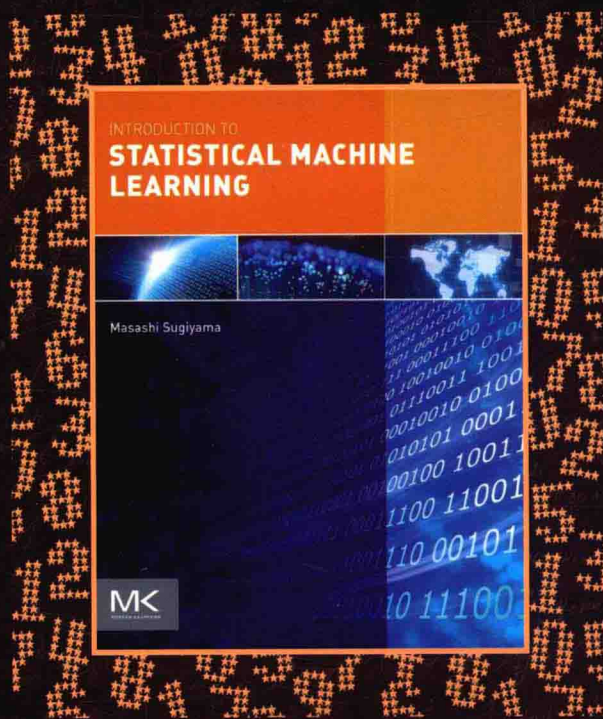


统计机器学习导论

[日] 杉山将 (Masashi Sugiyama) 著

谢宁 李柏杨 肖竹 罗宇轩 等译



INTRODUCTION TO
STATISTICAL MACHINE LEARNING



机械工业出版社
China Machine Press

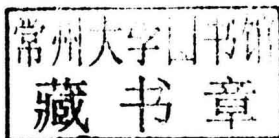
数据科学与工程丛书

INTRODUCTION TO
STATISTICAL MACHINE LEARNING

统计机器学习导论

[日] 杉山将 (Masashi Sugiyama) 著

谢宁 李柏杨 肖竹 罗宇轩 等译



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

统计机器学习导论 / (日) 杉山将 (Masashi Sugiyama) 著; 谢宁等译. —北京: 机械工业出版社, 2018.4

(数据科学与工程丛书)

书名原文: Introduction to Statistical Machine Learning

ISBN 978-7-111-59679-0

I. 统… II. ①杉… ②谢… III. 机器学习 IV. TP181

中国版本图书馆 CIP 数据核字 (2018) 第 067578 号

本书版权登记号: 图字 01-2016-4752

ELSEVIER

Elsevier (Singapore) Pte Ltd.

3 Killiney Road, #08-01 Winsland House I, Singapore 239519

Tel: (65) 6349-0200; Fax: (65) 6733-1817

Introduction to Statistical Machine Learning

Masashi Sugiyama

Copyright © 2016 by Elsevier Inc. All rights of reproduction in any form reserved.

ISBN-13: 9780128021217

This translation of Introduction to Statistical Machine Learning by Masashi Sugiyama was undertaken by China Machine Press and is published by arrangement with Elsevier (Singapore) Pte Ltd.

Introduction to Statistical Machine Learning by Masashi Sugiyama 由机械工业出版社进行翻译, 并根据机械工业出版社与爱思唯尔(新加坡)私人有限公司的协议约定出版。

《统计机器学习导论》(谢宁 李柏杨 肖竹 罗宇轩 等译)

ISBN: 978-7-111-59679-0

Copyright © 2018 by Elsevier (Singapore) Pte Ltd.

All rights reserved. No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from Elsevier (Singapore) Pte Ltd. Details on how to seek permission, further information about the Elsevier's permissions policies and arrangements with organizations such as the Copyright Clearance Center and the Copyright Licensing Agency, can be found at our website: www.elsevier.com/permissions.

This book and the individual contributions contained in it are protected under copyright by Elsevier (Singapore) Pte Ltd. and China Machine Press (other than as may be noted herein).

注 意

本译本由 Elsevier (Singapore) Pte Ltd. 和机械工业出版社完成。相关从业及研究人员必须凭借其自身经验和知识对文中描述的信息数据、方法策略、搭配组合、实验操作进行评估和使用。由于医学科学发展迅速, 临床诊断和给药剂量尤其需要经过独立验证。在法律允许的最大范围内, 爱思唯尔、译文的原文作者、原文编辑及原文内容提供者均不对译文或因产品责任、疏忽或其他操作造成的人身及 / 或财产伤害及 / 或损失承担责任, 亦不对由于使用文中提到的方法、产品、说明或思想而导致的人身及 / 或财产伤害及 / 或损失承担责任。

Printed in China by China Machine Press under special arrangement with Elsevier (Singapore) Pte Ltd. This edition is authorized for sale in the People's Republic of China only, excluding Hong Kong SAR, Macau SAR and Taiwan. Unauthorized export of this edition is a violation of the contract.

本书封底贴有 Elsevier 防伪标签, 无标签者不得销售。

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: 唐晓琳

责任校对: 殷虹

印刷: 三河市宏图印务有限公司

版次: 2018年5月第1版第1次印刷

开本: 185mm × 260mm 1/16

印张: 22

书号: ISBN 978-7-111-59679-0

定价: 89.00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: hzjsj@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

译者序

机器学习是计算机科学的重要分支之一，旨在研究原理、算法以及能够像人类一样学习的系统的应用。同时，其亦是一门交叉学科，涉及概率论、统计学、逼近论、凸分析、算法复杂度理论等多门学科。机器学习作为人工智能的核心部分，是计算机获得智能的根本途径，其应用遍及人工智能的各个领域，发挥着不可替代的重要作用。

本书是日本人工智能和机器学习领域的新一代领军人物杉山将(Masashi Sugiyama)的统计机器学习力作。本书致力于讲解数学背景及多种机器学习技术的实用化算法。其结构清晰，内容丰富，案例详实，系统地介绍了统计机器学习的概念、技术及应用。通过对本书的学习，读者可以了解统计机器学习的基本概念和知识，同时培养统计机器学习的基本技能。阅读本书需要了解计算机科学、概率论与统计学等相关基础知识。本书适用于计算机及相关专业的本科生、研究生以及相关领域的研究人员和专业技术人员。

本书翻译工作得到了课题组成员的鼎力支持和大力协作。

谢宁作为本次翻译活动的倡议者和联络人，负责和参与翻译了前言、作者简介及第1章，并帮助分析、修改各章中的疑难点。

第一部分(第1章)，由谢宁和李柏杨共同完成。

第二部分(第2~10章)，由李煜玮、周飞宇、苏秋霖和文洋负责翻译。

第三部分(第11~20章)，由王磊和徐颖负责翻译。

第四部分(第21~29章)，由罗宇轩负责翻译。

第五部分(第30~39章)，由李柏杨负责翻译。

此外，肖竹负责翻译第一到五部分的引言内容。在翻译组内部审校阶段，谢宁担任内部审校总负责人，李柏杨、肖竹和张帅担任内部审校主要负责人。

本书中文版能够出版发行，首先要感谢本书的作者杉山将教授，是他为我们著作了一本好书。其次要感谢机械工业出版社华章公司引进了本书的中文版权，使得我们能够获得为博士生导师杉山将教授翻译此书的机会，并实现将其介绍给国内广大读者的良好愿望。此外，特别感谢本书的编辑曲熠以及所有为此书的出版做出贡献

的排校人员，是他们的辛勤劳动才使本书能够付诸印刷和出版，在此表示深深的感谢和崇高的敬意！

本书对原著的错误之处做了一些修正，在原著难懂或需要提醒的地方添加了一些译者说明。尽管我们在翻译过程中力图做得更好，但因个人的业务水平、英文水平乃至中文文学水平的限制，以及翻译过程中的粗心和不够严谨，可能使得本书中文版中存在错误、不足和不当之处。热切期望读者对本书提出宝贵意见、建议和勘误，并欢迎与我们联系(seanxiening@gmail.com)。

2018年3月

前 言

机器学习是计算机领域的一个学科，旨在研究原理、算法以及能够像人类一样学习的系统的应用。近年来，计算机和传感器的发展使得我们能够访问不同领域的海量数据(如文本、音频、图片、电影、电子商务、电气、医学和生物学等)。在此类大数据的分析和利用方面，机器学习起到了核心的作用。

本书致力于讨论机器学习的数学背景及多种机器学习技术的实用化算法。目标读者定位于计算机和相关专业的本科生和研究生。在工作中应用机器学习技术的工程师和分析数据的科学家也会从本书中获益。

本书特色在于每章的主题简明扼要，给出具体机器学习技术的数学推导并附以简洁的 MATLAB 程序。由此，读者在学习数学概念的同时，可掌握多种机器学习技术的实用价值。全部 MATLAB 程序可以从如下网址获得：

<http://www.ms.k.u-tokyo.ac.jp/software/SMLbook.zip>

本书第一部分给出机器学习领域的简要概述。紧接着，第二部分介绍了概率和统计的基本概念，它们构成了统计机器学习的数学基础。第二部分的成文基于：

Sugiyama, M.
Probability and Statistics for Machine Learning,
Kodansha, Tokyo, Japan, 2015. (in Japanese).

第三部分和第四部分分别在生成和判别框架下，介绍了一系列实用机器学习算法。随后，第五部分介绍高级论题，进而处理更具挑战的机器学习任务。第三部分的成文基于：

Sugiyama, M.
Statistical Pattern Recognition: Pattern Recognition Based on Generative Models,
Ohmsha, Tokyo, Japan, 2009. (in Japanese),

第四部分和第五部分的成文基于：

Sugiyama, M.
An Illustrated Guide to Machine Learning,
Kodansha, Tokyo, Japan, 2013. (in Japanese).

在此感谢东京大学和东京工业大学相关研究组的研究员和学生针对本书早期手稿给出的有价值的反馈。

杉山将
东京大学

作者简介

杉山将 分别于1997、1999和2001年获得日本东京工业大学颁发的计算机专业工学学士、硕士和博士学位。2001年，他被东京工业大学聘为助理教授，随后于2003年晋升为副教授。2014年，他转到东京大学任教授。2003至2004年，他获得亚历山大·冯洪堡特基金会研究奖学金并前往位于德国柏林的弗朗霍夫研究所从事研究工作。2006年，他获得欧盟委员会项目 Erasmus Mundus 奖学金并前往英国爱丁堡大学从事研究工作。他曾被授予2007年度 IBM 学者奖(表彰他在机器学习领域中非平稳性方面的贡献)、2011年度日本信息处理学会颁发的 Nagao 特别研究员奖，以及由文部科学省颁发的科学技术领域的青年科学家奖(表彰他对机器学习领域密度比范式的贡献)。他的研究兴趣包括机器学习与数据挖掘的理论与算法，及其广泛的应用(如信号处理、图像处理和机器人控制等)。



目 录

译者序
前言
作者简介

第一部分 绪论

第 1 章 统计机器学习	2
1.1 学习的类型	2
1.2 机器学习任务举例	3
1.2.1 监督学习	3
1.2.2 非监督学习	4
1.2.3 进一步的主题	4
1.3 本书结构	5

第二部分 概率与统计

第 2 章 随机变量与概率分布	8
2.1 数学基础	8
2.2 概率	9
2.3 随机变量和概率分布	10
2.4 概率分布的性质	11
2.4.1 期望、中位数和众数	11
2.4.2 方差和标准差	13
2.4.3 偏度、峰度和矩	13
2.5 随便变量的变换	15
第 3 章 离散概率分布的实例	17
3.1 离散均匀分布	17

3.2 二项分布	17
3.3 超几何分布	18
3.4 泊松分布	21
3.5 负二项分布	23
3.6 几何分布	24

第 4 章 连续概率分布的实例	25
4.1 连续均匀分布	25
4.2 正态分布	25
4.3 伽马分布、指数分布和卡方分布	27
4.4 Beta 分布	29
4.5 柯西分布和拉普拉斯分布	31
4.6 t 分布和 F 分布	33

第 5 章 多维概率分布	35
5.1 联合概率分布	35
5.2 条件概率分布	36
5.3 列联表	36
5.4 贝叶斯定理	36
5.5 协方差与相关性	38
5.6 独立性	39

第 6 章 多维概率分布的实例	42
6.1 多项分布	42
6.2 多元正态分布	43
6.3 狄利克雷分布	45
6.4 威沙特分布	48

第 7 章 独立随机变量之和	50	第 10 章 假设检验	67
7.1 卷积	50	10.1 假设检验基础	67
7.2 再生性	50	10.2 正态样本期望的检验	68
7.3 大数定律	51	10.3 尼曼-皮尔森引理	68
7.4 中心极限定理	53	10.4 列联表检验	69
第 8 章 概率不等式	55	10.5 正态样本期望差值检验	70
8.1 联合界	55	10.5.1 无对应关系的两组 样本	70
8.2 概率不等式	55	10.5.2 有对应关系的两组 样本	71
8.2.1 马尔可夫不等式和切尔 诺夫不等式	55	10.6 秩的无参检验	72
8.2.2 坎泰利不等式和切比雪夫 不等式	56	10.6.1 无对应关系的两组 样本	72
8.3 期望不等式	57	10.6.2 有对应关系的两组 样本	73
8.3.1 琴生不等式	57	10.7 蒙特卡罗检验	74
8.3.2 赫尔德不等式和施瓦茨 不等式	57		
8.3.3 闵可夫斯基不等式	58		
8.3.4 康托洛维奇不等式	58		
8.4 独立随机变量和的不等式	59		
8.4.1 切比雪夫不等式和切尔 诺夫不等式	59		
8.4.2 霍夫丁不等式和伯恩 斯坦不等式	59		
8.4.3 贝内特不等式	60		
第 9 章 统计估计	62		
9.1 统计估计基础	62		
9.2 点估计	62		
9.2.1 参数密度估计	62		
9.2.2 非参数密度估计	63		
9.2.3 回归和分类	64		
9.2.4 模型选择	64		
9.3 区间估计	65		
9.3.1 基于正态样本期望的 区间估计	65		
9.3.2 bootstrap 置信区间	65		
9.3.3 贝叶斯置信区间	66		
		第三部分 统计模式识别的生成式方法	
		第 11 章 通过生成模型估计的模式 识别	76
		11.1 模式识别的公式化	76
		11.2 统计模式识别	77
		11.3 分类器训练的准则	79
		11.3.1 最大后验概率规则	79
		11.3.2 最小错误分类率 准则	79
		11.3.3 贝叶斯决策规则	80
		11.3.4 讨论	81
		11.4 生成式方法和判别式方法	81
		第 12 章 极大似然估计	83
		12.1 定义	83
		12.2 高斯模型	84
		12.3 类-后验概率的计算	86
		12.4 Fisher 线性判别分析	88
		12.5 手写数字识别	90
		12.5.1 预备知识	90

12.5.2	线性判别分析的 实现	90	第 17 章 贝叶斯推理	123
12.5.3	多分类器方法	91	17.1 贝叶斯预测分布	123
第 13 章 极大似然估计的性质		93	17.1.1 定义	123
13.1	一致性	93	17.1.2 与极大似然估计的 比较	124
13.2	渐近无偏性	93	17.1.3 计算问题	124
13.3	渐近有效性	94	17.2 共轭先验	125
13.3.1	一维的情况	94	17.3 最大后验估计	126
13.3.2	多维的情况	94	17.4 贝叶斯模型选择	128
13.4	渐近正态性	95	第 18 章 边缘相似的解析近似	131
13.5	总结	97	18.1 拉普拉斯近似	131
第 14 章 极大似然估计的模型 选择		98	18.1.1 高斯密度估计	131
14.1	模型选择	98	18.1.2 例证	132
14.2	KL 散度	99	18.1.3 应用于边际似然 逼近	133
14.3	AIC 信息论准则	100	18.1.4 贝叶斯信息准则	133
14.4	交叉检验	102	18.2 变分近似	134
14.5	讨论	103	18.2.1 变分贝叶斯最大 期望算法	134
第 15 章 高斯混合模型的极大似然 估计		104	18.2.2 与一般最大期望 法的关系	135
15.1	高斯混合模型	104	第 19 章 预测分布的数值近似	137
15.2	极大似然估计	105	19.1 蒙特卡罗积分	137
15.3	梯度上升算法	107	19.2 重要性采样	138
15.4	EM 算法	108	19.3 采样算法	139
第 16 章 非参数估计		112	19.3.1 逆变换采样	139
16.1	直方图方法	112	19.3.2 拒绝采样	141
16.2	问题描述	113	19.3.3 马尔可夫链蒙特 卡罗方法	142
16.3	核密度估计	115	第 20 章 贝叶斯混合模型	147
16.3.1	Parzen 窗法	115	20.1 高斯混合模型	147
16.3.2	利用核的平滑	116	20.1.1 贝叶斯公式化	147
16.3.3	带宽的选择	117	20.1.2 变分推断	148
16.4	最近邻密度估计	118	20.1.3 吉布斯采样	151
16.4.1	最近邻距离	118	20.2 隐狄利克雷分配模型	154
16.4.2	最近邻分类器	118	20.2.1 主题模型	154

20.2.2	贝叶斯公式化	154	25.2	ℓ_1 损失最小化	187
20.2.3	吉布斯采样	155	25.3	Huber 损失最小化	187
第四部分 统计机器学习的判别式方法					
第 21 章	学习模型	158	25.3.1	定义	188
21.1	线性参数模型	158	25.3.2	随机梯度算法	188
21.2	核模型	159	25.3.3	迭代加权最小二乘	188
21.3	层次模型	161	25.3.4	ℓ_1 约束 Huber 损失最小化	190
第 22 章	最小二乘回归	163	25.4	Tukey 损失最小化	193
22.1	最小二乘法	163	第 26 章	最小二乘分类器	195
22.2	线性参数模型的解决方案	163	26.1	基于最小二乘回归的分类器	195
22.3	最小二乘法的特性	166	26.2	0/1 损失和间隔	196
22.4	大规模数据的学习算法	167	26.3	多类分类器	198
22.5	层次模型的学习算法	168	第 27 章	支持向量分类	200
第 23 章	具有约束的最小二乘回归	171	27.1	最大间隔分类	200
23.1	子空间约束的最小二乘	171	27.1.1	硬间隔支持向量分类	200
23.2	ℓ_2 约束的最小二乘	172	27.1.2	软间隔支持向量分类	201
23.3	模型选择	175	27.2	支持向量分类的对偶最优化问题	201
第 24 章	稀疏回归	178	27.3	对偶解的稀疏性	203
24.1	ℓ_1 约束的最小二乘	178	27.4	使用核技巧的非线性模型	204
24.2	解决 ℓ_1 约束的最小二乘	179	27.5	多类扩展	206
24.3	稀疏学习的特征选择	181	27.6	损失最小化观点	207
24.4	若干扩展	181	27.6.1	Hinge 损失最小化	207
24.4.1	广义 ℓ_1 约束最小二乘	182	27.6.2	平方 Hinge 损失最小化	208
24.4.2	ℓ_p 约束最小二乘	182	27.6.3	Ramp 损失最小化	210
24.4.3	$\ell_1 + \ell_2$ 约束最小二乘	183	第 28 章	概率分类法	212
24.4.4	$\ell_{1,2}$ 约束最小二乘	184	28.1	Logistic 回归	212
24.4.5	迹范数约束最小二乘	184	28.1.1	Logistic 模型与极大似然估计	212
第 25 章	稳健回归	186	28.1.2	损失最小化的观点	214
25.1	ℓ_2 损失最小化的非稳健性	186	28.2	最小二乘概率分类	214

第 29 章 结构化分类	217	33.1.1 输入样本的流形 结构	248
29.1 序列分类器	217	33.1.2 计算解决方案	249
29.2 序列的概率分类	217	33.2 协变量移位的适应	251
29.2.1 条件随机场	218	33.2.1 重要度加权学习	251
29.2.2 极大似然估计	219	33.2.2 相对重要度加权 学习	252
29.2.3 递归计算	219	33.2.3 重要度加权交叉 检验	253
29.2.4 新样本预测	221	33.2.4 重要度估计	253
29.3 序列的确定性分类	222	33.3 类别平衡变化下的适应	255
第五部分 高级主题			
第 30 章 集成学习	226	33.3.1 类别平衡加权 学习	256
30.1 决策树桩分类器	226	33.3.2 类别平衡估计	256
30.2 bagging 算法	227	第 34 章 多任务学习	259
30.3 boosting 算法	228	34.1 任务相似度正则化	259
30.3.1 adaboost 算法	228	34.1.1 公式化	259
30.3.2 损失最小化观点	230	34.1.2 解析解	260
30.4 泛化集成学习	233	34.1.3 多任务的有效计算 方法	260
第 31 章 在线学习	234	34.2 多维函数学习	261
31.1 随机梯度下降法	234	34.2.1 公式化	261
31.2 被动攻击学习	235	34.2.2 有效的分析解决 方案	263
31.2.1 分类	235	34.3 矩阵正则化	263
31.2.2 回归	237	34.3.1 参数矩阵正则化	264
31.3 加权向量的自适应正则化	238	34.3.2 迹范数正则化的 梯度法	265
31.3.1 参数的不确定性	238	第 35 章 线性降维	268
31.3.2 分类	239	35.1 维度灾难	268
31.3.3 回归	240	35.2 无监督降维法	269
第 32 章 预测的置信度	241	35.2.1 主成分分析	270
32.1 ℓ_2 正则化最小二乘的预测 方差	241	35.2.2 局部保留投影	271
32.2 bootstrap 法置信区间估计	243	35.3 分类的线性判别分析	272
32.3 应用	244	35.3.1 Fisher 判别 分析法	273
32.3.1 时间序列预测	244		
32.3.2 调整参数的优化	245		
第 33 章 半监督学习	248		
33.1 流形正则化	248		

35.3.2	局部 Fisher 判别 分析法	274	第 37 章	聚类	297
35.3.3	半监督局部 Fisher 判别分析法	276	37.1	k 均值聚类	297
35.4	回归问题的充分降维	277	37.2	核 k 均值聚类	299
35.4.1	信息论公式化	278	37.3	谱聚类	299
35.4.2	直接导数估计	279	37.4	调谐参数的选择	299
35.5	矩阵插补	282	第 38 章	异常检测	304
第 36 章	非线性降维	285	38.1	密度估计和局部异常因子 ..	304
36.1	利用核技巧的降维	285	38.2	支持向量数据描述	305
36.1.1	核主成分分析	285	38.3	基于正常值的异常检测	308
36.1.2	拉普拉斯特征映射 ..	288	第 39 章	变化检测	312
36.2	通过神经网络的监督 降维法	289	39.1	基于分布模型的变化检测 ..	312
36.3	通过自编码器的非监督 降维法	290	39.1.1	KL 散度	312
36.3.1	自编码器	290	39.1.2	Pearson 散度	313
36.3.2	通过梯度下降法的 训练	290	39.1.3	L_2 距离	313
36.3.3	稀疏自编码器	292	39.1.4	L_1 距离	315
36.4	通过受限玻尔兹曼机的非监督 降维法	292	39.1.5	最大均值差异	317
36.4.1	模型	293	39.1.6	能量距离	317
36.4.2	通过梯度下降法的 训练	293	39.1.7	时序变化检测的 应用	317
36.5	深度学习	296	39.2	基于结构模型的变化检测 ..	318
			39.2.1	稀疏极大似然估计 ..	319
			39.2.2	稀疏密度比估计 ..	321
			参考文献	324
			索引	329

第一部分

绪 论

第 1 章 统计机器学习

1
>
2

第 1 章

统计机器学习

计算机与网络的近期发展使得我们能够即刻访问大量信息，如文字、声音、图像与影像。此外，日志、消费记录和病历等广泛的个人数据也日复一日地累积起来。如此大量的数据被称为大数据(big data)。并且，存在一股通过从数据中抽取有用的知识来创造新价值与商机的增长趋势。这一过程通常被称为数据挖掘(data mining)。机器学习是用于抽取有用知识的关键技术。

1.1 学习的类型

依据可用数据的类型，机器学习能够分为监督学习(supervised learning)、非监督学习(unsupervised learning)和强化学习(reinforcement learning)。

监督学习可谓机器学习最基本的类型。它被视为一种学生学习的过程，即向导师提问并回答。在机器学习情境中，学生对应于计算机，导师对应于计算机的用户；计算机从问与答的成对样本中学习一种从问题到其答案的映射。监督学习的目标在于获得泛化能力(generalization ability)。其指的是一种能够为从未被学习过的问题猜出恰当答案的能力。因此，用户不必再将每件事情都教给计算机，而计算机仅通过学习一小部分知识就能够自动地应对未知情况。监督学习已经被成功地应用到广泛的真实问题中，如：手写字母识别、语音识别、图像识别、垃圾邮件过滤、信息检索、在线广告、推荐系统、脑电波分析、基因分析、股票价格预测、天气预报和天文数据分析。监督学习的问题也会被具体称为回归(regression)，当答案是一个实数值(如：温度)；分类(classification)，如果答案是一个分类值(如：“是”或“否”)以及排序(ranking)，如果答案是一个数列值(如：“好”、“中”或“劣”)。

非监督学习被认为是，导师不存在并且学生自学。在机器学习情境中，计算机通过互联网自动地收集数据并且尝试在没有用户任何指导下抽取有用的知识。因此，非监督学习比监督学习更加自动化，尽管其目标不一定指定清楚。非监督学习的典型任务包括数据聚类(data clustering)和异常点检测(outlier detection)。同时，这些非监督学习技术已经就广泛的真实问题取得巨大成功，如系统诊断、安全、事件检测和社交网络分析。非监督学习也通常被用作监督学习的预处理过程。

强化学习与监督学习类似，也是以使计算机获得对没有学习过的问题做出正确解答的泛化能力为目标，但是在学习过程中，不设置导师提示对错、告知最终答案的环节。相

反, 导师评价(evaluate)学生的行为并给予其反馈。强化学习的目标是基于来自导师的反馈, 使得学生提高其行为, 从而最大化导师的评价。强化学习是一个人与机器人的行为的重要的模型。它已经被广泛地应用于诸多领域, 如: 自主机器人控制、电脑游戏和营销策略优化。在强化学习背后, 监督学习与非监督学习的方法, 如: 回归、分类与聚类, 都常被使用。

本书重点在于监督学习与非监督学习。对于强化学习, 请参考文献[99, 105]。

1.2 机器学习任务举例

在本节, 会详尽介绍不同的监督学习与非监督学习任务。

1.2.1 监督学习

回归(regression)的目标是为了从其样本近似估计一个实数值函数(如图 1.1 所示)。让我们不妨定义输入为 d 维实数向量 \mathbf{x} , 输出为实数标量 y , 学习目标函数为 $y=f(\mathbf{x})$ 。这个学习目标函数 f 被认为是未知的, 但是其输入输出样本对 $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ 是可以被观测的。在实践中, 被观测输出值 y_i 可能被某些误差 ϵ_i 所污染, 即

$$y_i = f(\mathbf{x}_i) + \epsilon_i$$

在这个设定中, x_i 对应学生向导师提出的一个问题; y_i 对应导师给学生的回答。噪声 ϵ_i 或许对应导师的口误或者学生的误解。学习目标函数 f 对应导师的知识, 这些知识使得他(她)可以回答任何问题。回归的目标在于让学生学习这个函数, 学生由此也可以回答任何问题。泛化的水平能够通过真实函数 f 与其近似 \hat{f} 的接近程度来测量。

另一个方面, 以监督学习方式, 分类则是一个模式识别(pattern recognition)问题(如图 1.2 所示)。不妨用 d 维向量 $\hat{\mathbf{x}}$ 来定义输入模式和由标量 $y \in \{1, \dots, c\}$ 定义其类别, 其中 c 定义类别数目。为了训练分类器, 输入-输出对样本 $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ 以与回归一样的方式提供。倘若真实分类规则被定义为 $y=f(\mathbf{x})$, 那么分类也能够被视为函数逼近问题。然而, 在回归情境下, 一个本质的不同在于在 y 中不存在接近程度这一概念 $y: y=2$ 比 $y=1$ 更接近 $y=3$, 但是“ y 和 y' 是否相同”是分类唯一关心的。

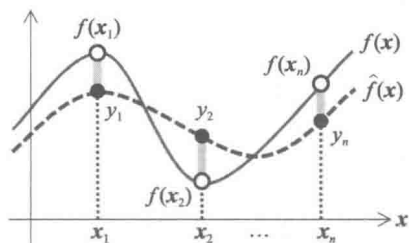


图 1.1 回归

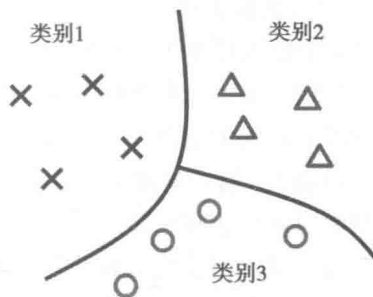


图 1.2 分类

在监督学习中, 排序(ranking)问题旨在学习样本 \mathbf{x} 的次序 y 。因为次序有顺序, 如: $1 < 2 < 3$, 排序更像是回归而非分类。为此, 排序问题也可以被视为有序回归(ordinal regression)。然而, 不同于回归, 提出输出值 y 是没有必要被预测的, 但也仅仅是其相对