



华章教育

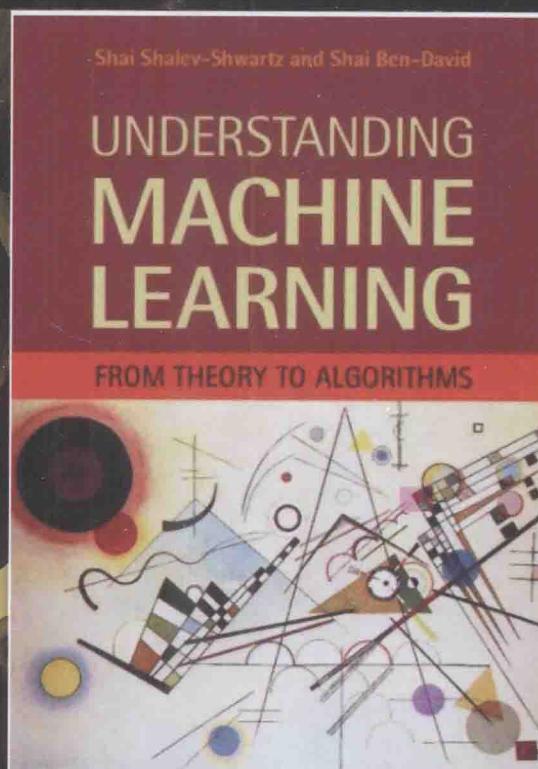
计 算 机 科 学 从 书

CAMBRIDGE

深入理解机器学习 从原理到算法

[以] 沙伊·沙莱夫-施瓦茨 (Shai Shalev-Shwartz) 著
[加] 沙伊·本-戴维 (Shai Ben-David) /
张文生 等译

Understanding Machine Learning
From Theory to Algorithms



机械工业出版社
China Machine Press

计 算 机 科 学 丛 书

深入理解机器学习

从原理到算法

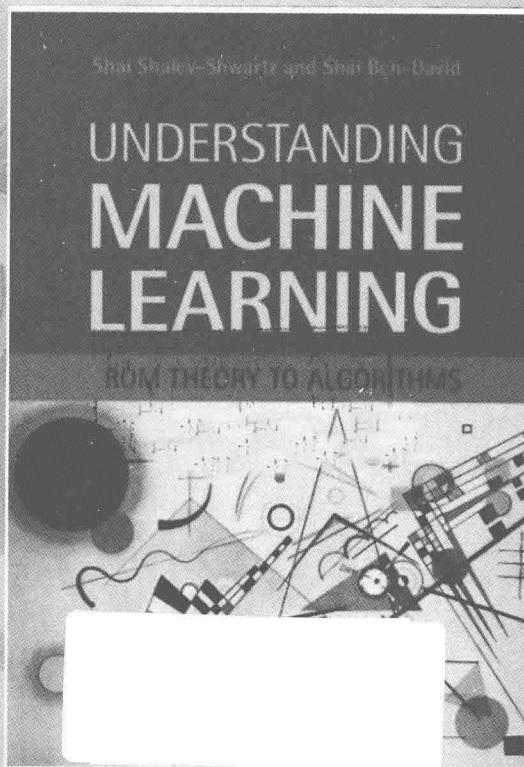
[以] 沙伊·沙莱夫-施瓦茨 (Shai Shalev-Shwartz) 著

[加] 沙伊·本-戴维 (Shai Ben-David)

张文生 等译

Understanding Machine Learning

From Theory to Algorithms



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

深入理解机器学习：从原理到算法 / (以) 沙伊·沙莱夫 - 施瓦茨 (Shai Shalev-Shwartz) 等著；张文生等译。—北京：机械工业出版社，2016.7 (2016.11 重印)
(计算机科学丛书)

书名原文：Understanding Machine Learning: From Theory to Algorithms

ISBN 978-7-111-54302-2

I. 深… II. ①沙… ②张… III. 机器学习 IV. TP181

中国版本图书馆 CIP 数据核字 (2016) 第 157549 号

本书版权登记号：图字：01-2016-3281

This is a Chinese Simplified edition of the following title published by Cambridge University Press: Shai Shalev-Shwartz and Shai Ben-David, Understanding Machine Learning: From Theory to Algorithms (ISBN 978-1-107-05713-5).

© Shai Shalev-Shwartz and Shai Ben-David 2014.

This Chinese Simplified edition for the People's Republic of China (excluding Hong Kong, Macau and Taiwan) is published by arrangement with the Press Syndicate of the University of Cambridge, Cambridge, United Kingdom.

© Cambridge University Press and China Machine Press in 2016.

This Chinese Simplified edition is authorized for sale in the People's Republic of China (excluding Hong Kong, Macau and Taiwan) only. Unauthorized export of this simplified Chinese is a violation of the Copyright Act. No part of this publication may be reproduced or distributed by any means, or stored in a database or retrieval system, without the prior written permission of Cambridge University Press and China Machine Press.

本书原版由剑桥大学出版社出版。

本书简体字中文版由剑桥大学出版社与机械工业出版社合作出版。未经出版者预先书面许可，不得以任何方式复制或抄袭本书的任何部分。

此版本仅限在中华人民共和国境内（不包括香港、澳门特别行政区及台湾地区）销售。

本书涵盖了机器学习领域中的严谨理论和实用方法，讨论了学习的计算复杂度、凸性和稳定性、PAC-贝叶斯方法、压缩界等概念，并介绍了一些重要的算法范式，包括随机梯度下降、神经元网络以及结构化输出。

全书讲解全面透彻，适合有一定基础的高年级本科生和研究生学习，也适合作为 IT 行业从事数据分析和挖掘的专业人员以及研究人员参考阅读。

出版发行：机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码：100037)

责任编辑：和 静

责任校对：董纪丽

印 刷：三河市宏图印务有限公司

版 次：2016 年 11 月第 1 版第 2 次印刷

开 本：185mm×260mm 1/16

印 张：20.25

书 号：ISBN 978-7-111-54302-2

定 价：79.00 元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

客服热线：(010) 88378991 88361066

投稿热线：(010) 88379604

购书热线：(010) 68326294 88379649 68995259

读者信箱：hzsj@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问：北京大成律师事务所 韩光 / 邹晓东

文艺复兴以来，源远流长的科学精神和逐步形成的学术规范，使西方国家在自然科学的各个领域取得了垄断性的优势；也正是这样的优势，使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中，美国的产业界与教育界越来越紧密地结合，计算机学科中的许多泰山北斗同时身处科研和教学的最前线，由此而产生的经典科学著作，不仅擘划了研究的范畴，还揭示了学术的源变，既遵循学术规范，又自有学者个性，其价值并不会因年月的流逝而减退。

近年，在全球信息化大潮的推动下，我国的计算机产业发展迅猛，对专业人才的需求日益迫切。这对计算机教育界和出版界都既是机遇，也是挑战；而专业教材的建设在教育战略上显得举足轻重。在我国信息技术发展时间较短的现状下，美国等发达国家在其计算机科学发展的几十年间积淀和发展的经典教材仍有许多值得借鉴之处。因此，引进一批国外优秀计算机教材将对我国计算机教育事业的发展起到积极的推动作用，也是与世界接轨、建设真正的世界一流大学的必由之路。

机械工业出版社华章公司较早意识到“出版要为教育服务”。自1998年开始，我们就将工作重点放在了遴选、移译国外优秀教材上。经过多年的不懈努力，我们与 Pearson, McGraw-Hill, Elsevier, MIT, John Wiley & Sons, Cengage 等世界著名出版公司建立了良好的合作关系，从他们现有的数百种教材中甄选出 Andrew S. Tanenbaum, Bjarne Stroustrup, Brian W. Kernighan, Dennis Ritchie, Jim Gray, Alfred V. Aho, John E. Hopcroft, Jeffrey D. Ullman, Abraham Silberschatz, William Stallings, Donald E. Knuth, John L. Hennessy, Larry L. Peterson 等大师名家的一批经典作品，以“计算机科学丛书”为总称出版，供读者学习、研究及珍藏。大理石纹理的封面，也正体现了这套丛书的品位和格调。

“计算机科学丛书”的出版工作得到了国内外学者的鼎力相助，国内的专家不仅提供了中肯的选题指导，还不辞劳苦地担任了翻译和审校的工作；而原书的作者也相当关注其作品在中国的传播，有的还专门为本书的中译本作序。迄今，“计算机科学丛书”已经出版了近两百个品种，这些书籍在读者中树立了良好的口碑，并被许多高校采用为正式教材和参考书籍。其影印版“经典原版书库”作为姊妹篇也被越来越多实施双语教学的学校所采用。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑，这些因素使我们的图书有了质量的保证。随着计算机科学与技术专业学科建设的不断完善和教材改革的逐渐深化，教育界对国外计算机教材的需求和应用都将步入一个新的阶段，我们的目标是尽善尽美，而反馈的意见正是我们达到这一终极目标的重要帮助。华章公司欢迎老师和读者对我们的工作提出建议或给予指正，我们的联系方法如下：

华章网站：www.hzbook.com

电子邮件：hzjsj@hzbook.com

联系电话：(010)88379604

联系地址：北京市西城区百万庄南街1号

邮政编码：100037



华章教育

华章科技图书出版中心

译者序 |

Understanding Machine Learning: From Theory to Algorithms

以色列希伯来大学副教授 Shai Shalev-Shwartz 和加拿大滑铁卢大学教授 Shai Ben-David 的专著《Understanding Machine Learning: From Theory to Algorithms》是机器学习领域一部具有里程碑意义的著作。

近几年，机器学习是人工智能研究领域中最活跃的分支之一，已成为信息科学领域解决实际问题的重要方法，它的应用已遍及人工智能的各个应用领域。机器学习又是一个多学科的交叉领域，涉及数学、自动化、计算机科学、应用心理学、生物学和神经生理学等。这种学科交叉融合带来的良性互动，无疑促进了包括机器学习在内的诸学科的发展与繁荣。

本书内容十分丰富，作者以前所未有的广度和深度，介绍了目前机器学习中重要的理论和关键的算法。本书没有陷入“科普”式的堆砌材料的写作方式，由于作者是该领域的权威专家，因此在介绍各种理论和算法时，时刻不忘将不同理论、算法的对比与作者自身的研究成果传授给读者，使读者不至于对如此丰富的理论和算法无所适从。另外，特别值得指出的是，本书第一部分非常有特色，也是非常重要的一部分。这部分内容从更高的观点和更深的层次探讨机器学习的许多理论基础，引入对指导理论研究和实际应用都至关重要的概率近似正确(Probably Approximately Correct, PAC)学习理论。该理论旨在回答由机器学习得到的结果到底有多高的可信度与推广能力，从某种意义上来说，只有懂得了该部分，才可能透彻地理解和更好地运用其他章节的内容。国内关于 PAC 学习的资料非常少，在翻译过程中团队成员碰到了极大的困难，我们人工智能与机器学习研究团队为此进行了多方论证并多次召开专题讨论会。

本书主要面向人工智能、机器学习、模式识别、数据挖掘、计算机应用、生物信息学、数学和统计学等领域的研究生和相关领域的科技人员。翻译出版中译本的目的，是希望能为国内广大从事相关研究的学者和研究生提供一本全面、系统、权威的教科书和参考书。如果能做到这一点，译者将感到十分欣慰。

必须说明的是，本书的翻译是中国科学院自动化研究所人工智能与机器学习研究团队集体努力的结果，团队的成员杨雪冰、匡秋明、蒋晓娟、薛伟、魏波、李思园、张似衡、曾凡霞、于廷照、王鑫、李涛、杨叶辉、胡文锐、张志忠、唐永强、陈东杰、何泽文、张英华、李悟、李硕等参与了本书的翻译工作，李思园老师参与了全书的审校与修正。感谢机械工业出版社华章分社的大力协助，倘若没有他们的热情支持，本书的中译版难以如此迅速地与大家见面。另外，本书的翻译得到了国家自然科学基金委重点项目和面上项目(61472423、U1135005、61432008、61532006、61305018、61402481 等)的资助，特此感谢。

在翻译过程中，我们力求准确地反映原著内容，同时保留原著的风格。但由于译者水平有限，书中难免有不妥之处，恳请读者批评指正。

最后，谨把本书的中译版献给我的博士生导师王珏研究员！王珏老师生前对机器学习理论、算法和应用非常关注，对于 PAC 可学习理论也有着独到而深刻的理解，他启发并引领了我们研究团队对机器学习理论和算法的研究工作，使我们终身受益。

中国科学院自动化研究所

张文生

2016 年 4 月于北京

前 言 |

Understanding Machine Learning: From Theory to Algorithms

“机器学习”旨在从数据中自动识别有意义的模式。过去几十年中，机器学习成为一项常用工具，几乎所有需要从大量数据集合中提取信息的任务都在使用它。我们身边的许多技术都以机器学习为基础：搜索引擎学习在带给我们最佳的搜索结果的同时，植入可以盈利的广告；屏蔽软件学习过滤垃圾邮件；用于保护信用卡业务的软件学习识别欺诈。数码相机学习人脸识别，智能电话上的个人智能助手学习识别语音命令。汽车配备了用机器学习算法搭建的交通事故预警系统。同时机器学习还被广泛应用于各个科学领域，例如生物信息学、医药以及天文学等。

这些应用领域的一个共同特点在于，与相对传统的计算机应用相比，所需识别的模式更复杂。在这些情景中，对于任务应该如何执行，人类程序员无法提供明确的、细节优化的具体指令。以智能生物为例，我们人类的许多技能都是通过从经验中学习而取得并逐步提高的（而非遵从别人给我们的具体指令）。机器学习工具关注的正是赋予程序“学习”和适应不同情况的能力。

本书的第一个目标是，提供一个准确而简明易懂的导论，介绍机器学习的基本概念：什么是学习？机器怎样学习？学习某概念时，如何量化所需资源？学习始终都是可能的吗？我们如何知道学习过程是成功或失败？

本书的第二个目标是，为机器学习提供几个关键的算法。我们提供的算法，一方面已经成功投入实际应用，另一方面广泛地考虑到不同的学习技术。此外，我们特别将注意力放到了大规模学习（即俗称的“大数据”）上，因为近几年来，世界越来越“数字化”，需要学习的数据总量也在急剧增加。所以在许多应用中，数据量是充足的，而计算时间是主要瓶颈。因此，学习某一概念时，我们会明确量化数据量和计算时间这两个数值。

本书分为四部分。第一部分对于“学习”的基础性问题给出初步而准确的定义。我们会介绍 Valiant 提出的“概率近似正确(PAC)”可学习模型的通用形式，它将是对“何为学习”这一问题的第一个有力回答。我们还会介绍“经验风险最小化(ERM)”“结构风险最小化(SRM)”和“最小描述长度(MDL)”这几个学习规则，展现“机器是如何学习的”。我们量化使用 ERM、SRM 和 MDL 规则学习时所需的数据总量，并用“没有免费的午餐”定理说明，什么情况下学习可能会失败。此外，我们还探讨了学习需要多少计算时间。本书第二部分介绍多种算法。对于一些算法，我们先说明其主要学习原则，再介绍该算法是如何依据其原则运作的。前两部分将重点放在 PAC 模型上，第三部分将范围扩展到更广、更丰富的学习模型。最后，第四部分讨论最前沿的理论。

我们尽量让本书能够自成一体，不过我们假设读者熟悉概率论、线性代数、数学分析和算法设计的基本概念。前三部分为计算机科学、工程学、数学和统计学研究生一年级学生设计，具有相关背景的本科生也可以使用。高级章节适用于想要对理论有更深入理解的研究者。

致 谢

Understanding Machine Learning: From Theory to Algorithms

本书以“机器学习入门”课程为蓝本，这门课程由 Shai Shalev-Shwartz 和 Shai Ben-David 分别在希伯来大学和滑铁卢大学讲授。本书的初稿由 Shai Shalev-Shwartz 在 2010 至 2013 年间在希伯来大学所开课程的教案整理而成。感谢 2010 年的助教 Ohad Shamir 和 2011 至 2013 年的助教 Alon Gonen 的帮助，他们为课堂准备了一些教案以及许多课后练习。特别感谢 Alon 在全书编写过程中所做出的贡献，此外他还撰写了一册习题答案。

我们由衷地感谢 Dana Rubinstein 的辛勤工作。Dana 从科学的角度校对了书稿，对原稿进行了编辑，将它从章节教案的形式转换成连贯流畅的文本。

特别感谢 Amit Daniely，他仔细阅读了本书的高级部分，并撰写了多分类可学习性的章节。我们还要感谢耶路撒冷的一个阅读俱乐部的成员们，他们认真阅读了原稿的每一页，并提出了建设性的意见。他们是：Maya Alroy, Yossi Arjevani, Aharon Birnbaum, Alon Cohen, Alon Gonen, Roi Livni, Ofer Meshi, Dan Rosenbaum, Dana Rubinstein, Shahar Somin, Alon Vinnikov 和 Yoav Wald。还要感谢 Gal Elidan, Amir Globerson, Nika Haghtalab, Shie Mannor, Amnon Shashua, Nati Srebro 和 Ruth Urner 参与的有益讨论。

目 录 |

Understanding Machine Learning: From Theory to Algorithms

出版者的话	4.3 小结	26
译者序	4.4 文献评注	27
前言	4.5 练习	27
致谢		
第1章 引论	第5章 偏差与复杂性权衡	28
1.1 什么是学习	5.1 “没有免费的午餐”定理	28
1.2 什么时候需要机器学习	5.2 误差分解	31
1.3 学习的种类	5.3 小结	31
1.4 与其他领域的关系	5.4 文献评注	32
1.5 如何阅读本书	5.5 练习	32
1.6 符号		
第一部分 理论基础	第6章 VC维	33
第2章 简易入门	6.1 无限的类也可学习	33
2.1 一般模型——统计学习理论	6.2 VC维概述	34
框架	6.3 实例	35
2.2 经验风险最小化	6.3.1 阈值函数	35
2.3 考虑归纳偏置的经验风险	6.3.2 区间	35
最小化	6.3.3 平行于轴的矩形	35
2.4 练习	6.3.4 有限类	36
	6.3.5 VC维与参数个数	36
第3章 一般学习模型	6.4 PAC学习的基本定理	36
3.1 PAC学习理论	6.5 定理6.7的证明	37
3.2 更常见的学习模型	6.5.1 Sauer引理及生长函数	37
3.2.1 放宽可实现假设——	6.5.2 有小的有效规模的类的一致收敛性	39
不可知PAC学习	6.6 小结	40
3.2.2 学习问题建模	6.7 文献评注	41
3.3 小结	6.8 练习	41
3.4 文献评注		
3.5 练习		
第4章 学习过程的一致收敛性	第7章 不一致可学习	44
4.1 一致收敛是可学习的充分条件	7.1 不一致可学习概述	44
4.2 有限类是不可知PAC	7.2 结构风险最小化	46
可学习的	7.3 最小描述长度和奥卡姆剃刀	48
	7.4 可学习的其他概念——一致收敛性	50
	7.5 探讨不同的可学习概念	51

7.6 小结	53	第 11 章 模型选择与验证	85
7.7 文献评注	53	11.1 用结构风险最小化进行模型 选择	85
7.8 练习	54	11.2 验证法	86
第 8 章 学习的运行时间	56	11.2.1 留出的样本集	86
8.1 机器学习的计算复杂度	56	11.2.2 模型选择的验证法	87
8.2 ERM 规则的实现	58	11.2.3 模型选择曲线	88
8.2.1 有限集	58	11.2.4 k 折交叉验证	88
8.2.2 轴对称矩形	59	11.2.5 训练-验证-测试拆分	89
8.2.3 布尔合取式	59	11.3 如果学习失败了应该做什么	89
8.2.4 学习三项析取范式	60	11.4 小结	92
8.3 高效学习，而不通过合适的 ERM	60	11.5 练习	92
8.4 学习的难度*	61	第 12 章 凸学习问题	93
8.5 小结	62	12.1 凸性、利普希茨性和光滑性	93
8.6 文献评注	62	12.1.1 凸性	93
8.7 练习	62	12.1.2 利普希茨性	96
第二部分 从理论到算法		12.1.3 光滑性	97
第 9 章 线性预测	66	12.2 凸学习问题概述	98
9.1 半空间	66	12.2.1 凸学习问题的可学习性	99
9.1.1 半空间类线性规划	67	12.2.2 凸利普希茨/光滑有界 学习问题	100
9.1.2 半空间感知器	68	12.3 替代损失函数	101
9.1.3 半空间的 VC 维	69	12.4 小结	102
9.2 线性回归	70	12.5 文献评注	102
9.2.1 最小平方	70	12.6 练习	102
9.2.2 多项式线性回归	71	第 13 章 正则化和稳定性	104
9.3 逻辑斯谛回归	72	13.1 正则损失最小化	104
9.4 小结	73	13.2 稳定规则不会过拟合	105
9.5 文献评注	73	13.3 Tikhonov 正则化作为 稳定剂	106
9.6 练习	73	13.3.1 利普希茨损失	108
第 10 章 boosting	75	13.3.2 光滑和非负损失	108
10.1 弱可学习	75	13.4 控制适合与稳定性的权衡	109
10.2 AdaBoost	78	13.5 小结	111
10.3 基础假设类的线性组合	80	13.6 文献评注	111
10.4 AdaBoost 用于人脸识别	82	13.7 练习	111
10.5 小结	83	第 14 章 随机梯度下降	114
10.6 文献评注	83	14.1 梯度下降法	114
10.7 练习	84		

第 14 章 梯度下降法	111
14.1 梯度下降法 111	
14.1.1 梯度下降法 111	
14.1.2 梯度下降法的收敛性 112	
14.2 次梯度 116	
14.2.1 计算次梯度 117	
14.2.2 利普希茨函数的次梯度 118	
14.2.3 次梯度下降 118	
14.3 随机梯度下降 118	
14.4 SGD 的变型 120	
14.4.1 增加一个投影步 120	
14.4.2 变步长 121	
14.4.3 其他平均技巧 121	
14.4.4 强凸函数* 121	
14.5 用 SGD 进行学习 123	
14.5.1 SGD 求解风险极小化 123	
14.5.2 SGD 求解凸光滑学习问题的分析 124	
14.5.3 SGD 求解正则化损失极小化 125	
14.6 小结 125	
14.7 文献评注 125	
14.8 练习 126	
第 15 章 支持向量机 127	
15.1 间隔与硬 SVM 127	
15.1.1 齐次情况 129	
15.1.2 硬 SVM 的样本复杂度 129	
15.2 软 SVM 与范数正则化 130	
15.2.1 软 SVM 的样本复杂度 131	
15.2.2 间隔、基于范数的界与维度 131	
15.2.3 斜坡损失* 132	
15.3 最优化条件与“支持向量”* 133	
15.4 对偶* 133	
15.5 用随机梯度下降法实现软 SVM 134	
15.6 小结 135	
15.7 文献评注 135	
15.8 练习 135	
第 16 章 核方法 136	
16.1 特征空间映射 136	
16.2 核技巧 137	
16.2.1 核作为表达先验的一种形式 140	
16.2.2 核函数的特征* 141	
16.3 软 SVM 应用核方法 141	
16.4 小结 142	
16.5 文献评注 143	
16.6 练习 143	
第 17 章 多分类、排序与复杂预测问题 145	
17.1 一对多和一对一 145	
17.2 线性多分类预测 147	
17.2.1 如何构建 Ψ 147	
17.2.2 对损失敏感的分类 148	
17.2.3 经验风险最小化 149	
17.2.4 泛化合页损失 149	
17.2.5 多分类 SVM 和 SGD 150	
17.3 结构化输出预测 151	
17.4 排序 153	
17.5 二分排序以及多变量性能测量 157	
17.6 小结 160	
17.7 文献评注 160	
17.8 练习 161	
第 18 章 决策树 162	
18.1 采样复杂度 162	
18.2 决策树算法 163	
18.2.1 增益测量的实现方式 164	
18.2.2 剪枝 165	
18.2.3 实值特征基于阈值的拆分规则 165	
18.3 随机森林 165	
18.4 小结 166	
18.5 文献评注 166	
18.6 练习 166	

第 19 章 最近邻	167	22.3.3 非归一化的谱聚类	207
19.1 k 近邻法	167	22.4 信息瓶颈*	208
19.2 分析	168	22.5 聚类的进阶观点	208
19.2.1 1-NN 准则的泛化界	168	22.6 小结	209
19.2.2 “维数灾难”	170	22.7 文献评注	210
19.3 效率实施*	171	22.8 练习	210
19.4 小结	171		
19.5 文献评注	171		
19.6 练习	171		
第 20 章 神经元网络	174	第 23 章 维度约简	212
20.1 前馈神经网络	174	23.1 主成分分析	212
20.2 神经网络学习	175	23.1.1 当 $d \gg m$ 时一种更加有效的求解方法	214
20.3 神经网络的表达力	176	23.1.2 应用与说明	214
20.4 神经网络样本复杂度	178	23.2 随机投影	216
20.5 学习神经网络的运行时	179	23.3 压缩感知	217
20.6 SGD 和反向传播	179	23.4 PCA 还是压缩感知	223
20.7 小结	182	23.5 小结	223
20.8 文献评注	183	23.6 文献评注	223
20.9 练习	183	23.7 练习	223
第三部分 其他学习模型		第 24 章 生成模型	226
第 21 章 在线学习	186	24.1 极大似然估计	226
21.1 可实现情况下的在线分类	186	24.1.1 连续随机变量的极大似然估计	227
21.2 不可实现情况下的在线识别	191	24.1.2 极大似然与经验风险最小化	228
21.3 在线凸优化	195	24.1.3 泛化分析	228
21.4 在线感知器算法	197	24.2 朴素贝叶斯	229
21.5 小结	199	24.3 线性判别分析	230
21.6 文献评注	199	24.4 隐变量与 EM 算法	230
21.7 练习	199	24.4.1 EM 是交替最大化算法	232
第 22 章 聚类	201	24.4.2 混合高斯模型参数估计的 EM 算法	233
22.1 基于链接的聚类算法	203	24.5 贝叶斯推理	233
22.2 k 均值算法和其他代价最小聚类	203	24.6 小结	235
22.3 谱聚类	206	24.7 文献评注	235
22.3.1 图割	206	24.8 练习	235
22.3.2 图拉普拉斯与松弛图割算法	206		
第 25 章 特征选择与特征生成	237		
25.1 特征选择	237		
25.1.1 滤波器	238		

25.1.2 贪婪选择方法	239	第 29 章 多分类可学习性	271
25.1.3 稀疏诱导范数	241	29.1 纳塔拉詹维	271
25.2 特征操作和归一化	242	29.2 多分类基本定理	271
25.3 特征学习	244	29.3 计算纳塔拉詹维	272
25.4 小结	246	29.3.1 基于类的一对多	272
25.5 文献评注	246	29.3.2 一般的多分类到二分类 约简	273
25.6 练习	246	29.3.3 线性多分类预测器	273
第四部分 高级理论			
第 26 章 拉德马赫复杂度	250	29.4 好的与坏的 ERM	274
26.1 拉德马赫复杂度概述	250	29.5 文献评注	275
26.2 线性类的拉德马赫复杂度	255	29.6 练习	276
26.3 SVM 的泛化误差界	256	第 30 章 压缩界	277
26.4 低 ℓ_1 范数预测器的泛化 误差界	258	30.1 压缩界概述	277
26.5 文献评注	259	30.2 例子	278
第 27 章 覆盖数	260	30.2.1 平行于轴的矩形	278
27.1 覆盖	260	30.2.2 半空间	279
27.2 通过链式反应从覆盖到 拉德马赫复杂度	261	30.2.3 可分多项式	279
27.3 文献评注	262	30.2.4 间隔可分的情况	279
第 28 章 学习理论基本定理的 证明	263	30.3 文献评注	280
28.1 不可知情况的上界	263	第 31 章 PAC-贝叶斯	281
28.2 不可知情况的下界	264	31.1 PAC-贝叶斯界	281
28.2.1 证明 $m(\epsilon, \delta) \geqslant$ $0.5 \log(1/(4\delta))/\epsilon^2$	264	31.2 文献评注	282
28.2.2 证明 $m(\epsilon, 1/8) \geqslant$ $8d/\epsilon^2$	265	31.3 练习	282
28.3 可实现情况的上界	267	附录 A 技术性引理	284
		附录 B 测度集中度	287
		附录 C 线性代数	294
		参考文献	297
		索引	305

引 论

本书的主题是“自动学习”，后文中我们更经常称之为“机器学习”。机器学习的含义是，希望通过计算机编程，使它能够根据已有的输入数据进行学习。粗略地说，学习是一个将经验转化为专业技能或知识的过程。输入学习算法的是代表经验的训练数据，而输出的则是知识。这种知识通常以一种可以被其他计算机程序执行任务时所用的形式存在。为寻求这一概念的形式化数学解释，我们必须更明确地了解其中涉及的每个术语的准确含义：程序获取的训练数据是什么？学习过程是如何自动进行的？如何评价这一学习过程的成败(即学习程序输出结果的质量)？

1.1 什么是学习

我们首先来看几个存在于大自然的动物学习的例子。从这些熟悉的例子中可以看出，机器学习的一些基本问题也存在于自然界。

怯饵效应——老鼠学习躲避毒饵：当老鼠遇到有新颖外观或气味的食物时，它们首先会少量进食，随后的进食量将取决于事物本身的风味及其生理作用。如果产生不良反应，那么新的食物往往与这种不良后果相关联，随之，老鼠不再进食这种食物。很显然，这里有一个学习机制在起作用——动物通过经验来获取判断食物安全性的技能。如果对一种食物过去的经验是负标记的，那么动物会预测在未来遇到它时也会产生负面影响。

前文的示例解释了什么是学习成功，下面我们再举例说明什么是典型的机器学习任务。假设我们想对一台机器进行编程，使其学会如何过滤垃圾邮件。一个最简单的解决方案是仿照老鼠学习躲避毒饵的过程。机器只须记住所有以前被用户标记为垃圾的邮件。当一封新邮件到达时，机器将在先前垃圾邮件库中进行搜索。如果匹配其中之一，它会被丢弃。否则，它将被移动到用户的收件箱文件夹。

虽然上述“通过记忆进行学习”的方法时常是有用的，但是它缺乏一个学习系统的重要特性——标记未见邮件的能力。一个成功的学习器应该能够从个别例子进行泛化，这也称为归纳推理。在前面提到的“怯饵效应”例子中，老鼠遇到一种特定类型的食物后，它们会对新的、没见过的、有相似气味和口味的食物采取同样的态度。为了实现垃圾邮件过滤任务的泛化，学习器可以扫描以前见过的电子邮件，并提取那些垃圾邮件的指示性的词集；然后，当新电子邮件到达时，这台机器可以检查它是否含有可疑的单词，并相应地预测它的标签。这种系统应该有能力正确预测未见电子邮件的标签。

但是，归纳推理有可能推导出错误的结论。为了说明这一点，我们再来思考一个动物学习的例子。

鸽子迷信：心理学家 B. F. Skinner 进行过一项实验，他在笼子里放了一群饥饿的鸽子。笼子上附加了一个自动装置，不管鸽子当时处于什么行为状态，都会以固定的时间间隔为它们提供食物。饥饿的鸽子在笼子里走来走去，当食物第一次送达时，每只鸽子都在进行某项活动(啄食、转动头部等)。食物的到来强化了它们各自特定的行为，此后，每只鸟都倾向于花费更多的时间重复这种行为。接下来，随机的食物送达又增加了

每只鸟做出这种行为的机会。结果是，不管第一次食物送达时，每只鸟处于什么行为状态，这一连串的事件都增强了食物送达和这种行为之间的关联。进而，鸽子们也更勤奋地做出这种行为^②。

有用的学习机制与形成迷信的学习机制有何差别？这个问题对自动学习器的发展至关重要。尽管人类可以依靠常识来滤除随机无意义的学习结论，但是一旦我们将学习任务付之于一台机器，就必须提供定义明确、清晰的规则，来防止程序得出无意义或无用的结论。发展这些规则是机器学习理论的一个核心目标。

是什么使老鼠的学习比鸽子更成功？作为回答这个问题的第一步，我们仔细看一下老鼠在“怯饵效应”实验中的心理现象。

重新审视“怯饵效应”——老鼠未能获得食物与电击或声音与反胃之间的关联：老鼠的怯饵效应机制可能比你想象中的更复杂。Garcia 进行的实验(Garcia & Koelling 1996)表明，当进食后伴随的是不愉快的刺激时，比如说电击(不是反胃反应)，那么关联没有出现。即使将进食后电击的机制重复多次，老鼠仍然倾向于进食。同样，食物引起的反胃(口味或气味)与声音之间的关联实验也失败了。老鼠似乎有一些“内置的”先验知识，告诉它们，虽然食物和反胃存在因果相关，但是食物与电击或声音与反胃之间不太可能存在因果关系。

由此我们得出结论，怯饵效应和鸽子迷信的一个关键区别点是先验知识的引入使学习机制产生偏差，也称为“归纳偏置”。在实验中，鸽子愿意采取任何食物送达时发生的行为。然而，老鼠“知道”食物不能导致电击，也知道与食物同现的噪音不可能影响这种食物的营养价值。老鼠的学习过程偏向于发现某种模式，而忽略其他的关联。

事实证明，引入先验知识导致学习过程产生偏差，这对于学习算法的成功必不可少(正式陈述与证明参见第5章中的“没有免费的午餐”定理)。这种方法的发展，即能够表示领域知识，将其转化为一个学习偏置，并量化偏置对学习成功的影响，是机器学习理论的一个核心主题。粗略地讲，具有的先验知识(先验假设)越强，越容易从样本实例中进行学习。但是，先验假设越强，学习越不灵活——受先验假设限制。第5章将详细讨论这些问题。

1.2 什么时候需要机器学习

什么时候需要机器学习，而不是直接动手编程完成任务？在指定问题中，程序能否在“经验”的基础上自我学习和提高，有两方面的考量：问题本身的复杂性和对自适应性的需要。

1. 过于复杂的编程任务

- **动物/人可执行的任务：**虽然人类可以习惯性地执行很多任务，但是反思我们如何完成任务的内省机制还不够精细，无法从中提取一个定义良好的程序。汽车驾驶、语音识别和图像识别都属于此类任务。面对此类任务，只要接触到足够多的训练样本，目前最先进的机器学习程序，即能“从经验中学习”的程序，就可以达到比较满意的效果。
- **超出人类能力的任务：**受益于机器学习技术，另一大系列任务都涉及对庞大且复杂的数据集进行分析：天文数据，医疗档案转化为医学知识，气象预报，基因组数据

^② <http://psychclassics.yorku.ca/Skinner/Pigeon>。

分析，网络搜索引擎和电子商务。随着越来越多的数字数据的出现，显而易见的是，隐含在数据里的有意义、有价值的信息过于庞大复杂，超出了人类的理解能力。学习在大量复杂数据中发现有意义的模式是一个有前途的领域，无限内存容量加上不断提高的处理速度，更为这一领域开辟了新的视野。

2. 自适应性

编程的局限之一是其刻板性——一旦程序的编写与安装完成，它将保持不变。但是，任务会随着时间的推移而改变，用户也会出现变更。机器学习方法——其行为自适应输入数据的程序——为这个难题提供了一个解决方案。机器学习方法天生具备自适应于互动环境变化的性质。机器学习典型的成功应用有：能够适应不同用户的手写体识别，自动适应变化的垃圾邮件检测，以及语音识别。

1.3 学习的种类

学习是一个非常广泛的领域。因此，机器学习根据学习任务的不同分为不同的子类。这里给出一个粗略的分类，旨在对本书中属于机器学习广泛领域的那部分内容提供一些视角。

下面给出四种分类方式。

监督与无监督：学习涉及学习器与环境之间的互动，那么可以根据这种互动的性质划分学习任务。首先需要关注的是监督学习与无监督学习之间的区别。下面以垃圾邮件检测和异常检测为例说明。对于垃圾邮件检测任务，学习器的训练数据是带标签的邮件(是/否垃圾邮件)。在这种训练的基础上，学习器应该找出标记新电子邮件的规则。相反，对于异常检测任务，学习器的训练数据是大量没有标签的电子邮件，学习器的任务是检测出“不寻常”的消息。

抽象一点来讲，如果我们把学习看做一个“利用经验获取技能”的过程，那么监督学习正是这样的一种场景：经验是包含显著信息(是/否垃圾邮件)的训练数据，“测试数据”缺少这些显著信息，但可从学到的“技能”中获取。此种情况下，获得的“技能”旨在预测测试数据的丢失信息，我们可以将环境看做通过提供额外信息(标签)来“监督”学习器的老师。然而，无监督学习的训练数据和测试数据之间没有区别。学习器处理输入数据的目标是提取概括信息(浓缩数据)。聚类(相似数据归为一类)是执行这样任务的一个典型例子。

还有一种中间情况，训练数据比测试数据包含更多的信息，也要求学习器预测更多信息。举个例子，当学习数值函数判断国际象棋游戏中白棋和黑棋谁更有利时，训练过程中提供给学习器的唯一信息是，谁在整个实际的棋牌类游戏中最终赢得那场比赛的标签。这种学习被称作“强化学习”。

主动学习器与被动学习器：学习可依据学习器扮演的角色不同分类为“主动”和“被动”学习器。主动学习器在训练时通过提问或实验的方式与环境交互，而被动学习器只观察环境(老师)所提供的信息而不影响或引导它。请注意，垃圾邮件过滤任务通常是被动学习——等待用户标记电子邮件。我们可以设想，在主动学习中，要求用户来标记学习器挑选的电子邮件，以提高学习器对“垃圾邮件是什么”的理解。

老师的帮助：人类的学习过程中(在家的幼儿或在校的学生)往往会有一个人良师，他向学习者传输最有用的信息以实现学习目标。相比之下，科学家研究自然时，环境起到了老师的作用。环境的作用是消极的——苹果坠落、星星闪烁、雨点下落从不考虑学习者的需

求。在对这种学习情境建模时，我们假定训练数据（学习者的经验）是由随机过程产生的，这是统计机器学习的一个基本构成单元。此外，学习也发生在学习者的输入是由对立“老师”提供的。垃圾邮件过滤任务（如果垃圾邮件制作者尽力误导垃圾邮件过滤器设计者）和检测欺诈学习任务就是这种情况。当不存在更好的假设时，我们也会使用对立老师这一最坏方案。如果学习器能够从对立老师中学习，那么遇到任何老师都可以成功。

在线与批量：在线响应还是处理大量数据后才获得技能，是对学习器的另一种分类方式。举个例子，股票经纪人必须基于当时的经验信息做出日常决策。随着时间推移，他或许会成为专家，但是也会犯错并付出高昂的代价。相比之下，在大量的数据挖掘任务中，学习器，也就是数据挖掘器，往往是在处理大量训练数据之后才输出结论。

在本书中，我们只选取一部分机器学习技术进行讨论。重点是被动的、有监督的、统计批量学习（例如，基于大量独立收集的且带有病人最终结果标记的诊断记录，学习如何预测病人结果）。另外，本书也对在线学习和无监督批量学习（尤其是聚类）做了介绍。

5

1.4 与其他领域的关系

作为一门交叉学科，机器学习与统计学、信息论、博弈论、最优化等众多数学分支有着共同点。我们的最终目标是在计算机上编写程序，所以机器学习自然也是计算机科学的一个分支。在某种意义上，机器学习可以视为人工智能的一个分支，毕竟，要将经验转变成专业知识或从复杂感知数据中发现有意义的模式的能力是人类和动物智能的基石。但是，应该注意的是，与传统人工智能不同，机器学习并不是试图自动模仿智能行为，而是利用计算机的优势和特长与人类的智慧相得益彰。机器学习常用于执行远远超出人类能力的任务。例如，机器学习程序通过浏览和处理大型数据，能够检测到超出人类感知范围的模式。

机器学习（的经验）训练涉及的数据往往是随机生成的。机器学习的任务就是处理这些背景下的随机生成样本，得出与背景相符的结论。这样的描述强调了机器学习与统计学的密切关系。两个学科之间确实有很多共同点，尤其表现在目标和技术方面。但是，两者之间仍然存在显著的差别：如果一个医生提出吸烟与心脏病之间存在关联这一假设，这时应该由统计学家去查看病人样本并检验假设的正确性（这是常见的统计任务——假设检验）。相比之下，机器学习的任务是利用患者样本数据找出心脏病的原因。我们希望自动化技术能够发现被人类忽略的、有意义的模式（或假设）。

与传统统计学不同，算法在机器学习中（尤其在本书里）扮演了重要的角色。机器学习算法要靠计算机来执行，因此算法问题是关键。我们开发算法完成学习任务，同时关心算法的计算效率。两者的另外一个区别是，统计关心算法的渐近性（如随着样本量增长至无穷大，统计估计的收敛问题），机器学习理论侧重于有限样本。也就是说，给定有限可用样本，机器学习理论旨在分析学习器可达到的准确度。

6

机器学习与统计学之间还有很多差异，我们在此仅提到了少数。比如，在统计学中，常首先提出数据模型假设（生成数据呈正态分布或依赖函数为线性）；在机器学习中常考虑“非参数”背景，对数据分布的性质假设尽可能地少，学习算法自己找出最接近数据生成过程的模型。深入讨论这个问题需要更多的技术基础，详见第5章。

1.5 如何阅读本书

本书第一部分是机器学习的基本理论知识，从某种意义上讲，这是本书其余部分的基础。这部分应该作为机器学习理论入门课程的基础。