

采菊东篱下，悠然见南山

机器学习的实与境

机器学习是经典而又现代的学科，它的发展过程交织着理想和务实。机器学习期待着机器具有人一样的自主学习能力，其名称本身就充满着理想主义色彩；许多相关想法蕴含着超前的理念。但作为一门学科，它务实地从数据、模型、算法、系统、实现等着手，结合数学和工程手段解决实际问题；它的许多模型或算法朴素、简洁而又优美。一个超前的理念往往能起牵引作用，使现在的理想在将来可能变为现实。

机器学习经历了以专家系统和句法模式识别为代表的符号主义，以神经网络为模型的连接主义等早期发展阶段。符号主义试图建立一套形式化体系来表示人类对目标对象的感知和认知，以形成知识库，从而自然地进行规则推理。连接主义受人脑神经元机制的启发，模仿提出神经网络模型以通向人工智能。但由于时代的局限，比如计算资源和算法性能所限，这些思路没有能有效地解决实际问题，使得理想未能成为现实。

上世纪中期，机器学习和统计学深度融合，置关注点于探索算法的构造和机理，以期从数据或经验中学习，改善学习系统预测或决策性能。这成就了机器学习其后十多年的黄金发展期。主要标志是学术界涌现出一批重要成果，比如，基于统计学习理论的支持向量机，随机森林和Boosting等集成分类方法，概率图模型，基于再生核理论的非线性数据分析与处理方法，非参数贝叶斯方法，稀疏学习模型及应用等等。这一期间，机器学习界没有过多渲染“智能”或者“认知”，而是务实地关注引入统计学等来建立学科的理论基础，面向数据分析与处理，以无监督学习和有监督学习为两大主要研究问题，提出和开发了一系列模型、方法和算法等，切实解决了工业界所面临的一些实际问题。

深度学习的概念在2006年被提出来，但直到2012年由于AlexNet模型在ImageNet数据集上分类性能突破性进展才被广泛认同。由此，开启了深度神经网络蓬勃发展，给学术界和工业界带来了深刻的变革和机遇。深度学习的成功不是由于脑科学或认知科学的进展，而是源自大数据的驱动，计算能力的极大提升，和一批面向机器学习的底层架构，比如PyTorch和TensorFlow等的先后被开发出来。但是脑科学和认知科学新进展势必给人工智能的探索以灵感和想象。

深度学习促进人工智能产生了突破性发展,特别是为计算机视觉、语音识别、自然语言处理以及游戏类等应用领域带来了颠覆性的进展;同时也给数学、统计学和理论计算机科学等带来了新的学科方向和发展机遇。深度学习激励了无监督生成模型、强化学习和因果学习等的发展,并催生了新的研究领域,比如自动机器学习、迁移学习、对抗学习、元学习等,给机器学习注入了新的活力和开辟了新的方向。

经过30年左右时间的发展,机器学习的触角已经渗入许多学科和领域,涉及统计学、概率论、控制论、信息论、运筹学、数值计算、并行计算、计算理论等多种学科。设计和分析机器学习算法将涉及可学习性、泛化性、稳定性、可扩展性和计算复杂性等一些基础问题。层次化、正则化、平均化和自适应化等是求解这些问题的有效技术途径。它们诠释和展现了数学与工程的美妙结合。特别地,卷积结构、参数共享、Dropout、Batch Normalization, 和Robbins-Monro算法等技术手段蕴含简约之美。这些让机器学习这个学科生机盎然、绚丽多姿。

王国维先生提出“境界”说来评价诗词,先生认为诗词的创作手法分为“造境”和“写境”,分别对应于“理想”和“写实”两派,然而这两者有时颇难区分开。造、写二境或偏于想象和虚拟,或侧重模仿和写实。“境界”说也适用于评价学术研究,且造境和写境是我们从事学术研究的要旨。

有境界则格局自高,自为名作。机器学习所以脱颖此是重要一端。机器学习的发展,比如深度学习、强化学习、元学习和自动机器学习等的提出和研究,是境与实之合,是造境和写境二者的相辅相成乃至高度融合。一个好的理论、模型、算法或思路无不蕴含造境和写境之意味。

所造之境,必合乎自然;所写之境,亦必邻于理想。并以造境和写境两者之完美融合使学术研究臻于胜境、化境,而入无我之境。

致谢:该文是为华为诺亚方舟实验室撰写的《机器学习的基础与前沿》白皮书所做的跋。文章得到了我的一位不愿署名老师的精心修改和润色,一些同事在读初稿时也提出了许多宝贵意见,在此表达我诚挚的感谢!

张志华

2019年11月8日于承泽园