

TURING

图灵数学·统计学丛书

统计机器学习界泰斗作品

Statistical Learning with Sparsity:
The Lasso and Generalizations

稀疏统计学习 及其应用

[美] Trevor Hastie
Robert Tibshirani 著
Martin Wainwright
刘波 景鹏杰 译

全面介绍稀疏统计模型及其研究成果
用lasso模型解决大数据挖掘、机器学习等热点问题

CRC Press
Taylor & Francis Group

中国工信出版集团

人民邮电出版社
POSTS & TELECOM PRESS

非
外
借

TURING

图灵数学·统计学丛书

Statistical Learning with Sparsity:
The Lasso and Generalizations

稀疏统计学习 及其应用

[美] | Trevor Hastie
Robert Tibshirani | 著
Martin Wainwright

刘波 景鹏杰 译

人民邮电出版社

北京

图书在版编目(CIP)数据

稀疏统计学习及其应用/(美)特里瓦·哈斯蒂
(Trevor Hastie), (美)罗伯特·蒂伯沙拉尼
(Robert Tibshirani), (美)马丁·韦恩怀特
(Martin Wainwright)著; 刘波, 景鹏杰译—北京:
人民邮电出版社, 2018. 1

(图灵数学·统计学丛书)

ISBN 978-7-115-47261-8

I. ①稀… II. ①特… ②罗… ③马… ④刘… ⑤景
… III. ①统计学 IV. ①C8

中国版本图书馆CIP数据核字(2017)第282216号

内 容 提 要

稀疏统计模型只具有少数非零参数或权重, 经典地体现了化繁为简的理念, 因而广泛应用于诸多领域。本书就稀疏性统计学习做出总结, 以lasso方法为中心, 层层推进, 逐渐囊括其他方法, 深入探讨诸多稀疏性问题的求解和应用; 不仅包含大量的例子和清晰的图表, 还附有文献注释和课后练习, 是深入学习统计学知识的极佳参考。

本书适合算法、统计学和机器学习专业人士。

◆ 著 [美] Trevor Hastie, Robert Tibshirani,
Martin Wainwright

译 刘波 景鹏杰

责任编辑 朱巍

责任印制 彭志环

◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路11号

邮编 100164 电子邮件 315@ptpress.com.cn

网址 <http://www.ptpress.com.cn>

北京隆昌伟业印刷有限公司印刷

◆ 开本: 700×1000 1/16

印张: 18.75

字数: 368千字

印数: 1-3 000册

彩插: 4

2018年1月第1版

2018年1月北京第1次印刷

著作权合同登记号 图字: 01-2015-8299号

定价: 89.00元

读者服务热线: (010)51095186 转 600 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147号

站在巨人的肩上
Standing on Shoulders of Giants



iTuring.cn

站在巨人的肩上
Standing on Shoulders of Giants



iTuring.cn

版 权 声 明

Statistical Learning with Sparsity: The Lasso and Generalizations / by Trevor Hastie, Robert Tibshirani, Martin Wainwright / ISBN: 978-1-4987-1216-3.

© 2015 by Taylor & Francis Group. LLC

Authorized translation form English language edition published by CRC Press, an imprint of Taylor & Francis Group LCC. All rights reserved.

Post & Telecom Press is authorized to publish and distribute exclusively the Chinese (Simplified Characters) language edition. This edition is authorized for sale throughout Mainland of China. No part of the publication may be reproduced or distributed by any means, or stored in a database or retrieval system, without the prior written permission of the publisher.

Copies of this book sold without a Taylor & Francis sticker on the cover are unauthorized and illegal.

本书原版由 Taylor & Francis 出版集团旗下 CRC 出版公司出版, 并经其授权翻译出版。版权所有, 侵权必究。

本书中文简体翻译版授权由人民邮电出版社独家出版并仅限在中国大陆地区销售。未经出版者书面许可, 不得以任何方式复制或发行本书的任何部分。

本书封面贴有 Taylor & Francis 公司防伪标签, 无标签者不得销售。

谨以此书献给我们的父母亲：

Valerie Hastie、Patrick Hastie

Vera Tibshirani、Sami Tibshirani

Patricia Wainwright、John Wainwright

也献给我们的亲人：

Samantha、Timothy、Lynda

Charlie、Ryan、Jesss、Julie、Cheryl

Haruko、Hana

译者序

我们怀着无比敬仰的心情译完了这本书，因为本书的三位作者都是统计机器学习界的泰斗，他们在基于稀疏的统计学习理论方面的造诣有目共睹。能翻译此书是我们的荣幸。

随着大数据时代的到来，稀疏性成为研究大数据的重要手段。斯坦福大学的 Robert Tibshirani 于 1996 年首次提出将 l_1 范数作为普通最小二乘的正则项，从而得到了著名的 lasso 模型。lasso 在拟和数据的同时，也利用 l_1 范数让系数具有稀疏性来选择特征。这种稀疏性可增加模型的可解释性，并提高计算效率，这些特性对高维数据尤其重要。因此， l_1 范数的这一特性很快就吸引了大批学者进行研究，并在各种机器学习模型中广泛使用，这些模型不仅仅包括广义线性模型，还包括图模型、信号处理中的压缩感知模型、矩阵填充模型等。 l_1 范数是一个凸函数，但不可微，因此，在求解以 l_1 范数作为惩罚项的模型时通常比较麻烦。用来求解这类模型的常用算法是次梯度方法，但这种方法的效率通常不高，尤其在大规模的高维数据上更是如此。最近又出现了大量求解该问题的高效求解算法，比如坐标下降法等。

本书从最基本的 lasso 模型出发介绍了 lasso 模型的物理意义和相应的求解算法，然后介绍了 lasso 模型的推广和最新的研究成果。本书的特点是：(1) 采用深入浅出、图文并茂的方式介绍各种抽象理论；(2) 内容新颖，书中大量内容都是最近几年的研究成果，比如筛选规则等；(3) 涉及面广，书中的内容包含了与稀疏性相关的各个重要的研究领域，比如多元统计中的稀疏性问题等。因此，本书是广大机器学习研究人员、工程人员不可多得的参考书，也可用于研究生的机器学习课程的教材。

本书的第 1 章、第 2 章、第 5 章、第 7 章、第 8 章、第 10 章由重庆工商大学计算机科学与信息工程学院的刘波博士翻译，第 3 章、第 4 章、第 6 章、第 9 章、第 11 章由上海期货信息技术有限公司的景鹏杰翻译。

由于本书涉及的知识面广，内容也很新颖，因此许多术语尚无固定译法。我们虽然经过反复推敲和讨论，但仍然可能出现词不达意的情况。同时由于时间和精力有限，书中内容难免不出差错。若有问题或建议，读者可通过电子邮件 liubo7971@163.com 或 pjjing@foxmail.com 与我们联系，欢迎大家对本书的翻译进行指正或提出宝贵的建议。本书翻译的勘误信息会发布在 <http://www.cnblogs.com/ml-cv/> 上，欢迎关注。

本书翻译过程得到如下项目资助：(1) 重庆市教委研究项目“多核正则化机器学习理论研究”，项目号为 KJ130709；(2) 重庆工商大学研究项目“基于多核学习的高维数据分析研究”，项目号为 2013-56-09；(3) 大数据稀疏表示判别字典学习及其应用技术研究，项目号为 KJ1400612；(4) 重庆工商大学研究生院教改项目“基于二维码的研究生互动教学改革”；(5) 电子商务及供应链系统重庆市重点实验室项目“基于迹比率的特征选择及关键技术研究”，项目编号为 1456025。

感谢翻译过程中图灵公司的朱巍编辑给予的帮助和支持；感谢高敬雅老师，她给予了我们在统计学上的帮助；感谢重庆工商大学计信学院金融信息化专业的曾芳同学，她帮助我们录入了本书的公式。感谢刘波的妻子杨雪莉的支持，感谢刘波两个小女儿刘典、刘恩丫（此书完成时她才 7 个月大）对他的忍耐。

前 言

在这本专著中，我们将概述基于稀疏性的统计学习的最新研究。稀疏的统计模型仅具有少数非零参数或权值。它代表了“少即是多”的经典情形：与稠密模型相比，稀疏模型更容易估计和解释。在这个大数据的时代，对一个人或目标进行度量的特征数量可能更多，而且有可能比观测样本数更多。借助稀疏性假设，我们能够解决这些问题，并从大数据集中提取有用的、可重复性模式。

这里所陈述的观点代表了统计学和机器学习方面全体研究人员的工作，我们感谢大家对这个令人激动的领域所作出的贡献。我们要特别感谢斯坦福大学和加州大学伯克利分校的同事，感谢我们的合作者，以及现在及过去在这个领域工作的同学们。他们分别是：Alekh Agarwal, Arash Amini, Francis Bach, Jacob Bien, Stephen Boyd, Andreas Buja, Emmanuel Candes, Alexandra Chouldechova, David Donoho, John Duchi, Brad Efron, Will Fithian, Jerome Friedman, Max G'Sell, Iain Johnstone, Michael Jordan, Ping Li, Po-Ling Loh, Michael Lim, Jason Lee, Richard Lockhart, Rahul Mazumder, Balasubramanian Narashimhan, Sahand Negahban, Guillaume Obozinski, Mee-Young Park, Junyang Qian, Garvesh Raskutti, Pradeep Ravikumar, Saharon Rosset, Prasad Santhanam, Noah Simon, Dennis Sun, Yukai Sun, Jonathan Taylor, Ryan Tibshirani^①, Stefan Wager, Daniela Witten, Bin Yu, Yuchen Zhang, Ji Zhou, and Hui Zou。感谢我们的编辑 John Kimmel 对本书的建议和支持。

Trevor Hastie, Robert Tibshirani, 斯坦福大学
Martin Wainwright, 加州大学伯克利分校

^①本书文献中提到的 Tibshirani₂ 是 Ryan Tibshirani，而 Tibshirani 则是其父 Robert Tibshirani。

目 录

第 1 章	引言	1
第 2 章	lasso 线性模型	6
2.1	引言	6
2.2	lasso 估计	7
2.3	交叉验证和推断	10
2.4	lasso 解的计算	12
2.4.1	基于单变量的软阈值法	12
2.4.2	基于多变量的循环坐标下降法	13
2.4.3	软阈值与正交基	15
2.5	自由度	15
2.6	lasso 解的唯一性	16
2.7	理论概述	17
2.8	非负 garrote	17
2.9	ℓ_q 惩罚和贝叶斯估计	19
2.10	一些观点	20
	习题	21
第 3 章	广义线性模型	24
3.1	引言	24
3.2	逻辑斯蒂回归模型	26
3.2.1	示例: 文本分类	27
3.2.2	算法	29
3.3	多分类逻辑斯蒂回归	30
3.3.1	示例: 手写数字	31
3.3.2	算法	32
3.3.3	组 lasso 多分类	33
3.4	对数线性模型及泊松广义线性模型	33
3.5	Cox 比例风险模型	35
3.5.1	交叉验证	37
3.5.2	预验证	38
3.6	支持向量机	39

3.7 计算细节及 glmnet	43
参考文献注释	44
习题	45
第 4 章 广义 lasso 惩罚	47
4.1 引言	47
4.2 弹性网惩罚	47
4.3 组 lasso	50
4.3.1 组 lasso 计算	53
4.3.2 稀疏组 lasso	54
4.3.3 重叠组 lasso	56
4.4 稀疏加法模型和组 lasso	59
4.4.1 加法模型和 backfitting	59
4.4.2 稀疏加法模型和 backfitting	60
4.4.3 优化方法与组 lasso	61
4.4.4 稀疏加法模型的多重惩罚	64
4.5 融合 lasso	65
4.5.1 拟合融合 lasso	66
4.5.2 趋势滤波	69
4.5.3 近保序回归	70
4.6 非凸惩罚	72
参考文献注释	74
习题	75
第 5 章 优化方法	80
5.1 引言	80
5.2 凸优化条件	80
5.2.1 优化可微问题	80
5.2.2 非可微函数和次梯度	83
5.3 梯度下降	84
5.3.1 无约束的梯度下降	84
5.3.2 投影梯度法	86
5.3.3 近点梯度法	87
5.3.4 加速梯度方法	90
5.4 坐标下降	92
5.4.1 可分性和坐标下降	93
5.4.2 线性回归和 lasso	94

5.4.3	逻辑斯蒂回归和广义线性模型	97
5.5	仿真研究	99
5.6	最小角回归	100
5.7	交替方向乘子法	103
5.8	优化-最小化算法	104
5.9	双凸问题和交替最小化	105
5.10	筛选规则	108
	参考文献注释	111
	附录 A lasso 的对偶	112
	附录 B DPP 规则的推导	113
	习题	114
第 6 章	统计推断	118
6.1	贝叶斯 lasso	118
6.2	自助法	121
6.3	lasso 法的后选择推断	125
6.3.1	协方差检验	125
6.3.2	选择后推断的更广方案	128
6.3.3	检验何种假设	133
6.3.4	回到向前逐步回归	134
6.4	通过去偏 lasso 推断	134
6.5	后选择推断的其他建议	136
	参考文献注释	137
	习题	138
第 7 章	矩阵的分解、近似及填充	141
7.1	引言	141
7.2	奇异值分解	142
7.3	缺失数据和矩阵填充	143
7.3.1	Netflix 电影挑战赛	144
7.3.2	基于原子范数的矩阵填充	146
7.3.3	矩阵填充的理论结果	149
7.3.4	最大间隔分解及相关方法	153
7.4	减秩回归	154
7.5	通用矩阵回归框架	156
7.6	惩罚矩阵分解	157
7.7	矩阵分解的相加形式	160

参考文献注释	164
习题	165
第 8 章 稀疏多元方法	169
8.1 引言	169
8.2 稀疏组成成分分析	169
8.2.1 背景	169
8.2.2 稀疏主成分	171
8.2.3 秩大于 1 的解	174
8.2.4 基于 Fantope 投影的稀疏 PCA	176
8.2.5 稀疏自编码和深度学习	176
8.2.6 稀疏 PCA 的一些理论	178
8.3 稀疏典型相关分析	179
8.4 稀疏线性判别分析	182
8.4.1 标准理论和贝叶斯规则	182
8.4.2 最近收缩中心	183
8.4.3 Fisher 线性判别分析	184
8.4.4 最佳评分	188
8.5 稀疏聚类	190
8.5.1 聚类的一些背景知识	191
8.5.2 稀疏层次聚类	191
8.5.3 稀疏 K 均值聚类	192
8.5.4 凸聚类	193
参考文献注释	195
习题	196
第 9 章 图和模型选择	202
9.1 引言	202
9.2 图模型基础	202
9.2.1 分解和马尔可夫特性	202
9.2.2 几个例子	204
9.3 基于惩罚似然的图选择	206
9.3.1 高斯模型的全局似然性	207
9.3.2 图 lasso 算法	208
9.3.3 利用块对角化结构	210
9.3.4 图 lasso 的理论保证	211
9.3.5 离散模型的全局似然性	212

9.4 基于条件推断的图选择	213
9.4.1 高斯分布下基于近邻的似然概率	214
9.4.2 离散模型下基于近邻的似然概率	214
9.4.3 混合模型下的伪似然概率	217
9.5 带隐变量的图模型	218
参考文献注释	219
习题	221
第 10 章 信号近似与压缩感知	225
10.1 引言	225
10.2 信号与稀疏表示	225
10.2.1 正交基	225
10.2.2 用正交基逼近	228
10.2.3 用过完备基来重构	229
10.3 随机投影与近似	231
10.3.1 Johnson-Lindenstrauss 近似	231
10.3.2 压缩感知	232
10.4 ℓ_0 恢复与 ℓ_1 恢复之间的等价性	234
10.4.1 受限零空间性质	235
10.4.2 受限零空间的充分条件	235
10.4.3 证明	237
参考文献注释	238
习题	239
第 11 章 lasso 的理论结果	242
11.1 引言	242
11.1.1 损失函数类型	242
11.1.2 稀疏模型类型	243
11.2 lasso ℓ_2 误差的界限	244
11.2.1 经典情形中的强凸性	244
11.2.2 回归受限特征值	245
11.2.3 基本一致性结果	246
11.3 预测误差的界	250
11.4 线性回归中的支持恢复	252
11.4.1 lasso 的变量选择一致性	252
11.4.2 定理 11.3 的证明	256
11.5 超越基础 lasso	259

参考文献注释	260
习题	261
参考文献	264

第1章 引言

“我从来都不记记分卡或击球率。我讨厌统计。我会把必须知道的东西记在脑海里。”

这段话是棒球投手 Dizzy Dean 说的，他曾在 1930 ~1947 年参加美国职业棒球大联盟的比赛。

一晃 75 年过去了，世界发生了很大的变化！如今，人们在科学、娱乐、商业和工业各领域收集和挖掘大量数据，并对其进行研究和应用。医学家们通过研究患者的基因组选择最佳的治疗方法，并由此了解这些疾病产生的根本原因。在线电影和网上书店会研究客户的评价，以便向他们推荐新的电影或书籍。社交网络会研究其会员及好友的资料，优化在线体验。而且，现在多数大联盟棒球队都有统计员收集和分析击球手和投手的详细信息，帮助球队经理和队员做出更好的决策。

由此可知，这个世界淹没在了数据中。而 Rutherford D. Roger 等人则说：

“我们淹没在了信息的海洋里，却渴求着知识。”

海量信息亟待整理，取其精华去其糟粕。为了成功完成这项工作，人们期望真实情况得以简化：也许人体内大约 30 000 个基因并非都与癌症的发展过程直接相关；也许只需要客户对 50 或 100 部电影做出评价就足以揭示他们的爱好；也许左撇子投手对付左撇子击球手会比较轻松。

这些情形背后都有简单性假设。稀疏性 (sparsity) 是简单性的一种形式，这也是本书的中心主题。简而言之，在一个稀疏统计模型中，仅有较少参数（也称预测子，predictor）在发挥重要作用。本书将介绍如何利用稀疏性来恢复一组数据中的基础信号。

最典型的例子是线性回归，即有 N 组观测值，每组观测值由一个输出变量 y_i 和 p 个相关预测子变量（也称特征） $x_i = (x_{i1}, \dots, x_{ip})^T$ 所组成。线性回归的目标是通过预测子来预测输出值，既要正确预测将来的数据，又要找出哪些预测子在起重要作用。一个线性回归模型可设为：

$$y_i = \beta_0 + \sum_{j=1}^p x_{ij}\beta_j + e_i \quad (1.1)$$

其中， β_0 和 $\beta = (\beta_1, \beta_2, \dots, \beta_p)$ 是未知参数， e_i 为误差项。这些参数可用最小二乘法来估计，即最小化最小二乘目标函数：

$$\underset{\beta_0, \beta}{\text{minimize}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j \right)^2 \quad (1.2)$$