

# 大数据、数据挖掘 与智慧运营

梁栋 张兆静 彭木根◎编著



BIG DATA MINING  
AND INTELLIGENT OPERATION

清华大学出版社

BIG DATA MINING  
AND INTELLIGENT OPERATION

# 大数据、数据挖掘 与智慧运营

梁栋 张兆静 彭木根〇编著

清华大学出版社  
北京

## 内 容 简 介

本书系统地介绍了大数据挖掘的基本概念、经典挖掘算法、挖掘工具和企业智慧运营应用案例。

全书分为 9 章，内容包括：大数据挖掘与智慧运营的概念，数据预处理，数据挖掘中的四种主流算法：聚类分析、分类分析、回归分析、关联分析，增强型数据挖掘算法，数据挖掘在运营商智慧运营中的应用案例，未来大数据挖掘的发展趋势等。

全书以运用大数据挖掘方法提升企业运营业绩与效率为主线，从运营商实际工作中选取了大量运营和销售案例，详细讲述了数据采集、挖掘建模、模型落地与精准营销的全部过程。书中大部分案例的代码、软件操作流程和微课视频可以通过扫描本书封底的二维码下载。

本书主要面向运营商及其他高科技企业员工、高等院校相关专业本科生和研究生，以及其他对数据挖掘与精准营销感兴趣的读者。

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

大数据、数据挖掘与智慧运营 / 梁栋, 张兆静, 彭木根编著. — 北京: 清华大学出版社, 2017

ISBN 978-7-302-48337-3

I. ①大… II. ①梁… ②张… ③彭… III. ①数据处理 IV. ①TP274

中国版本图书馆 CIP 数据核字 (2017) 第 218392 号

责任编辑：刘 洋

封面设计：李召霞

版式设计：方加青

责任校对：王荣静

责任印制：杨 艳

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](mailto:c-service@tup.tsinghua.edu.cn)

质 量 反 馈：010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印 刷 者：北京富博印刷有限公司

装 订 者：北京市密云县京文制本装订厂

经 销：全国新华书店

开 本：187mm×235mm 印 张：26.25 字 数：494 千字

版 次：2017 年 11 月第 1 版 印 次：2017 年 11 月第 1 次印刷

印 数：1 ~ 4000

定 价：99.00 元

---

产品编号：075552-01

# 前言

数据挖掘（Data Mining），是指从数据中发现知识的过程（Knowledge Discovery in Databases，KDD）。狭义的数据挖掘一般指从大量的、不完全的、有噪声的、模糊的、随机的实际应用数据中，提取隐含其中的、人们事先不知道的、但又是潜在有用知识的过程。自从计算机发明之后，科学家们先后提出了许多优秀的数据挖掘算法。2006年12月，在数据挖掘领域的权威学术会议 *the IEEE International Conference on Data Mining (ICDM)* 上，科学家们评选出了该领域的十大经典算法：C4.5、K-Means、SVM、Apriori、EM、PageRank、AdaBoost、kNN、Naive Bayes 和 CART。这是数据挖掘学科的一个重要里程碑，从此数据挖掘在理论研究和实际应用两方面均进入飞速发展时期，并得到广泛关注。

在实际生产活动中，许多问题都可以用数据挖掘方法来建模，从而提升运营效率。例如，某企业在其移动终端应用（App）上售卖各种商品，它希望向不同的客户群体精准推送差异化的产品和服务，从而提升销售业绩。在这个案例中，如何将千万量级的客户划分为不同的客户群体，可以由数据挖掘中的聚类分析算法来完成；针对某个客户群体，如何判断某个产品是否是他们感兴趣的，可以由数据挖掘中的分类分析算法来完成；如何发现某个客户群体感兴趣的各种产品之间的关联性，应该把哪些产品打包为套餐，可以由数据挖掘中的关联分析算法来完成；如何发现某个客户群体的兴趣爱好的长期趋势，可以由数据挖掘中的回归算法来完成；如何综合考虑公司的 KPI 指标、营销政策和 App 页面限制等条件，制订最终的落地营销方案，可以基于数据挖掘中的 ROC 曲线建立数学模型求得最优解来解决。

当前，许多企业正面临前所未有的竞争压力。以运营商企业为例，从政策层面看，国家提出了“提速降费”的战略指示：一方面要提高网络连接速度、提供更好的服务，这意味着公司成本的提高；另一方面要降低资费标准，这意味着单个产品收入的下降，运营商该如何化解这对矛盾？从运营商内部数据统计看，传统的语音和短信、彩信业

务收入占比正不断下降，传统的利润点已经风光不再；流量收入目前已占据主要位置并保持上涨趋势，但单纯的流量经营又将面临“管道化”压力；未来的利润增长点要让位于被称为“第三条曲线”的数字化服务。运营商该如何经营这一新鲜事物？从外部环境看，互联网和电子商务企业借助其在各方面的优势，已经对运营商形成了巨大的压力，特别是在数字化服务营销领域，传统运营商企业已经不再具备优势，又该如何应对互联网企业的全面竞争？

随着移动互联网和物联网时代的来临，人和万事万物被广泛地联系在一起。人们在联系的过程产生了大量的数据，例如用户基础信息、网页浏览记录、历史消费记录、视频监控影像，等等。据此，以 Google 为首的互联网公司提出了“大数据”（Big Data）的概念，并声称人类已经脱离了信息时代（Information Time, IT），进入了大数据时代（Data Time, DT）。显然，海量数据包含了非常丰富的浅层次信息和深层次知识。对于同一竞争领域的企业，谁能获取最大量的数据，展开最精准的数据挖掘与建模分析，并加以精细化的落地实施，谁便能在行业竞争中取得优势。对于运营商企业而言，其具备的一个显著优势便是手握海量数据资源。如果能运用先进的数据挖掘技术找出客户的行为规律，从传统的经验式、粗放式、“一刀切”式的运营决策向数据化、精细化、个性化的运营决策转型，运营商将迎来新的腾飞。上述运营模式转型的目标，便是所谓的“智慧运营”。

目前，人类对大数据尚没有统一的、公认的定义，但几乎所有学者和企业都认同大数据具备四大特征（四大挑战）：体量巨大（Volume）、类型繁多（Variety）、价值密度低（Value）、需要实时处理（Velocity）。这其中最重要的一点是类型繁多，即过去人类的数据储备以结构化数据为主，而未来将以非结构化数据为主。回到之前提到的 App 营销案例，企业基于用户的基础信息、历史消费信息、简单的网络行为信息等结构化数据展开挖掘建模，被认为是传统的“基于数据挖掘的智慧运营”。随着时代的发展，企业还掌握了用户观看在线视频的内容数据、在营业网点接受营业员推荐的表情信息和语言交流数据、用户在客服热线中的语音咨询数据等。这些数据被统称为非结构化数据，随着语音识别、人脸识别、语义识别等新技术的发展成熟，对非结构化数据的分析挖掘已成为可能，并将获得广阔的商业应用空间。基于非结构化数据的挖掘建模又被称为“基于人工智能的智慧运营”。考虑当前大部分企业的实际运营现状，本书将主要围绕“基于数据挖掘的智慧运营”展开讨论，“基于人工智能的智慧运营”将在后续书籍中展开讨论。

本书共分为九章：第 1 章大数据、数据挖掘与智慧运营综述，讲述数据挖掘的基本概念和发展史、大数据的时代特征、当前结构化数据挖掘进展、非结构化数据挖掘

与人工智能进展、数据挖掘的主流软件等；第2章数据统计与数据预处理，讲述在数据挖掘之前的数据集成、数据清洗、数据衍生、数据统计等；第3章聚类分析，重点讲述K-means、BIRCH、DBSCAN、CLIQUE等几种主流经典聚类算法；第4章分类分析，重点讲述决策树、KNN、贝叶斯、神经网络、SVM等几种主流分类算法；第5章回归分析，重点讲述线性回归、非线性回归、逻辑回归等几种主流回归算法；第6章关联分析，重点讲述Apriori、FP-tree等几种主流关联算法；第7章增强型数据挖掘算法，重点讲述随机森林、Bagging、Boosting等几种主流增强算法；第8章数据挖掘在运营商智慧运营中的应用，展开讲述数据挖掘方法在外呼营销、精准推送、套餐适配、客户保有、投诉预警、网络质量监控、室内定位中的应用；第9章面向未来大数据的数据挖掘与机器学习发展趋势，简要讲述数据挖掘领域的前沿研究进展。

全书以运用大数据挖掘方法提升企业运营业绩与效率为主线。第3章至第7章组成本书的理论知识部分，在讲述理论知识的同时，这部分每章都配套列举了大量实际应用案例，及其在SPSS等分析软件中的具体操作流程。此外，第8章从运营商实际工作中选取了大量运营和销售案例，详细讲述了数据采集、挖掘建模、模型落地与精准营销的全部过程。书中大部分案例的代码、软件操作流程和微课视频可以通过扫描本书封底的二维码下载。

本书基于作者所带领的研究团队多年研究积累和在运营商企业广泛落地应用的基础上提炼而成。全书由曾丽丽博士组织并统稿，梁栋、张兆静和彭木根撰写了主要章节，研究团队中的谢花花、柯联兴、张笑凯、鲁晨、李子凡等在读研究生参与了部分章节的写作，胡林、唐糖等团队外专家参与了部分章节的写作并给出了宝贵的意见。在本书写作过程中，中国移动及许多省市分公司（特别是广西分公司）给予了大力支持。在本书出版前，许多素材被中国移动广西分公司选为教材并展开了广泛落地应用，获得了2016年中国移动集团公司颁发的“培训案例最佳实践奖”。在本书出版过程中，得到了深圳市傲举企业管理顾问有限公司的大力支持。在此对有关人员一并表示诚挚的感谢！

由于作者能力所限，疏漏之处在所难免，希望各位读者海涵，并批评指正。

作 者

2017年9月于北京邮电大学

# 目 录

<b>第1章 大数据、数据挖掘与智慧运营综述</b>	<b>1</b>
1.1 数据挖掘的发展史	2
1.1.1 数据挖掘的定义与起源	2
1.1.2 数据挖掘的早期发展	3
1.1.3 数据挖掘的算法前传	4
1.1.4 数据挖掘的第一个里程碑	8
1.1.5 最近十年的发展与应用	11
1.2 数据挖掘的主要流程与金字塔模型	13
1.2.1 数据挖掘的任务	14
1.2.2 数据挖掘的基本步骤	16
1.2.3 数据挖掘的架构——云计算	17
1.2.4 “金字塔”模型	20
1.3 数据挖掘对智慧运营的意义	22
1.3.1 “互联网+”时代的来临及其对运营商的冲击和挑战	22
1.3.2 大数据时代的来临及其对运营商的挑战和机遇	24
1.3.3 电信运营商运营发展面临的主要瓶颈	26
1.3.4 电信运营商发展的“三条曲线”	27
1.3.5 智慧运营与大数据变现	29
1.3.6 数据挖掘对于提升智慧运营效率的意义	30
1.4 大数据时代已经来临	31
1.4.1 大数据的定义	31

1.4.2 大数据的“4V”特征.....	32
1.4.3 结构化数据与非结构化数据.....	33
1.5 非结构化数据挖掘的研究进展.....	34
1.5.1 文本挖掘.....	34
1.5.2 模式识别.....	36
1.5.3 语音识别.....	40
1.5.4 视频识别.....	44
1.5.5 其他非结构化数据挖掘.....	48
1.6 数据挖掘与机器学习、深度学习、人工智能及云计算.....	50
1.6.1 机器学习.....	51
1.6.2 深度学习.....	53
1.6.3 人工智能.....	55
1.6.4 云计算.....	56
1.7 现有数据挖掘的主要分析软件与系统.....	61
1.7.1 Hadoop.....	61
1.7.2 Storm.....	63
1.7.3 Spark.....	65
1.7.4 SPASS (SPSS) .....	66
1.7.5 SAS .....	68
参考文献 .....	70
<b>第2章 数据统计与数据预处理.....</b>	<b>73</b>
2.1 数据属性类型 .....	74
2.1.1 数据属性定义 .....	74
2.1.2 离散属性 .....	74
2.1.3 连续属性 .....	75
2.2 数据的统计特性 .....	77
2.2.1 中心趋势度量 .....	77
2.2.2 数据散布度量 .....	78
2.2.3 数据相关性 .....	82
2.3 数据预处理 .....	87
2.3.1 数据预处理概述 .....	87

2.3.2 数据预处理的主要任务.....	88
2.3.3 数据清理.....	89
2.3.4 数据集成.....	91
2.3.5 数据规约.....	94
2.3.6 数据变换和离散化.....	97
2.4 数据字段的衍生.....	100
2.4.1 数据字段的拆分.....	101
2.4.2 统计特征的构造.....	101
2.4.3 数据域的变换.....	102
2.5 SPSS 软件中的数据预处理案例 .....	103
2.5.1 缺失值的实操处理.....	103
2.5.2 噪声数据的实操处理.....	104
2.5.3 主成分分析的实操处理.....	105
参考文献 .....	107
<b>第3章 聚类分析 .....</b>	<b>109</b>
3.1 概述 .....	110
3.2 聚类算法的评估 .....	113
3.3 基于划分的聚类: K-means.....	117
3.3.1 基于划分的聚类算法概述.....	117
3.3.2 K-means聚类算法原理 .....	118
3.3.3 K-means算法的优势与劣势 .....	120
3.3.4 K-means算法优化 .....	121
3.3.5 SPSS软件中的K-means算法应用案例 .....	123
3.4 基于层次化的聚类: BIRCH .....	134
3.4.1 基于层次化的聚类算法概述.....	134
3.4.2 BIRCH算法的基本原理 .....	135
3.4.3 BIRCH算法的优势与劣势 .....	136
3.5 基于密度的聚类: DBSCAN .....	137
3.5.1 基于密度的聚类算法概述.....	137
3.5.2 DBSCAN算法的基本原理 .....	137
3.5.3 DBSCAN算法的优势与劣势 .....	140

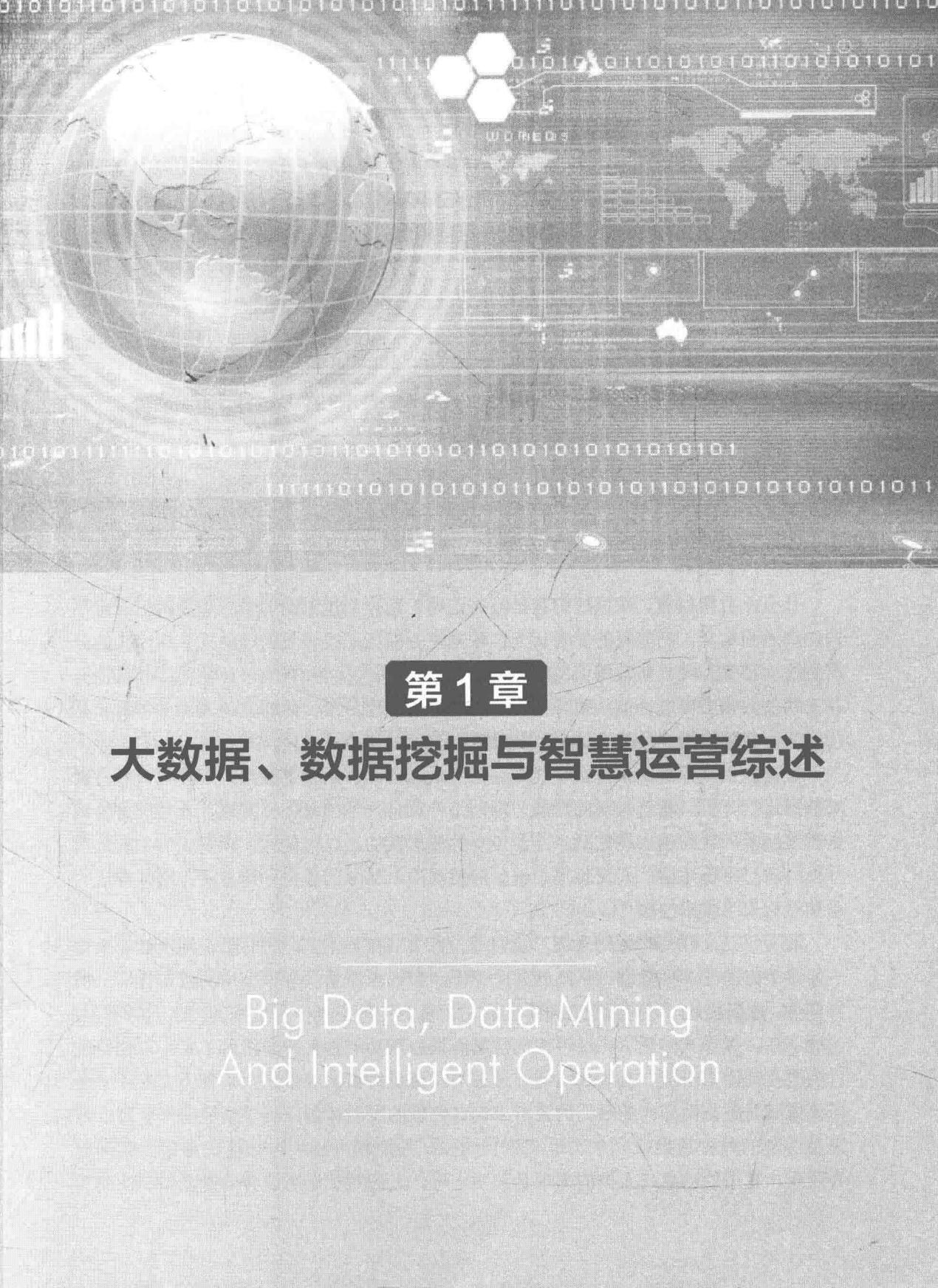
3.6 基于网格的聚类: CLIQUE.....	140
3.6.1 基于网格的聚类算法概述.....	140
3.6.2 CLIQUE算法的基本原理.....	141
3.6.3 CLIQUE算法的优势与劣势.....	142
参考文献 .....	143
<b>第4章 分类分析 .....</b>	<b>145</b>
4.1 分类分析概述 .....	146
4.2 分类分析的评估 .....	148
4.3 决策树分析 .....	152
4.3.1 决策树算法的基本原理.....	152
4.3.2 CHAID决策树 .....	160
4.3.3 ID3决策树 .....	167
4.3.4 C4.5决策树 .....	171
4.3.5 CART决策树 .....	175
4.3.6 决策树中的剪枝问题.....	179
4.3.7 决策树在SPSS中的应用 .....	180
4.4 最近邻分析 (KNN) .....	185
4.4.1 KNN算法的基本原理 .....	185
4.4.2 KNN算法流程 .....	186
4.4.3 KNN算法的若干问题 .....	187
4.4.4 KNN分类器的特征 .....	188
4.4.5 KNN算法在SPSS中的应用 .....	188
4.5 贝叶斯分析 .....	191
4.5.1 贝叶斯定理 .....	191
4.5.2 朴素贝叶斯分类 .....	192
4.5.3 贝叶斯网络 .....	195
4.6 神经网络 .....	199
4.6.1 感知器 .....	200
4.6.2 多重人工神经网络 .....	201
4.6.3 人工神经网络的特点 .....	203
4.7 支持向量机 .....	204

4.7.1 支持向量机简介.....	204
4.7.2 最大边缘超平面.....	205
4.7.3 数据线性可分的情况.....	206
4.7.4 数据非线性可分的情况.....	209
4.7.5 支持向量机的特征.....	210
参考文献 .....	210
<b>第5章 回归分析 .....</b>	<b>213</b>
5.1 回归分析概述 .....	214
5.2 一元线性回归 .....	214
5.2.1 一元线性回归的基本原理.....	215
5.2.2 一元线性回归性能评估.....	216
5.2.3 SPSS软件中一元线性回归应用案例.....	216
5.3 多元线性回归 .....	222
5.3.1 多元线性回归基本原理.....	223
5.3.2 自变量选择方法.....	223
5.3.3 SPSS软件中的多元线性回归应用案例.....	224
5.4 非线性回归 .....	230
5.4.1 非线性回归基本原理.....	231
5.4.2 幂函数回归分析.....	231
5.4.3 指数回归分析.....	232
5.4.4 对数回归分析.....	232
5.4.5 多项式回归分析.....	233
5.4.6 非线性模型线性和曲线回归.....	234
5.5 逻辑回归 .....	235
5.5.1 逻辑回归基本原理.....	235
5.5.2 二元逻辑回归.....	238
5.5.3 多元逻辑回归.....	238
5.5.4 SPSS软件中的逻辑回归应用案例.....	238
参考文献 .....	242

<b>第6章 关联分析</b>	245
6.1 关联分析概述	246
6.2 关联分析的评估指标	247
6.2.1 支持度	247
6.2.2 置信度	248
6.2.3 算法复杂度	248
6.3 Apriori 算法	249
6.3.1 频繁项集的定义与产生	249
6.3.2 先验原理	251
6.3.3 基于支持度的计数与剪枝	252
6.3.4 候选项集生成	253
6.3.5 基于置信度的剪枝	259
6.3.6 Apriori算法规则生成	259
6.4 FP-tree 算法	261
6.4.1 频繁模式树	261
6.4.2 FP-tree算法频繁项集的产生	263
6.4.3 FP-tree算法规则生成	263
6.4.4 算法性能对比与评估	264
6.5 SPSS Modeler 关联分析实例	265
参考文献	269
<b>第7章 增强型数据挖掘算法</b>	271
7.1 增强型数据挖掘算法概述	272
7.1.1 组合方法的优势	272
7.1.2 构建组合分类器的方法	272
7.2 随机森林	273
7.2.1 随机森林的原理	273
7.2.2 随机森林的优缺点	276
7.2.3 随机森林的泛化误差	276
7.2.4 输入特征的选择方法	277
7.3 Bagging 算法	277

7.4 AdaBoost 算法.....	280
7.4.1 AdaBoost算法简介.....	281
7.4.2 AdaBoost算法原理.....	281
7.4.3 AdaBoost算法的优缺点.....	285
7.5 提高不平衡数据的分类准确率 .....	286
7.5.1 不平衡数据.....	286
7.5.2 不平衡数据的处理方法——数据层面.....	288
7.5.3 不平衡数据的处理方法——算法层面.....	292
7.6 迁移学习 .....	295
7.6.1 迁移学习的基本原理.....	295
7.6.2 迁移学习的分类.....	296
7.6.3 迁移学习与数据挖掘.....	298
7.6.4 迁移学习的发展.....	301
参考文献 .....	302
<b>第8章 数据挖掘在运营商智慧运营中的应用 .....</b>	<b>305</b>
8.1 概述 .....	306
8.2 单个业务的精准营销——合约机外呼营销 .....	307
8.2.1 总结历史营销规律.....	308
8.2.2 预测潜在客户群体.....	321
8.2.3 客户群体细分.....	327
8.2.4 制定层次化、个性化精准营销方案.....	328
8.3 多种互联网业务的精准推送 .....	329
8.3.1 根据历史营销规律总结单个业务的历史营销规律.....	329
8.3.2 预测潜在客户群体、预测单个业务的潜在客户群体及多个业务的联合建模.....	341
8.3.3 制定多业务层次化个性化联合精准营销方案.....	342
8.3.4 落地效果评估.....	343
8.4 套餐精准适配 .....	344
8.4.1 痛点.....	344
8.4.2 挖掘潜在客户群体.....	346
8.4.3 探寻强相关字段.....	349

8.4.4 多元线性回归建模	350
8.4.5 制定层次化、个性化精准营销方案	351
8.4.6 落地效果评估与模型调优	352
8.5 客户保有	353
8.5.1 总结客户流失的历史规律	354
8.5.2 细分潜在流失客户群体	357
8.5.3 客户保有效益建模与最优决策	359
8.5.4 落地效果评估	361
8.6 投诉预警	363
8.6.1 客户投诉现象分析	363
8.6.2 挖掘潜在客户群体	366
8.6.3 制定个性化关怀方案	368
8.7 网络质量栅格化呈现	368
8.7.1 栅格化呈现的基本原理	369
8.7.2 覆盖栅格化	370
8.7.3 基于流量聚簇的网络优化策略	372
8.8 无线室内定位	376
8.8.1 传统室内定位方法	376
8.8.2 基于Wi-Fi信号的指纹定位算法	378
8.8.3 基于数据挖掘算法的改进定位方法	379
参考文献	383
<b>第9章 面向未来大数据的数据挖掘与机器学习发展趋势</b>	<b>385</b>
9.1 大数据时代数据挖掘与机器学习面临的新挑战	386
9.2 IEEE ICDM会议数据挖掘与机器学习的最新研究进展	395
9.3 “计算机奥运会”——Sort Benchmark	400
参考文献	402



## 第1章

# 大数据、数据挖掘与智慧运营综述

Big Data, Data Mining  
And Intelligent Operation

近年来，大数据、数据挖掘、机器学习、云计算和人工智能等词语日渐为人们所熟悉。本章将围绕上述基本概念和话题展开讨论。本章 1.1 节介绍数据挖掘的概念和发展史，1.2 节介绍数据挖掘的主要流程和金字塔模型，1.3 节介绍数据挖掘对企业智慧运营的重要意义，1.4 节介绍大数据的基本概念、特征和挑战，1.5 节介绍非结构化数据挖掘的概念和研究进展，1.6 节介绍结构化数据挖掘与机器学习、深度学习和人工智能之间的关联关系，1.7 节介绍常见的数据挖掘分析软件与系统。

## 1.1 数据挖掘的发展史

### 1.1.1 数据挖掘的定义与起源

什么是数据挖掘，数据挖掘包括哪些范畴？迄今为止不同的学者和公司仍有有着不同的理解和定义。例如有的学者认为：数据挖掘即指摆脱传统的经验式、规律式的分析方法，转变为纯粹从数据出发来探索问题的本质。又例如有的公司认为：数据挖掘是一种从数据中榨取价值，提升公司运营效率的重要手段。然而，绝大部分学者和公司都认同数据挖掘的最基本定义：从数据中获取知识。

数据挖掘具体起源于什么年代现在已无从考证。自从有了数据，人类就开始尝试对数据进行分析。随着时代的发展，特别是计算机技术的诞生和发展，人类拥有的数据越来越多，种类越来越复杂，之前传统的浅层次的、以经验式、观察式为主的数据分析方法已不再适用，人类急需一整套深层次的、科学的数据分析方法，这些方法的总和被称为“数据挖掘”。

随着移动互联网时代的来临，我们每天都生活在数据中，时时刻刻都接触着来自生活各个方面各种数据：早高峰各个十字路口的车流量，各个公司的股市行情、销售票务、产品描述、用户反馈，科学实验记录着的种种信息……数据的产生无时不在，无处不在。爆炸式增长、广泛可用的巨量数据急需功能强大和通用的工具，以便发现它们潜在的巨大价值。交警部门需要通过对车流量数据的观察来决定警力支配；公司需要通过对方方面面商业数据的分析来制订合理的发展计划；科学的研究工作者需要对来自实验的种种数据研究来实现实验目的……人们越来越关注如何把海量的数据变为直观、有用的信息。人类的需求是发明之母，人们对数据所蕴含的潜在知识的需求

促使了数据挖掘的诞生。

近年来，数据挖掘引起了信息产业界的极大关注，其主要原因是存在大量数据可以广泛使用，并且迫切需要将这些数据转换成有用的信息和知识。获取的信息和知识可以被广泛用于各种应用，包括商务管理、生产控制、市场分析、工程设计和科学探索等。

数据挖掘利用了来自如下领域的思想：

- (1) 来自统计学的抽样、估计和假设检验。
- (2) 人工智能、模式识别和机器学习的搜索算法、建模技术和学习理论。

数据挖掘也迅速地接纳了来自其他领域的思想，这些领域包括最优化、进化计算、信息论、信号处理、可视化和信息检索。一些其他领域也起到重要的支撑作用。特别的，需要数据库系统提供有效的存储、索引和查询处理支持。源于高性能（并行）计算的技术在处理海量数据集方面常常是重要的。分布式技术也能帮助处理海量数据，并且当数据不能集中到一起处理时更是至关重要。

## 1.1.2 数据挖掘的早期发展

数据挖掘起始于 20 世纪下半叶，是在多个学科发展的基础上逐步发展起来的。随着大数据与数据库技术的发展应用，数据量不断积累与膨胀，这导致基础的查询和统计操作已经无法满足企业的商业需求。如何挖掘出数据隐含的信息是当前亟须解决的难题。与此同时，计算机领域的人工智能（Artificial Intelligence）方向也取得了巨大进展，进入了机器学习的阶段。因此，人们将两者结合起来，用数据库管理系统存储数据，用计算机分析数据，并且尝试挖掘数据背后的信息。这两者的结合促生了一门新的学科，即数据库中的知识发现（Knowledge Discovery in Databases，KDD）。1989 年 8 月召开的第 11 届国际人工智能联合会议的专题讨论会上首次出现了“知识发现”这个术语，到目前为止，知识发现的重点已经从发现方法转向了实践应用。

数据挖掘（Data Mining）则是 KDD 的核心部分，它指的是从数据集合中自动抽取隐藏在数据中那些有用信息的非平凡过程，这些信息的表现形式为：规则、概念、规律及模式等。进入 21 世纪，数据挖掘已经成为一门比较成熟的交叉学科，并且数据挖掘技术也伴随着信息技术的发展日益成熟起来。总体来说，数据挖掘融合了数据库、人工智能、机器学习、统计学、高性能计算、模式识别、神经网络、数据可视化、信息检索和空间数据分析等多个领域的理论和技术，是 21 世纪初期对人类产生重大影响的十大新兴技术之一。