

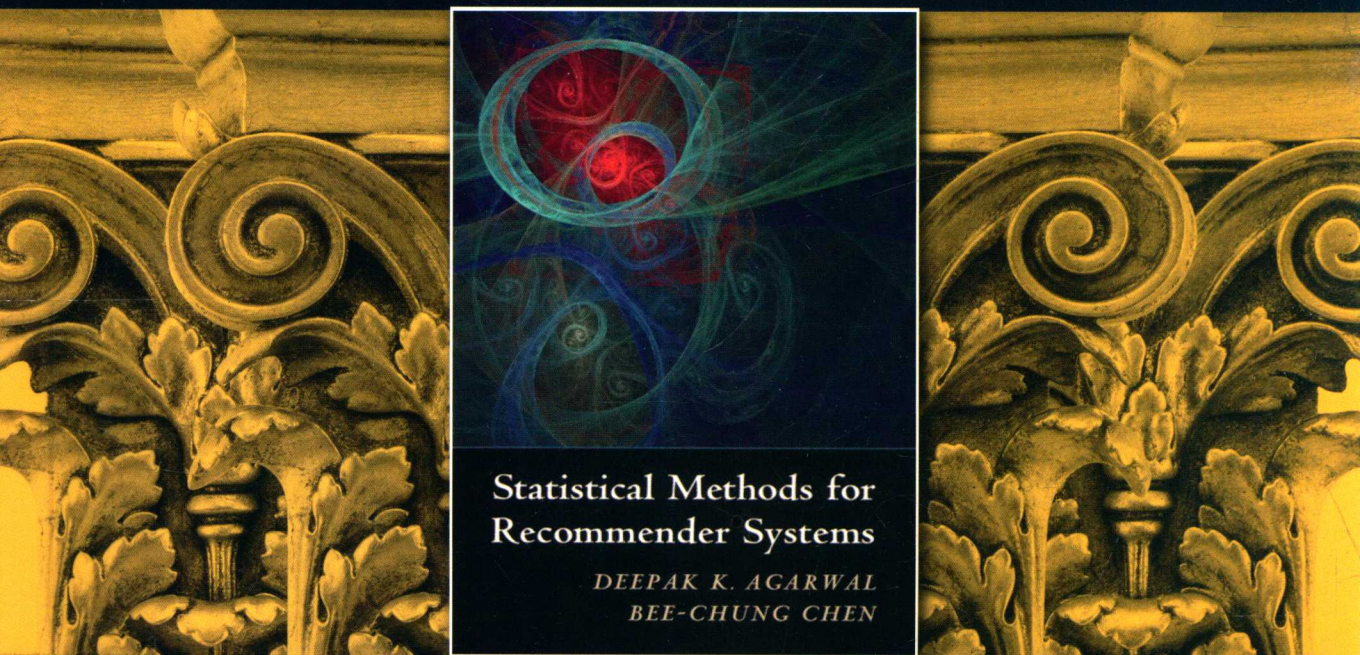
统计推荐系统

迪帕克·K.阿加瓦尔 陈必衷
(Deepak K. Agarwal) (Bee-Chung Chen) 著

LinkedIn公司

戴薇 潘微科 明仲 译
深圳大学

涵盖算法理论、实验分析和结果展示，分享大规模推荐系统的开发经验



Statistical Methods for Recommender Systems



机械工业出版社
China Machine Press

计

算

书

统计推荐系统

迪帕克·K.阿加瓦尔 陈必衷

(Deepak K. Agarwal) (Bee-Chung Chen)

[美]

著

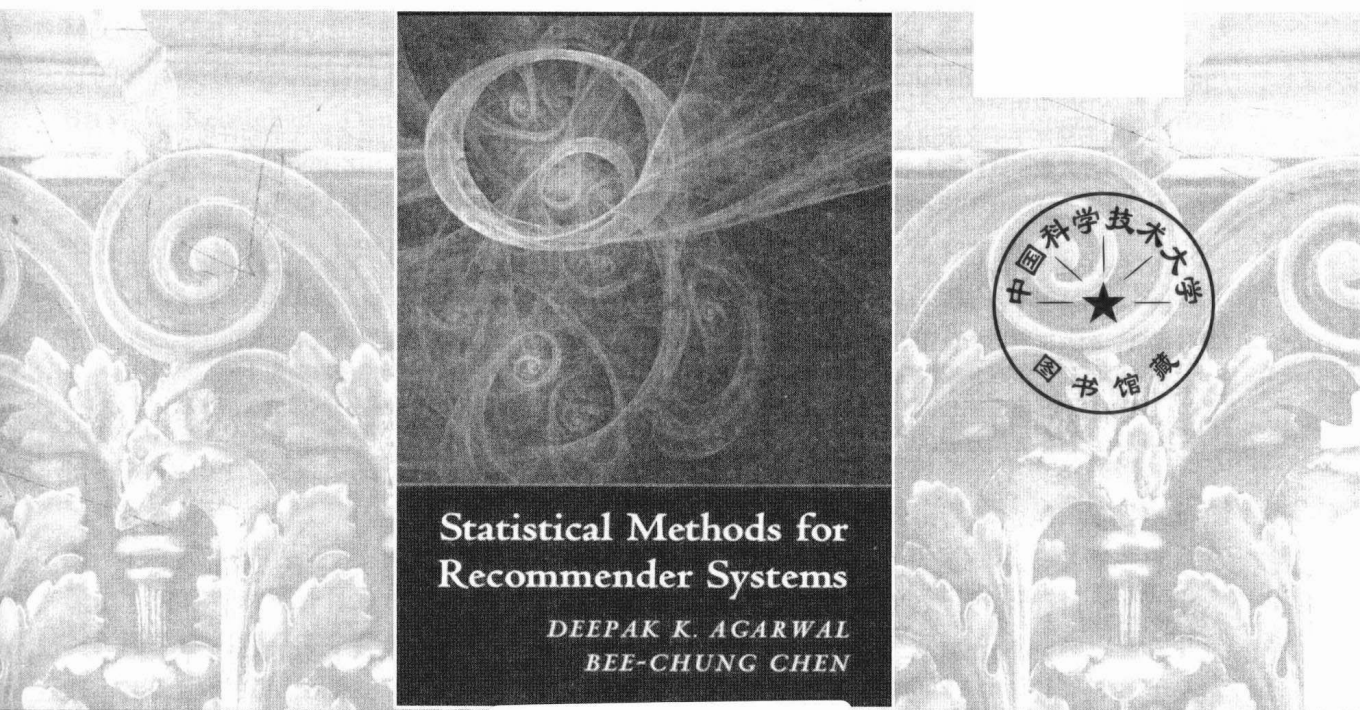
LinkedIn公司

戴薇 潘微科 明仲

译

深圳大学

Statistical Methods for Recommender Systems



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

统计推荐系统 / (美) 迪帕克·K. 阿加瓦尔 (Deepak K. Agarwal) 等著; 戴薇, 潘微科, 明仲译. —北京: 机械工业出版社, 2019.9

(计算机科学丛书)

书名原文: Statistical Methods for Recommender Systems

ISBN 978-7-111-63573-4

I. 统… II. ①迪… ②戴… ③潘… ④明… III. 统计程序 IV. TP319

中国版本图书馆 CIP 数据核字 (2019) 第 188819 号

本书版权登记号: 图字 01-2019-0736

This is a Simplified-Chinese edition of the following title published by Cambridge University Press:

Deepak K. Agarwal, Bee-Chung Chen: Statistical Methods for Recommender Systems (ISBN 978-1-107-03607-9).

© Deepak K. Agarwal and Bee-Chung Chen 2016.

This Simplified-Chinese 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 2019.

This Simplified-Chinese 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。

本书由 LinkedIn 公司的技术专家撰写, 着眼于推荐系统的核心——统计方法, 不仅讲解理论知识, 而且分享了作者在 LinkedIn 和 Yahoo! 的实践经验。全书分为三部分: 第一部分介绍推荐系统的组成、经典推荐方法及评估方法, 并引出了探索与利用问题; 第二部分围绕点击通过率 (CTR) 预估这一重要问题, 重点介绍快速在线双线性因子模型和面向回归的隐因子模型, 为热门推荐和个性化推荐提供解决方案; 第三部分讨论进阶主题, 涵盖分解的隐含狄利克雷分布模型、张量分解模型、层次收缩模型以及多目标优化方法。

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: 赵 静

责任校对: 李秋荣

印刷: 北京文昌阁彩色印刷有限责任公司

版次: 2019 年 9 月第 1 版第 1 次印刷

开本: 185mm×260mm 1/16

印张: 14.5

书号: ISBN 978-7-111-63573-4

定价: 89.00 元

客服电话: (010) 88361066 88379833 68326294

投稿热线: (010) 88379604

华章网站: www.hzbook.com

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

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，从它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、Afred 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，有的还专门为其书的中译本作序。迄今，“计算机科学丛书”已经出版了近500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，我们的联系方式如下：

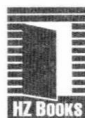
华章网站：www.hzbook.com

电子邮件：hzsj@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

近几年，推荐系统发展之迅猛超乎想象，正如作者所说，“推荐系统无处不在，已然成为我们日常生活的一部分。”无论是在工业界还是在学术界，人们对探索推荐系统的热情都有增无减。本书着眼于推荐系统的核心部分——统计方法，虽然是介绍算法理论，但作者曾领导团队开发过雅虎和领英的多个推荐系统，对于算法在实际系统中的应用也有着独到的见解。因此，本书不仅包含对统计方法的详解，还包含详尽的实验分析和丰富的结果展示，这无疑为需求不一的读者提供了极大的便利。

在得知有机会翻译这本专业性极强的书籍后，激动之情难以言表：能够为那些因语言受限而在推荐系统的研究中苦苦摸索的同仁带去一丝光亮，这是一件多么有意义的事情！激动之余也意识到，正因为意义非凡，我们更应该以严肃认真的态度对待它。推荐系统涉及的专业知识范围广、难度大，对于较生疏的内容，我们也借此机会填补空缺，拓展知识面。尽管如此，也难免存在不足之处，敬请读者批评指正。

在此感谢机械工业出版社华章公司的曲熠编辑和赵静编辑，以及参与部分章节审校的陈子翔、陈宪聪、黄云峰、林晶、刘基雄、廖道虢、梁锋、马万绮、章湘鑫、钟柳兰和庄燊等学生。另外，感谢国家自然科学基金项目的支持（NO. 61872249, NO. 61836005）。

戴薇 潘微科 明仲

2019年5月于深圳大学大数据系统计算技术国家工程实验室

这本书讲什么

推荐系统是一类自动化的计算机程序，能够在不同场景下将物品和用户进行匹配。推荐系统无处不在，已然成为我们日常生活的一部分。例如，亚马逊购物网站上的产品推荐，雅虎上的内容推荐，Netflix 上的电影推荐，领英上的工作推荐等。匹配算法的构建需要用到大量高频数据，它们来源于用户与物品的历史交互行为。从本质上来看，推荐算法属于统计学范畴，在序贯决策过程、高维类别数据的建模以及开发可伸缩的统计方法等领域都面临着挑战。在推荐系统领域，算法的推陈出新依赖于计算机科学家、机器学习专家、统计学家、优化专家、系统专家，当然还有领域专家之间的密切合作。可以说，推荐系统是大数据领域最振奋人心的应用之一。

我们为什么写这本书

虽然计算机科学、机器学习和统计学等领域已有大量关于推荐系统的书籍，但它们仅针对问题的某些特定方面，没有综合考虑所有的统计问题，也没有分析这些统计问题是如何相互关联的。而我们也是在雅虎和领英部署推荐系统时才意识到这个问题，例如，统计学和机器学习的重点在于最小化样本外的预测误差，但达成这个目标并不意味着实践中的所有重要问题都得到了解决。从统计学意义上来说，推荐系统是一个高维序贯过程，研究实验设计类问题与开发精密的统计模型一样重要。事实上，这两者关系密切，高效的实验设计需要借助模型克服维数灾难。此外，大多数现有工作倾向于对单一反馈建模，例如电影评分、购买和点击率。但随着 Facebook、领英和推特等社交媒体的兴起，多种反馈随之而来，例如，一个新闻推荐应用可能需要同时对用户的点击率、分享率和发文率这三类数据建模。这种面向多种反馈的建模是很有挑战性的。最后的问题是，即便我们获得了能够实现这种多变量预测的方法，又该如何构建效用函数去完成推荐呢？优化分享率比优化点击率更重要吗？关于这些问题的解答，我们可以与多目标优化领域的专家密切合作，利用多目标优化来获得一些效用参数。

本书的目的是对推荐系统中的问题进行全面讨论，另外，也对当前最先进的统计方法，如自适应序贯设计（多臂赌博机方法）、双线性随机效应模型（矩阵分解）以及现代的基于分布式计算框架的可伸缩模型，进行详细且深入的探讨。我们希望通过本书分享我们在工业界开发大规模推荐系统的丰富经验，也希望能够引起统计学、机器学习和计算机科学等领域相关人士的关注。我们相信，这对许多方面都是有益的。本书有助于推进高维大数据统计的研究，这类研究尤其有利于 Web 应用的发展。此类学术研究离不开处理海量数据的软件，为此，我们将本书用到的隐因子模型的代码公布在以下网址：<https://github.com/beechnung/Latent-Factor-Models>。我们也相信本书能够成为连接理论研究与实际应用的桥梁。一方面，本书可以帮助对推荐有疑惑的学者理解推荐系统中的统计知识；另一方面，如果建模人员在实际应用中遇到复杂的统计问题，本书也能提供深入的解答。

章节组织结构

本书共分为三个部分。

在第一部分中，我们将介绍推荐系统问题、存在的挑战、应对挑战的主要思路以及所需的背景知识。在第 2 章中，我们将概述几种开发推荐系统的经典方法。这些方法将用户和物品表示为特征向量，然后通过一些相似度计算函数、标准监督学习或协同过滤来预测用户 - 物品的评分。这些经典方法通常会忽略推荐问题中探索与利用之间的权衡。因此，我们将在第 3 章论述在推荐系统中权衡探索与利用的重要性，并介绍用它解决后面章节中问题的主要思路。在深入研究技术性方案之前，我们将在第 4 章回顾一些用于评估不同推荐算法性能的方法。

在第二部分中，我们将提供针对常见问题设置的详细解决方案。在第 5 章中，我们将介绍不同的问题设置，并展示一个系统架构案例。接下来的三章分别对应三个常见的问题设置。第 6 章将为热门推荐问题提供几种解决方案，尤其注重探索和利用之间的权衡。第 7 章将基于特征回归解决个性化推荐问题，重点在于如何利用最新的用户 - 物品交互数据不断更新模型，使其快速收敛至最优。第 8 章将第 7 章中基于特征的回归模型扩展成因子模型（矩阵分解），同时还将为因子模型中的冷启动问题提供一个合适的解决方案。

在第三部分中，我们将讨论三个进阶主题。在第 9 章中，我们将介绍一个结合隐含狄利克雷分布（LDA）主题模型的矩阵分解模型，该模型可以同时确定物品蕴涵的主题和用户对不同主题的偏好度。在第 10 章中，我们将研究上下文相关推荐问题，即物品不仅需要与用户具有高度的关联性，还必须与上下文相关（例如，推荐与用户正在阅读的新闻相关的物品）。在第 11 章中，我们将讨论一个基于约束优化方法的多目标优化框架，试图在其他目标的有界损失范围内（例如，点击损失不超过 5%）最大化某一特定目标（例如，收入）。

缺点

与其他书籍一样，本书也难免存在不足。首先，我们没有深入涉及现代计算框架，比如可以用来拟合一定规模模型的 Spark 框架。其次，如果用户构成了一个社交网络，那么传统的实验设计方法无法用于模型的在线评估，这就需要我们开发适用于社交图谱推理的新技术。以上这些进阶主题都不在本书的范围内。全书从始至终都将基于回归的响应预测方法作为主要工具来解决推荐问题，主要是因为这些模型的输出很容易被后续方法所使用。所以，我们也没有详细讨论直接优化排序损失函数的方法。当然，这两种方法的对比也是一个值得探讨的话题。

致谢

特别感谢 Raghu Ramakrishnan、Liang Zhang、Xuanhui Wang、Pradheep Elango、Bo Long、Bo Pang、Rajiv Khanna、Nitin Motgi、Seung-Taek Park、Scott Roy、Joe Zachariah，我们多次与他们合作，进行了深入的探讨。我们还要感谢雅虎和领英的同事们的鼓励和支持，没有他们，我们的许多想法将难以付诸实现。

出版者的话
译者序
前言

第一部分 基础知识

第 1 章 简介	2	2.5 混合方法	27
1.1 面向网络应用的推荐系统概述	3	2.6 小结	28
1.1.1 算法	3	2.7 练习	28
1.1.2 优化指标	5	第 3 章 面向推荐问题的探索与利用	29
1.1.3 探索与利用之间的权衡	5	3.1 探索与利用之间的权衡简介	30
1.1.4 推荐系统的评估	5	3.2 多臂赌博机问题	31
1.1.5 推荐和搜索：推送与拉取	6	3.2.1 贝叶斯方法	31
1.2 一个简单的评分模型：热门推荐	7	3.2.2 极小化极大方法	34
1.3 练习	10	3.2.3 启发式赌博方案	35
第 2 章 经典推荐方法	11	3.2.4 方法评价	36
2.1 物品特征	11	3.3 推荐系统中的探索与利用	36
2.1.1 分类	12	3.3.1 热门推荐	36
2.1.2 词袋模型	13	3.3.2 个性化推荐	36
2.1.3 主题建模	15	3.3.3 数据稀疏性的挑战	37
2.1.4 其他物品特征	16	3.4 处理数据稀疏性的探索与利用	37
2.2 用户特征	16	3.4.1 降维方法	37
2.2.1 声明的个人信息	17	3.4.2 降维中的探索与利用	39
2.2.2 基于内容的画像	17	3.4.3 在线模型	39
2.2.3 其他用户特征	18	3.5 小结	40
2.3 基于特征的方法	18	3.6 练习	40
2.3.1 无监督方法	18	第 4 章 评估方法	41
2.3.2 有监督方法	19	4.1 传统的离线评估方法	41
2.3.3 上下文信息	22	4.1.1 数据划分方法	42
2.4 协同过滤	22	4.1.2 准确度指标	44
2.4.1 基于用户 - 用户相似度的方法	23	4.1.3 排序指标	45
2.4.2 基于物品 - 物品相似度的方法	24	4.2 在线分桶测试	49
2.4.3 矩阵分解	24	4.2.1 设置分桶测试	49
		4.2.2 在线性能指标	50
		4.2.3 测试结果分析	51
		4.3 离线模拟	52
		4.4 离线回放	54

4.4.1 基本回放估计	55	7.1.2 FOBFM 详解	91
4.4.2 回放的扩展	57	7.2 离线训练	93
4.5 小结	58	7.2.1 EM 算法	94
4.6 练习	58	7.2.2 E 步骤	95
		7.2.3 M 步骤	96
		7.2.4 可扩展性	97
第二部分 常见问题设置		7.3 在线学习	97
第 5 章 问题设置与系统架构	60	7.3.1 在线高斯模型	97
5.1 问题设置	60	7.3.2 在线逻辑模型	98
5.1.1 常见的推荐模块	60	7.3.3 探索与利用方案	99
5.1.2 应用设置	63	7.3.4 在线模型选择	99
5.1.3 常见的统计方法	65	7.4 雅虎数据集上的效果展示	100
5.2 系统架构	66	7.4.1 My Yahoo! 数据集	101
5.2.1 主要组件	66	7.4.2 雅虎首页数据集	103
5.2.2 示例系统	67	7.4.3 不包含离线双线性项的 FOBFM	105
第 6 章 热门推荐	69	7.5 小结	105
6.1 应用案例：雅虎“今日”模块	69	7.6 练习	106
6.2 问题定义	71	第 8 章 基于因子模型的个性化	107
6.3 贝叶斯方案	72	8.1 面向回归的隐因子模型	107
6.3.1 2×2 案例：两件物品， 两个间隔	73	8.1.1 从矩阵分解到 RLFM	108
6.3.2 $K \times 2$ 案例： K 件物品， 两个间隔	75	8.1.2 模型详解	109
6.3.3 一般解	77	8.1.3 RLFM 的随机过程	112
6.4 非贝叶斯方案	79	8.2 拟合算法	113
6.5 实验评估	81	8.2.1 适用于高斯响应的 EM 算法	114
6.5.1 比较分析	81	8.2.2 适用于逻辑响应的基于 ARS 的 EM 算法	118
6.5.2 方案刻画	83	8.2.3 适用于逻辑响应的变分 EM 算法	121
6.5.3 分段分析	85	8.3 冷启动效果展示	124
6.5.4 桶测试结果	86	8.4 时间敏感物品的大规模推荐	127
6.6 大规模内容池	87	8.4.1 在线学习	127
6.7 小结	87	8.4.2 并行拟合算法	128
6.8 练习	88	8.5 大规模问题效果展示	130
第 7 章 基于特征回归的个性化	89	8.5.1 MovieLens-1M 数据	131
7.1 快速在线双线性因子模型	90	8.5.2 小规模雅虎首页数据	132
7.1.1 FOBFM 概述	90	8.5.3 大规模雅虎首页数据	134

8.5.4 结果讨论	137
8.6 小结	138
8.7 练习	138

第三部分 进阶主题

第9章 基于隐含狄利克雷分布的分解

9.1 简介	140
9.2 模型	141
9.2.1 模型概述	141
9.2.2 模型详解	142
9.3 训练和预测	145
9.3.1 模型拟合	145
9.3.2 预测	150
9.4 实验	150
9.4.1 MovieLens 数据	150
9.4.2 Yahoo! Buzz 应用	151
9.4.3 BookCrossing 数据集	153
9.5 相关工作	154
9.6 小结	155

第10章 上下文相关推荐

10.1 张量分解模型	157
10.1.1 建模	157
10.1.2 模型拟合	158
10.1.3 讨论	159
10.2 层次收缩模型	160
10.2.1 建模	160

10.2.2 模型拟合	161
10.2.3 局部增强张量模型	164
10.3 多角度新闻文章推荐	165
10.3.1 探索性数据分析	166
10.3.2 实验评估	171
10.4 相关物品推荐	176
10.4.1 语义相关性	177
10.4.2 响应预测	177
10.4.3 预测响应和预测相关性的结合	178
10.5 小结	178

第11章 多目标优化

11.1 应用设置	179
11.2 分段方法	180
11.2.1 问题设置	180
11.2.2 目标优化	181
11.3 个性化方法	183
11.3.1 原始表示	184
11.3.2 拉格朗日对偶	185
11.4 近似方法	188
11.4.1 聚类	188
11.4.2 采样	189
11.5 实验	189
11.5.1 实验设置	190
11.5.2 实验结果	191
11.6 相关工作	197
11.7 小结	198

参考文献

参考文献	199
索引	205

| 第一部分 |

Statistical Methods for Recommender Systems

基础知识

简介

推荐系统是一类在不同的上下文中为用户推荐“最佳”物品的计算机程序。“最佳”匹配通常可以通过优化一些特定目标而得到，如总点击数、总收入、总销售额等。推荐系统在网络上无处不在，已经成为我们日常生活的组成部分。例如：电商网站为了最大化销售额，会向用户推荐商品；新闻网站为了最大化总点击数，会向访问的用户推荐新闻内容；视频网站为了最大化用户参与度，同时提高订阅量，会向用户推荐电影；求职网站为了最大化工作申请数，会向用户推荐工作。以上这些算法的输入通常包含与用户、物品、上下文有关的信息以及用户与物品发生交互时获取的反馈信息。

图 1-1 展示了一个典型的网络推荐系统示例。首先，用户通过浏览器访问某网页，然后浏览器向网站服务器提交 HTTP 请求。为了在页面上进行推荐（如新闻门户页面上的热门新闻报道），网站服务器会调用推荐服务，推荐服务会检索出一组物品，并将其展示在网页上。这样一项推荐服务往往需要完成大量不同类型的运算才能挑选出最佳物品。这些运算通常混合了离线运算和实时运算，并且为了确保页面加载足够迅速（通常为几百毫秒），它们必须严格符合效率要求。一旦网页加载成功，用户就能与物品进行交互，如点击、喜欢或分享。从交互行为中获得的数据反过来又用于更新底层推荐算法的参数，以便为未来访问网站的用户提供更精准的推荐服务。参数更新的频率与应用有关，以新闻推荐为例，新闻报道对时间敏感，且生存期短暂，必须经常更新参数（例如每隔几分钟）；而对于生存期较长的应用（如电影推荐），参数更新不频繁（如一天更新一次）也不会对系统的整体推荐效果造成太大影响。

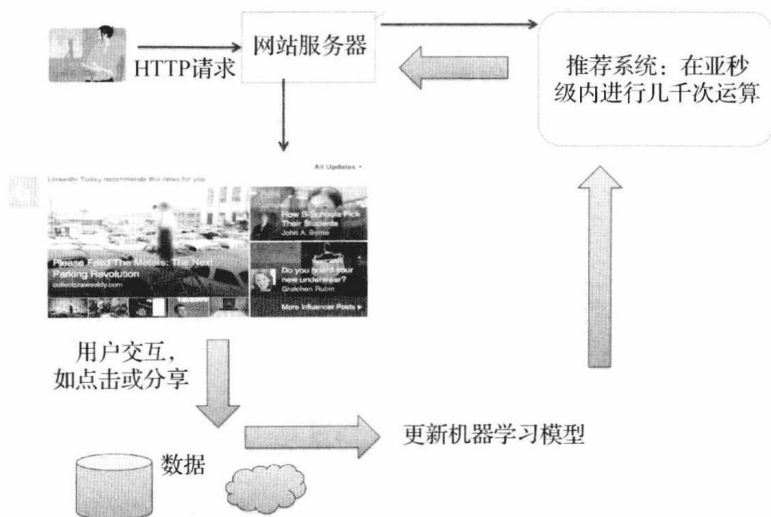


图 1-1 典型的推荐系统

在成功的推荐系统的底层，一定有一个挑选最佳物品的好算法。本书全面介绍了这类算法涉及的统计和机器学习方法。简单起见，我们在本书中粗略地将这些算法称为推荐系统，但请注意，它们仅仅是在客户端与服务端交互过程中为用户推荐物品的一个可伸缩组件。不过，我们并不能否认它们的重要性。

1.1 面向网络应用的推荐系统概述

在开发推荐系统之前，我们先考虑以下几个问题：

- 可用的输入信息有哪些？在构建用于预测用户在给定的上下文中可能与哪些物品发生交互的机器学习模型时，我们可以利用很多信息，包括：每件物品的内容和来源；用户的兴趣画像（既反映了用户的历史访问数据中隐含的长期兴趣，也反映了用户在当前会话中表现出的短期兴趣）；用户已声明的信息，如人口统计信息；还有“流行度”指标，例如观测到的点击通过率（即 CTR，表示物品被点击的次数与物品展示给用户的次数之比）；以及社交分享度，如物品被转推、分享或喜欢的次数。
- 可优化的目标有哪些？供网站选择的优化目标有很多，可分为短期目标和长期目标。短期目标如点击数、收入或用户的正向显式评分；长期目标如在网站上花费的时间的延长、用户回头率和留存率的提高、社交行为的增加、订阅量的增长等。

各种不同的推荐算法便是基于以上问题的答案开发出来的。

1.1.1 算法

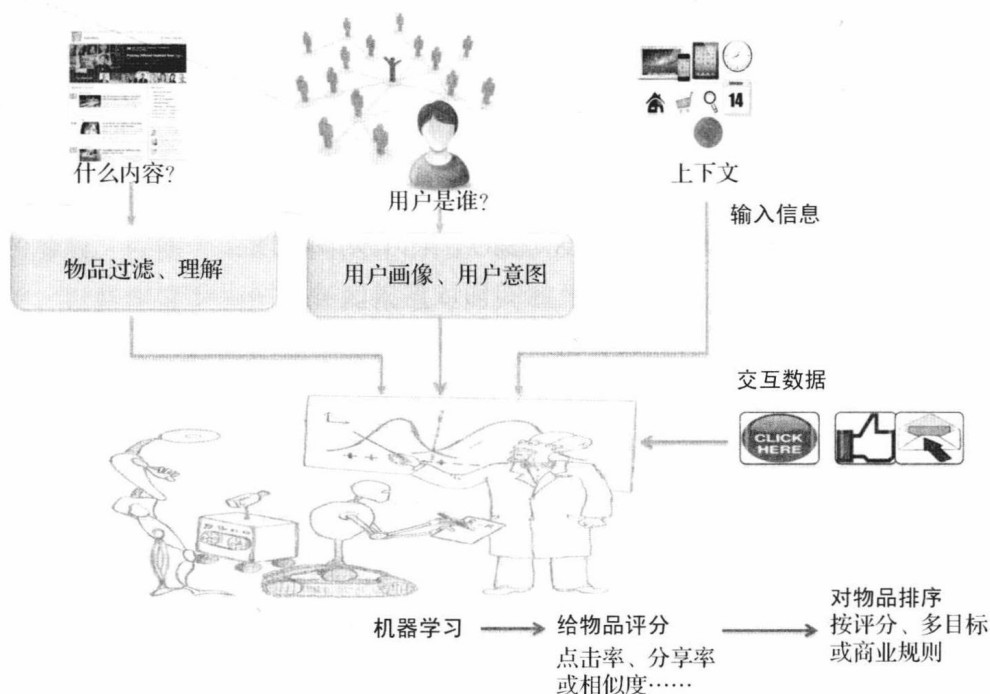
通常，推荐系统中的算法需要完成以下四项任务：

- 内容过滤和理解。我们需要一个高效的算法来过滤掉物品池（候选物品集）中的低质量内容。因为推荐低质量内容不仅会降低用户体验度，还会破坏网站的品牌形象。不同的应用对低质量内容的定义不同：在新闻网站中，知名出版商认为色情内容是低质量内容；电商网站不会代售信誉分过低的店家的商品。大多数情况下，确定并标记低质量内容是一项复杂的任务，需要运用一系列不同的方法才能解决，比如（编辑）打标签、众包或者分类等机器学习方法。除了过滤低质量内容之外，分析和理解质量达标的物品内容也很重要。构建能够精准捕捉内容的物品画像（如特征向量）是一种高效的方法。特征的构建可以借助词袋模型、短语提取、实体提取和主题提取等方法。
- 用户画像建模。除了物品画像，我们还需要构建用户画像，它能反映用户可能会购买哪些物品。用户画像可以根据人口统计信息、用户注册时提交的身份信息、

5 社交网络信息或用户的行为信息来构建。

- 评分。有了用户画像和物品画像，接下来要设计评分函数。评分函数用来估计在给定的上下文（可能是用户正在浏览的网页、正在使用的设备或当前所处的地点）中，将一个物品展示给用户后产生的未来“价值”（如 CTR、与用户当前目标的语义相关性，或期望的收入）。
- 排序。最后，为了最大化目标函数的期望值，我们需要一种排序机制来筛选出一个有序的推荐物品列表。最简单的无非是根据单一的分数对物品进行排序，如每件物品的 CTR。但在实际情况中，排序比想象的更复杂，因为要综合考虑各种不同的因素，如语义相关性、量化不同效用方法的分数，或者为确保良好的用户体验的多样性要求以及为维护品牌形象而设定的商业规则。

图 1-2 将前面介绍的不同算法组件关联了起来。从最上面开始，将用户信息、物品信息和用户 - 物品的历史交互数据输入机器学习统计模型中，然后，模型输出用于衡量用户与物品关联度的评分。最后，排序模块结合评分和单个或多个优化目标，生成优先级从高到低排列的物品列表。



6 图 1-2 推荐系统概览

内容过滤和理解技术在很大程度上由推荐物品的类型决定，例如处理文本的技术与处理图片的技术是截然不同的。我们不打算全面地介绍内容过滤和理解技术，但我们将在第 2 章做简要回顾。我们也不打算介绍各种各样生成用户画像的技术，但会介绍一些可以从用户 - 物品的历史交互数据中自动“学习”用户画像和物品画像的技术。通过这

些技术学到的画像也可以与其他技术生成的画像完美结合。

1.1.2 优化指标

为了制定适用于网站推荐问题的解决方案，我们需要考虑众多重要事项，首要的是确定优化指标。大部分应用程序只有一个优化指标，例如，最大化给定时间段内的总点击数、总收入或总销售额。但有些应用程序会要求同时优化多个指标，例如，在满足一些后续参与度约束的条件下，最大化内容链接的总点击数。参与度约束可能是确保跳离点击数（点击了但未实际阅读的点击数）小于某个阈值。当然，我们可能也想平衡其他因素，例如多样性（随着时间的推移，确保用户能看到不同的主题）和惊喜度（确保不会过度推荐物品给用户，从而限制了新兴趣的发现），这些因素都有利于优化长期用户体验。

优化指标确定后，接下来需要定义优化问题的输入，即分数。如果目标是最大化总点击数，则 CTR 可以很好地衡量一个物品对用户的价值；如果目标有多个，可能需要用到多个分数，如 CTR 和期望时间花销。不得不说，这是一项非常重要的任务，因为要求我们能开发出一种准确估计分数的统计方法。一旦分数估计完成，我们便能根据考虑的优化问题将其运用于排序模块。

1.1.3 探索与利用之间的权衡

可靠地估计分数是推荐系统中一项基本的统计学挑战。分数具体化到应用，可能是期望的正向响应率，例如点击率、显式评分、分享率（分享物品的概率）或喜欢率（用户点击“喜欢”物品这一按钮的概率）。期望响应率可以根据每个可能响应的效用（或价值）进行加权，这是一种基于预期效用对物品进行排序的常用方法。因此，响应率（适当加权）是我们在本书中考虑的主要评分函数。

7

为了准确估计每件候选物品的响应率，我们将每件候选物品展示给一部分用户，并及时收集物品的响应数据，以此方式对候选物品进行探索。之后，利用响应率估值高的物品来优化目标。然而，看似完美的探索过程也存在机会成本。如果仅根据当前收集的数据估算物品的响应率，那么，实际响应效果更好的物品可能没有机会展示给用户。因此，对候选物品的探索和利用便构成了探索与利用的权衡问题。

探索与利用是本书的主题之一，我们会在第3章中介绍，并在第6章中讨论具体的技术细节。第7章和第8章中的方法也是为解决这个问题而提出的。

1.1.4 推荐系统的评估

要了解推荐系统是否能完成目标，需要在开发周期的不同阶段评估其性能。从评估的角度来看，我们将推荐算法的开发分为两个阶段：

- 预部署阶段：处于在线部署算法为网站的部分用户提供推荐服务之前。在此阶段，我们用历史数据评估算法的性能。这种评估存在一定的局限性，因为它是离线的，而我们也没有用户对算法推荐的物品的响应数据。
- 后部署阶段：从在线部署算法服务用户时开始。它主要包括在线分桶测试（也称为 A/B 测试），用于测量合适的指标。虽然在线分桶测试很接近实际情况，但进行此类测试也有一定的代价。常见的解决方法是，在预部署阶段根据离线评估结果，把性能较差的算法过滤掉。

推荐系统的各种组件要用不同的评估方法进行评估：

- 分数的评估。分数通常由预测用户会对物品做何种响应的统计方法给出，我们一般用预测准确度来衡量这类统计方法的性能。举个例子，某统计方法预测出用户给物品的分数，那么该统计方法的误差可以是所有用户的预测分数与实际分数的绝对误差的均值，误差的倒数便是准确度。其他测量准确度的方法将在 4.1.2 节介绍。
- 排序的评估。排序的目的是优化推荐系统的目标。在后部署阶段，为了评估推荐算法，我们利用在线测试收集的数据直接计算兴趣指标（如 CTR，或用户在推荐物品上的时间花销）。在 4.2 节中，我们将讨论如何设置实验，以及如何合理地分析实验结果。但在预部署阶段，因为没有算法为用户服务的数据，很难用离线估计的算法性能来模拟其在线行为。在 4.3 节和 4.4 节中，我们将介绍几种解决这一问题的方法。

8

1.1.5 推荐和搜索：推送与拉取

在界定本书的内容范围时，我们注意到用户意图是区分不同网站应用程序的重要因素。如果用户意图是明确且强烈的（例如，在搜索引擎中查询关键词），那么寻找或推荐与用户意图匹配的物品的的问题通过拉取模型就能解决——检索与明确的用户需求信息相关的物品。但是，在许多推荐场景中，无法获取这种明确的意图信息，最多可以从某种程度上推断出来。因此，系统常采用推送模型，直接将信息推送给用户，目的是提供可能吸引用户的物品。

真实场景中的推荐问题总会归结为一系列的推送与拉取过程。例如，新闻网站难以获取明确的用户意图，因此主要通过推送模型推荐文章。一旦用户开始阅读文章，表现出明确的意图，系统就可以推荐与用户正在阅读的文章主题相关的新闻报道。这种相关新闻推荐系统通常由推送和拉取模型混合而成，先检索与用户当前正在阅读的文章主题相关的文章，然后将它们排序，达到最大化用户参与度的目的。

我们关注的不是搜索网站这类主要靠拉取模型且严重依赖于估计查询语句和物品之间语义相似度的计算方法的的应用程序。我们的重点更多地放在用户意图较弱的应用程序

序上。因此对每个用户来说，基于该用户与物品的历史交互数据给物品评分变得尤为重要。

1.2 一个简单的评分模型：热门推荐

为了说明评分的基本思路，我们考虑推荐热门物品的问题，在网页的单个槽位上为所有用户推荐热门物品（CTR 最高的物品）以最大化总点击数。热门推荐问题虽然简单，但它涵盖了物品推荐的基本要素，也为后续章节介绍的更复杂的技术提供了强有力的基准线。我们假设物品池中的物品数相对于访问数和点击数而言较小。对于物品池的组成，我们不做任何假设，物品池中可能会有新物品加入，随着时间的推移，旧物品也可能会消失。

我们的示例应用是在雅虎首页的今日模块上推荐新闻报道（图 1-3 为模块截图），为了便于说明，这个应用会贯穿整本书。该模块是一个多槽位面板，每个槽位展示一个从物品池中挑选的物品（即新闻报道），物品池中的物品都是在编辑的监督下创建的。为了简便和易于说明，我们只最大化最显眼的槽位中物品的点击数，因为该槽位获得了绝大部分点击数。



图 1-3 雅虎首页的今日模块

令 p_{it} 为物品 i 在 t 时刻的瞬时 CTR。如果每个候选物品的 p_{it} 已知，那么问题就简单了，只要在 t 时刻将瞬时 CTR 最高的物品展示给所有用户即可。换言之，在 t 时刻，我们选择物品 $i_t^* = \arg \max_i p_{it}$ 供用户访问。然而，瞬时 CTR 是未知的，需要从数据中估算。令 \hat{p}_{it} 为从数据中估算的 CTR，那么将瞬时 CTR 最高的物品推荐给用户是不是就够了呢？从数学的角度来说， $\hat{i}_t^* = \arg \max_i \hat{p}_{it}$ 是一个好的 i_t^* 的估计吗？显然不一定，因为估计值的统计方差因物品而异。举个例子，假设有两个物品，且 $\hat{p}_{1t} \sim \mathcal{D}(\text{mean} = 0.01, \text{var} = .005)$ ， $\hat{p}_{2t} \sim \mathcal{D}(\text{mean} = 0.015, \text{var} = .001)$ ， \mathcal{D} 为近似正态的概率分布。那么， $P(\hat{p}_{1t} > \hat{p}_{2t}) = .47$ ，也就是说第一个物品有 47% 的概率被选中，即使它比第二个物品差。出现这种情况的原因是，在样本数较少的情况下，第一个物品 CTR 估计值的方差比第二个物品大很多。因此，在实际应用中，简单粗暴地选择 CTR 估计值最高的物品很有可能会产生假正样本（选中的不是真正的最佳物品）。那么，有可以减少平均假正样本数的其他方案吗？答案是肯定的，那便是我们在前面提到过的探索与利