
 - ▶ Independence;
 - ▶ Normality.
 - ▶ Variance equality(homogeneity);

One-sided Two-sample T-Test

- ▶ $H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$.
- ▶ Rejection field: $\{T > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha)$.
- ▶ p-value: $1 - F_{t(n_1+n_2-2)}(T)$.
- ▶ Prerequisites for two-sample t-test:
 - ▶ Independence;
 - ▶ Normality.
 - ▶ Variance equality(homogeneity);

One-sided Two-sample T-Test

- ▶ $H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$.
- ▶ Rejection field: $\{T > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha)$.
- ▶ p-value: $1 - F_{t(n_1+n_2-2)}(T)$.
- ▶ Prerequisites for two-sample t-test:
 - ▶ Independence;
 - ▶ Normality.
 - ▶ Variance equality(homogeneity);

One-sided Two-sample T-Test

- ▶ $H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$.
- ▶ Rejection field: $\{T > \lambda\}$, $\lambda = F_{t(n_1+n_2-2)}^{-1}(1 - \alpha)$.
- ▶ p-value: $1 - F_{t(n_1+n_2-2)}(T)$.
- ▶ Prerequisites for two-sample t-test:
 - ▶ Independence;
 - ▶ Normality.
 - ▶ Variance equality(homogeneity);

Example: Two-sided Two-sample T-Test

To compare mean height of girls and boys.

```
proc ttest data=sasuser.class;  
  class sex;  
  var height;  
run;
```

Statistics											
Variable	sex	N	Lower CL Mean	Mean	Upper CL Mean	Lower CL Std Dev	Std Dev	Upper CL Std Dev	Std Err	Min.	Max.
height	F	9	56.731	60.589	64.446	3.3897	5.0183	9.614	1.6728	51.3	66.5
height	M	10	60.378	63.91	67.442	3.3965	4.9379	9.0147	1.5615	57.3	72
height	Diff (1-2)		-8.145	-3.321	1.5025	3.7339	4.9759	7.4596	2.2863		

T-Tests					
Variable	Method	Variances	DF	t Value	Pr > t
height	Pooled	Equal	17	-1.45	0.1645
height	Satterthwaite	Unequal	16.7	-1.45	0.1652

Equality of Variances					
Variable	Method	Num DF	Den DF	F Value	Pr > F
height	Folded F	8	9	1.03	0.9527

- ▶ Female mean height 60.598, male mean height 63.91.
- ▶ Two-sided t-test $p\text{-value}=0.1652$. No significant difference between height of the two groups.
- ▶ Equal variance? There is a table for equality of variance test. P-value 0.9527 means that we can assume the two group have equal variance.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

- ▶ Female mean height 60.598, male mean height 63.91.
- ▶ Two-sided t-test $p\text{-value}=0.1652$. No significant difference between height of the two groups.
- ▶ Equal variance? There is a table for equality of variance test. P-value 0.9527 means that we can assume the two group have equal variance.

- ▶ Female mean height 60.598, male mean height 63.91.
- ▶ Two-sided t-test $p\text{-value}=0.1652$. No significant difference between height of the two groups.
- ▶ Equal variance? There is a table for equality of variance test. P-value 0.9527 means that we can assume the two group have equal variance.

Satterthwaite Test

- ▶ Satterthwaite test do not assume equal variance.



$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{S_x^2}{n_1} + \frac{S_y^2}{n_2}}} \sim t(\text{df}_{\text{Satterthwaite}}) \text{ (Approximate)}$$

$$\text{df}_{\text{Satterthwaite}} = \frac{(n_1 - 1)(n_2 - 1)}{(n_1 - 1)(1 - c)^2 + (n_2 - 1)c^2}$$

$$c = \frac{S_x^2/n_1}{\frac{S_x^2}{n_1} + \frac{S_y^2}{n_2}}$$

Satterthwaite Test

- ▶ Satterthwaite test do not assume equal variance.



$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{S_x^2}{n_1} + \frac{S_y^2}{n_2}}} \sim t(\text{df}_{\text{Satterthwaite}}) \text{ (Approximately, under } H_0)$$

$$\text{df}_{\text{Satterthwaite}} = \frac{(n_1 - 1)(n_2 - 1)}{(n_1 - 1)(1 - c)^2 + (n_2 - 1)c^2}$$

$$c = \frac{S_x^2/n_1}{\frac{S_x^2}{n_1} + \frac{S_y^2}{n_2}}$$

One-sided Two-sample T-test

- ▶ H_a should be sensible. If not, just accept H_0 .
- ▶ Now that female mean height is 60.598, male mean height 63.91, we could only test if male mean height(denote μ_M) is higher than female mean height(denote μ_F):

$$H_0 : \mu_F \geq \mu_M \longleftrightarrow H_a : \mu_F < \mu_M$$

- ▶ P-value is half the two-sided p-value, $0.1652/2 = 0.0826$. Male height is not significantly higher than female height at significance level 0.05.

One-sided Two-sample T-test

- ▶ H_a should be sensible. If not, just accept H_0 .
- ▶ Now that female mean height is 60.598, male mean height 63.91, we could only test if male mean height(denote μ_M) is higher than female mean height(denote μ_F):

$$H_0 : \mu_F \geq \mu_M \longleftrightarrow H_a : \mu_F < \mu_M$$

- ▶ P-value is half the two-sided p-value, $0.1652/2 = 0.0826$. Male height is not significantly higher than female height at significance level 0.05.

One-sided Two-sample T-test

- ▶ H_a should be sensible. If not, just accept H_0 .
- ▶ Now that female mean height is 60.598, male mean height 63.91, we could only test if male mean height(denote μ_M) is higher than female mean height(denote μ_F):

$$H_0 : \mu_F \geq \mu_M \longleftrightarrow H_a : \mu_F < \mu_M$$

- ▶ P-value is half the two-sided p-value, $0.1652/2 = 0.0826$. Male height is not significantly higher than female height at significance level 0.05.

Wilcoxon Rank Sum Test

- ▶ What if the two populations are not normal? Use the nonparametric Wilcoxon rank sum test.
- ▶ **Rank**: just like the ranks of grade scores, but rank 1 corresponds the smallest value.
- ▶ Compare the mean rank of two independent groups to see which group has higher value.

Wilcoxon Rank Sum Test

- ▶ What if the two populations are not normal? Use the nonparametric Wilcoxon rank sum test.
- ▶ **Rank**: just like the ranks of grade scores, but rank 1 corresponds the smallest value.
- ▶ Compare the mean rank of two independent groups to see which group has higher value.

Wilcoxon Rank Sum Test

- ▶ What if the two populations are not normal? Use the nonparametric Wilcoxon rank sum test.
- ▶ **Rank**: just like the ranks of grade scores, but rank 1 corresponds the smallest value.
- ▶ Compare the mean rank of two independent groups to see which group has higher value.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

Example: Two-sided Two-sample T-Test

Compare the GPA score of boys and girls.

```
proc npar1way data=sasuser.gpa  
    wilcoxon;  
    class sex;  
    var gpa;  
run;
```

Wilcoxon Scores (Rank Sums) for Variable gpa Classified by Variable sex					
sex	N	Sum of Scores	Expected Under H0	Std Dev Under H0	Mean Score
female	145	16067.50	16312.50	463.429146	110.810345
male	79	9132.50	8887.50	463.429146	115.601266
Average scores were used for ties.					

Dongfeng Li

One Sample Test of the Mean

Comparing Two Groups

Paired Comparison

One Sample Test of the Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

- ▶ Female mean score 110.8, male mean score 115.6, male scores better. Is it significant?
- ▶ For two-sided test, using normal approximation for the test statistic, p-value is 0.5978, no significant difference between the GPA scores of female and male students.
- ▶ For one sided test, we can only test H_a : male scores better. P-value is 0.2989, male is not significantly better than female regarding GPA scores.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

- ▶ Female mean score 110.8, male mean score 115.6, male scores better. Is it significant?
- ▶ For two-sided test, using normal approximation for the test statistic, p-value is 0.5978, no significant difference between the GPA scores of female and male students.
- ▶ For one sided test, we can only test H_a : male scores better. P-value is 0.2989, male is not significantly better than female regarding GPA scores.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

- ▶ Female mean score 110.8, male mean score 115.6, male scores better. Is it significant?
- ▶ For two-sided test, using normal approximation for the test statistic, p-value is 0.5978, no significant difference between the GPA scores of female and male students.
- ▶ For one sided test, we can only test H_a : male scores better. P-value is 0.2989, male is not significantly better than female regarding GPA scores.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

Paired Comparison

- ▶ Comparing two measurements of the same subject, instead of comparing the same measurement of two groups of subjects, is a different problem from two-sample test. It is called **paired-comparison**, we use paired t-test to solve the problem.
- ▶ Example: Comparing the blood pressure of the same subject before and after treatment by some drugs; in a fitness program, comparing the heart rate at entrance and at end of the program, etc.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

Paired Comparison

- ▶ Comparing two measurements of the same subject, instead of comparing the same measurement of two groups of subjects, is a different problem from two-sample test. It is called **paired-comparison**, we use paired t-test to solve the problem.
- ▶ Example: Comparing the blood pressure of the same subject before and after treatment by some drugs; in a fitness program, comparing the heart rate at entrance and at end of the program, etc.

Paired T-test

- ▶ Let X be the “before” measurement, Y be the “after” measurement, both belong to the same subject, to compare the mean of X and Y , let $Z = X - Y$, we simply compare μ_Z with 0.
- ▶ Program solution 1: use PROC TTEST with PAIRED statement. Need Z normal assumption.
- ▶ Program solution 2: first compute Z , then do one sample test $H_0 : \mu_Z = 0 \longleftrightarrow H_a : \mu_Z \neq 0$ using PROC UNIVARIATE. This could also give the **signed rank test** and the **sign test**.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

Paired T-test

- ▶ Let X be the “before” measurement, Y be the “after” measurement, both belong to the same subject, to compare the mean of X and Y , let $Z = X - Y$, we simply compare μ_Z with 0.
- ▶ Program solution 1: use PROC TTEST with PAIRED statement. Need Z normal assumption.
- ▶ Program solution 2: first compute Z , then do one sample test $H_0 : \mu_Z = 0 \longleftrightarrow H_a : \mu_Z \neq 0$ using PROC UNIVARIATE. This could also give the signed rank test and the sign test.

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

Paired T-test

- ▶ Let X be the “before” measurement, Y be the “after” measurement, both belong to the same subject, to compare the mean of X and Y , let $Z = X - Y$, we simply compare μ_Z with 0.
- ▶ Program solution 1: use PROC TTEST with PAIRED statement. Need Z normal assumption.
- ▶ Program solution 2: first compute Z , then do one sample test $H_0 : \mu_Z = 0 \longleftrightarrow H_a : \mu_Z \neq 0$ using PROC UNIVARIATE. This could also give the **signed rank test** and the **sign test**.

Example: Paired T-Test Using PROC TTEST

A stimulus is being examined to determine its effect on systolic blood pressure. Twelve men participate in the study. Their systolic blood pressure is measured both before and after the stimulus is applied. Program:

```
title 'Paired Comparison';
data pressure;
    input SBPbefore SBPafter @@;
    datalines;
120 128    124 131    130 131    118 127
140 132    128 125    140 141    135 137
126 118    130 132    126 129    127 135
;
run;
proc ttest;
    paired SBPbefore*SBPafter;
run;
```

Paired T-Test: result

Statistics										
Difference	N	Lower CL Mean	Mean	Upper CL Mean	Lower CL Std Dev	Std Dev	Upper CL Std Dev	Std Err	Minimum	Maximum
SBPbefore - SBPafter	12	-5.536	-1.833	1.8698	4.1288	5.8284	9.8958	1.6825	-9	8

T-Tests			
Difference	DF	t Value	Pr > t
SBPbefore - SBPafter	11	-1.09	0.2992

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

One Sample Test of the
Mean

Comparing Two Groups

Paired Comparison

Comparing Proportions

Analysis of
Variance

Example: Paired T-Test Using PROC UNIVARIATE

```
data _tmp_;  
    set pressure;  
    diff = SBPbefore - SBPafter;  
run;  
proc univariate;  
    var diff;  
run;  
proc datasets library=work nolist;  
    delete _tmp_;  
quit;
```


PROC UNIVARIATE Result

Tests for Location: $\mu_0=0$				
Test	Statistic		p Value	
Student's t	t	-1.08965	Pr > t	0.2992
Sign	M	-3	Pr >= M	0.1460
Signed Rank	S	-14.5	Pr >= S	0.2700

Comparing One Proportion

- ▶ $X \sim B(1, p)$. Sample X_1, X_2, \dots, X_n .
- ▶ $H_0 : p = p_0 \longleftrightarrow H_a : p \neq p_0$.
- ▶ When $np \geq 5, n(1 - p) \geq 5$, use the approximate Z test:

$$Z = \frac{\bar{X} - p_0}{\sqrt{p_0(1 - p_0)/n}} \stackrel{H_0}{\sim} N(0,1) \text{ (approximately)}$$

- ▶ Two-sided test p-value: $2(1 - \Phi(|Z|))$.
- ▶ For $H_0 : p \leq p_0 \longleftrightarrow H_a : p > p_0$, p-value is $1 - \Phi(Z)$.
Note that if $\bar{X} < p_0$ then p-value is greater than 0.5.
- ▶ For $H_0 : p \geq p_0 \longleftrightarrow H_a : p < p_0$, p-value is $\Phi(Z)$.
Note that if $\bar{X} > p_0$ then p-value is greater than 0.5.

Comparing One Proportion

- ▶ $X \sim B(1, p)$. Sample X_1, X_2, \dots, X_n .
- ▶ $H_0 : p = p_0 \longleftrightarrow H_a : p \neq p_0$.
- ▶ When $np \geq 5, n(1 - p) \geq 5$, use the approximate Z test:

$$Z = \frac{\bar{X} - p_0}{\sqrt{p_0(1 - p_0)/n}} \stackrel{H_0}{\sim} N(0,1) \text{ (approximately)}$$

- ▶ Two-sided test p-value: $2(1 - \Phi(|Z|))$.
- ▶ For $H_0 : p \leq p_0 \longleftrightarrow H_a : p > p_0$, p-value is $1 - \Phi(Z)$.
Note that if $\bar{X} < p_0$ then p-value is greater than 0.5.
- ▶ For $H_0 : p \geq p_0 \longleftrightarrow H_a : p < p_0$, p-value is $\Phi(Z)$.
Note that if $\bar{X} > p_0$ then p-value is greater than 0.5.

Comparing One Proportion

- ▶ $X \sim B(1, p)$. Sample X_1, X_2, \dots, X_n .
- ▶ $H_0 : p = p_0 \longleftrightarrow H_a : p \neq p_0$.
- ▶ When $np \geq 5, n(1 - p) \geq 5$, use the approximate Z test:

$$Z = \frac{\bar{X} - p_0}{\sqrt{p_0(1 - p_0)/n}} \stackrel{H_0}{\sim} N(0,1) \text{ (approximately)}$$

- ▶ Two-sided test p-value: $2(1 - \Phi(|Z|))$.
- ▶ For $H_0 : p \leq p_0 \longleftrightarrow H_a : p > p_0$, p-value is $1 - \Phi(Z)$.
Note that if $\bar{X} < p_0$ then p-value is greater than 0.5.
- ▶ For $H_0 : p \geq p_0 \longleftrightarrow H_a : p < p_0$, p-value is $\Phi(Z)$.
Note that if $\bar{X} > p_0$ then p-value is greater than 0.5.

Comparing One Proportion

- ▶ $X \sim B(1, p)$. Sample X_1, X_2, \dots, X_n .
- ▶ $H_0 : p = p_0 \longleftrightarrow H_a : p \neq p_0$.
- ▶ When $np \geq 5, n(1 - p) \geq 5$, use the approximate Z test:

$$Z = \frac{\bar{X} - p_0}{\sqrt{p_0(1 - p_0)/n}} \stackrel{H_0}{\sim} N(0,1) \text{ (approximately)}$$

- ▶ Two-sided test p-value: $2(1 - \Phi(|Z|))$.
- ▶ For $H_0 : p \leq p_0 \longleftrightarrow H_a : p > p_0$, p-value is $1 - \Phi(Z)$.
Note that if $\bar{X} < p_0$ then p-value is greater than 0.5.
- ▶ For $H_0 : p \geq p_0 \longleftrightarrow H_a : p < p_0$, p-value is $\Phi(Z)$.
Note that if $\bar{X} > p_0$ then p-value is greater than 0.5.

Comparing One Proportion

- ▶ $X \sim B(1, p)$. Sample X_1, X_2, \dots, X_n .
- ▶ $H_0 : p = p_0 \longleftrightarrow H_a : p \neq p_0$.
- ▶ When $np \geq 5, n(1 - p) \geq 5$, use the approximate Z test:

$$Z = \frac{\bar{X} - p_0}{\sqrt{p_0(1 - p_0)/n}} \stackrel{H_0}{\sim} N(0,1) \text{ (approximately)}$$

- ▶ Two-sided test p-value: $2(1 - \Phi(|Z|))$.
- ▶ For $H_0 : p \leq p_0 \longleftrightarrow H_a : p > p_0$, p-value is $1 - \Phi(Z)$.
Note that if $\bar{X} < p_0$ then p-value is greater than 0.5.
- ▶ For $H_0 : p \geq p_0 \longleftrightarrow H_a : p < p_0$, p-value is $\Phi(Z)$.
Note that if $\bar{X} > p_0$ then p-value is greater than 0.5.

Comparing One Proportion

- ▶ $X \sim B(1, p)$. Sample X_1, X_2, \dots, X_n .
- ▶ $H_0 : p = p_0 \longleftrightarrow H_a : p \neq p_0$.
- ▶ When $np \geq 5, n(1 - p) \geq 5$, use the approximate Z test:

$$Z = \frac{\bar{X} - p_0}{\sqrt{p_0(1 - p_0)/n}} \stackrel{H_0}{\sim} N(0,1) \text{ (approximately)}$$

- ▶ Two-sided test p-value: $2(1 - \Phi(|Z|))$.
- ▶ For $H_0 : p \leq p_0 \longleftrightarrow H_a : p > p_0$, p-value is $1 - \Phi(Z)$.
Note that if $\bar{X} < p_0$ then p-value is greater than 0.5.
- ▶ For $H_0 : p \geq p_0 \longleftrightarrow H_a : p < p_0$, p-value is $\Phi(Z)$.
Note that if $\bar{X} > p_0$ then p-value is greater than 0.5.

Example: Hepatitis B Percent

- ▶ Is China's hepatitis B infected percent 8%?
- ▶ A random sample of 100 persons show 5 infected with hepatitis B.

```
%MACRO percentzt (n,n1,p0);  
data _null_;  
    file print;  
    p0 = &p0.; n = &n.; n1 = &n1.;  
    xbar = n1/n;  
    Z = (xbar - p0)/sqrt(p0 * (1-p0)/n);  
    ptwosided = 2*(1 - probnorm(abs(Z)));  
    prightsided = 1 - probnorm(Z);  
    pleftsided = probnorm(Z);  
    put '==== Test for percent =====';  
    put 'n = ' n ' p = ' xbar;  
    put 'p0 = ' p0;  
    put 'Z = ' Z;  
    put 'Pr > |Z|: ' ptwosided pvalue.;  
    put 'Pr > Z: ' prightsided pvalue.;  
    put 'Pr < Z: ' pleftsided pvalue.;  
run;  
%MEND percentzt;  
%percentzt (100,5,0.08);
```

- ▶ Conclusion: no significant difference with 8%.

Example: Hepatitis B Percent

- ▶ Is China's hepatitis B infected percent 8%?
- ▶ A random sample of 100 persons show 5 infected with hepatitis B.

```
%MACRO percentzt (n,n1,p0);  
data _null_;  
    file print;  
    p0 = &p0.; n = &n.; n1 = &n1.;  
    xbar = n1/n;  
    Z = (xbar - p0)/sqrt(p0 * (1-p0)/n);  
    ptwosided = 2*(1 - probnorm(abs(Z)));  
    prightsided = 1 - probnorm(Z);  
    pleftsided = probnorm(Z);  
    put '==== Test for percent =====';  
    put 'n = ' n ' p = ' xbar;  
    put 'p0 = ' p0;  
    put 'Z = ' Z;  
    put 'Pr > |Z|: ' ptwosided pvalue.;  
    put 'Pr > Z: ' prightsided pvalue.;  
    put 'Pr < Z: ' pleftsided pvalue.;  
run;  
%MEND percentzt;  
%percentzt (100,5,0.08);
```

- ▶ Conclusion: no significant difference with 8%.

Example: Hepatitis B Percent

- ▶ Is China's hepatitis B infected percent 8%?
- ▶ A random sample of 100 persons show 5 infected with hepatitis B.

```
%MACRO percentzt (n,n1,p0);  
data _null_;  
    file print;  
    p0 = &p0.; n = &n.; n1 = &n1.;  
    xbar = n1/n;  
    Z = (xbar - p0)/sqrt(p0 * (1-p0)/n);  
    ptwosided = 2*(1 - probnorm(abs(Z)));  
    prightsided = 1 - probnorm(Z);  
    pleftsided = probnorm(Z);  
    put '==== Test for percent =====';  
    put 'n = ' n ' p = ' xbar;  
    put 'p0 = ' p0;  
    put 'Z = ' Z;  
    put 'Pr > |Z|: ' ptwosided pvalue.;  
    put 'Pr > Z: ' prightsided pvalue.;  
    put 'Pr < Z: ' pleftsided pvalue.;  
run;  
%MEND percentzt;  
%percentzt (100,5,0.08);
```

- ▶ Conclusion: no significant difference with 8%.

Comparing Two Proportions—Two-Sided

- ▶ $X \sim B(1, p_1)$, $Y \sim B(1, p_2)$, independent, test $H_0 : p_1 = p_2 \longleftrightarrow H_a : p_1 \neq p_2$.
- ▶ n_1 trials of X gives s_1 successes, n_2 trials of Y gives s_2 successes.
- ▶ When n_1 and n_2 large, let $W = \frac{s_1}{n_1} - \frac{s_2}{n_2}$, then $SE^2(W) \approx \hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$, $W/SE(W)$ is asymptotically distributed $N(0,1)$. $\hat{p} = \frac{s_1 + s_2}{n_1 + n_2}$.
- ▶ Let $Z = W/SE(W)$, p-value is $2(1 - \Phi(|Z|))$.

Comparing Two Proportions—Two-Sided

- ▶ $X \sim B(1, p_1)$, $Y \sim B(1, p_2)$, independent, test $H_0 : p_1 = p_2 \longleftrightarrow H_a : p_1 \neq p_2$.
- ▶ n_1 trials of X gives s_1 successes, n_2 trials of Y gives s_2 successes.
- ▶ When n_1 and n_2 large, let $W = \frac{s_1}{n_1} - \frac{s_2}{n_2}$, then $SE^2(W) \approx \hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$, $W/SE(W)$ is asymptotically distributed $N(0,1)$. $\hat{p} = \frac{s_1 + s_2}{n_1 + n_2}$.
- ▶ Let $Z = W/SE(W)$, p-value is $2(1 - \Phi(|Z|))$.

Comparing Two Proportions—Two-Sided

- ▶ $X \sim B(1, p_1)$, $Y \sim B(1, p_2)$, independent, test $H_0 : p_1 = p_2 \longleftrightarrow H_a : p_1 \neq p_2$.
- ▶ n_1 trials of X gives s_1 successes, n_2 trials of Y gives s_2 successes.
- ▶ When n_1 and n_2 large, let $W = \frac{s_1}{n_1} - \frac{s_2}{n_2}$, then $SE^2(W) \approx \hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$, $W/SE(W)$ is asymptotically distributed $N(0,1)$. $\hat{p} = \frac{s_1 + s_2}{n_1 + n_2}$.
- ▶ Let $Z = W/SE(W)$, p-value is $2(1 - \Phi(|Z|))$.

Comparing Two Proportions—Two-Sided

- ▶ $X \sim B(1, p_1)$, $Y \sim B(1, p_2)$, independent, test $H_0 : p_1 = p_2 \longleftrightarrow H_a : p_1 \neq p_2$.
- ▶ n_1 trials of X gives s_1 successes, n_2 trials of Y gives s_2 successes.
- ▶ When n_1 and n_2 large, let $W = \frac{s_1}{n_1} - \frac{s_2}{n_2}$, then $SE^2(W) \approx \hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$, $W/SE(W)$ is asymptotically distributed $N(0,1)$. $\hat{p} = \frac{s_1 + s_2}{n_1 + n_2}$.
- ▶ Let $Z = W/SE(W)$, p-value is $2(1 - \Phi(|Z|))$.

Comparing Two Proportions—One-Sided

- ▶ When n_1 and n_2 large, let $W = \hat{p}_1 - \hat{p}_2$, then $SE^2(W) \approx \frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}$, $W/SE(W)$ is asymptotically distributed $N(0,1)$ when $p_1 = p_2$.
 $\hat{p}_1 = \frac{s_1}{n_1}$, $\hat{p}_2 = \frac{s_2}{n_2}$.
- ▶ Let $Z = W/SE(W)$.
- ▶ For right-sided alternative, p-value is $1 - \Phi(Z)$.
- ▶ For left-sided alternative, p-value is $\Phi(Z)$.

Comparing Two Proportions—One-Sided

- ▶ When n_1 and n_2 large, let $W = \hat{p}_1 - \hat{p}_2$, then $SE^2(W) \approx \frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}$, $W/SE(W)$ is asymptotically distributed $N(0,1)$ when $p_1 = p_2$.
 $\hat{p}_1 = \frac{s_1}{n_1}$, $\hat{p}_2 = \frac{s_2}{n_2}$.
- ▶ Let $Z = W/SE(W)$.
- ▶ For right-sided alternative, p-value is $1 - \Phi(Z)$.
- ▶ For left-sided alternative, p-value is $\Phi(Z)$.

Comparing Two Proportions—One-Sided

- ▶ When n_1 and n_2 large, let $W = \hat{p}_1 - \hat{p}_2$, then $SE^2(W) \approx \frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}$, $W/SE(W)$ is asymptotically distributed $N(0,1)$ when $p_1 = p_2$.
 $\hat{p}_1 = \frac{s_1}{n_1}$, $\hat{p}_2 = \frac{s_2}{n_2}$.
- ▶ Let $Z = W/SE(W)$.
- ▶ For right-sided alternative, p-value is $1 - \Phi(Z)$.
- ▶ For left-sided alternative, p-value is $\Phi(Z)$.

Comparing Two Proportions—One-Sided

- ▶ When n_1 and n_2 large, let $W = \hat{p}_1 - \hat{p}_2$, then $SE^2(W) \approx \frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}$, $W/SE(W)$ is asymptotically distributed $N(0,1)$ when $p_1 = p_2$.
 $\hat{p}_1 = \frac{s_1}{n_1}$, $\hat{p}_2 = \frac{s_2}{n_2}$.
- ▶ Let $Z = W/SE(W)$.
- ▶ For right-sided alternative, p-value is $1 - \Phi(Z)$.
- ▶ For left-sided alternative, p-value is $\Phi(Z)$.

Example: Hepatitis B Proportions of Male and Female

- ▶ Do male and female have the same proportion of Hepatitis B infected?
- ▶ A random sample of 100 males, 100 females show 10 males infected with hepatitis B, 8 females infected.
- ▶ Test result:

$$\Pr > |Z| : 0.6212$$

$$\Pr > Z : 0.3105$$

$$\Pr < Z : 0.6895$$

- ▶ Conclusion: no significant difference between male and female.

Example: Hepatitis B Proportions of Male and Female

- ▶ Do male and female have the same proportion of Hepatitis B infected?
- ▶ A random sample of 100 males, 100 females show 10 males infected with hepatitis B, 8 females infected.
- ▶ Test result:

$$\Pr > |Z| : 0.6212$$

$$\Pr > Z : 0.3105$$

$$\Pr < Z : 0.6895$$

- ▶ Conclusion: no significant difference between male and female.

Example: Hepatitis B Proportions of Male and Female

- ▶ Do male and female have the same proportion of Hepatitis B infected?
- ▶ A random sample of 100 males, 100 females show 10 males infected with hepatitis B, 8 females infected.
- ▶ Test result:

$$\Pr > |Z| : 0.6212$$

$$\Pr > Z : 0.3105$$

$$\Pr < Z : 0.6895$$

- ▶ Conclusion: no significant difference between male and female.

Example: Hepatitis B Proportions of Male and Female

- ▶ Do male and female have the same proportion of Hepatitis B infected?
- ▶ A random sample of 100 males, 100 females show 10 males infected with hepatitis B, 8 females infected.
- ▶ Test result:

$$\Pr > |Z| : 0.6212$$

$$\Pr > Z : 0.3105$$

$$\Pr < Z : 0.6895$$

- ▶ Conclusion: no significant difference between male and female.

```
%MACRO percent2z (n1,s1,n2,s2);
data _null_;
    file print;
    n1=&n1; s1=&s1; n2=&n2; s2=&s2;
    hatp = (s1+s2)/(n1+n2);
    hatp1 = s1/n1; hatp2 = s2/n2;
    Z2s = (hatp1 - hatp2) /
        sqrt(hatp*(1-hatp)*(1/n1 + 1/n2));
    Z1s = (hatp1 - hatp2) /
        sqrt(hatp1*(1-hatp1)/n1
            +hatp2*(1-hatp2)/n2);
    ptwosided = 2*(1 - probnorm(abs(Z2s)));
    prightsided = 1 - probnorm(Z1s);
    pleftsided = probnorm(Z1s);
    put '==== Test for percent =====';
    put 'n1 = ' n1 ' s1 = ' s1 ' p1=' hatp1;
    put 'n2 = ' n2 ' s2 = ' s2 ' p2=' hatp2;
    put 'Pr > |Z|: ' ptwosided pvalue.;
    put 'Pr > Z: ' prightsided pvalue.;
    put 'Pr < Z: ' pleftsided pvalue.;
run;
%MEND percent2z;

%percent2z(100,10,100,8);
```

Analysis of Variance

- ▶ Test for significant effect of some **factors** on a **response**.
- ▶ One-way ANOVA.
- ▶ Nonparametric Kruskal-Wallis test.
- ▶ Two-way ANOVA, additive models, interactions.

Analysis of Variance

- ▶ Test for significant effect of some **factors** on a **response**.
- ▶ One-way ANOVA.
- ▶ Nonparametric Kruskal-Wallis test.
- ▶ Two-way ANOVA, additive models, interactions.

Analysis of Variance

- ▶ Test for significant effect of some **factors** on a **response**.
- ▶ One-way ANOVA.
- ▶ Nonparametric Kruskal-Wallis test.
- ▶ Two-way ANOVA, additive models, interactions.

Analysis of Variance

- ▶ Test for significant effect of some **factors** on a **response**.
- ▶ One-way ANOVA.
- ▶ Nonparametric Kruskal-Wallis test.
- ▶ Two-way ANOVA, additive models, interactions.

One-Way ANOVA

- ▶ Two sample t-test → one-way anova:
- ▶ Response Y , group C (factor).
- ▶ Two sample t-test: Compare the mean of Y in the two groups of C .
- ▶ One-way ANOVA: Compare the mean of Y in more than two groups of C .
- ▶ Question:
 - ▶ Is there any significant difference among the means of Y of different groups? Equivalently, does factor C have significant effect on the mean level of Y ?
 - ▶ Which pair of groups have significant mean difference? This is the multiple comparison problem.

One-Way ANOVA

- ▶ Two sample t-test \longrightarrow one-way anova:
- ▶ Response Y , group C (factor).
- ▶ Two sample t-test: Compare the mean of Y in the two groups of C .
- ▶ One-way ANOVA: Compare the mean of Y in more than two groups of C .
- ▶ Question:
 - ▶ Is there any significant difference among the means of Y of different groups? Equivalently, does factor C have significant effect on the mean level of Y ?
 - ▶ Which pair of groups have significant mean difference? This is the multiple comparison problem.

One-Way ANOVA

- ▶ Two sample t-test \longrightarrow one-way anova:
- ▶ Response Y , group C (factor).
- ▶ Two sample t-test: Compare the mean of Y in the two groups of C .
- ▶ One-way ANOVA: Compare the mean of Y in more than two groups of C .
- ▶ Question:
 - ▶ Is there any significant difference among the means of Y of different groups? Equivalently, does factor C have significant effect on the mean level of Y ?
 - ▶ Which pair of groups have significant mean difference? This is the multiple comparison problem.

One-Way ANOVA

- ▶ Two sample t-test \longrightarrow one-way anova:
- ▶ Response Y , group C (factor).
- ▶ Two sample t-test: Compare the mean of Y in the two groups of C .
- ▶ One-way ANOVA: Compare the mean of Y in more than two groups of C .
- ▶ Question:
 - ▶ Is there any significant difference among the means of Y of different groups? Equivalently, does factor C have significant effect on the mean level of Y ?
 - ▶ Which pair of groups have significant mean difference? This is the multiple comparison problem.

One-Way ANOVA

- ▶ Two sample t-test \longrightarrow one-way anova:
- ▶ Response Y , group C (factor).
- ▶ Two sample t-test: Compare the mean of Y in the two groups of C .
- ▶ One-way ANOVA: Compare the mean of Y in more than two groups of C .
- ▶ Question:
 - ▶ Is there any significant difference among the means of Y of different groups? Equivalently, does factor C have significant effect on the mean level of Y ?
 - ▶ Which pair of groups have significant mean difference? This is the multiple comparison problem.

One-Way ANOVA

- ▶ Two sample t-test \longrightarrow one-way anova:
- ▶ Response Y , group C (factor).
- ▶ Two sample t-test: Compare the mean of Y in the two groups of C .
- ▶ One-way ANOVA: Compare the mean of Y in more than two groups of C .
- ▶ Question:
 - ▶ Is there any significant difference among the means of Y of different groups? Equivalently, does factor C have significant effect on the mean level of Y ?
 - ▶ Which pair of groups have significant mean difference? This is the multiple comparison problem.

One-Way ANOVA

- ▶ Two sample t-test \longrightarrow one-way anova:
- ▶ Response Y , group C (factor).
- ▶ Two sample t-test: Compare the mean of Y in the two groups of C .
- ▶ One-way ANOVA: Compare the mean of Y in more than two groups of C .
- ▶ Question:
 - ▶ Is there any significant difference among the means of Y of different groups? Equivalently, does factor C have significant effect on the mean level of Y ?
 - ▶ Which pair of groups have significant mean difference? This is the multiple comparison problem.

Example: The Comparison of Veneer Brands

- ▶ Compare WEAR of 5 BRANDS. Each brand has 4 samples.
- ▶ Use PROC ANOVA or PROC GLM. PROC GLM should be used for unbalanced design.
- ▶ Theoretical prerequisites: Independence, normal distribution, equal variances.
- ▶ Data:

```
data veneer;  
  input brand $ wear @@;  
  datalines;  
  
ACME      2.3  ACME      2.1  ACME      2.4  ACME      2.5  
CHAMP     2.2  CHAMP     2.3  CHAMP     2.4  CHAMP     2.6  
AJAX      2.2  AJAX      2.0  AJAX      1.9  AJAX      2.1  
TUFFY     2.4  TUFFY     2.7  TUFFY     2.6  TUFFY     2.7  
XTRA      2.3  XTRA      2.5  XTRA      2.3  XTRA      2.4  
  
;  
run;
```

Example: The Comparison of Veneer Brands

- ▶ Compare WEAR of 5 BRANDS. Each brand has 4 samples.
- ▶ Use PROC ANOVA or PROC GLM. PROC GLM should be used for unbalanced design.
- ▶ Theoretical prerequisites: Independence, normal distribution, equal variances.
- ▶ Data:

```
data veneer;  
  input brand $ wear @@;  
  datalines;  
ACME      2.3  ACME      2.1  ACME      2.4  ACME      2.5  
CHAMP     2.2  CHAMP     2.3  CHAMP     2.4  CHAMP     2.6  
AJAX      2.2  AJAX      2.0  AJAX      1.9  AJAX      2.1  
TUFFY     2.4  TUFFY     2.7  TUFFY     2.6  TUFFY     2.7  
XTRA      2.3  XTRA      2.5  XTRA      2.3  XTRA      2.4  
;  
run;
```

Example: The Comparison of Veneer Brands

- ▶ Compare WEAR of 5 BRANDS. Each brand has 4 samples.
- ▶ Use PROC ANOVA or PROC GLM. PROC GLM should be used for unbalanced design.
- ▶ Theoretical prerequisites: Independence, normal distribution, equal variances.
- ▶ Data:

```
data veneer;  
  input brand $ wear @@;  
  datalines;  
  
ACME      2.3  ACME      2.1  ACME      2.4  ACME      2.5  
CHAMP     2.2  CHAMP     2.3  CHAMP     2.4  CHAMP     2.6  
AJAX      2.2  AJAX      2.0  AJAX      1.9  AJAX      2.1  
TUFFY     2.4  TUFFY     2.7  TUFFY     2.6  TUFFY     2.7  
XTRA      2.3  XTRA      2.5  XTRA      2.3  XTRA      2.4  
  
;  
run;
```

Example: The Comparison of Veneer Brands

- ▶ Compare WEAR of 5 BRANDS. Each brand has 4 samples.
- ▶ Use PROC ANOVA or PROC GLM. PROC GLM should be used for unbalanced design.
- ▶ Theoretical prerequisites: Independence, normal distribution, equal variances.
- ▶ Data:

```
data veneer;  
  input brand $ wear @@;  
  datalines;  
ACME    2.3  ACME    2.1  ACME    2.4  ACME    2.5  
CHAMP    2.2  CHAMP    2.3  CHAMP    2.4  CHAMP    2.6  
AJAX     2.2  AJAX     2.0  AJAX     1.9  AJAX     2.1  
TUFFY    2.4  TUFFY    2.7  TUFFY    2.6  TUFFY    2.7  
XTRA     2.3  XTRA     2.5  XTRA     2.3  XTRA     2.4  
;  
run;
```

► Example using PROC ANOVA:

```
proc anova data=veneer;  
  class brand;  
  model wear = brand;  
quit;
```

► Example using PROC GLM:

```
proc glm data=samp.veneer;  
  class brand;  
  model wear = brand;  
quit;
```

► Example using PROC ANOVA:

```
proc anova data=veneer;  
  class brand;  
  model wear = brand;  
quit;
```

► Example using PROC GLM:

```
proc glm data=samp.veneer;  
  class brand;  
  model wear = brand;  
quit;
```


Dongfeng Li

One-Way ANOVA

Class Level Information		
Class	Levels	Values
Brand	5	ACME AJAX CHAMP TUFFY XTRA

Number of Observations Read	20
Number of Observations Used	20

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	4	0.61700000	0.15425000	7.40	0.0017
Error	15	0.31250000	0.02083333		
Corrected Total	19	0.92950000			

R-Square	Coeff Var	Root MSE	Wear Mean
0.663798	6.155120	0.144338	2.345000

Source	DF	Anova SS	Mean Square	F Value	Pr > F
Brand	4	0.61700000	0.15425000	7.40	0.0017

PROC GLM Result

The data information, model analysis of variance, model fit statistics part are very similar to PROC ANOVA results.

Different results:

Source	DF	Type I SS	Mean Square	F Value	Pr > F
Brand	4	0.61700000	0.15425000	7.40	0.0017

Source	DF	Type III SS	Mean Square	F Value	Pr > F
Brand	4	0.61700000	0.15425000	7.40	0.0017

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

Multiple Comparison

- ▶ Two test which pairs have significant differences, which are not significantly different.
- ▶ For 5 groups, we have $\binom{5}{2} = 10$ pairs to compare.
- ▶ Two kinds of type I errors:
 - ▶ Comparisonwise error rate(CER), only controls the type I error rate of each comparison;
 - ▶ Experimentwise error rate(FWR, Familywise Error Rate), controls the type I error rate of all of the comparisons.
- ▶ Controlling experimental error rate is more cautious.
- ▶ Control CER to discovery more possible differences.
- ▶ Control FWR to make sound inference.

- ▶ Fisher's LSD test controls the comparison error rate.
- ▶ Like repeated two-sample tests, but use a common SE estimate of the numerator.
- ▶ Let y_{ij} be the observation of j 'th individual of the i 'th category, $j = 1, \dots, r$, $i = 1, \dots, g$. \bar{y}_i is the sample mean of the i 'th group.
- ▶ Let

$$T_{ik} = \frac{\bar{y}_{i.} - \bar{y}_{k.}}{\text{SE}(\bar{y}_{i.} - \bar{y}_{k.})}$$

Fisher's Least Significant Difference Test

- ▶ Fisher's LSD test controls the comparison error rate.
- ▶ Like repeated two-sample tests, but use a common SE estimate of the numerator.
- ▶ Let y_{ij} be the observation of j 'th individual of the i 'th category, $j = 1, \dots, r$, $i = 1, \dots, g$. $\bar{y}_{i\cdot}$ is the sample mean of the i 'th group.
- ▶ Let

$$SE^2(\bar{y}_{i\cdot} - \bar{y}_{k\cdot}) \equiv \frac{1}{g(r-1)} \sum_i \sum_j (y_{ij} - \bar{y}_{i\cdot})^2$$

▶

$$T_{ik} = \frac{\bar{y}_{i\cdot} - \bar{y}_{k\cdot}}{SE(\bar{y}_{i\cdot} - \bar{y}_{k\cdot})}$$

Fisher's Least Significant Difference Test

- ▶ Fisher's LSD test controls the comparison error rate.
- ▶ Like repeated two-sample tests, but use a common SE estimate of the numerator.
- ▶ Let y_{ij} be the observation of j 'th individual of the i 'th category, $j = 1, \dots, r$, $i = 1, \dots, g$. $\bar{y}_{i\cdot}$ is the sample mean of the i 'th group.
- ▶ Let

$$SE^2(\bar{y}_{i\cdot} - \bar{y}_{k\cdot}) \equiv \frac{1}{g(r-1)} \sum_i \sum_j (y_{ij} - \bar{y}_{i\cdot})^2$$

▶

$$T_{ik} = \frac{\bar{y}_{i\cdot} - \bar{y}_{k\cdot}}{SE(\bar{y}_{i\cdot} - \bar{y}_{k\cdot})}$$

Fisher's Least Significant Difference Test

- ▶ Fisher's LSD test controls the comparison error rate.
- ▶ Like repeated two-sample tests, but use a common SE estimate of the numerator.
- ▶ Let y_{ij} be the observation of j 'th individual of the i 'th category, $j = 1, \dots, r$, $i = 1, \dots, g$. $\bar{y}_{i\cdot}$ is the sample mean of the i 'th group.
- ▶ Let

$$SE^2(\bar{y}_{i\cdot} - \bar{y}_{k\cdot}) \equiv \frac{1}{g(r-1)} \sum_i \sum_j (y_{ij} - \bar{y}_{i\cdot})^2$$

▶

$$T_{ik} = \frac{\bar{y}_{i\cdot} - \bar{y}_{k\cdot}}{SE(\bar{y}_{i\cdot} - \bar{y}_{k\cdot})}$$

Fisher's Least Significant Difference Test

- ▶ Fisher's LSD test controls the comparison error rate.
- ▶ Like repeated two-sample tests, but use a common SE estimate of the numerator.
- ▶ Let y_{ij} be the observation of j 'th individual of the i 'th category, $j = 1, \dots, r$, $i = 1, \dots, g$. $\bar{y}_{i\cdot}$ is the sample mean of the i 'th group.
- ▶ Let

$$SE^2(\bar{y}_{i\cdot} - \bar{y}_{k\cdot}) \equiv \frac{1}{g(r-1)} \sum_i \sum_j (y_{ij} - \bar{y}_{i\cdot})^2$$

▶

$$T_{ik} = \frac{\bar{y}_{i\cdot} - \bar{y}_{k\cdot}}{SE(\bar{y}_{i\cdot} - \bar{y}_{k\cdot})}$$

Example of Fisher's LSD Test

```
proc anova data=samp.veneer;  
  class brand;  
  model wear = brand;  
  means brand / t;  
quit;
```

t Tests (LSD) for Wear

Note

This test controls the Type I comparisonwise error rate, not the experimentwise error rate.

Alpha	0.05
Error Degrees of Freedom	15
Error Mean Square	0.020833
Critical Value of t	2.13145
Least Significant Difference	0.2175

Means with the same letter are not significantly different.

t Grouping	Mean	N	Brand
A	2.6000	4	TUFFY
B	2.3750	4	XTRA
B			
B	2.3750	4	CHAMP
B			
B	2.3250	4	ACME
C	2.0500	4	AJAX

REGWQ Test

- ▶ REGWQ test is one of the multiple testing methods which controls the experiment wise type I error rate.
- ▶ Example:

```
proc anova data=samp.veneer;  
  class brand;  
  model wear = brand;  
  means brand / regwq;  
quit;
```

- ▶ REGWQ test is one of the multiple testing methods which controls the experiment wise type I error rate.
- ▶ Example:

```
proc anova data=samp.veneer;  
  class brand;  
  model wear = brand;  
  means brand / regwq;  
quit;
```

SAS Programming
in Clinical Trials
Chapter 3. SAS
STAT

Dongfeng Li

One-Way ANOVA

Two-Way ANOVA

The Kruskal-Wallis Test

- ▶ What if the samples are not normally distributed? Use the nonparametric Kruskal-Wallis test, similar to the Wilcoxon rank sum test.
- ▶ Compute ranks of Y observations. Compare the mean rank of the groups.
- ▶ Use PROC NPAR1WAY with the WILCOXON option.

The Kruskal-Wallis Test

- ▶ What if the samples are not normally distributed? Use the nonparametric Kruskal-Wallis test, similar to the Wilcoxon rank sum test.
- ▶ Compute ranks of Y observations. Compare the mean rank of the groups.
- ▶ Use PROC NPAR1WAY with the WILCOXON option.

The Kruskal-Wallis Test

- ▶ What if the samples are not normally distributed? Use the nonparametric Kruskal-Wallis test, similar to the Wilcoxon rank sum test.
- ▶ Compute ranks of Y observations. Compare the mean rank of the groups.
- ▶ Use PROC NPAR1WAY with the WILCOXON option.

Example: Kruskal-Wallis Test

- ▶ Compare wear of different brands.

- ▶ Code:

```
proc npar1way data=samp.veneer wilcoxon;  
  class brand;  
  var wear;
```

Example: Kruskal-Wallis Test

- ▶ Compare wear of different brands.
- ▶ Code:

```
proc npar1way data=samp.veneer wilcoxon;  
  class brand;  
  var wear;
```

Dongfeng Li

Two-Way ANOVA

◀ ◻ ▶ ◀ ◼ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

Two-Way ANOVA—Main Effects

- ▶ Consider the simple case of balanced complete design of two factors, with repetition.
- ▶ Main effects model(additive model)

$$y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$$
$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 1)$$

- ▶ y_{ijk} is the response of factor A at level i , factor B at level j , repetition k .
- ▶ μ is the over all average.
- ▶ α_i is the **main effect** of factor A at level i . $\sum_i \alpha_i = 0$.
- ▶ β_j is the main effect of factor B at level j . $\sum_j \beta_j = 0$.
- ▶ e_{ijk} i.i.d. $N(0, \sigma^2)$.

Two-Way ANOVA—Main Effects

- ▶ Consider the simple case of balanced complete design of two factors, with repetition.
- ▶ Main effects model(additive model)

$$y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$$

$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 1)$$

- ▶ y_{ijk} is the response of factor A at level i , factor B at level j , repetition k .
- ▶ μ is the over all average.
- ▶ α_i is the **main effect** of factor A at level i . $\sum_i \alpha_i = 0$.
- ▶ β_j is the main effect of factor B at level j . $\sum_j \beta_j = 0$.
- ▶ e_{ijk} i.i.d. $N(0, \sigma^2)$.

Two-Way ANOVA—Main Effects

- ▶ Consider the simple case of balanced complete design of two factors, with repetition.
- ▶ Main effects model(additive model)

$$y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$$

$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 1)$$

- ▶ y_{ijk} is the response of factor A at level i , factor B at level j , repetition k .
- ▶ μ is the over all average.
- ▶ α_i is the **main effect** of factor A at level i . $\sum_i \alpha_i = 0$.
- ▶ β_j is the main effect of factor B at level j . $\sum_j \beta_j = 0$.
- ▶ e_{ijk} i.i.d. $N(0, \sigma^2)$.

Two-Way ANOVA—Main Effects

- ▶ Consider the simple case of balanced complete design of two factors, with repetition.
- ▶ Main effects model(additive model)

$$y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$$

$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 1)$$

- ▶ y_{ijk} is the response of factor A at level i , factor B at level j , repetition k .
- ▶ μ is the over all average.
- ▶ α_i is the **main effect** of factor A at level i . $\sum_i \alpha_i = 0$.
- ▶ β_j is the main effect of factor B at level j . $\sum_j \beta_j = 0$.
- ▶ e_{ijk} i.i.d. $N(0, \sigma^2)$.

Two-Way ANOVA—Main Effects

- ▶ Consider the simple case of balanced complete design of two factors, with repetition.
- ▶ Main effects model(additive model)

$$y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$$

$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 1)$$

- ▶ y_{ijk} is the response of factor A at level i , factor B at level j , repetition k .
- ▶ μ is the over all average.
- ▶ α_i is the **main effect** of factor A at level i . $\sum_i \alpha_i = 0$.
- ▶ β_j is the main effect of factor B at level j . $\sum_j \beta_j = 0$.
- ▶ e_{ijk} i.i.d. $N(0, \sigma^2)$.

Two-Way ANOVA—Main Effects

- ▶ Consider the simple case of balanced complete design of two factors, with repetition.
- ▶ Main effects model(additive model)

$$y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$$

$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 1)$$

- ▶ y_{ijk} is the response of factor A at level i , factor B at level j , repetition k .
- ▶ μ is the over all average.
- ▶ α_i is the **main effect** of factor A at level i . $\sum_i \alpha_i = 0$.
- ▶ β_j is the main effect of factor B at level j . $\sum_j \beta_j = 0$.
- ▶ e_{ijk} i.i.d. $N(0, \sigma^2)$.

Two-Way ANOVA—Main Effects

- ▶ Consider the simple case of balanced complete design of two factors, with repetition.
- ▶ Main effects model(additive model)

$$y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$$

$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 1)$$

- ▶ y_{ijk} is the response of factor A at level i , factor B at level j , repetition k .
- ▶ μ is the over all average.
- ▶ α_i is the **main effect** of factor A at level i . $\sum_i \alpha_i = 0$.
- ▶ β_j is the main effect of factor B at level j . $\sum_j \beta_j = 0$.
- ▶ e_{ijk} i.i.d. $N(0, \sigma^2)$.

Two-Way ANOVA—Interactions

- ▶ Interaction effect model

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{i,j} + e_{ijk}$$
$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 2)$$

- ▶ $\gamma_{i,j}$ is called an **interaction effect**.
 $\sum_i \gamma_{i,j} = 0, \sum_j \gamma_{i,j} = 0.$
- ▶ Use PROC ANOVA or PROC GLM. For unbalanced design, use PROC GLM.

Two-Way ANOVA—Interactions

- ▶ Interaction effect model

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{i,j} + e_{ijk}$$
$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 2)$$

- ▶ $\gamma_{i,j}$ is called an **interaction effect**.
 $\sum_i \gamma_{i,j} = 0, \sum_j \gamma_{i,j} = 0.$
- ▶ Use PROC ANOVA or PROC GLM. For unbalanced design, use PROC GLM.

Two-Way ANOVA—Interactions

- ▶ Interaction effect model

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{i,j} + e_{ijk}$$
$$i = 1, \dots, n, j = 1, \dots, m, k = 1, \dots, r (r \geq 2)$$

- ▶ $\gamma_{i,j}$ is called an **interaction effect**.
 $\sum_i \gamma_{i,j} = 0, \sum_j \gamma_{i,j} = 0.$
- ▶ Use PROC ANOVA or PROC GLM. For unbalanced design, use PROC GLM.

Example of Interactions

- ▶ Example with only main effects: $n = m = 2$, $r \geq 2$, $\mu = 10$, $\alpha_1 = 4$, $\alpha_2 = -4$, $\beta_1 = 2$, $\beta_2 = -2$. If interaction effect does not exist ($\gamma_{i,j} \equiv 0$), then

$$Ey_{11k} = 10 + 4 + 2 = 16$$

$$Ey_{12k} = 10 + 4 - 2 = 12$$

$$Ey_{21k} = 10 - 4 + 2 = 8$$

$$Ey_{22k} = 10 - 4 - 2 = 4$$

- ▶ Example with interactions: Suppose $\gamma_{1,1} = 1$, then $\gamma_{1,2} = -1$, $\gamma_{2,1} = -1$, $\gamma_{2,2} = 1$, so

$$Ey_{11k} = 16 + 1 = 17$$

$$Ey_{12k} = 12 - 1 = 11$$

$$Ey_{21k} = 8 - 1 = 7$$

$$Ey_{22k} = 4 + 1 = 5$$

the interaction increases the expectation when A and B are both 1 or both 2, decreases the expectation other wise.

Example of Interactions

- ▶ Example with only main effects: $n = m = 2$, $r \geq 2$, $\mu = 10$, $\alpha_1 = 4$, $\alpha_2 = -4$, $\beta_1 = 2$, $\beta_2 = -2$. If interaction effect does not exist ($\gamma_{i,j} \equiv 0$), then

$$Ey_{11k} = 10 + 4 + 2 = 16$$

$$Ey_{12k} = 10 + 4 - 2 = 12$$

$$Ey_{21k} = 10 - 4 + 2 = 8$$

$$Ey_{22k} = 10 - 4 - 2 = 4$$

- ▶ Example with interactions: Suppose $\gamma_{1,1} = 1$, then $\gamma_{1,2} = -1$, $\gamma_{2,1} = -1$, $\gamma_{2,2} = 1$, so

$$Ey_{11k} = 16 + 1 = 17$$

$$Ey_{12k} = 12 - 1 = 11$$

$$Ey_{21k} = 8 - 1 = 7$$

$$Ey_{22k} = 4 + 1 = 5$$

the interaction increases the expectation when A and B are both 1 or both 2, decreases the expectation other wise.

Two-way ANOVA Example

- ▶ To study the effects of different production factors on the strength of some rubber product, consider 3 levels of factor A, 4 levels of factor B, complete experiment with $3 \times 4 = 12$ combinations, each repeated 2 times, so we have $n = 24$ experiments.
- ▶ Data:

```
data rubber;  
  do a=1 to 3; do b=1 to 4; do r=1 to 2;  
    input stren @@;  
    output;  
  end; end; end;  
  cards;  
31 33 34 36 35 36 39 38  
33 34 36 37 37 39 38 41  
35 37 37 38 39 40 42 44  
;  
run;
```


Two-way ANOVA Example

- ▶ To study the effects of different production factors on the strength of some rubber product, consider 3 levels of factor A, 4 levels of factor B, complete experiment with $3 \times 4 = 12$ combinations, each repeated 2 times, so we have $n = 24$ experiments.
- ▶ Data:

```
data rubber;  
  do a=1 to 3; do b=1 to 4; do r=1 to 2;  
    input stren @@;  
    output;  
  end; end; end;  
  cards;  
31 33 34 36 35 36 39 38  
33 34 36 37 37 39 38 41  
35 37 37 38 39 40 42 44  
;  
run;
```

► Interaction model:

```
proc anova data=rubber;  
  class a b;  
  model  stren = a b a*b;  
run;
```

► Main effects(additive) model:

```
proc anova data=rubber;  
  class a b;  
  model  stren = a b;  
run;
```

► Interaction model:

```
proc anova data=rubber;  
  class a b;  
  model  stren = a b a*b;  
run;
```

► Main effects(additive) model:

```
proc anova data=rubber;  
  class a b;  
  model  stren = a b;  
run;
```

Result

Class Level Information		
Class	Levels	Values
a	3	1 2 3
b	4	1 2 3 4

Number of Observations Read	24
Number of Observations Used	24

Some Statistics
Background

Descriptive
statistics: concepts
and programs

Comparing the
Mean

Analysis of
Variance

One-Way ANOVA

Two-Way ANOVA

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	11	193.4583333	17.5871212	12.06	<.0001
Error	12	17.5000000	1.4583333		
Corrected Total	23	210.9583333			

R-Square	Coeff Var	Root MSE	stren Mean
0.917045	3.260152	1.207615	37.04167

Source	DF	Anova SS	Mean Square	F Value	Pr > F
a	2	56.5833333	28.2916667	19.40	0.0002
b	3	132.1250000	44.0416667	30.20	<.0001
a*b	6	4.7500000	0.7916667	0.54	0.7665

Additive model result

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	5	188.7083333	37.7416667	30.53	<.0001
Error	18	22.2500000	1.2361111		
Corrected Total	23	210.9583333			

R-Square	Coeff Var	Root MSE	stren Mean
0.894529	3.001499	1.111805	37.04167

Source	DF	Anova SS	Mean Square	F Value	Pr > F
a	2	56.5833333	28.2916667	22.89	<.0001
b	3	132.1250000	44.0416667	35.63	<.0001