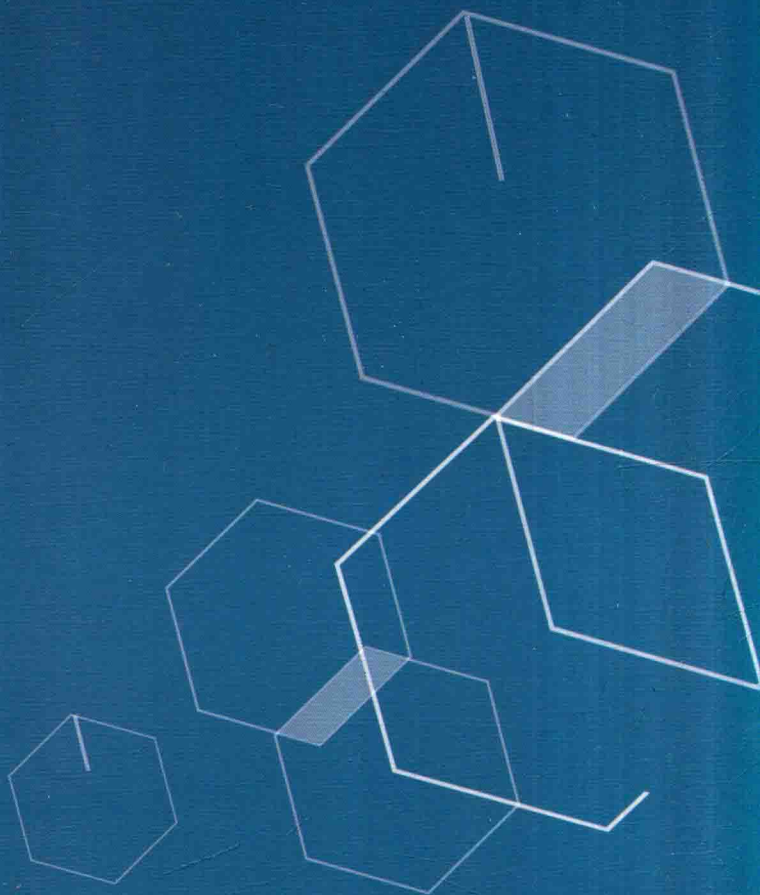
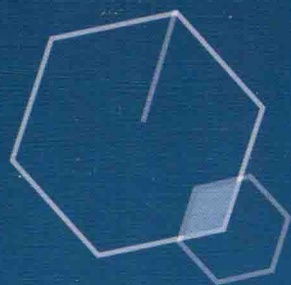


统计学习方法

(第2版)

李航著



清华大学出版社

统计学习方法

(第2版)

李航 著

清华大学出版社
北京

内 容 简 介

统计学习方法即机器学习方法,是计算机及其应用领域的一门重要学科。本书分为监督学习和无监督学习两篇,全面系统地介绍了统计学习的主要方法。包括感知机、 k 近邻法、朴素贝叶斯法、决策树、逻辑斯谛回归与最大熵模型、支持向量机、提升方法、EM算法、隐马尔可夫模型和条件随机场,以及聚类方法、奇异值分解、主成分分析、潜在语义分析、概率潜在语义分析、马尔可夫链蒙特卡罗法、潜在狄利克雷分配和PageRank算法等。

本书是统计机器学习及相关课程的教学参考书,适用于高等院校文本数据挖掘、信息检索及自然语言处理等专业的大学生、研究生,也可供计算机应用等专业的研发人员参考。

本书封面贴有清华大学出版社防伪标签,无标签者不得销售。

版权所有,侵权必究。侵权举报电话:010-62782989 13701121933

图书在版编目(CIP)数据

统计学习方法/李航著.—2版.—北京:清华大学出版社,2019
ISBN 978-7-302-51727-6

I. ①统… II. ①李… III. ①统计学 IV. ①C8

中国版本图书馆CIP数据核字(2018)第267477号

责任编辑:薛慧

封面设计:李祥榕

责任校对:刘玉霞

责任印制:宋林

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦A座 邮 编:100084

社总机:010-62770175 邮 购:010-62786544

投稿与读者服务:010-62776969, c-service@tup.tsinghua.edu.cn

质量反馈:010-62772015, zhiliang@tup.tsinghua.edu.cn

印装者:三河市龙大印装有限公司

经 销:全国新华书店

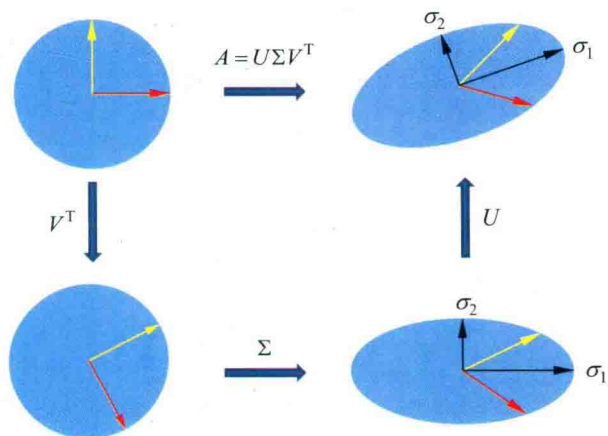
开 本:170mm×240mm 印 张:30.25 插 页:1 字 数:593千字

版 次:2012年3月第1版 2019年5月第2版 印 次:2019年5月第1次印刷

印 数:1~20000

定 价:98.00元

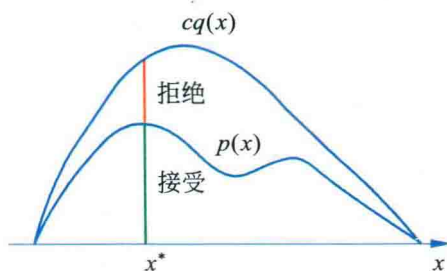
产品编号:081329-01



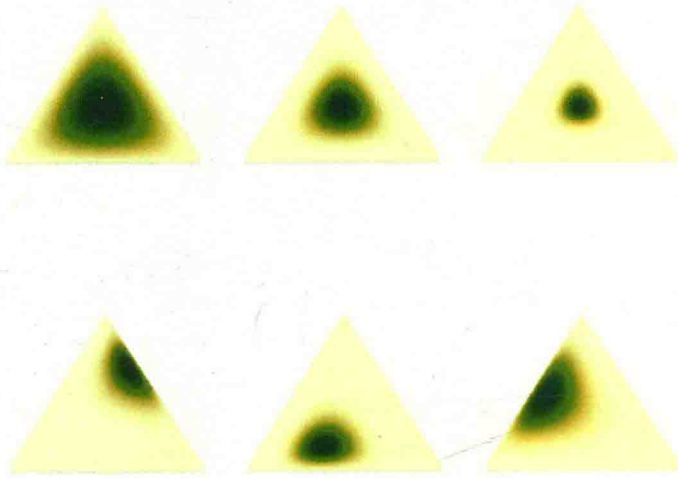
彩图 15.1 奇异值分解的几何解释

	doc 1	doc 2	doc 3	doc 4
word 1	2	2	4	3
word 2	2	1	5	3
word 3	1	1	2	0
word 4	0	1	2	1

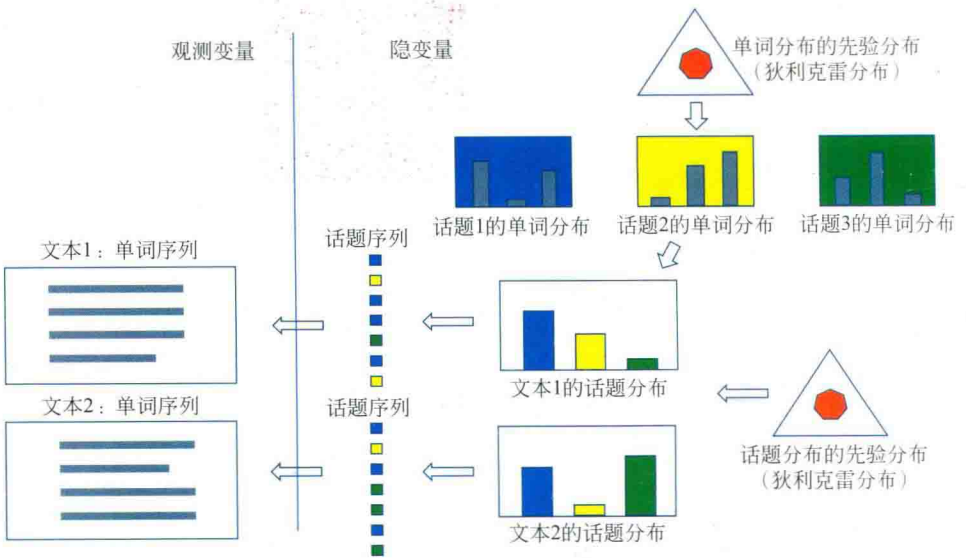
彩图 18.1 概率潜在语义分析的直观解释



彩图 19.1 接受-拒绝抽样法



彩图 20.1 狄利克雷分布例



彩图 20.3 LDA 的文本生成过程

献给我的母亲

第 2 版序言

《统计学习方法》第 1 版于 2012 年出版，讲述了统计机器学习方法，主要是一些常用的监督学习方法。第 2 版增加了一些常用的无监督学习方法，由此本书涵盖了传统统计机器学习方法的主要内容。

在撰写《统计学习方法》伊始，对全书内容做了初步规划。第 1 版出版之后，即着手无监督学习方法的写作。由于写作是在业余时间进行，常常被主要工作打断，历经六年多时间才使这部分工作得以完成。犹未能加入深度学习和强化学习等重要内容，希望今后能够增补，完成整本书的写作计划。

《统计学习方法》第 1 版的出版正值大数据和人工智能的热潮，生逢其时，截至 2019 年 4 月本书共印刷 25 次，152000 册，得到了广大读者的欢迎和支持。有许多读者指出本书对学习和掌握机器学习技术有极大的帮助，也有许多读者通过电子邮件、微博等方式指出书中的错误，提出改进的建议和意见。一些高校将本书作为机器学习课程的教材或参考书。有的同学在网上发表了读书笔记，有的同学将本书介绍的方法在计算机上实现。清华大学深圳研究生院袁春老师精心制作了第 1 版十二章的课件，在网上公布，为大家提供教学之便。众多老师、同学、读者的支持和鼓励，让作者深受感动和鼓舞。在这里向所有的老师、同学、读者致以诚挚的谢意！

能为中国的计算机科学、人工智能领域做出一点微薄的贡献，感到由衷的欣慰，同时也感受到作为知识传播者的重大责任，让作者决意把本书写好。也希望大家今后不吝指教，多提宝贵意见，以帮助继续提高本书的质量。在写作中作者也深切体会到教学相长的道理，经常发现自己对基础知识的掌握不够扎实，通过写作得以对相关知识进行了深入的学习，受益匪浅。

本书是一部机器学习的基本读物，要求读者拥有高等数学、线性代数和概率统计的基础知识。书中主要讲述统计机器学习的方法，力求系统全面又简明扼要地阐述这些方法的理论、算法和应用，使读者能对这些机器学习的基本技术有很好的掌握。针对每个方法，详细介绍其基本原理、基础理论、实际算法，给出细致的数学推导和具体实例，既帮助读者理解，也便于日后复习。

第 2 版增加的无监督学习方法,王泉、陈嘉怡、柴琛林、赵程绮等帮助做了认真细致的校阅,提出了许多宝贵意见,在此谨对他们表示衷心的感谢。清华大学出版社的薛慧编辑一直对本书的写作给予非常专业的指导和帮助,在此对她表示衷心的感谢!

由于本人水平有限,本书一定存在不少错误,恳请各位专家、老师和同学批评指正。

李 航

2019 年 4 月

第 1 版序言

计算机与网络已经融入人们的日常学习、工作和生活之中，成为人们不可或缺的助手和伙伴。计算机与网络的飞速发展完全改变了人们的学习、工作和生活的方式。智能化是计算机研究与开发的一个主要目标。近几十年来的实践表明，统计机器学习方法是实现这一目标的最有效手段，尽管它还存在着一定的局限性。

本人一直从事利用统计学习方法对文本数据进行各种智能性处理的研究，包括自然语言处理、信息检索、文本数据挖掘。近 20 年来，这些领域发展之快，应用之广，实在令人惊叹！可以说，统计机器学习是这些领域的核心技术，在这些领域的发展及应用中起着决定性的作用。

本人在日常的研究工作中经常指导学生，并在国内外一些大学及讲习班上多次做关于统计学习的报告和演讲。在这一过程中，同学们学习热情很高，希望得到指导，这使作者产生了撰写本书的想法。

国内外已出版了多本关于统计机器学习的书籍，比如，Hastie 等人的《统计学习基础》，该书对统计学习的诸多问题有非常精辟的论述，但对初学者来说显得有些深奥。统计学习范围甚广，一两本书很难覆盖所有问题。本书主要是面向将统计学习方法作为工具的科研人员与学生，特别是从事信息检索、自然语言处理、文本数据挖掘及相关领域的研究与开发的科研人员与学生。

本书力求系统而详细地介绍统计学习的方法。在内容选取上，侧重介绍那些最重要、最常用的方法，特别是关于分类与标注问题的方法。对其他问题及方法，如聚类等，计划在今后的写作中再加以介绍。在叙述方式上，每一章讲述一种方法，各章内容相对独立、完整；同时力图用统一框架来论述所有方法，使全书整体不失系统性，读者可以从头到尾通读，也可以选择单个章节细读。对每一种方法的讲述力求深入浅出，给出必要的推导证明，提供简单的实例，使初学者易于掌握该方法的基本内容，领会方法的本质，并准确地使用方法。对相关的深层理论，则予以简述。在每章后面，给出一些习题，介绍一些相关的研究动向和阅读材料，列出参考文献，以满足读者进一步学习的需求。本书第 1 章简要叙述统计学习方法的基本概念，最后一章对统计学习方

法进行比较与总结。此外,在附录中简要介绍一些共用的最优化理论与方法。

本书可以作为统计机器学习及相关课程的教学参考书,适用于信息检索及自然语言处理等专业的大学生、研究生。

本书初稿完成后,田飞、王佳磊、武威、陈凯、伍浩铖、曹正、陶宇等人分别审阅了全部或部分章节,提出了许多宝贵意见,对本书质量的提高有很大帮助,在此向他们表示衷心的感谢。在本书写作和出版过程中,清华大学出版社的责任编辑薛慧给予了很多帮助,在此特向她致谢。

由于本人水平所限,书中难免有错误和不当之处,欢迎各位专家和读者给予批评指正。

李 航

2011 年 4 月 23 日

目 录

第 1 篇 监督学习

第 1 章 统计学习及监督学习概论	3
1.1 统计学习	3
1.2 统计学习的分类	5
1.2.1 基本分类	6
1.2.2 按模型分类	11
1.2.3 按算法分类	13
1.2.4 按技巧分类	13
1.3 统计学习方法三要素	15
1.3.1 模型	15
1.3.2 策略	16
1.3.3 算法	19
1.4 模型评估与模型选择	19
1.4.1 训练误差与测试误差	19
1.4.2 过拟合与模型选择	20
1.5 正则化与交叉验证	23
1.5.1 正则化	23
1.5.2 交叉验证	24
1.6 泛化能力	24
1.6.1 泛化误差	24
1.6.2 泛化误差上界	25
1.7 生成模型与判别模型	27
1.8 监督学习应用	28
1.8.1 分类问题	28

1.8.2	标注问题	30
1.8.3	回归问题	32
	本章概要	33
	继续阅读	33
	习题	33
	参考文献	34
第 2 章	感知机	35
2.1	感知机模型	35
2.2	感知机学习策略	36
2.2.1	数据集的线性可分性	36
2.2.2	感知机学习策略	37
2.3	感知机学习算法	38
2.3.1	感知机学习算法的原始形式	38
2.3.2	算法的收敛性	41
2.3.3	感知机学习算法的对偶形式	43
	本章概要	46
	继续阅读	46
	习题	46
	参考文献	47
第 3 章	k 近邻法	49
3.1	k 近邻算法	49
3.2	k 近邻模型	50
3.2.1	模型	50
3.2.2	距离度量	50
3.2.3	k 值的选择	52
3.2.4	分类决策规则	52
3.3	k 近邻法的实现: kd 树	53
3.3.1	构造 kd 树	53
3.3.2	搜索 kd 树	55
	本章概要	57
	继续阅读	57

习题	58
参考文献	58
第 4 章 朴素贝叶斯法	59
4.1 朴素贝叶斯法的学习与分类	59
4.1.1 基本方法	59
4.1.2 后验概率最大化的含义	61
4.2 朴素贝叶斯法的参数估计	62
4.2.1 极大似然估计	62
4.2.2 学习与分类算法	62
4.2.3 贝叶斯估计	64
本章概要	65
继续阅读	66
习题	66
参考文献	66
第 5 章 决策树	67
5.1 决策树模型与学习	67
5.1.1 决策树模型	67
5.1.2 决策树与 if-then 规则	68
5.1.3 决策树与条件概率分布	68
5.1.4 决策树学习	69
5.2 特征选择	71
5.2.1 特征选择问题	71
5.2.2 信息增益	72
5.2.3 信息增益比	76
5.3 决策树的生成	76
5.3.1 ID3 算法	76
5.3.2 C4.5 的生成算法	78
5.4 决策树的剪枝	78
5.5 CART 算法	80
5.5.1 CART 生成	81
5.5.2 CART 剪枝	85

本章概要	87
继续阅读	88
习题	89
参考文献	89
第 6 章 逻辑斯谛回归与最大熵模型	91
6.1 逻辑斯谛回归模型	91
6.1.1 逻辑斯谛分布	91
6.1.2 二项逻辑斯谛回归模型	92
6.1.3 模型参数估计	93
6.1.4 多项逻辑斯谛回归	94
6.2 最大熵模型	94
6.2.1 最大熵原理	94
6.2.2 最大熵模型的定义	96
6.2.3 最大熵模型的学习	98
6.2.4 极大似然估计	102
6.3 模型学习的最优化算法	103
6.3.1 改进的迭代尺度法	103
6.3.2 拟牛顿法	107
本章概要	108
继续阅读	109
习题	109
参考文献	109
第 7 章 支持向量机	111
7.1 线性可分支持向量机与硬间隔最大化	112
7.1.1 线性可分支持向量机	112
7.1.2 函数间隔和几何间隔	113
7.1.3 间隔最大化	115
7.1.4 学习的对偶算法	120
7.2 线性支持向量机与软间隔最大化	125
7.2.1 线性支持向量机	125
7.2.2 学习的对偶算法	127

7.2.3 支持向量	130
7.2.4 合页损失函数	131
7.3 非线性支持向量机与核函数	133
7.3.1 核技巧	133
7.3.2 正定核	136
7.3.3 常用核函数	140
7.3.4 非线性支持向量分类机	141
7.4 序列最小最优化算法	142
7.4.1 两个变量二次规划的求解方法	143
7.4.2 变量的选择方法	147
7.4.3 SMO 算法	149
本章概要	149
继续阅读	152
习题	152
参考文献	153
第 8 章 提升方法	155
8.1 提升方法 AdaBoost 算法	155
8.1.1 提升方法的基本思路	155
8.1.2 AdaBoost 算法	156
8.1.3 AdaBoost 的例子	158
8.2 AdaBoost 算法的训练误差分析	160
8.3 AdaBoost 算法的解释	162
8.3.1 前向分步算法	162
8.3.2 前向分步算法与 AdaBoost	164
8.4 提升树	166
8.4.1 提升树模型	166
8.4.2 提升树算法	166
8.4.3 梯度提升	170
本章概要	172
继续阅读	172
习题	173
参考文献	173

第 9 章 EM 算法及其推广	175
9.1 EM 算法的引入	175
9.1.1 EM 算法	175
9.1.2 EM 算法的导出	179
9.1.3 EM 算法在无监督学习中的应用	181
9.2 EM 算法的收敛性	181
9.3 EM 算法在高斯混合模型学习中的应用	183
9.3.1 高斯混合模型	183
9.3.2 高斯混合模型参数估计的 EM 算法	183
9.4 EM 算法的推广	187
9.4.1 F 函数的极大-极大算法	187
9.4.2 GEM 算法	189
本章概要	191
继续阅读	192
习题	192
参考文献	192
第 10 章 隐马尔可夫模型	193
10.1 隐马尔可夫模型的基本概念	193
10.1.1 隐马尔可夫模型的定义	193
10.1.2 观测序列的生成过程	196
10.1.3 隐马尔可夫模型的 3 个基本问题	196
10.2 概率计算算法	197
10.2.1 直接计算法	197
10.2.2 前向算法	198
10.2.3 后向算法	201
10.2.4 一些概率与期望值的计算	202
10.3 学习算法	203
10.3.1 监督学习方法	203
10.3.2 Baum-Welch 算法	204
10.3.3 Baum-Welch 模型参数估计公式	206
10.4 预测算法	207
10.4.1 近似算法	208
10.4.2 维特比算法	208

本章概要	212
继续阅读	212
习题	213
参考文献	213
第 11 章 条件随机场	215
11.1 概率无向图模型	215
11.1.1 模型定义	215
11.1.2 概率无向图模型的因子分解	217
11.2 条件随机场的定义与形式	218
11.2.1 条件随机场的定义	218
11.2.2 条件随机场的参数化形式	220
11.2.3 条件随机场的简化形式	221
11.2.4 条件随机场的矩阵形式	223
11.3 条件随机场的概率计算问题	224
11.3.1 前向-后向算法	225
11.3.2 概率计算	225
11.3.3 期望值的计算	226
11.4 条件随机场的学习算法	227
11.4.1 改进的迭代尺度法	227
11.4.2 拟牛顿法	230
11.5 条件随机场的预测算法	231
本章概要	235
继续阅读	235
习题	236
参考文献	236
第 12 章 监督学习方法总结	237
第 2 篇 无监督学习	
第 13 章 无监督学习概论	245
13.1 无监督学习基本原理	245
13.2 基本问题	246