

中外学者论

AI



机器学习

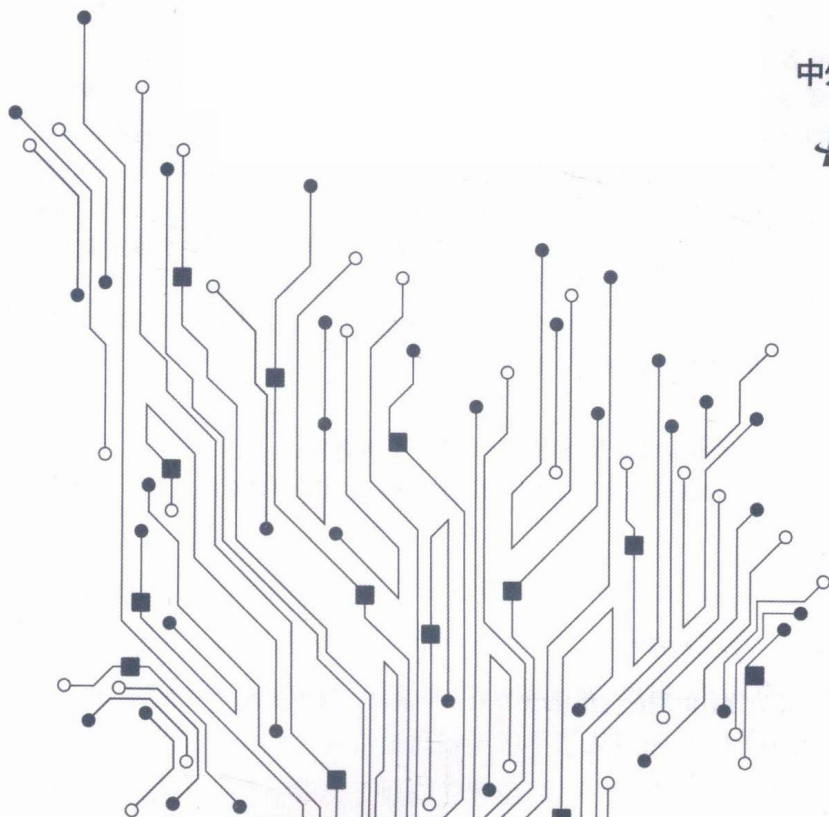
算法背后的理论与优化

◎ 史春奇 卜晶祎 施智平 著

清华大学出版社



中外学者论



机器学习

算法背后的理论与优化

◎史春奇 卜晶祎 施智平 著



清华大学出版社

北京

内 容 简 介

以机器学习为核心的人工智能已经成为新一代生产力发展的主要驱动因素。新的技术正在向各行各业渗透,大有变革各个领域的趋势。传统产业向智慧产业的升级迫使原行业从业人员逐渐转型,市场上对相关学习材料的需求也日益高涨。帮助广大学习者更好地理解 and 掌握机器学习,是编写本书的目的。

本书针对机器学习领域中最常见的一类问题——有监督学习,从入门、进阶、深化三个层面由浅入深地进行了讲解。三个层面包括基础入门算法、核心理论及理论背后的数学优化。入门部分用以逻辑回归为代表的广义线性模型为出发点,引入书中所有涉及的知识;进阶部分的核心理论涵盖了经验风险最小、结构风险最小、正则化及统一的分类边界理论;深化部分的数学优化则主要包括最大熵原理、拉格朗日对偶等理论在数学上的推导,以及对模型求解的主流最优化方法的探讨等。

本书由浅入深,从个别到普遍,从自然算法到优化算法,从各个角度深入剖析了机器学习,力求帮助读者循序渐进地掌握机器学习的概念、算法和优化理论。

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

机器学习:算法背后的理论与优化/史春奇,卜晶祎,施智平著. —北京:清华大学出版社,2019
(中外学者论 AI)

ISBN 978-7-302-51718-4

I. ①机… II. ①史… ②卜… ③施… III. ①机器学习-算法 IV. ①TP181

中国版本图书馆 CIP 数据核字(2018)第 267470 号

责任编辑:王 芳 王冰飞

封面设计:台禹微

责任校对:焦丽丽

责任印制:丛怀宇

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

课 件 下 载: <http://www.tup.com.cn>, 010-62770175 转 4506

印 装 者:北京嘉实印刷有限公司

经 销:全国新华书店

开 本:185mm×260mm 印 张:12.75

字 数:277千字

版 次:2019年7月第1版

印 次:2019年7月第1次印刷

定 价:69.00元

产品编号:078279-01

史春奇博士 毕业于日本京都大学，美国 Brandeis University 博士后，现为港辉金融信息 Vice President。曾任通用电气（中国）有限公司资深数据科学家。

卜晶祎 毕业于上海交通大学，现为友邦保险集团人工智能主管。曾就职于通用电气（中国）研究开发中心有限公司，任资深数据科学家；曾任飞利浦亚洲研究院高级研究员。

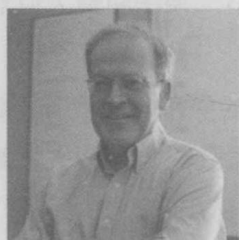
施智平博士 首都师范大学信息工程学院教授，院长，毕业于中科院计算技术研究所。于 2012 年和 2017 年获北京市科学技术奖二等奖两次，中国计算机学会高级会员，形式化方法专委会委员，人工智能学会会员，智能服务专委会委员，IEEE/ACM 会员。

P 前言

preface



在当今的人工智能领域中，最热门的技术毫无疑问当属深度学习。深度学习在 Geoffrey Hinton、Yoshua Bengio、Yann LeCun 和 Juergen Schmidhuber 等巨擘们持续不断的贡献下，在文本、图像、自然语言等领域均取得了革命性的进展。当然，深度学习只是机器学习的一个分支，能取得当前的成就也是建立在机器学习不断发展的基础之上。在机器学习领域，很多著名科学家（如图 1 所示）提出了他们的理论，做出了他们的贡献。Leslie Valiant 提出的概率近似正确学习 (Probably Approximately Correct Learning, PAC) 理论打下了计算学习理论的基石，并在此后提出了自举 (Bootstrapping) 思想。Vladimir Vapnik 提出的支持向量机 (Support Vector Machine, SVM) 是一个理论和应用都十分强大的算法。与此同时他所提出的经验风险最小与结构风险最小理论，以及背后更深层次的 VC 维 (Vapnik-Chervonenkis dimension) 理论，为部分统一分类问题提供了理论基础。Judea Pearl 提出



(a) Leslie Valiant



(b) Vladimir Vapnik



(c) Judea Pearl



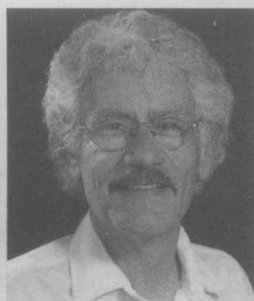
(d) Michael I. Jordan



(e) Leo Breiman



(f) Robert Schapire



(g) Jerome H. Friedman

图 1 机器学习领域 (支持向量机、集成学习、概率图模型) 的著名科学家



了贝叶斯网络，而 Michael I. Jordan 则在此基础上发展了概率图模型。Leo Breiman 在集成 (Ensemble) 学习的思想下设计了随机森林 (Random Forest) 算法，Robert Schapire 和 Jerome H. Friedman 则基于 Boosting 分别发明了 AdaBoost 和 Gradient Boosting 算法。至此，机器学习中最耀眼的算法——支持向量机、集成学习和概率图模型交相辉映，为整个机器学习理论的发展奠定了深厚的基础。

本书首先尝试把机器学习的经典算法，包括逻辑回归 (Logistic Regression)、支持向量机和 AdaBoost 等，在经验风险最小和结构风险最小的框架下进行统一，并且借助 Softmax 模型和概率图模型中的 Log-Linear 模型阐述它们的内在联系；其次从熵的角度解读概率分布、最大似然估计、指数分布族、广义线性模型等概念；最后深入剖析用于求解的最优化算法及其背后的数学理论。

本书的主要内容

全书分为 9 个章节，从单一算法到统一框架，再到一致最优化求解，各章节的设置如下。

第 1 章，首先提出并探讨几个基本问题，包括回归思想、最优模型评价标准、数理统计与机器学习的关系等。然后介绍两个最简单、最常见的有监督学习算法——线性回归和逻辑回归，并从计算的角度分析两种模型内在的关联，从而为学习“广义线性模型”打下基础。在本章的最后部分初步讲解两个模型的求解方法——最小二乘法和最大似然估计。

第 2 章，主要内容是线性回归的泛化形式——广义线性模型。本章详细介绍广义线性模型，并在第 1 章的基础上从 Fisher 信息、KL 散度、Bregman 距离的角度深入讲解最大似然估计。本章可以看作是第 3 章的基础引入。

第 3 章，在前两章的基础上提出泛化误差和经验风险最小等概念，并且将最小二乘和最大似然并入损失函数的范畴。在此基础之上，我们便将逻辑回归、支持向量机和 Ada Boost 算法统一到分类界面的框架下。至此，我们会看到不同的算法只是分别对应了不同的损失函数。

第 4 章，介绍经验风险最小的不足与过拟合的概念，之后引出正则化。紧接着介绍有监督学习算法中的常见正则化方法，包括 L_1 和 L_2 正则化 XG Boost 和树。本章从两个角度对 L_1 和 L_2 正则化进行深入讲解——贝叶斯和距离空间。这两个观点分别对应本书后续的两大部分——熵和最优化。

第 5 章，介绍贝叶斯统计和熵之间的关系，并且基于熵重新解读了最大似然估计、指

数分布族等概念。本章可以看作是前四章中出现的内容在熵概念下的再定义。同时也为下一章的 Log-Linear 模型作出铺垫。

第 6 章, 介绍 Softmax 和 Log-Linear 的变化, 并且将第 3 章的二分类界面泛化到多分类界面, 把分类问题的思路扩展到了多分类和结构分类。在本章中通过 Log-Linear 关联了概率图模型, 通过 Softmax 关联了深度学习。

第 7 章, 承接第 4 章中 L_1 和 L_2 正则化在最优化角度的解释, 从凸共轭开始递进地推导出拉格朗日对偶、Fenchel 对偶、增广拉格朗日乘法、交替方向乘法。

第 8 章, 介绍有监督学习模型在机器学习场景下的统一求解方法——随机梯度下降法及其改进算法。本章对随机梯度下降法进行了收敛性分析, 并根据分析结果针对其缺点着重介绍了两类改进策略——方差缩减和加速与适应。

第 9 章, 主要对数学意义上的最优化方法进行探讨, 可以看作是连接第 7 章和第 8 章的桥梁。第 7 章的内容是本章的理论部分, 而第 8 章的内容则是本章介绍的算法应用在机器学习场景中的特例, 主要内容包括一阶、二阶最优化算法及其收敛性分析。

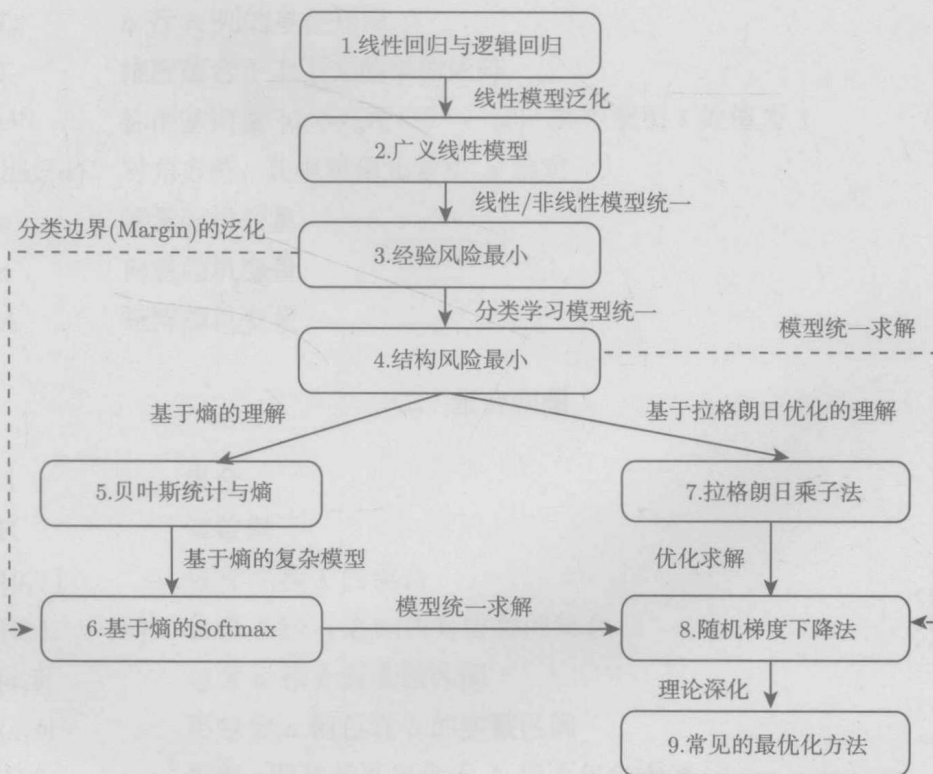


图 2 章节结构关系图



史春奇与卜晶祎共同为本书的第一作者。书中第 3~6 章主要由史春奇博士撰写，第 1、2、7~9 章主要由卜晶祎撰写，施智平教授参与了本书的组织结构设计并提出了很多宝贵意见。由于作者的能力与水平有限，本书对机器学习的探讨难免会有不全面、不深刻等不足之处，敬请各位读者批评指正，如蒙赐教将不胜感激。

各个章节结构之间的关系如图 2 所示。对于基础稍浅的读者，可以按照图示循序渐进地阅读；对于有一定基础的读者，可以跳过部分章节直接阅读感兴趣的章节。

作 者

2019 年 1 月

M 数学符号

mathematical symbol



本部分介绍本书所使用的数学符号。

一、数和数组

a	标量 (整数或实数)
\boldsymbol{a}	向量
\boldsymbol{A}	矩阵
\mathcal{A}	张量
\boldsymbol{I}_n	n 行 n 列的单位矩阵
\boldsymbol{I}	维度蕴含于上下文的单位矩阵
$\boldsymbol{e}^{(i)}$	标准基向量 $[0, \dots, 0, 1, 0, \dots, 0]$, 其中索引 i 处值为 1
$\text{diag}(\boldsymbol{a})$	对角方阵, 其中对角元素由 \boldsymbol{a} 给定
a	标量随机变量
\boldsymbol{a}	向量随机变量
\boldsymbol{A}	矩阵随机变量

二、集合和图

\mathbb{A}	集合
\mathbb{R}	实数集
$\{0, 1\}$	包含 0 和 1 的集合
$\{0, 1, \dots, n\}$	包含 0 和 n 之间所有整数的集合
$[a, b]$	包含 a 和 b 的实数区间
$(a, b]$	不包含 a 但包含 b 的实数区间
$\mathbb{A} \setminus \mathbb{B}$	差集, 即其元素包含于 \mathbb{A} 但不包含于 \mathbb{B}
\mathcal{G}	图
$\text{Pa}_{\mathcal{G}}(x_i)$	图 \mathcal{G} 中 x_i 的父节点



三、索引

a_i	向量 a 的第 i 个元素, 其中索引从 1 开始
a_{-i}	除了第 i 个元素, a 的所有元素
$A_{i,j}$	矩阵 A 的 i, j 元素
$A_{i,:}$	矩阵 A 的第 i 行
$A_{:,i}$	矩阵 A 的第 i 列
$A_{i,j,k}$	三维张量 A 的 (i, j, k) 元素
$A_{:,:,i}$	三维张量的二维切片
a_i	随机向量 a 的第 i 个元素

四、线性代数中的操作

A^\top	矩阵 A 的转置
A^+	A 的 Moore-Penrose 伪逆
$A \odot B$	A 和 B 的逐元素乘积 (Hadamard 乘积)
$\det(A)$	A 的行列式

五、微积分

$\frac{dy}{dx}$	y 关于 x 的导数
$\frac{\partial y}{\partial x}$	y 关于 x 的偏导
$\nabla_x y$	y 关于 x 的梯度
$\nabla_X y$	y 关于 X 的矩阵导数
$\nabla_X y$	y 关于 X 求导后的张量
$\frac{\partial f}{\partial x}$	$f: \mathbf{R}^n \rightarrow \mathbf{R}^m$ 的 Jacobian 矩阵 $J \in \mathbf{R}^{m \times n}$
$\nabla_x^2 f(x)$ or $H(f)(x)$	f 在点 x 处的 Hessian 矩阵
$\int f(x) dx$	x 整个域上的定积分
$\int_S f(x) dx$	集合 S 上关于 x 的定积分



六、概率和信息论

$a \perp b$	a 和 b 相互独立的随机变量
$a \perp b \mid c$	给定 c 后条件独立
$P(a)$	离散变量上的概率分布
$p(a)$	连续变量 (或变量类型未指定时) 上的概率分布
$a \sim P$	具有分布 P 的随机变量 a
$\mathbb{E}_{x \sim P}[f(x)]$ 或 $\mathbb{E}f(x)$	$f(x)$ 关于 $P(x)$ 的期望
$\text{Var}(f(x))$	$f(x)$ 在分布 $P(x)$ 下的方差
$\text{Cov}(f(x), g(x))$	$f(x)$ 和 $g(x)$ 在分布 $P(x)$ 下的协方差
$H(x)$	随机变量 x 的香农熵
$D_{\text{KL}}(P \parallel Q)$	P 和 Q 的 KL 散度
$\mathcal{N}(x; \mu, \Sigma)$	均值为 μ , 协方差为 Σ , x 上的高斯分布

七、函数

$f: \mathbb{A} \rightarrow \mathbb{B}$	定义域为 \mathbb{A} 、值域为 \mathbb{B} 的函数 f
$f \circ g$	f 和 g 的组合
$f(x; \theta)$	由 θ 参数化, 关于 x 的函数 (有时为简化表示, 忽略 θ 记为 $f(x)$)
$\ln x$	x 的自然对数
$\sigma(x)$	Logistic sigmoid, $\frac{1}{1 + \exp(-x)}$
$\zeta(x)$	Softplus, $\ln(1 + \exp(x))$
$\ \mathbf{x}\ _p$	\mathbf{x} 的 L^p 范数
$\ \mathbf{x}\ $	\mathbf{x} 的 L^2 范数
x^+	x 的正数部分, 即 $\max(0, x)$
$\mathbf{1}_{\text{condition}}$	如果条件为真则为 1, 否则为 0

有时候使用函数 f , 它的参数是一个标量, 但应用到一个向量、矩阵或张量: $f(\mathbf{x})$ 、 $f(\mathbf{X})$ 、 $f(X)$ 。这表示逐元素地将 f 应用于数组。例如, $C = \sigma(X)$, 则对于所有合法的 i 、 j 和 k , $C_{i,j,k} = \sigma(X_{i,j,k})$ 。



八、数据集和分布

p_{data}	数据生成分布
\hat{p}_{train}	由训练集定义的经验分布
\mathcal{X}	训练样本的集合
$\mathbf{x}^{(i)}, \mathbf{x}_i$	数据集的第 i 个样本 (输入)
$y^{(i)}, \mathbf{y}^{(i)}, y_i$ 或 \mathbf{y}_i	监督学习中与 $\mathbf{x}^{(i)}$ 关联的目标
\mathbf{X}	$m \times n$ 的矩阵, 其中行 $\mathbf{X}_{i,:}$ 为输入样本 $\mathbf{x}^{(i)}$

本套丛书从人工智能科普、人工智能相关计算、人工智能前沿技术以及人工智能相关创新创业等方面出发，内容涵盖人工智能的各个方面，可以使读者全面了解人工智能。

《人工智能：来路与前程》

《人工智能世代的创新与创业》

《计算人生》

《强化学习》

C 目 录 ontents



第 1 章 线性回归与逻辑回归	1
1.1 线性回归	1
1.1.1 函数关系与统计关系	1
1.1.2 统计与机器学习	2
1.2 最小二乘法与高斯-马尔可夫定理	5
1.2.1 最小二乘法	5
1.2.2 高斯-马尔可夫定理	6
1.3 从线性回归到逻辑回归	8
1.4 最大似然估计求解逻辑回归	9
1.5 最小二乘与最大似然	11
1.5.1 逻辑回归与伯努利分布	11
1.5.2 线性回归与正态分布	12
1.6 小结	13
参考文献	13
第 2 章 广义线性模型	15
2.1 广义线性模型概述	15
2.1.1 广义线性模型的定义	15
2.1.2 链接函数与指数分布簇	17
2.2 广义线性模型求解	20
2.3 最大似然估计 I: Fisher 信息	21
2.4 最大似然估计 II: KL 散度与 Bregman 散度	23
2.4.1 KL 散度	23
2.4.2 Bregman 散度	25
2.5 小结	26

参考文献	26
第 3 章 经验风险最小	28
3.1 经验风险与泛化误差概述	28
3.1.1 经验风险	30
3.1.2 泛化误差	30
3.1.3 欠拟合和过拟合	34
3.1.4 VC 维	37
3.2 经验风险最小的算法	40
3.3 分类边界	42
3.3.1 分类算法的损失函数	42
3.3.2 分类算法的边界	45
3.4 小结	48
参考文献	48
第 4 章 结构风险最小	49
4.1 经验风险最小和过拟合	49
4.2 结构风险最小和正则化	51
4.2.1 从空间角度理解 SRM	52
4.2.2 从贝叶斯观点理解 SRM	54
4.3 回归的正则化	55
4.3.1 L_2 正则化和岭回归	56
4.3.2 L_1 正则化和 Lasso 回归	57
4.3.3 L_1 、 L_2 组合正则化和 ElasticNet 回归	58
4.4 分类的正则化	60
4.4.1 支持向量机和 L_2 正则化	60
4.4.2 XGBoost 和树正则化	62
4.4.3 神经网络和 DropOut 正则化	65
4.4.4 正则化的优缺点	66
4.5 小结	67
参考文献	67
第 5 章 贝叶斯统计与熵	68
5.1 统计学习的基础：参数估计	68
5.1.1 矩估计	68



5.1.2	最大似然估计	69
5.1.3	最小二乘法	71
5.2	概率分布与三大统计思维	72
5.2.1	频率派和正态分布	72
5.2.2	经验派和正态分布	75
5.2.3	贝叶斯派和正态分布	76
5.2.4	贝叶斯统计和熵的关系	79
5.3	信息熵的理解	79
5.3.1	信息熵简史	79
5.3.2	信息熵定义	80
5.3.3	期望编码长度解释	81
5.3.4	不确定性公理化解释	81
5.3.5	基于熵的度量	84
5.4	最大熵原理	86
5.4.1	最大熵的直观理解	86
5.4.2	最大熵解释自然指数分布簇	87
5.4.3	最大熵解释最大似然估计	89
5.5	小结	90
	参考文献	91
第 6 章	基于熵的 Softmax	92
6.1	二项分布和多项分布	92
6.2	Logistic 回归和 Softmax 回归	93
6.2.1	广义线性模型的解释	93
6.2.2	Softmax 回归	94
6.2.3	最大熵原理与 Softmax 回归的等价性	96
6.3	最大熵条件下的 Log-Linear	101
6.4	多分类界面	103
6.4.1	感知机和多分类感知机	104
6.4.2	多分类感知机和结构感知机	105
6.5	概率图模型里面的 Log-Linear	106
6.6	深度学习里面的 Softmax 层	108
6.7	小结	109

参考文献	109
第 7 章 拉格朗日乘子法	111
7.1 凸共轭	111
7.1.1 凸共轭的定义	111
7.1.2 凸共轭定理	113
7.2 拉格朗日对偶	114
7.2.1 拉格朗日对偶概述	115
7.2.2 Salter 条件	117
7.2.3 KKT 条件	118
7.3 Fenchel 对偶	120
7.4 增广拉格朗日乘子法	123
7.4.1 近端	123
7.4.2 增广拉格朗日乘子法和对偶上升算法	126
7.5 交替方向乘子法	129
7.5.1 对偶分解	130
7.5.2 交替方向乘子法概述	131
7.6 小结	131
参考文献	132
第 8 章 随机梯度下降法	134
8.1 随机梯度下降法概述	134
8.1.1 机器学习场景	134
8.1.2 随机梯度下降法的定义	135
8.1.3 随机梯度下降法收敛性分析	136
8.1.4 收敛性证明	139
8.2 随机梯度下降法进阶 I：方差缩减	140
8.2.1 方差缩减的效果	141
8.2.2 方差缩减的实现	143
8.3 随机梯度下降法进阶 II：加速与适应	145
8.3.1 加速	146
8.3.2 适应	148
8.3.3 加速 × 适应	151
8.4 随机梯度下降法的并行实现	156