

普通高等教育统计与大数据专业“十三五”规划教材

数据挖掘方法与应用

徐雪琪 编著

清华大学出版社

北 京

内容简介

本书以应用为导向介绍数据挖掘的相关工具、理论和方法,包括数据挖掘概述、数据挖掘工具、数据与数据平台、数据预处理、关联分析、决策树、贝叶斯分类和神经网络。通过循序渐进地讲解数据挖掘可使用的工具、数据存储及分析环境、原始数据可能存在的问题及相应的预处理方法、数据挖掘经典算法等相关知识,使读者对数据挖掘有整体的认识 and 了解。此外,本书以解决问题为目的,结合实例阐述了使用 IBM SPSS Modeler 和 R 软件进行数据挖掘的方法与步骤,便于读者更好地理解和掌握。

本书可作为统计学、大数据等相关专业高年级本科生及硕士研究生数据挖掘课程的教材,也可作为其他数据挖掘爱好者的参考用书。

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

数据挖掘方法与应用 / 徐雪琪 编著. —北京:清华大学出版社, 2020.6

普通高等教育统计与大数据专业“十三五”规划教材

ISBN 978-7-302-55062-4

I. ①数… II. ①徐… III. ①数据采集—高等学校—教材 IV. ①TP274

中国版本图书馆 CIP 数据核字(2020)第 039854 号

责任编辑:崔 伟

装帧设计:马筱琨

责任校对:牛艳敏

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

质 量 反 馈:010-62772015, zhiliang@tup.tsinghua.edu.cn

印 装 者:北京国马印刷厂

经 销:全国新华书店

开 本:185mm×260mm 印 张:17.5 字 数:426 千字

版 次:2020 年 8 月第 1 版 印 次:2020 年 8 月第 1 次印刷

定 价:49.00 元

产品编号:083786-01

前 言

数据挖掘出现于 20 世纪 80 年代后期,随着信息化技术的持续发展,它不断汲取统计学、机器学习、数据库技术、人工智能、模式识别和数据可视化等多学科领域的知识,无可争议地成为当今利用大数据分析获取知识的核心利器。

本教材是浙江省“十三五”优势专业(经济统计学)、浙江省一流学科(统计学)、浙江省优势特色学科(统计学)的建设成果之一,具有以下显著特点:

(1) 重视数据挖掘项目实现的整个流程,除了包含数据挖掘的经典理论与方法,还详细介绍了数据挖掘工具、挖掘的数据类型和存储环境、大数据平台及数据预处理方法。

(2) 重视数据挖掘理论和方法的基本思想,在不失严谨的前提下,略过了一些复杂程度高,但又不影响理解的数学推导,将各个知识点言简意赅地阐述透彻。

(3) 重视实际案例应用及实现,每类方法结合多个案例,以运用恰当的方法解决实际问题为导向,以培养分析问题能力为重点,详细介绍 IBM SPSS Modeler 和 R 软件的实现过程。

本教材共分为 8 章:第 1 章为数据挖掘概述,主要介绍数据挖掘的发展历程、相关技术与发展趋势等;第 2~3 章主要介绍数据挖掘工具、数据类型及数据平台;第 4 章介绍数据预处理相关技术;第 5~8 章介绍了各种数据挖掘经典算法原理、案例应用及实现。

本教材主要针对统计学、大数据相关专业的高年级本科生和硕士研究生编写,以学生深入理解并掌握数据挖掘的基本方法、了解相关的应用环境、熟练运用相关软件进行数据挖掘为目标,也可作为其他各专业读者学习数据挖掘方法与应用的教材或参考书。

本教材教学资源丰富,除了教学课件之外,还提供了每章的案例数据,读者可以通过 <http://www.tupwk.com.cn> 下载使用。

本教材由浙江工商大学徐雪琪副教授结合十多年的教学工作经验编写而成。结合笔者的教学实践,以 48 学时为例(一学期 16 周,每周 3 学时),本教材的理论教学内容可安排 33 学时,第 5~8 章的应用部分可安排 15 学时实验教学。在编写过程中,笔者参考了国内外数据挖掘领域许多学者的研究成果,在此深表谢意!

笔者虽已尽心竭力,但限于水平和时间仓促,书中谬误之处在所难免,敬请读者批评指正。

徐雪琪
2020 年 4 月

目 录

| | | | |
|-----------------------------------|----|----------------------------|----|
| 第 1 章 数据挖掘概述 | 1 | 2.3.3 R 语言与数据挖掘 | 42 |
| 1.1 数据挖掘的产生与发展 | 1 | 2.4 Python 语言 | 45 |
| 1.1.1 数据挖掘概念的提出 | 2 | 2.4.1 Python 语言简述 | 45 |
| 1.1.2 数据挖掘系统的发展 | 3 | 2.4.2 Python 与数据分析 | 46 |
| 1.1.3 当前热点和未来趋势 | 5 | 2.4.3 Anaconda | 46 |
| 1.2 数据挖掘过程 | 10 | 2.5 练习与拓展 | 50 |
| 1.2.1 Fayyad 过程模型 | 10 | 第 3 章 数据与数据平台 | 51 |
| 1.2.2 CRISP-DM 过程模型 | 11 | 3.1 数据类型 | 51 |
| 1.3 数据挖掘功能与使用技术 | 21 | 3.1.1 数据形态与数据类型 | 51 |
| 1.3.1 数据挖掘功能 | 21 | 3.1.2 数据环境与数据类型 | 54 |
| 1.3.2 数据挖掘使用技术 | 22 | 3.2 关系型数据库 | 55 |
| 1.4 数据挖掘应用 | 26 | 3.2.1 关系型数据库概述 | 55 |
| 1.4.1 金融领域的数据挖掘 | 26 | 3.2.2 关系型数据库管理系统 | 56 |
| 1.4.2 电信领域的数据挖掘 | 26 | 3.3 NoSQL 数据库 | 57 |
| 1.4.3 零售与电子商务领域的 数据挖掘 | 27 | 3.3.1 键值数据库 | 57 |
| 1.4.4 政府政务领域的数据挖掘 | 27 | 3.3.2 文档数据库 | 58 |
| 1.4.5 医疗领域的数据挖掘 | 28 | 3.3.3 列族数据库 | 60 |
| 1.4.6 科学领域的数据挖掘 | 28 | 3.3.4 图数据库 | 61 |
| 1.5 练习与拓展 | 28 | 3.4 数据仓库与大数据平台 | 63 |
| 第 2 章 数据挖掘工具 | 30 | 3.4.1 数据仓库 | 63 |
| 2.1 Weka | 30 | 3.4.2 大数据平台 | 68 |
| 2.1.1 Weka 简述 | 30 | 3.5 练习与拓展 | 74 |
| 2.1.2 Weka 运行界面 | 31 | 第 4 章 数据预处理 | 75 |
| 2.2 IBM SPSS Modeler | 34 | 4.1 数据预处理概述 | 75 |
| 2.2.1 IBM SPSS Modeler 简述 | 34 | 4.1.1 原始数据中存在的问题 | 75 |
| 2.2.2 IBM SPSS Modeler 主界面 及功能 | 35 | 4.1.2 数据预处理的主要任务 | 77 |
| 2.3 R 语言 | 41 | 4.2 数据清洗 | 77 |
| 2.3.1 R 语言简述 | 41 | 4.2.1 缺失数据处理 | 77 |
| 2.3.2 RStudio | 42 | 4.2.2 异常数据处理 | 78 |
| | | 4.3 数据集成 | 80 |
| | | 4.3.1 模式匹配及数值一致化 | 80 |

| | | | |
|--------------------------------------|------------|--------------------------------------|------------|
| 4.3.2 删除冗余数据····· | 81 | 6.1.1 决策树分析相关概念····· | 137 |
| 4.4 数据变换····· | 82 | 6.1.2 决策树分析核心问题····· | 138 |
| 4.4.1 定性数据数值化····· | 82 | 6.2 ID3 算法····· | 138 |
| 4.4.2 定量数据离散化和规范化····· | 83 | 6.2.1 信息论的基本概念····· | 138 |
| 4.4.3 不平衡数据处理····· | 84 | 6.2.2 ID3 算法基本原理····· | 139 |
| 4.5 数据归约····· | 85 | 6.2.3 使用 ID3 算法建立决策树····· | 141 |
| 4.5.1 属性的归约····· | 85 | 6.3 C5.0 算法····· | 143 |
| 4.5.2 记录的归约····· | 87 | 6.3.1 C5.0 算法的决策树生长····· | 144 |
| 4.5.3 数值的归约····· | 88 | 6.3.2 C5.0 算法的决策树修剪····· | 149 |
| 4.6 练习与拓展····· | 89 | 6.4 基于 IBM SPSS Modeler 的 应用····· | 151 |
| 第 5 章 关联分析····· | 90 | 6.4.1 数据读取与审核····· | 152 |
| 5.1 关联分析概述····· | 90 | 6.4.2 探索性分析····· | 153 |
| 5.1.1 关联分析基本概念····· | 91 | 6.4.3 数据预处理····· | 158 |
| 5.1.2 关联规则挖掘的基本过程····· | 93 | 6.4.4 决策树模型构建与评估： 基于 C5.0 算法····· | 163 |
| 5.2 Apriori 算法····· | 94 | 6.4.5 预测结果····· | 170 |
| 5.2.1 Apriori 性质····· | 94 | 6.5 基于 R 语言的应用····· | 171 |
| 5.2.2 Apriori 算法的频繁项集 产生····· | 95 | 6.5.1 数据探索····· | 172 |
| 5.3 强关联规则的悖论····· | 99 | 6.5.2 数据分区····· | 177 |
| 5.3.1 强关联规则不一定是有趣的 规则····· | 99 | 6.5.3 模型训练与评估····· | 178 |
| 5.3.2 基于提升度过滤无趣的 强关联规则····· | 100 | 6.5.4 使用 boosting 和代价矩阵 调整模型····· | 181 |
| 5.3.3 基于支持度、置信度及 提升度的关联规则发现····· | 100 | 6.6 练习与拓展····· | 184 |
| 5.4 基于 IBM SPSS Modeler 的 应用····· | 103 | 第 7 章 贝叶斯分类····· | 185 |
| 5.4.1 事实表数据的应用示例····· | 103 | 7.1 贝叶斯分类概述····· | 185 |
| 5.4.2 事务表数据的应用示例····· | 113 | 7.1.1 贝叶斯定理····· | 186 |
| 5.5 基于 R 语言的应用····· | 123 | 7.1.2 贝叶斯信念网络····· | 186 |
| 5.5.1 数据初探····· | 123 | 7.2 朴素贝叶斯分类····· | 188 |
| 5.5.2 可视化交易数据····· | 125 | 7.2.1 朴素贝叶斯分类原理····· | 188 |
| 5.5.3 挖掘关联规则····· | 127 | 7.2.2 朴素贝叶斯分类计算示例····· | 191 |
| 5.5.4 可视化关联规则····· | 130 | 7.2.3 零概率问题：拉普拉斯 平滑····· | 193 |
| 5.6 练习与拓展····· | 134 | 7.3 TAN 贝叶斯分类····· | 194 |
| 第 6 章 决策树····· | 136 | 7.3.1 TAN 贝叶斯网络结构····· | 194 |
| 6.1 决策树概述····· | 136 | 7.3.2 TAN 贝叶斯分类过程····· | 195 |
| | | 7.4 基于 IBM SPSS Modeler 的 应用····· | 196 |

| | | | |
|---------------------------------|------------|--------------------------------------|------------|
| 7.4.1 数据读取与审核····· | 198 | 8.2.3 前馈神经网络计算示例····· | 238 |
| 7.4.2 探索性分析····· | 199 | 8.3 卷积神经网络····· | 240 |
| 7.4.3 数据预处理····· | 208 | 8.3.1 卷积层····· | 240 |
| 7.4.4 TAN 贝叶斯分类模型构建 与评估····· | 210 | 8.3.2 激活层····· | 243 |
| 7.5 基于 R 语言的应用····· | 214 | 8.3.3 池化层····· | 244 |
| 7.5.1 数据探索····· | 214 | 8.3.4 全连接层····· | 244 |
| 7.5.2 文本数据预处理····· | 215 | 8.4 基于 IBM SPSS Modeler 的 应用····· | 245 |
| 7.5.3 划分数据集····· | 219 | 8.4.1 数据读取····· | 246 |
| 7.5.4 词云分析····· | 221 | 8.4.2 “数据审核”节点预处理····· | 247 |
| 7.5.5 模型训练与评估····· | 223 | 8.4.3 探索性分析····· | 250 |
| 7.6 练习与拓展····· | 225 | 8.4.4 分区与平衡····· | 251 |
| 第 8 章 神经网络····· | 226 | 8.4.5 模型构建与评价····· | 252 |
| 8.1 神经网络概述····· | 226 | 8.5 基于 R 语言的应用····· | 260 |
| 8.1.1 生物神经元与人工神经元····· | 226 | 8.5.1 数据初探····· | 260 |
| 8.1.2 激活函数····· | 227 | 8.5.2 数据转换与分区····· | 263 |
| 8.1.3 神经网络的拓扑结构····· | 230 | 8.5.3 模型构建与评价····· | 263 |
| 8.2 BP 神经网络····· | 232 | 8.6 练习与拓展····· | 268 |
| 8.2.1 BP 神经网络的学习过程····· | 232 | 参考文献····· | 270 |
| 8.2.2 BP 算法描述····· | 237 | | |

第 1 章

数据挖掘概述

本章内容

- 数据挖掘的产生与发展
- 数据挖掘过程
- 数据挖掘功能与使用技术
- 数据挖掘应用

绕不开的传说：啤酒和尿布

不管真相如何，啤酒和尿布的传说已被写入数据挖掘发展的历史。过去几年，在介绍数据挖掘的不同场合及许多出版物上，不断有人以此为例。

据说在 20 世纪 90 年代，沃尔玛对其在美国本土超市的销售数据展开研究，结果发现，和尿布一起购买次数最多的商品竟然是啤酒！啤酒和尿布，似乎风马牛不相及，沃尔玛管理层对这个结果产生了疑问：真是这样吗？为什么？于是决定对同时购买过啤酒和尿布的顾客进行电话回访，询问其为什么会同时购买。答案是一些年轻的爸爸在下班途中经常会接到妻子的电话，要求其回家时顺便购买孩子的尿布，有 30%~40% 的爸爸会顺便买点啤酒犒劳自己。证实了这个规律后，管理层就把啤酒和尿布摆放在一起进行销售，不出意料，销售量双双增加。

1.1 数据挖掘的产生与发展

自 20 世纪 60 年代以来，随着信息技术的飞速发展，数据库及数据仓库技术被广泛应用，遍及超级销售市场、银行、天文学研究、医学研究以及政府部门等各个领域。以全球最大的零售企业沃尔玛为例，其创始人山姆·沃尔顿非常重视信息的沟通和信息系统的建设，早在 1969 年，便购买第一台计算机用来支持公司日常业务。20 世纪 70 年代，沃尔玛建立了物流的管理信息系统(management information system, MIS)。20 世纪 80 年代初，沃尔玛与休斯公司合作发射物流通信卫星，实现了全球联网；1983 年开始使用 POS 机；1985 年建立了电子数据交换系统(electronic data interchange, EDI)，开始无纸化作业，所有信息

全部在电脑上运作；1986 年建立了快速反应系统(quick response, QR)用于订货业务和付款通知业务。20 世纪 90 年代, 沃尔玛开始采用全球领先的卫星定位系统(GPS), 控制公司物流。由此, 沃尔玛成为全球第一个实现集团内部 24 小时计算机物流网络化监控的企业, 使采购库存、订货、配送和销售一体化。信息化建设, 使沃尔玛积累了大量的各类业务数据, 但是我们知道, 数据作为一种资源, 本身并没有什么直接的价值, 有价值的是从中所能获得的信息和知识。数据挖掘正是基于这种需要而产生、发展起来的, 如图 1.1 所示。

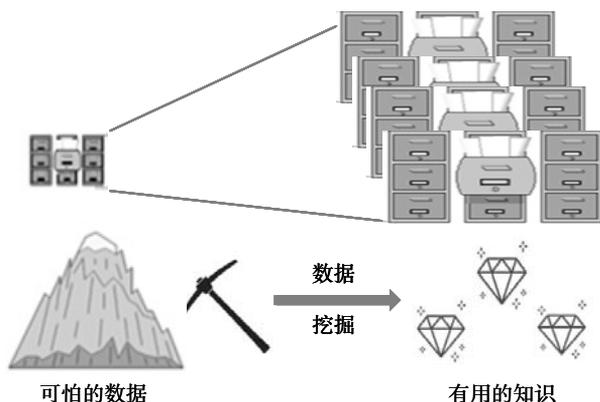


图 1.1 数据挖掘产生示意图

1.1.1 数据挖掘概念的提出

1. KDD 国际学术会议

1989 年 8 月在美国底特律召开的第 11 届国际联合人工智能学术会议(IJCAI-89)上, Gregory Piatetsky-Shapiro 组织了“数据库中的知识发现”(KDD: Knowledge Discovery in Database)专题讨论会。该讨论会的重点是强调发现的方法以及发现的知识两个方面, 这是基于数据挖掘概念的首次国际学术会议。

随后在 1991、1993 和 1994 年都举行了 KDD 专题讨论会, 来自各个领域的研究人员和应用开发者集中讨论了数据统计、海量数据分析算法、知识表示和知识运用等问题。随着参与科研和开发人员的不断增加, 国际 KDD 组委会于 1995 年把专题讨论会发展成为国际年会。在加拿大的蒙特利尔市召开了第一届 KDD 国际学术会议, 会议全称为 ACM SIGKDD(Special Interested Group on Knowledge Discovery in Databases)International Conference on Knowledge Discovery and Data Mining, 是世界数据挖掘领域的顶级学术会议。在这次会议上, “数据挖掘”(data mining)概念第一次由 Usama M. Fayyad 提出。Fayyad 同时界定了数据挖掘的内涵, 指出数据挖掘是从大量的、不完全的、有噪声的、模糊的、随机的数据中, 提取隐含在其中的、有效的、新颖的、潜在有用的并且最终可理解的模式的非平凡过程。以后每年召开一次, 参加人数由几十人发展到上千人, 研究重点也逐渐从发现方法转向系统应用, 并且注重多种发现策略和技术的集成, 以及多种学科之间的相互渗透。其中, 1997 年第三届 KDD 国际学术大会上进行的数据挖掘工具的竞赛评奖活动, 就

是一个生动的证明。1998年，在美国纽约举行的第四届 KDD 国际学术会议上，与会者不仅进行了学术讨论，而且领略了 30 多家软件公司展示的数据挖掘软件产品。最近一届于 2019 年 8 月 4 日至 8 日在美国阿拉斯加安克雷奇市举行。

2. 其他国际性数据挖掘年会

除了美国人工智能协会主办的 KDD 年会外，还有许多国际性数据挖掘年会，包括 ICDM、SDM、PAKDD、PKDD 等。ICDM(IEEE International Conference on Data Mining)是由 IEEE 组织主办的国际数据挖掘会议，会议涉及数据挖掘的所有内容，包括算法、软件、系统及应用，从 2001 年开始，每年召开一次，第 19 届会议于 2019 年 11 月 8 日至 11 日在中国北京举行。SDM(SIAM International Conference on Data Mining)是 SIAM(Society for Industrial and Applied Mathematics)组织召开的数据挖掘讨论会，2001 年 4 月召开第一届讨论会，专注于科学数据的数据挖掘，以后每年召开一次，最近一届于 2019 年 5 月 2 日至 4 日在加拿大艾伯塔省的卡尔加里市举行。PAKDD(Pacific-Asia Conference on Knowledge Discovery and Data Mining)是亚太平洋地区数据挖掘年会，从 1997 年开始，每年召开一次，第 23 届会议于 2019 年 4 月 14 日至 17 日在中国澳门举行。PKDD(European Symposium on Principles of Data Mining and Knowledge Discovery)是欧洲数据挖掘会议，也是从 1997 年开始，每年召开一次。但是从 2008 年开始，PKDD 已和欧洲机器学习会议(European Conference on Machine Learning, ECML)合并，称为 ECML PKDD，最近一届于 2019 年 9 月 16 日至 20 日在德国维尔茨堡举行。

1.1.2 数据挖掘系统的发展

数据挖掘技术所表现出的广阔应用前景及其所蕴含的巨大商业价值，吸引了国内外众多研究人员和商业机构从事数据挖掘系统的理论研究和原型开发。

1. 四代数据挖掘系统：基于技术角度的划分

从数据挖掘系统研究的技术角度看，早在 1998 年，Grossman 就提出把数据挖掘系统发展划分为四代的观点^[1]，如表 1.1 所示。

表 1.1 四代数据挖掘系统

| 代 | 特征 | 数据挖掘算法 | 集成 | 计算模型 分布形式 | 支持数据 类型 |
|-----|-----------------|------------------------------|---------------------------|-------------------|---------------------------------|
| 第一代 | 独立应用程序 | 一个或少数几个 算法 | 独立的系统 | 单台机器 | 向量数据 |
| 第二代 | 与数据库和数据 仓库集成 | 多个算法；能够 挖掘一次不能放 进内存的数据 | 数据管理系统， 包括数据库与数 据仓库 | 同质、局部区域的 计算机集群 | 一些系统支 持对象、文本 和连续的媒 体数据 |

(续表)

| 代 | 特征 | 数据挖掘算法 | 集成 | 计算模型分布形式 | 支持数据类型 |
|-----|----------------------|--------|--------------------|-----------------|----------------|
| 第三代 | 与预言模型系统集成 | 多个算法 | 数据管理系统和预言模型系统 | 内部/外部网络计算 | 半结构化数据和 Web 数据 |
| 第四代 | 与移动设备及各种计算设备结合(普适计算) | 多个算法 | 数据管理系统、预言模型系统、移动系统 | 移动和各种计算设备(普适计算) | 普遍存在的各种类型数据 |

(1) 第一代数据挖掘系统

第一代数据挖掘系统支持一个或少数几个数据挖掘算法，这些算法用来支持挖掘向量数据，作为一个独立的系统在单台机器上运行，数据一般一次性调进内存进行处理。这类工具要求用户对具体的算法和数据挖掘技术有相当的了解，还要预先完成大量的数据预处理工作。典型的系统有 Salford Systems 公司早期推出的 CART 系统等。

(2) 第二代数据挖掘系统

如果数据量非常大，需要利用数据库与数据仓库技术进行管理，第一代数据挖掘系统显然不能满足需求。第二代数据挖掘系统的主要特点是能够与数据库管理系统(DBMS)集成，支持数据库和数据仓库系统，与它们具有高性能的接口，具有高的可扩展性，支持多个算法，能够挖掘一次不能放进内存的数据，而且有些系统还能够支持挖掘对象、文本和连续的媒体数据。典型的系统如 DBMiner^[2]，能通过 DMQL 挖掘语言进行挖掘操作。

(3) 第三代数据挖掘系统

第三代数据挖掘系统除了可以与数据管理系统集成外，一个重要的优点是由数据挖掘系统产生的预言模型能够自动地被操作型系统吸收，从而与操作型系统中的预言模型相联合，提供决策支持的功能。另一个特点是支持半结构化数据和 Web 数据，能够挖掘网络环境下的分布式和高度异质的数据，并且能够有效地与操作型系统集成。典型的系统如早期被 SPSS 公司收购的 Clementine，以 PMML 格式提供与预言模型系统的接口。该系统现在被命名为 IBM SPSS Modeler，是 IBM 公司的数据挖掘工具之一。

PMML(predictive model markup language)是一种与平台无关的统计和数据挖掘模型表示标准，由数据挖掘协会(the Data Mining Group, DMG)开发，已经被 W3C(万维网联盟)接受，成为对数据挖掘模型进行描述和定义的国际标准。PMML 通过定义规范化的数据挖掘建模过程以及统一的模型表达，使得模型构造和基于模型的预测功能得以分离并可模块化实现，使得不同平台、不同数据挖掘产品之间能够共享所获得的数据挖掘模型，并为基于模型的可视化提供了条件。

(4) 第四代数据挖掘系统

第四代数据挖掘系统旨在挖掘嵌入式系统、移动系统及各种普适计算设备产生的各种类型数据。普适计算(ubiquitous computing)是软件工程和计算机科学中的概念，指计算可以使用任何设备，在任何位置，以任何格式进行。用户与计算机交互，计算机可以许多不同的形式存在，包括膝上型计算机、平板电脑和日常生活中的终端，例如汽车、冰箱或一副

眼镜。支持普适计算的基础技术包括 Internet、高级中间件、操作系统、移动代码、传感器、微处理器、新的输入输出(I/O)和用户界面、网络、移动协议、位置和定位以及新材料。物联网的不断发展,云计算、雾计算技术的广泛应用,将会进一步推动第四代数据挖掘系统的研究与发展。

2. 数据挖掘系统发展的三个阶段:基于应用角度的划分

从应用的角度,朱建秋将数据挖掘系统的发展归纳为三个阶段^[3]。

(1) 独立的数据挖掘系统

独立的数据挖掘系统对应第一代数据挖掘系统,出现在数据挖掘技术发展早期。一般研究人员开发出一种新型的数据挖掘算法,就会形成一个软件。如1993年Quinlan提出的C4.5决策树算法,1994年Agrawal和Srikant提出的Apriori关联挖掘算法等。

(2) 横向的数据挖掘工具

随着数据量的增大,数据库与数据仓库技术广泛应用于数据管理,数据挖掘系统与数据库和数据仓库的结合成为必然的选择;现实领域问题的多样性,导致一种或少数几种数据挖掘算法难以解决所有的问题;用于挖掘的数据通常不符合算法的要求,需要有数据清洗、转换等预处理的配合,才能得出有价值的模型。由于以上三方面的原因,人们认识到数据挖掘软件迫切需要结合数据库和数据仓库、多种类型的数据挖掘算法以及数据清洗、转换等预处理功能。1995年前后,软件开发商开始提供称之为“工具集”的数据挖掘系统。此类系统的特点是提供多种数据挖掘算法(通常包含分类、聚类和关联等),同时也包括数据的预处理与可视化,它是通用算法的集合,而非针对特定的应用,所以称为横向的数据挖掘工具。典型的横向工具有IBM公司的IBM Intelligent Miner、IBM SPSS Modeler和SAS公司的Enterprise Miner等。

(3) 纵向的数据挖掘解决方案

分析人员使用横向数据挖掘工具不仅需要熟悉分析的业务问题,还要精通数据挖掘算法。如果对业务或者算法不了解,就难以获得有效的模型用于决策。从1999年开始,国外大量的数据挖掘工具研制者开始提供纵向的数据挖掘解决方案,即针对特定的应用提供完整的数据挖掘方案。如在客户关系管理系统中嵌入基于神经网络的客户流失分析功能;在欺诈防护系统中嵌入基于贝叶斯的欺诈行为预测功能;在零售管理系统中嵌入客户行为分析功能,预测客户购买情况,并发送相应的优惠;在机场管理系统中嵌入旅客人数预测功能;在生产制造系统中嵌入质量控制功能等。

1.1.3 当前热点和未来趋势

1. 云计算与大数据

2006年,谷歌首席执行官埃里克·施密特推出了“Google 101计划”,正式提出“云”的概念和理论。2008年2月,美国《商业周刊》发表了一篇题为“Google及其云智慧”的文章,开篇就宣称:“这项全新的战略旨在把强大得超乎想象的计算能力分布到众人手中。”随后各大IT公司相继推出了自己的“云计划”。在中国,2009年以来,胡锦涛等前国家领

导人也先后在不同场合多次谈到“云计算”“云服务”，并把“云计算”“云服务”提到生产方式的高度。国内各大电信企业、地方政府和相关企业先后启动了云计算项目。所有这一切，预示着云计算和大数据时代的到来。

(1) 云计算

2006年，云计算创始人谷歌工程师克里斯托夫·比希利亚向首席执行官埃里克·施密特提出以谷歌设备为核心的“云计算”的想法。谷歌提供在线的网页创建、文档处理、电子表格处理等服务，用户只需要通过网络连接到谷歌的计算“云”，就可以进行相应的操作，而且能实现多人协同工作。自此，业界展开了“什么是云”“什么是云计算”“什么是云服务”的热烈讨论。

Mather等基于五个特性来定义云计算：多重租赁(分享资源)、大规模可扩展性、弹性、随用随付以及自行配置资源。^[4] Vaquero等分析了已有关于云计算的定义，认为现有定义都较多地体现某一项技术，缺乏全面性和综合性，其通过界定“云”将云计算定义为：云是一个具有大量易得易用的虚拟资源(如硬件、开发平台或服务)的资源池，这些资源可以根据不同的需求规模进行动态的重新分配，以提高资源的利用率，并实行按使用量付费的支付模式。^[5] Wang等从云计算系统功能的角度给出了云计算系统的定义，指出云计算系统不仅能够向用户提供硬件即服务(hardware as a service, HaaS)、软件即服务(software as a service, SaaS)、数据资源即服务(data as a service, DaaS)，而且还能够向用户提供能够配置的平台即服务(platform as a service, PaaS)，因此用户可以按需向计算平台提交自己的硬件配置、软件安装、数据访问需求。^[6] Fingar认为“云”包含三个层面：云计算，即一种设计模式，可实现自助服务自动化、可扩展、灵活、费用机动、数据分析方法丰富多样；云平台，即各种工具、编程与信息模型、辅助性的软件运行时组件及相关技术；云服务，一种用于信息服务的分发模型。^[7] Armbrust等认为云计算既指在互联网上以服务形式提供的应用，也指在数据中心里提供这些服务的硬件和软件，而这些硬件和软件则被称为“云”。^[8] 姚宏宇和田溯宁认为云计算应该包括两方面内容：服务和平台，云计算既是商业模式，也是技术。^[9]

基于以上不同学者的分析，本书认为云计算不仅是技术，更是一种全新的商业服务模式。云计算服务，以云资源为实现基础，以云计算技术为实现保障，以低成本、按需付费的形式，向用户提供软(硬)件基础设施、计算平台和软件服务，使用户在无基础投入的前提下直接实现数据的存储、管理和分析，也可利用提供的云服务平台创建和开发应用程序，或直接使用云服务平台提供的各类服务软件。

(2) 大数据

对于大数据，虽然众说纷纭，但有一个相对一致的说法是：大数据是超出了典型(传统、常用)硬件环境和软件工具收集、存储、管理和分析能力的数据集。由此可知，“大数据”是一个动态发展的、相对的概念。随着软(硬)件技术的发展，大数据的内涵会发生相应的变化。结合目前常用的软(硬)件技术，当下的“大数据”可以具体理解为日常关系型数据库无法收集、存储和管理的数据集。关系型数据库适合管理结构化数据，所以，当下的“大数据”除了数据量庞大(一般指PB量级)，数据形式还复杂、多样，不仅有大量的结构化数

据, 还有大量半结构化及非结构化的数据。社交网站、智能化移动设备及传感器的大规模使用, 促使数据产生的速度越来越快, 半结构化和非结构化的数据已占有绝对比重。虽然因为数据量大, 数据的价值密度相对低, 但从绝对数来看, 大数据中蕴含着大量有价值的信息。

正是因为大数据中蕴含着大量有价值的信息, 大数据被人们认为是下一个社会发展阶段的石油和金矿。各个国家把大数据当作一种全新的社会资源, 并把大数据产业的发展提升到国家战略发展的高度。类比于石油资源, 从石油的勘探、开采、运输、提炼到石油产品的生产与销售等多个环节形成了石油产业, 对于大数据的生产、采集、传输、存储、分析及应用则形成了大数据产业。在大数据产业链中, 大数据分析环节非常重要。它既是前几个环节的成果体现, 又是大数据应用及创新的基础。大数据分析的需要促进了大数据挖掘的发展, 与传统的数据挖掘相比较, 大数据挖掘将更多依赖于云计算技术, 虚拟化、可扩展的分布式数据存储模式使数据存储不仅在量上没有了限制, 而且数据形式也更为复杂, 不仅包含了大量半结构化及非结构化的数据, 还包括大量流数据。大数据挖掘将面临更海量的数据, 更复杂的数据预处理过程, 更多变的挖掘环境。

2. 从数据角度看当前热点

(1) Web 数据挖掘

Web(world wide web)是万维网的简称, 包括 Web 客户端和 Web 服务器程序。万维网可以让用户通过 Web 客户端(常用浏览器)访问 Web 服务器上带有超文本结构和多媒体的网页内容。Web 数据挖掘是数据挖掘技术在 Web 上的应用, 是从 Web 的网页内容、超链接结构和用户使用日志中获取有用知识的过程, 包括 Web 内容挖掘、Web 结构挖掘和 Web 使用挖掘。

Web 内容挖掘是指从 Web 上的文本、图像、音频、视频等其他各种类型数据及通过 Web 可以访问的数据库中的数据中发现有用知识的过程。如利用 Web 文本内容挖掘, 可以对文档进行分类、聚类。

Web 结构挖掘是从 Web 的组织结构和链接关系中发现有用知识的过程。它不仅包括文档之间的超链接结构, 还包括文档内容的结构。如利用 Web 结构挖掘, 分析网页间的超链接关系、被链接的次数及被链接的网页内容的质量, 可以对搜索引擎的文档进行重新排序。

Web 使用挖掘是从保留在 Web 服务器日志的用户访问和交互数据中提取有用信息的过程。如搜索公司利用 Web 使用挖掘探索用户的搜索模式, 预测其搜索趋势, 并向其进行相应推荐, 来提升服务质量。

(2) 文本数据挖掘

文本数据挖掘是从半结构化或非结构化文本中获取用户有用信息的过程。这一过程主要包括文本数据的获取、文本预处理、挖掘分析和结果可视化四个步骤。其中, 文本预处理一般包含数据清洗、分词(适用于中文文本)、词性标注(可选)、去停用词、特征构造、特征提取等环节; 挖掘分析主要有文本结构分析、关键词提取、文本摘要、文本分类、文本聚类、文本关联分析、文本主题模型、观点抽取、文本情感分析等。

从当前来看，数据挖掘与区块链技术还是相对独立的，但未来数据挖掘与区块链技术很可能形成技术互补。例如，区块链的可信任性、安全性和不可篡改性等，可以用来解决数据挖掘领域一直伴随的隐私与安全问题。

1.2 数据挖掘过程

从工程学的角度来看，数据挖掘是一个多环节、多处理阶段的闭环过程。如同软件工程中的软件过程模型在软件开发中的作用，数据挖掘过程模型为数据挖掘提供了宏观指导和工程方法。早期人们进行数据挖掘研究是为了将知识发现的研究成果应用于实际数据处理中，为科学决策提供支持。因此，大多数研究人员只着眼于数据挖掘的算法和应用层面，而忽视了其他方面。事实上，数据挖掘首先是一个处理过程，如果我们仅仅着重于挖掘，可能就看不到实际工作中数据处理过程的数据提取、组织和显示方式的难度。合理的数据挖掘过程模型能将各个处理阶段有机地结合在一起，指导人们更好地开发、使用数据挖掘系统和实施数据挖掘项目。从数据挖掘进入工程应用领域起，就有人对数据挖掘的过程进行归纳和总结，以便于人们开发及使用数据挖掘应用系统。目前被业界广泛认可并已应用于商用软件的数据挖掘过程模型主要有两种：一种是 Fayyad 等人总结的过程模型，另一种是遵循 CRISP-DM 标准的过程模型。

1.2.1 Fayyad 过程模型

Fayyad 等将知识发现过程定义为：从数据中鉴别出有效模式的非平凡过程，该模式是新颖的、可能有用的和最终可理解的。^[10]图 1.5 是 Fayyad 给出的过程模型。早期开发的大部分数据挖掘系统都是遵循 Fayyad 过程模型，例如 IBM Intelligent Miner 和 SAS Enterprise Miner 等。

如图 1.5 所示，Fayyad 过程模型包括数据选择(data selection)、数据预处理(data preprocessing)、数据转换(data transformation)、数据挖掘(data mining)、模式解释与评价(pattern interpretation)。

1. 数据选择

数据选择是指根据分析任务的要求从原始数据中提取与挖掘目标相关的数据，并将不同数据源中的数据集成在一起，形成本次数据挖掘任务的数据集。在此过程中，会利用一些数据库操作对数据进行处理。

2. 数据预处理

数据预处理是指对数据选择阶段产生的数据进行再加工，检查数据的完整性及数据的一致性，对其中的噪音数据进行处理，对缺失的数据进行填补等。

3. 数据转换

数据转换是指对经过预处理的数据，根据挖掘事务的任务对数据进行再处理，主要转

换成数据挖掘算法所需要的形式，如将连续型数据转换成离散型等。

4. 数据挖掘

数据挖掘是指运用合适的数据挖掘算法，从数据中提取出用户所需要的知识，这些知识可以用一种特定的方式表示或使用一些常用的表示方式，如产生规则等。

5. 模式解释与评价

模式解释与评价是指根据分析目的，对发现的模式进行解释，并评价模式的有效性。在此过程中，为了取得更有效的模式，可能会返回前面处理步骤中的某些步骤，从而提取出更有用的知识。

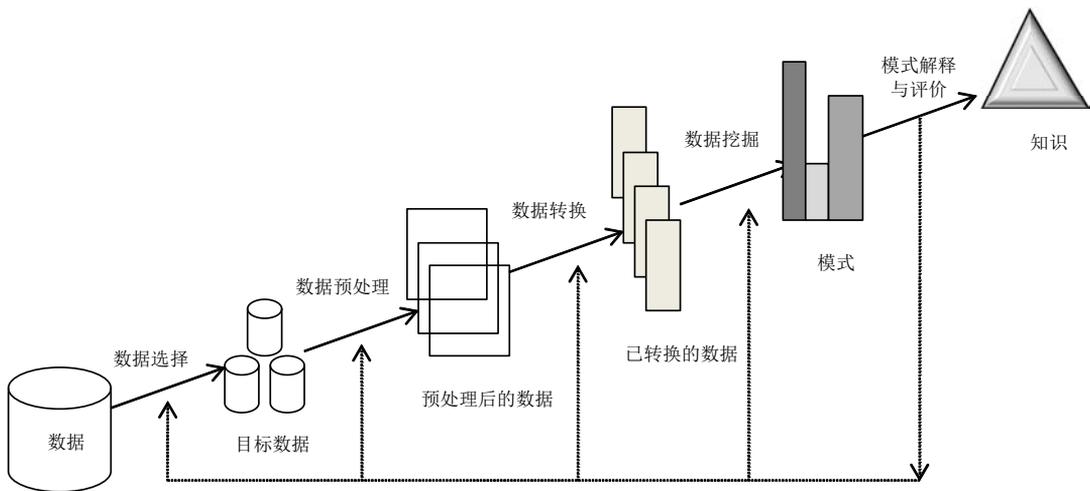


图 1.5 Fayyad 过程模型

从上述 Fayyad 过程模型来看，这个过程已经包括了数据挖掘过程中必要的各个处理阶段，并且也形成了一个可以根据各个处理阶段的结果来决定是否返回以前的阶段进行再处理的闭环过程。但是，Fayyad 过程模型从数据入手，到知识结束，过多地偏重于从技术的角度来理解数据挖掘过程。在实际使用过程中会存在两个问题：第一，数据选择对于整个分析至关重要，但是该如何选择，选择哪些数据呢？这是由具体的商业问题决定的，需要领域专家、数据管理员与数据挖掘专家一起讨论确定。如何明确商业问题，并把商业问题和数据相关联，这在 Fayyad 过程模型中没有反映。第二，数据挖掘一般在分析型环境中获得知识，获得的知识只有返回到操作型环境中使用，才能产生真正的价值。在 Fayyad 过程模型中，模式评价阶段结束后，对于挖掘到的知识应该如何使用，这方面也没有体现。

1.2.2 CRISP-DM 过程模型

CRISP-DM(cross-industry standard process for data mining)即跨行业数据挖掘过程标准，它由 SPSS、NCR 以及当时的戴姆勒-克莱斯勒等公司在 1996 年提出，后来得到欧洲共同体

研究基金的资助。2000年8月, CRISP-DM 1.0版正式推出^[1]。CRISP-DM强调, 数据挖掘不单是数据的组织或者呈现, 也不仅是数据分析和统计建模, 而是一个从理解业务需求、寻求解决方案到接受实践检验的完整过程。如图1.6所示, CRISP-DM过程模型包括商业理解(business understanding)、数据理解(data understanding)、数据准备(data preparation)、建模(modeling)、评价(evaluation)和部署(deployment)六个阶段。图1.6的外圈形象地表达了数据挖掘过程的循环特性。一个数据挖掘项目并不是一次部署完就结束, 在挖掘的过程中或部署过程中获得的经验可能会触发新的商业问题。后续的挖掘过程将从前一次的经验中受益。内部的箭头表示阶段之间最重要和最频繁发生的关联关系。阶段间的顺序不是严格不变的, 可以根据具体的任务需要进行来回选择。

CRISP-DM过程模型标准不仅被许多数据挖掘软件商用来指导开发数据挖掘软件, 如IBM公司的IBM SPSS Modeler就是遵循了CRISP-DM过程标准。同时, 该标准也被广泛用来指导数据挖掘项目的实施。

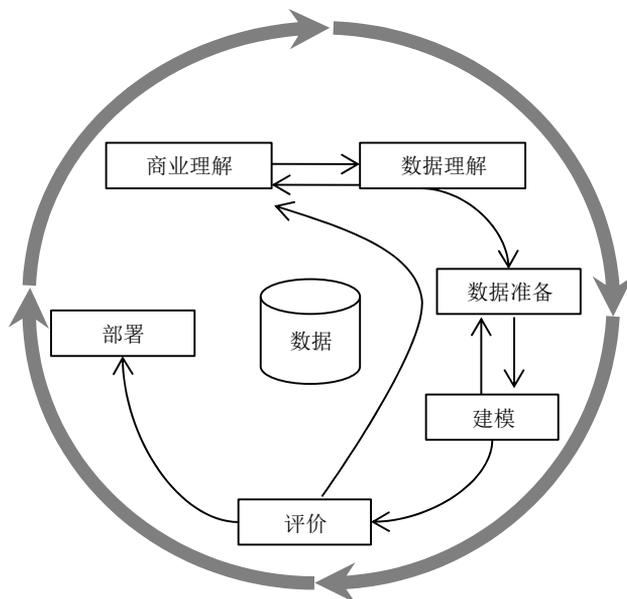


图 1.6 CRISP-DM 过程模型

1. 商业理解

商业理解是对企业运作、业务流程和行业背景进行了解, 专注于从商业的角度理解项目目标和需求, 然后将这种目标和需求转换成一个数据挖掘的问题定义及相应的项目计划, 其一般任务和输出内容如图1.7所示。

(1) 确定商业目标

数据分析师最重要的能力是对业务的理解和把握, 没有正确的业务理解, 再好的理论, 再强的工具, 都只会徒劳无益。所以, 一个数据挖掘项目的实施, 其首要任务就是从业务的角度真正理解所要解决的问题和所要实现的目标。完成确定商业目标这一任务, 其相应的输出文档内容一般包括背景、商业目标和商业成功标准三个方面。

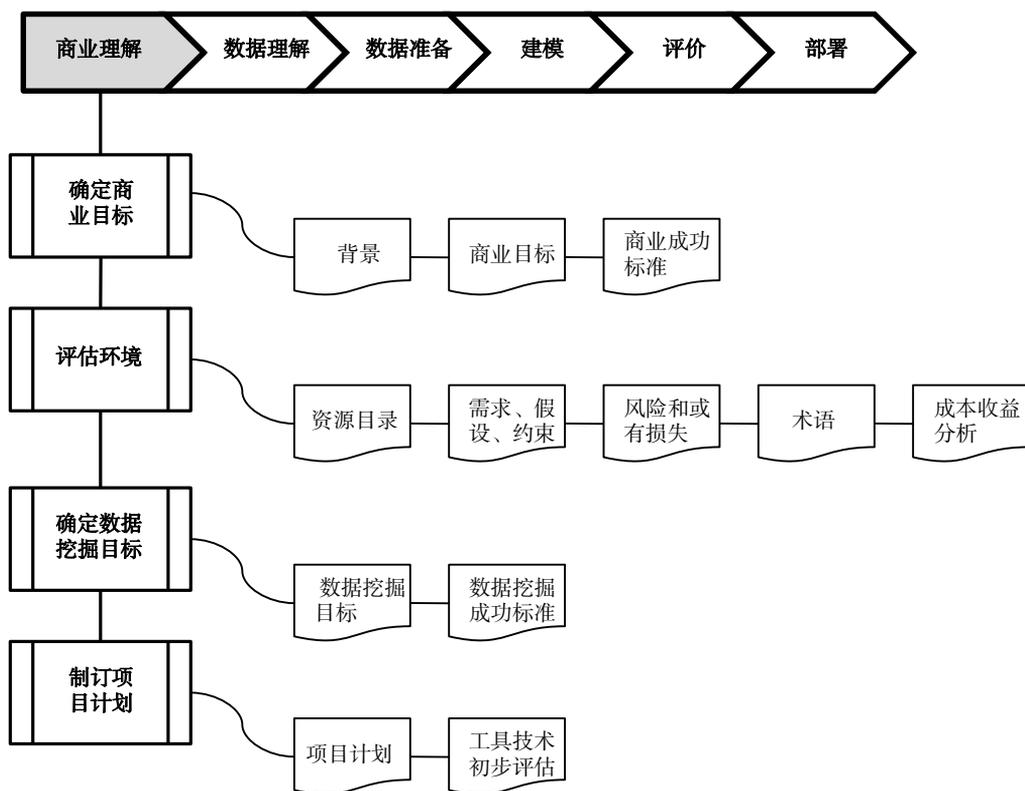


图 1.7 商业理解的一般任务(加粗显示部分)和输出

① 背景包括项目的商业环境，问题涉及的范围，项目的前提(如现有解决方案的优缺点、项目的动机、是否已经使用数据挖掘等)，项目需要的人力和物质，项目将会影响到的部门和使用项目结果的目标群体等。

② 商业目标是从商业的角度来描述打算用数据挖掘来解决的问题。尽可能准确地分析所有相关的商业问题，分清主要的商业目标及其他次要目标，制订尽可能实现的目标，并使用商业术语，详细说明期望收益。

③ 商业成功标准是从商业角度衡量项目结果成功的度量标准，包括客观度量标准(如投诉率下降 15%、下单转换率增加 20%等)和主观度量标准。主观度量标准要明确主观的主体，即是谁给出的主观判断。

(2) 评估环境

评估环境任务主要围绕已确定的商定目标和初步计划细化各种影响因素，其相应的输出文档内容一般包括资源目录，需求、假设和约束，风险和或有损失，术语及成本收益分析五个方面。

资源目录文档需要列出项目可用的各类资源，包括参与人员(项目发起人、相关商业领域专家、数据库管理员、市场分析师、数据挖掘专家及其他技术支持人员)，数据(企业内部固定抽取的数据、访问内部数据库或数据仓库的数据、外部调查或购买的数据等)，计算资源(硬件平台)和软件(数据挖掘工具及其他相关软件)。

需求、假设和约束文档要求列出项目执行的全部需求、围绕项目整个过程的各方面假设及约束。全部需求可包括：项目完成的时间进度表及相应进度的需求，项目和模型的可理解性、准确性、可部署性、可维护性和可重复性等方面的需求，安全、隐私及法律限制等方面的需求。假设包括对外部因素(如商业环境、经济问题、技术因素等)的假设，数据质量(如可用性、准确度等)的假设，模型理解、解释与评估时可能的假设等。约束包括一般性约束(如法律问题、经费、时间及其他所需资源)，数据源访问权利，数据访问时的技术性问题等。

风险和或有损失(contingencies)文档要求列出可能导致项目延期或失败的风险、可能的损失和为避免这些风险可采取的相应措施。确定每个风险可能发生的条件，如法律风险、商业风险、组织风险、经济风险、技术风险及与数据或数据源有关的风险(数据质量相关问题)等，并计算相应的可能损失，制订损失计划。

术语文档要求编辑一个与项目有关的术语表。术语表至少包括与商业问题有关的术语和与数据挖掘有关的术语两部分内容，以帮助不同专业背景的项目参与人员更好地理解项目。

成本和收益文档要求分析项目执行的成本和项目部署后可能产生的收益(如投资回报率、客户满意度等)。除了数据收集、项目开发和运行等成本，还必须考虑数据重复抽取和准备、工作流程的改变等隐含成本。

(3) 确定数据挖掘目标

确定数据挖掘目标这一任务就是要根据已确定的商业目标，从数据挖掘的角度，用数据挖掘技术术语来描述项目目标和项目成功的标准。其相应的输出文档内容一般包括数据挖掘目标和数据挖掘成功标准两个方面。

数据挖掘目标要求把商业问题转换成数据挖掘问题，也即确定业务问题需要什么类型的挖掘模型加以解决。如商业目标是要确定哪些客户会流失，则数据挖掘目标是构建一个客户流失预测模型，可以是客户是否流失的分类预测，也可以是客户流失概率预测。

数据挖掘成功标准指模型评估的标准。如对于客户是否流失的分类预测模型，可以使用准确率、精准率和召回率等评价指标来评估模型。如果是主观评价标准，和商业成功主观标准一样，需要明确这个标准是由哪个人或哪些人作出的主观判断。

(4) 制订项目计划

为达到数据挖掘目标进而实现商业目标，需要制订详细的项目计划。该计划要求详细列出项目需要完成的一系列步骤，包括对工具和技术的选择。其相应的输出文档内容一般包括项目计划及工具和技术的初步评估。

项目计划需要列出每个阶段的详细计划，包括持续的时间、需要的资源、输入、输出、可能的风险及关联性。在计划中要交代清楚可能的重复步骤及所需的时间。在估计项目时间进度时可以参考他人的经验，如数据理解和数据准备通常需要占用60%~80%的时间。分析时间进度和可能的风险之间的关联性，尽可能避免风险。

工具和技术的选择可能影响整个项目，所以要尽早列出工具和技术的选择标准，评估技术的合适程度，选择最合适的工具和技术。

2. 数据理解

数据理解是对企业现有应用系统进行了解，对数据挖掘所需数据进行全面调查以获取

完成挖掘目标所需的初步数据，然后从总体上对获得的数据的属性进行描述，包括数据格式、数据量、一致性、数据出处、收集时间频度等多个方面，并检查数据是否能够满足相关的要求，探测数据和检验数据质量等。其一般任务和输出文档内容如图 1.8 所示。

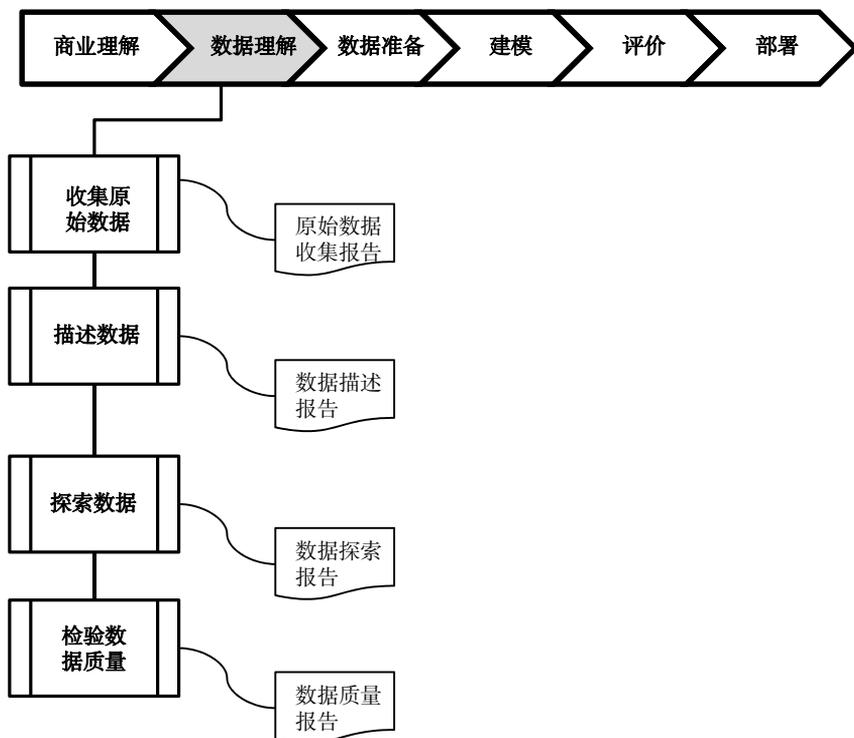


图 1.8 数据理解的一般任务(加粗显示部分)和输出文档

(1) 收集原始数据

收集原始数据任务是根据资源目录列出的数据资源选择感兴趣的表或文件，并选择表或文件中感兴趣的数据。完成这一任务要求产生相应的输出文档——原始数据收集报告。该报告应包括以下内容：数据来源(内部数据库或数据仓库、外部提供者)，负责维护、收集或购买此数据的人，调查或购买数据需要的费用，数据存储方式，安全和隐私需求、使用限制等。

(2) 描述数据

描述数据任务要求描述所获得的数据，包括数据数量(表、各个表的字段数和记录总数)，数据类型，编码方案，计量单位，取值范围或个数，属性和属性值的含义，主键和外键的关系，缺失数据占比等。该任务相应的输出文档是数据描述报告。

(3) 探索数据

探索数据任务是根据数据挖掘目标，结合数据描述报告，采用表格、图形和其他可视化技术细致探索数据，包括关键属性的分布、属性间的关系及一些简单的统计分析。这些分析丰富或细化了数据描述，可以作为后续数据准备工作的输入，或可能直接达到某个数据挖掘目标。这一任务将产生相应的输出文档——数据探索报告。

(4) 检验数据质量

检验数据质量任务需要对收集的数据从是否完整、是否缺失、是否一致、有无异常等方面进行检查，并产生该任务相应的输出文档——数据质量报告。该报告要求列出数据质量检验的结果，对于存在的质量问题，列出可能的解决方法。质量问题的解决方法很大程度上依赖于数据和商业知识。

3. 数据准备

数据准备是数据挖掘过程中最重要的一个环节，通常需要花费大量的时间，一般占用整个数据挖掘项目 50%~70%的时间和 workload。数据准备需要从所收集的大量原始数据中取出一个与业务目标相关的样本数据集，对该数据集进行描述，在此基础上，将该数据集转化为适合数据挖掘工具处理的最终目标数据，包括选择数据、清洗数据、构造数据、集成数据和格式化数据。其一般任务和输出文档内容如图 1.9 所示。

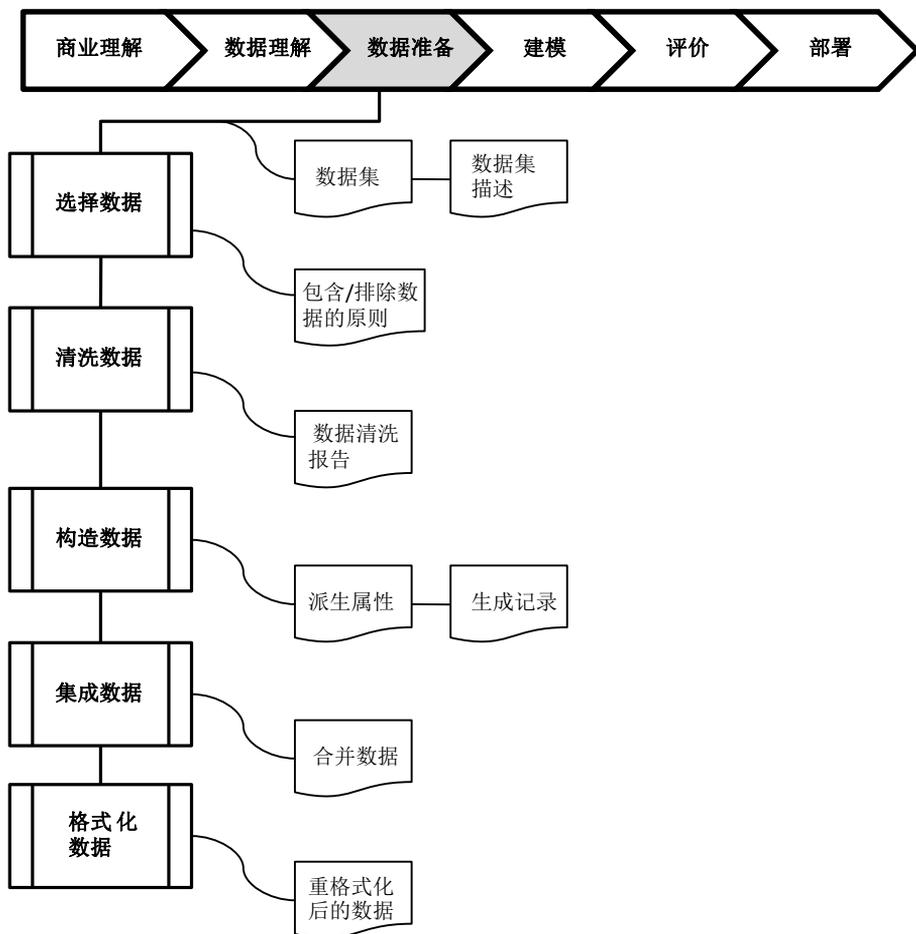


图 1.9 数据准备的一般任务(加粗显示部分)和输出文档

(1) 选择数据

选择数据需要确定用于分析的数据，包括对样本的选择和对属性或特征的选择。选择

的标准直接影响用于分析的数据质量，所以选择标准的确定至关重要，可以从与数据挖掘目标的相关性角度考虑，进行显著性检验或相关性分析作为属性或特征的选择标准，也可以从数据质量、容量与类型等方面限制作为选择的标准。其相应的输出文档为包含/排除数据的原则，需要列出被包含进来的和被排除出去的数据，并给出理由。

(2) 清洗数据

清洗数据主要是基于已选择的数据，选择合适的方法处理噪声、填补缺失值等，保证数据的正确性和一致性，提升数据质量。其相应的输出文档为数据清洗报告。该报告不仅要描述清洗的策略和行为，还要指出清洗后的数据用于挖掘时仍然可能存在的质量问题以及对挖掘结果的潜在影响。

(3) 构造数据

构造数据主要指派生属性(列或特征)、生成全新的记录(行)及对现有属性值进行转换等。派生属性是在一个或多个现有属性基础上构造符合挖掘目标需要的属性，例如为了预测客户是否会流失，通过对客户消费行为的分析，界定流失的内涵，构造新的属性“是否流失”，作为目标变量用于预测。该任务相应的输出文档即为构造的结果——派生属性和生成记录。

(4) 集成数据

集成数据是指把来自不同数据源的数据整合在一起，可以合并多个表，也可以通过数据合并构造新的记录和属性。例如，一家电子商务公司有两张客户信息表：一张为客户基本信息表，包括客户 ID 号、姓名、年龄、性别等客户基本信息；另一张为客户购买信息表，包括客户 ID 号、客户近一个月的购买明细记录，每一条记录对应每笔购买。对这两张表进行集成，可以先根据客户购买信息表生成一个新表，其中每条记录对应每个客户，属性则为客户 ID 号、购买次数、平均购买额、购买促销商品的比例等，再利用客户 ID 号，集成新表和客户信息表。该任务相应的输出文档即为集成的结果——合并数据。

(5) 格式化数据

格式化数据作为建模前的最后一个步骤，主要针对某些建模对数据的特殊格式要求进行改变。例如有些建模算法要求记录按某个属性值排序，有些建模算法又要求记录是随机排列的。对于文本数据，某些建模算法要求去掉文本字段内的标点符号，或者规定每个字段的值所允许的最大字符数。该任务相应的输出文档即为格式化后的结果——重格式化后的数据。

4. 建模

建模是根据对业务目标的理解，在数据准备的基础上，选择和应用多种不同的建模技术，调整它们的参数使其达到最优值，包括选择建模技术、产生测试设计、构建模型和评估模型。其一般任务和输出文档内容如图 1.10 所示。

(1) 选择建模技术

选择建模技术是结合数据挖掘目标确定实际所要使用的建模技术，可以是一种技术，也可以是多种技术，或者是基于多种技术的集成。确定了相应的技术后，需要了解所选技术对数据的假定要求，并产生相应的输出文档——建模技术和建模假设。

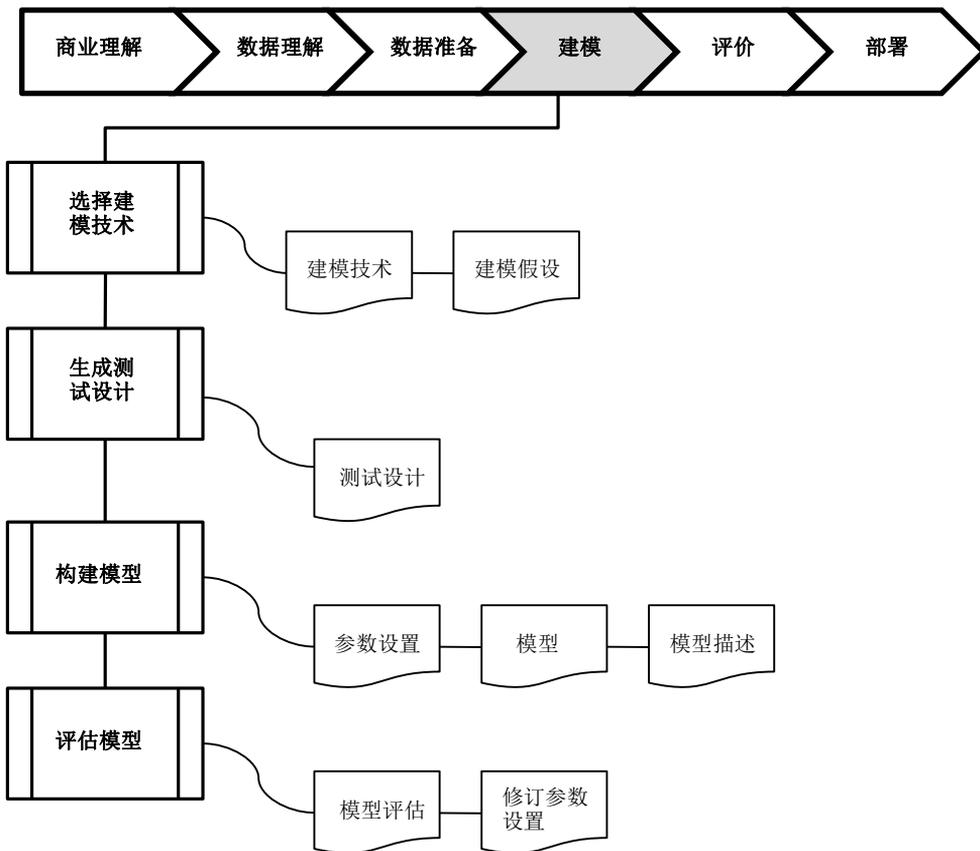


图 1.10 建模一般任务(加粗显示部分)和输出文档

(2) 生成测试设计

生成测试设计是指在实际构建模型前，建立一个用来测试模型质量和有效性的机制，包括数据集如何划分、划分成几部分(如训练集和测试集)、如何验证模型质量。其相应的输出文档是测试设计。

(3) 构建模型

构建模型指在准备好的数据集上使用建模工具，创建一个或多个模型。相应的输出文档为参数设置、模型和模型描述。参数设置列出模型需要调整的参数、相应的设置值及选择设置值的基本原则。模型指产生的实际模型，如决策树模型、神经网络模型。模型描述指描述模型的特征，生成解释模型的报告。

(4) 评估模型

评估模型是指数据挖掘工程师根据领域知识、数据挖掘目标成功标准和已生成的测试设计来解释模型。这一任务仅考虑模型，对后续的评价阶段会产生影响。评价阶段需要数据挖掘工程师和领域专家、业务分析人员一起考虑项目实施过程中产生的所有结果。相应的输出文档是模型评估和修订参数设置。模型评估列出全部建成的模型及其评估结果，如按准确率比较建成模型的优劣。根据模型评估结果，重新修订参数设置，并调整其值建立新的模型，直到数据挖掘工程师确信已找到最优模型为止。修订参数设置指记录所有这些修订和评估。

5. 评价

评价是由分析人员和领域专家一起从业务目标的角度全面地评价得到的模型，以确定它是否完全达到了业务目标，最终做出是否应用数据挖掘结果的决策。其一般任务和输出文档内容如图 1.11 所示。

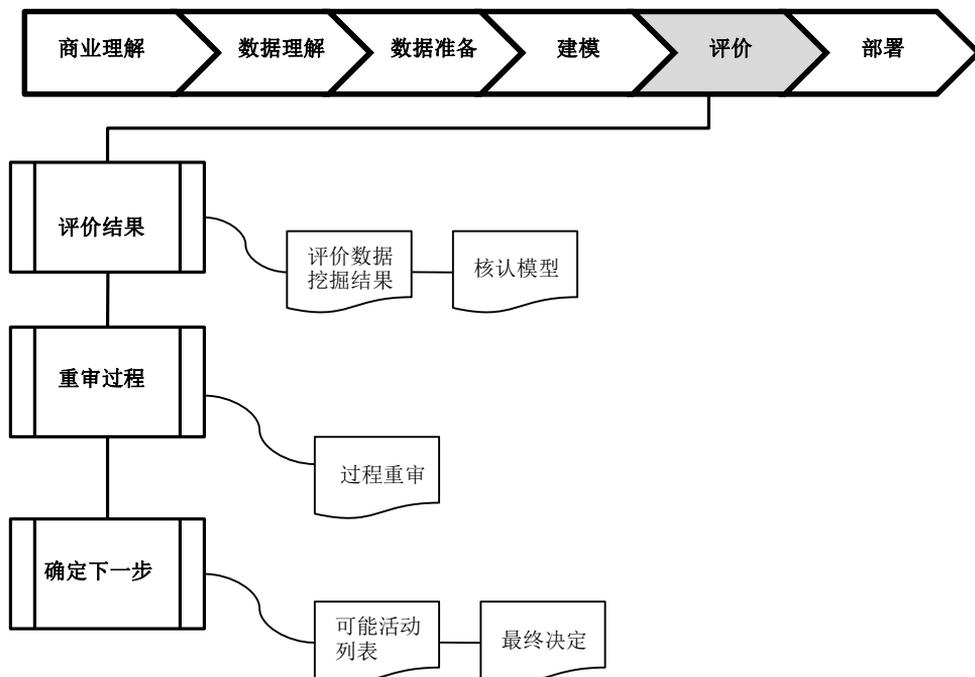


图 1.11 评价一般任务(加粗显示部分)和输出文档

(1) 评价结果

评价结果是评价模型是否符合商业目标，若存在不足，说明其商业理由，相应的输出文档为评价数据挖掘结果和核认模型。评价数据挖掘结果是指使用商业成功标准术语概述模型评价的结果，包括是否已满足既定商业目标的最终声明。核认模型是核准认可满足既定商业成功标准的模型。

(2) 重审过程

重审过程是指对数据挖掘项目实施的整个过程进行重新审核，用来确定是否忽略了某些重要的因素或任务，或者是否存在某些质量问题。其相应的输出文档为过程重审，即概述重审过程，并特别注明被忽略的因素或应该重复的环节。

(3) 确定下一步

确定下一步是指根据评价结果和重审过程，来分析项目该如何推进，需要确定是进入部署阶段还是继续重复前面步骤或者创建新的数据挖掘项目，同时，还要分析剩余的资源 and 预算。其相应的输出文档是可能活动列表和最终决定。可能活动列表列出潜在的进一步活动，并给出支持和反对每个结果的理由。最终决定描述如何合理推进。

6. 部署

部署是数据挖掘的最终目的，是将数据挖掘结果部署到商业环境中，成为日常商业运作的一部分，并生成一份基于项目整个过程的最终报告。其一般任务和输出文档内容如图 1.12 所示。

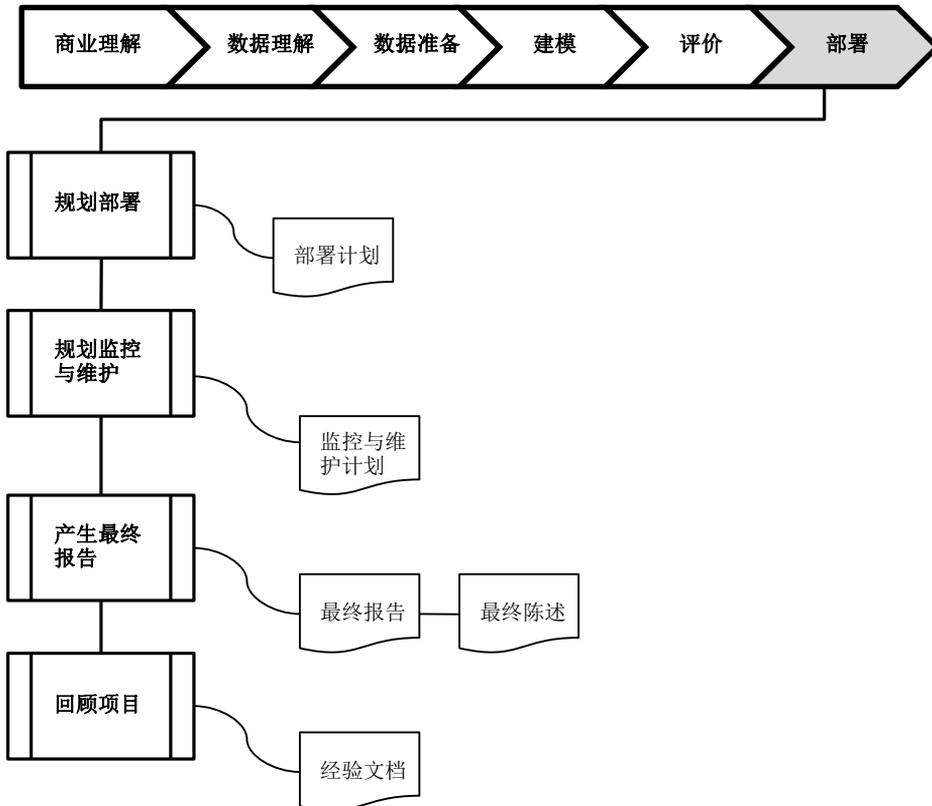


图 1.12 部署一般任务(加粗显示部分)和输出文档

(1) 规划部署

规划部署是指为了把数据挖掘结果部署到商业环境中，利用评估的结果给出部署的策略。其相应的输出文档是部署计划，即概述部署策略，包括必要的步骤和如何执行这些步骤。

(2) 规划监控与维护

数据挖掘结果成为日常商业运作的一部分时，监控和维护就成为重要问题。规划详细有效的监控和维护策略有助于避免长期错误应用数据挖掘结果。其相应的输出文档是监控与维护计划，即概述监控和维护策略，包括必要的步骤和如何执行这些步骤。

(3) 产生最终报告

项目的结束需要项目成员撰写一份最终报告，这份报告可能仅对项目和其经历进行概述，也可能对数据挖掘结果进行全面展示。其相应的输出文档是最终报告和最终陈述。最终报告可以描述全部过程并标明全部取得的结果，说明与原始计划的偏差，并给出将来工作的建议。其具体内容和形式很大程度依赖于报告的接受者。最终陈述一般只包括最终报

告的一部分内容，可以不同于报告的形式呈现。

(4) 回顾项目

回顾项目指总结经验，评论成功与失败之处，并指出如何改进。其相应的输出文档为经验文档，即描述项目期间获得的重要经验。

1.3 数据挖掘功能与使用技术

数据挖掘功能用于指定数据挖掘任务发现的模式。一般而言，这些任务可以分为两类：描述性的和预测性的。描述性数据挖掘任务是刻画目标数据中数据的一般性质。预测性数据挖掘任务是在当前数据上进行归纳，以便作出预测。^[12]随着信息技术的持续发展，数据挖掘吸纳了统计学、机器学习、模式识别、数据库与数据仓库、信息检索、可视化、分布式并行计算等更多领域的大量技术。

1.3.1 数据挖掘功能

常见的数据挖掘功能可以概括为六个方面：数据描述、聚类、偏差检测(孤立点检测)、关联分析、预测和分类，如图 1.13 所示。其中，数据描述、聚类、偏差检测和关联分析可以认为是描述性任务，分类和预测可以认为是预测性任务。

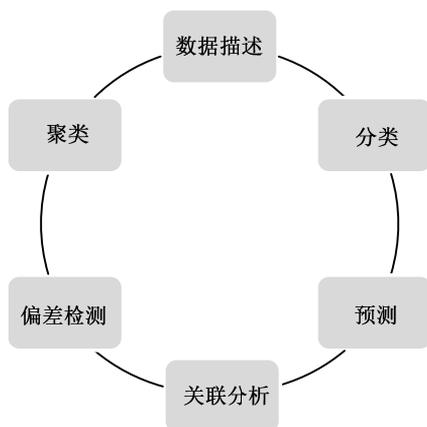


图 1.13 数据挖掘的主要功能

1. 数据描述

数据描述可以分为特征性描述和区别性描述。特征性描述用来反映目标数据的一般特征；区别性描述用来比较目标数据与一个或多个类比数据的不同特征。数据描述通常以图形、二维或多维表呈现描述结果，也可以以规则的形式呈现。

2. 聚类

聚类指按照尽量使同一个类(簇)中的数据之间具有较高的相似性，而不同类(簇)中的数

据之间具有较大的差异性的原则将数据划分成有意义或有用的类(簇)。数据事先不存在类标号。

3. 偏差检测(孤立点检测)

偏差检测也称孤立点检测，指通过发现数据集中特殊的变化，寻找孤立点，并对其进行分析，探究原因，以确定是否是事物发生的突变。

4. 关联分析

关联分析指通过挖掘频繁模式来发现大量数据中有趣的关联或相关联系。例如通过购物篮分析，确定哪些商品通常会被一起购买，从而制订交叉销售等相应的营销策略。

5. 预测

预测指用过去和现在的数据去拟合模型，并使用模型预测未来。广义上，预测包含分类，是对类别变量的预测，狭义的预测仅指对连续型变量的预测。

6. 分类

分类指找出描述并区分数据类或概念的模型(或函数)，从而使用该模型(或函数)来预测类标记未知的对象类。用于寻找模型的数据存在类标号，分类一般针对类别变量。

1.3.2 数据挖掘使用技术

数据挖掘的产生和发展一直受应用驱动。随着应用不断拓宽，其所使用的技术也越来越丰富，而且还将持续发展，如图 1.14 所示。

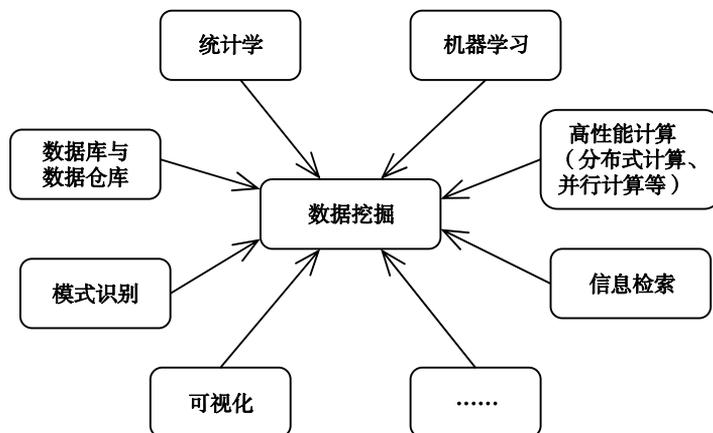


图 1.14 数据挖掘使用技术

1. 统计学

从统计学的发展过程看，统计学一直面临着来自自然科学、工业及商业等各个领域的挑战，从而不断得到充实和发展。随着计算机软硬件技术的飞速发展，数据存储能力无限量提高，面对海量且形式多样的数据，传统统计学方法在应用时遇到了新的挑战。数据

挖掘正是统计学适应这一变化的新的发展方向。数据挖掘并不是为了替代传统的统计分析技术，而是统计分析方法的延伸和扩展。Ganesh(2002)认为，从统计学的视角看，数据挖掘可以被看成是对大容量复杂数据的计算机自动化的探索和分析，可以被认为是“智能化统计”^[13]。因此，统计方法自然成为数据挖掘的一大技术支撑。

传统的统计方法可以分为描述统计和推断统计。描述统计主要对观察到的数据进行汇总、分类和计算，并用表格、图形和指标的形式来反映现象的数量特征。推断统计则以已知的数据(部分的或过去的)去推断未知的数据(整体的或未来的)。这两类方法正好符合数据挖掘两大任务(描述和预测)的需要，数据挖掘把统计学技术与计算机技术相结合，从数据中发现有用的知识。

2. 机器学习

机器学习是指计算机利用各种学习算法，从输入的数据中学习，识别复杂的模式，从而作出智能的决断。因为学习算法中涉及了大量的统计学理论，机器学习与推断统计学的联系尤为密切，也被称为统计学习理论。

机器学习的基础是数据，核心是各种学习算法，只有通过这些算法，机器才能识别分析这些数据，获得知识，从而不断提升自身性能。机器学习的算法很多，根据学习方式不同，可以分为有监督学习(supervised learning)、无监督学习(unsupervised learning)、半监督学习(semi-supervised learning)和强化学习(reinforcement learning)。

(1) 有监督学习

用于有监督学习训练的数据集包含输入(特征)和输出(目标)，也称为有标记的数据集。从有标记数据集中根据输入和输出学习出一个模型，即为有监督学习。当新的数据输入时，可以根据这个模型预测结果。由于训练集中存在目标，因此学习得到的模型可以使用历史数据进行验证，从而起到监督的作用。有监督学习算法主要应用于分类和回归，如决策树、朴素贝叶斯、Logistic 回归等。

(2) 无监督学习

用于无监督学习训练的数据集只包含输入(特征)，而没有输出(目标)，也称为无标记数据集。从无标记数据集中通过学习进行归纳，获得数据分布特征或数据与数据之间的关系，即为无监督学习。由于训练数据不存在目标，因此学习得到的模型不能使用历史数据进行验证，从而无法监督。无监督学习算法主要应用于聚类和关联分析，如 K-均值聚类、Apriori 算法等。

(3) 半监督学习

有两个数据集用于半监督学习，一个为有标记的数据集，一个为无标记的数据集，通常无标记数据集的数据量要远远大于有标记数据集的数据量。如上所述，如果单独使用有标记数据集，我们能够生成有监督模型；单独使用无标记数据集，我们能够生成无监督模型。但为了最大限度利用现有数据的信息，我们希望使用两个数据集进行学习。用户可以在有标记数据集中加入无标记数据，增强有监督学习的效果，如半监督支持向量机；也可以在没有标记数据集中加入有标记数据，增强无监督学习的效果，如半监督聚类。一般而言，半监督学习侧重于在有标记数据集中加入无标记数据来增强学习效果。

(4) 强化学习

强化学习是智能体(agent)在尝试的过程中学习在特定的环境下选择哪种行动可以得到最大的回报。如图 1.15 所示,智能体在学习的过程中选择一个动作,环境接受该动作后状态发生变化,同时产生一个强化信号(奖励或惩罚),反馈给智能体,智能体根据强化信号和环境当前状态再选择下一个动作,选择的原理是使受到的正强化(奖励)最大。智能体当下选择的动作不仅影响当下的强化值,而且影响环境下一时刻的状态及最终的强化值。

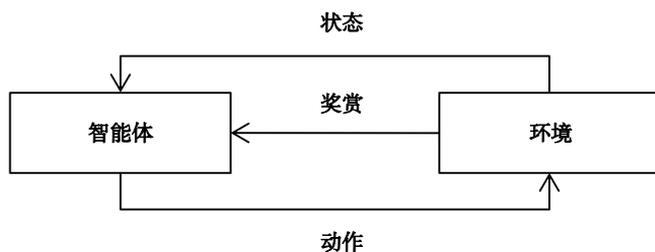


图 1.15 强化学习示意图

3. 数据库与数据仓库

(1) 数据库

数据库指的是以一定方式储存在一起、能为多个用户共享、具有尽可能小的冗余度、与应用程序彼此独立的数据集合。简单来说,可视为电子化的文件柜——存储电子文件的场所,用户可以对文件中的数据执行新增、截取、更新、删除等操作。数据库管理系统(database management system, DBMS)是管理数据库的大型计算机软件系统,用于建立、使用和维护数据库。它的主要功能包括:创建数据库,创建表,创建支撑结构(索引),读取数据库数据,修改(插入、更新、删除)数据库数据,维护数据库结构,执行规则,并发性控制,安全性控制,备份和恢复。其并发性控制功能确保一个用户的工作不会不适当地影响其他用户的工作,保证多个用户在同一时刻对同一数据进行读、写等操作时数据的一致性。

利用可伸缩的数据库技术,数据挖掘可以在大型数据集上获得高效率 and 可伸缩性。同时,数据挖掘技术也有利于扩充数据库系统的能力,满足高端用户复杂的数据分析需求,实现商务智能,如图 1.16 所示。

(2) 数据仓库

尽管对于小型数据库或者在线处理任务不多的数据库,直接从日常数据库中读取数据用于数据挖掘是可行的,但是对于更大的数据库或是要满足更多在线处理任务的数据库,直接从日常数据库中读取数据用于数据挖掘是不可行的。原因如下:第一,联机事务处理系统强调的是数据处理性能和系统的安全与可靠性,并不关心数据查询的方便与快捷,直接从日常数据库中读取数据用于数据挖掘会给日常 DBMS 带来很大负担,影响其运行性能;第二,用于数据挖掘的数据通常来源于不同的事务数据库,不同事务数据库数据的模式是针对具体事务处理而设计的,可能存在不一致等问题,不适合直接用于数据挖掘;第三,

联机事务处理(OLTP)系统可能缺少数据挖掘需要的大量历史数据。因为这些原因,很多企业选择使用数据仓库来进行数据挖掘。

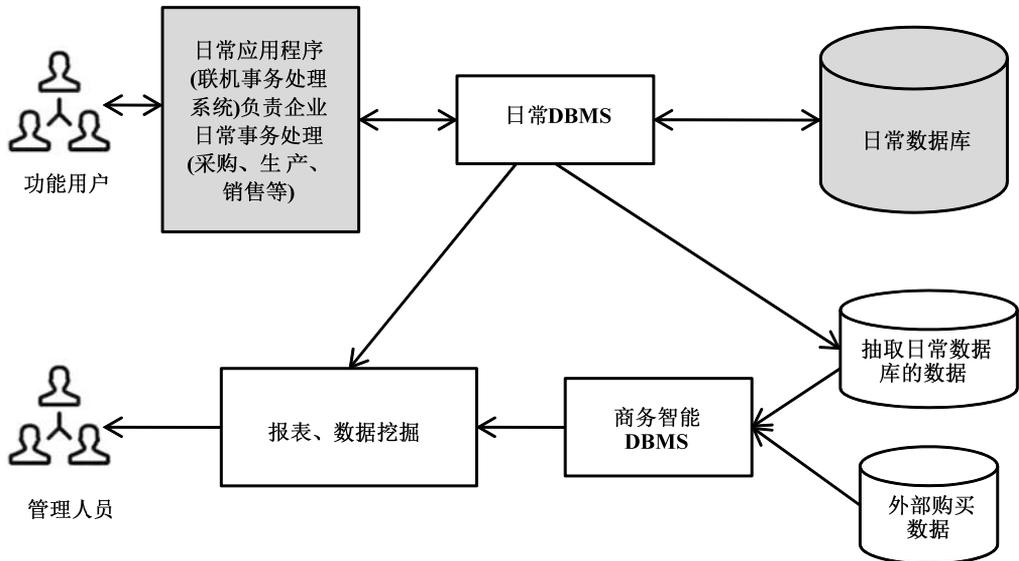


图 1.16 数据库与数据挖掘

数据仓库的数据是为了分析需要,按分析主题组织的(如客户、商品、供应商等),是集成的、稳定的,除了随时间批量载入外,是不能更改的,只能查询。如图 1.17 所示,数据仓库为数据挖掘提供了更好的、更广泛的数据源,为更好地实施数据挖掘提供了方便。同时,数据挖掘也为数据仓库提供了更复杂的数据分析,更有效的决策支持。

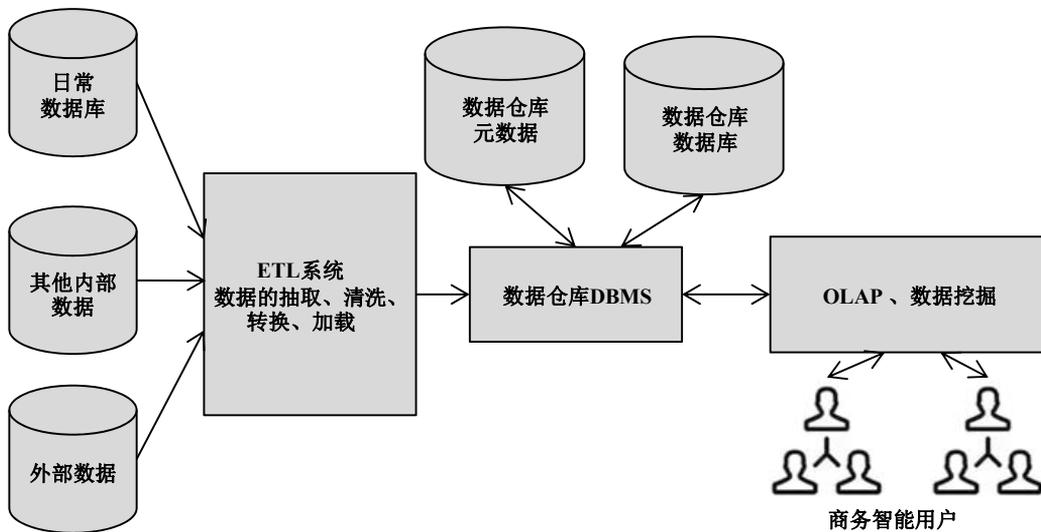


图 1.17 数据仓库与数据挖掘

1.4 数据挖掘应用

数据挖掘从一开始就是面向应用的，随着各行各业信息化的持续发展，数据挖掘应用领域也在不断发展和深化。目前数据挖掘在金融、电信、零售与电子商务、政府政务、医疗、科学等领域都有应用。

1.4.1 金融领域的数据挖掘

银行、证券和保险等金融领域，信息化建设较早，积累了大量的数据，是数据挖掘的重要应用领域。典型的应用有：金融风险分析、金融产品交叉销售、客户管理分析、洗黑钱等金融犯罪识别等。金融交易活动过程很可能存在洗黑钱等犯罪行为，把可能与侦破有关的数据集成(如金融机构交易数据库、犯罪历史数据库等)，运用合适的数据挖掘方法(数据可视化、孤立点分析等)，检测异常模式，可以为犯罪行为识别提供快速准确的参考。

银行业利用数据挖掘技术最集中的两个方面是风险管理和客户管理。风险管理，如信用风险评估，银行可通过建立信用风险模型，评估贷款申请人或信用卡申请人的风险，根据评估结果来决定是否接受申请，并确定贷款额度或信用额度。客户管理体现在客户生命周期的各个阶段，包括客户获取阶段的客户画像，客户保留阶段的客户细分、客户价值分析及客户流失分析等。在客户保留阶段，根据银行大量的客户基本属性数据、客户存款、贷款、金融产品使用等数据，利用聚类的方法，实现客户细分，将客户有效地划分为不同的类，从而针对每一类客户的特征设计出相应的产品组合、服务模式，以提高客户忠诚度。

证券业利用数据挖掘技术最集中的两个方面是客户管理和量化交易。证券公司可以利用客户个人基本信息、客户交易操作行为数据、软件使用习惯、自选股、常用分析指标等对客户的理财需求进行挖掘，实现精准营销。量化交易可借助数据挖掘方法，对证券期货市场的海量数据进行分析和挖掘，获得证券期货产品的价格变化规律，得到能带来超额收益的交易策略模型，然后通过分析结果来指导投资，以获得可持续、稳定且高于平均的超额回报。

保险业利用关联挖掘或各种推荐算法可以发现客户购买保险产品的关联与偏好，从而实现交叉销售。保险公司标的受损时，通过挖掘已有标的定损数据，可以对现有标的损失进行精确估计和预测，从而实现保险智能定损。随着保险业的发展，保险欺诈问题也日益突出，给保险公司和社会带来了极大危害。利用数据挖掘方法，分析并识别欺诈行为的特征，可以对保险欺诈行为进行实时监测与预警，从而促进保险业健康有序发展。

1.4.2 电信领域的数据挖掘

随着信息技术的迅速发展，电信业即将从 4G 时代进入 5G 时代，在电信业务迅速发

展的同时，电信行业的竞争也日益激烈。面对国内、国际电信业激烈的竞争态势，各大电信运营商纷纷使用数据挖掘技术了解行业动向、分析业务模式、洞察客户需求，实现精细化的管理和精准营销，提升自身服务质量，从而提高客户的满意度和忠诚度，增加竞争优势。

在客户关系管理方面，运营商使用数据挖掘可以对客户进行画像以提供个性化的业务推荐，可以对客户进行细分以发现不同价值的客户群体特征，可以通过客户流失分析制订相应的挽留策略，可以对客户之间的社会关系进行社交网络分析以获取潜在客户和保持现有客户，可以对客户流量使用进行异常识别，挖掘导致其流量异常的恶意程序和恶意 APP，以减少用户不必要的损失，并防止其他用户遭受同样的恶意攻击。在市场营销方面，可以使用关联挖掘进行电信业务的交叉销售。

运营商对网络信令数据进行挖掘，可以预测网络流量峰值，预警异常流量，防止网络堵塞和宕机，从而提高网络服务质量，提升用户体验。对移动用户的位置信息进行挖掘，与相关企业合作，可以提供基于位置的相关服务，如餐饮推荐、优惠券推送，这将改变运营商的盈利模式，而且具有非常广阔的应用前景。

1.4.3 零售与电子商务领域的数据挖掘

零售业的发展经历了从百货商店到超级市场、连锁商店、电子商务，再到如今线上线下相结合的“新零售”，积累了大量关于采购、销售、客户、物流等方面的数据。数据挖掘在零售与电子商务领域的应用非常广泛，如用户行为分析、个性化推荐、产品分析、广告追踪与优化、精准营销等。如顾客去商场购物，商场基于移动手机与 Wi-Fi 结合的数据，根据顾客所有的行动轨迹，分析顾客光顾的时间和频率、行径路线、驻留时间和地点，实现精准营销。

随着新零售业态的发展，线上线下系统对接和数据融合，零售企业借助数据挖掘技术可以对消费者全过程数据进行描述和产业链营销重构，实现数据化运营，探索新商业模式，建立新市场增长点。

1.4.4 政府政务领域的数据挖掘

政府信息化经过多年建设，已经有效实现了信息化办公。从 2015 年国家发布《促进大数据发展行动纲要》(国发〔2015〕50 号)开始，我国政府已将政务信息系统整合及共享提升到国家战略层面，对互联网+政务服务体系的建设和发展给出了明确指导意见和时间点要求。

国防、教育、公安、民政、司法、财政、交通运输、农业、商务、文化和旅游等政务部门信息系统的整合与共享，使数据挖掘的应用更加广泛。结合数据挖掘技术，政府加强统筹规划，实现智慧交通、智慧安防、智慧旅游等，加强智慧城市建设，使政务工作更高效、更开放、更透明。

1.4.5 医疗领域的数据挖掘

医疗领域积累了大量数据，尤其是海量的非格式化数据。数据挖掘在医疗领域的应用，主要集中在药品研发、疾病治疗、公共卫生管理、居民健康管理和健康影响因素分析等方面。

在药品研发方面，医药公司可以借助数据挖掘，在研发初期通过建模确定最有效率的投入产出比，配备最佳资源；在药物临床试验阶段，及时预测临床结果，选择最优药物。在疾病治疗方面，医生可以结合病人体征数据、费用数据和疗效数据进行挖掘，以确定在临床上对病人最有效和最具有成本效益的治疗方案。而且，对于医疗影像数据的分析和挖掘，会极大减轻医生的工作量，提高医疗效率。

在公共卫生管理、居民健康管理方面，卫生部门基于覆盖全国的电子病例数据进行挖掘，可以快速检测传染病，有效监测疫情，并提供有针对性的公众健康咨询，提高公众健康风险意识，降低传染病感染风险。

1.4.6 科学领域的数据挖掘

天文学、气象学、地质学、生物学等各科学领域使用全球定位系统、卫星遥感器及新一代生物学数据采集技术，收集了海量的包含时间和空间信息的高维数据、流数据和异构数据。早期，数据挖掘应用于天文学，在短短 4 个小时内发现的行星超过 20 多位天文学家 4 年的研究成果。

人类拥有 23 对染色体，约含有 30 亿对 DNA 碱基。1975 年，英国科学家 Frederick Sanger 发明了 Sanger 测序技术，由此开启了基因测序的新篇章。1990 年，由全球多个国家共同参与的人类基因组计划正式启动，被称为人类三大科学计划之一，旨在为这 30 亿对碱基构成的人类基因测序。数据挖掘技术应用于基因测序后，极大降低了测序成本，提升了测序速度。得益于此，从疾病的筛查、诊断到治疗，越来越多的临床基因检测项目落地，如新生儿疾病筛查、遗传病筛查、肿瘤易感基因筛查和肿瘤个性化用药等。

1.5 练习与拓展

1. 什么是数据挖掘？请结合实例加以说明。
2. 检索近几年数据挖掘国际学术会议的入选论文，分析数据挖掘研究现状及热点问题。
3. 查找相关资料，分析第四代数据挖掘系统的特点。
4. 什么是云计算？
5. 什么是大数据？大数据是否等于大数据分析？
6. 如何理解大数据被认为是下一个社会发展阶段的石油和金矿？
7. 什么是 Web 数据挖掘？

8. 什么是文本数据挖掘?
9. 分析说明 Fayyad 过程模型。
10. 分析说明 CRISP-DM 过程模型。
11. 结合 CRISP-DM 过程模型, 自选一个感兴趣的商业问题, 以小组为单位, 制订一份数据挖掘项目计划。
12. 数据挖掘的功能有哪些?
13. 结合数据挖掘使用技术, 分析其与相关学科之间的关系。
14. 什么是机器学习? 按学习方式不同, 机器学习可以分成哪几种? 分别具有什么特点?
15. 为什么说对于更大的数据库或是要满足更多在线处理任务的数据库, 直接从日常数据库中读取数据用于数据挖掘是不可行的?
16. 查阅相关资料, 说明什么是模式识别。
17. 数据挖掘可视化包含哪些方面?
18. 查阅相关资料, 说明什么是分布式计算, 什么是并行计算, 两者有什么关系。
19. 结合教材中提到的数据挖掘应用领域, 请举例说明, 并谈谈你的理解。
20. 除了教材中提到的数据挖掘应用领域, 请思考还有哪些应用领域, 并举例说明。