

普通高等院校数据科学与大数据技术专业“十三五”规划教材

数据仓库 与 数据挖掘

SHUJU
CANGKU

SHUJU
WAJUE

龙军 章成源 ◎ 编著



中南大学出版社
www.csupress.com.cn

17 18 19 22 23 24 25 26 27 28

普通高等院校数据科学与大数据技术专业“十三五”规划教材

数据仓库

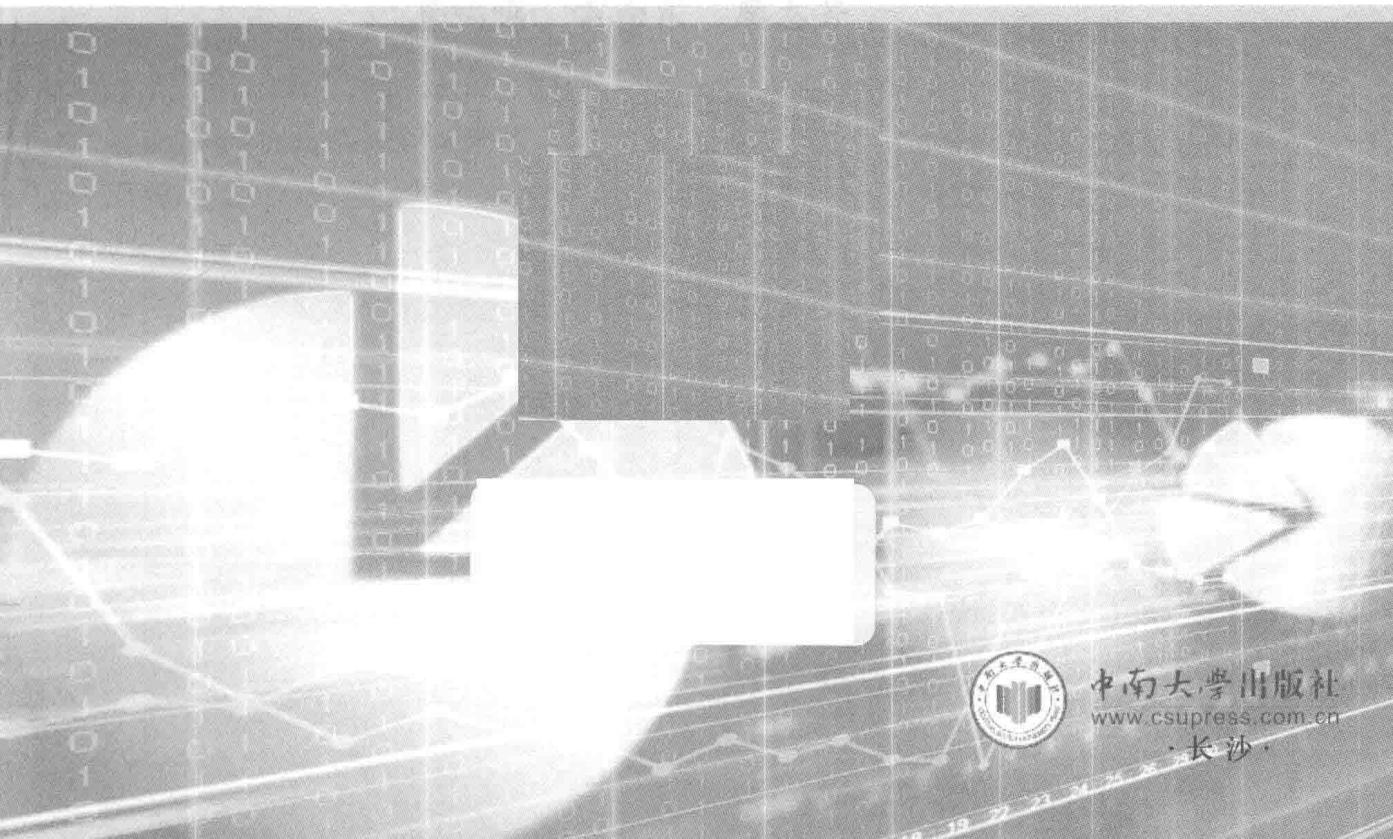
SHUJU
CANGKU

与

SHUJU
WAJUE

数据挖掘

龙军 章成源 ◎ 编著



中南大学出版社
www.csypress.com.cn
·长沙·

图书在版编目 (C I P) 数据

数据仓库与数据挖掘 / 龙军, 章成源编著. --长沙: 中南大学出版社, 2018. 12

ISBN 978 - 7 - 5487 - 3171 - 9

I . ①数… II . ①龙… ②章… III . ①数据库系统 ②数据采集
IV . ①TP311. 13 ②TP274

中国版本图书馆 CIP 数据核字(2018)第 071236 号

数据仓库与数据挖掘

龙军 章成源 编著

责任编辑 韩 雪

责任印制 易红卫

出版发行 中南大学出版社

社址: 长沙市麓山南路 邮编: 410083

发行科电话: 0731 - 88876770 传真: 0731 - 88710482

印 装 长沙德三印刷有限公司

开 本 787 × 1092 1/16 印张 17.25 字数 435 千字

版 次 2018 年 12 月第 1 版 印次 2018 年 12 月第 1 次印刷

书 号 ISBN 978 - 7 - 5487 - 3171 - 9

定 价 45.00 元

图书出现印装问题, 请与经销商调换

普通高等院校数据科学与大数据技术专业“十三五”规划教材

编委会

主任 桂卫华

副主任 邹北骥 吴湘华

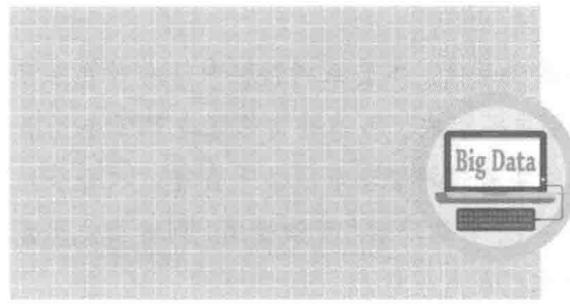
执行主编 郭克华 张祖平

委员 (按姓氏笔画排序)

龙军 刘丽敏 余腊生 周韵

高琰 桂劲松 高建良 章成源

鲁鸣鸣 雷向东 廖志芳



总序

Preface

随着移动互联网的兴起，全球数据呈爆炸性增长，目前 90% 以上的数据是近年产生的，数据规模大约每两年翻一番；而随着人工智能下物联网生态圈的形成，数据的采集、存储及分析处理、融合共享等技术需求都能得到响应，各行各业都在体验大数据带来的革命，“大数据时代”真正来临。这是一个产生大数据的时代，更是需要大数据力量的时代。

大数据具有体量巨大、速度极快、类型众多、价值巨大的特点，对数据从产生、分析到利用提出了前所未有的新要求。高等教育只有转变观念，更新方法与手段，寻求变革与突破，才能在大数据与人工智能的信息大潮面前立于不败之地。据预测，中国近年来大数据相关人才缺口达 200 万人，全世界相关人才缺口更超过 1000 万人之多。我国教育部门为了响应社会发展需要，率先于 2016 年开始正式开设“数据科学与大数据技术”本科专业及“大数据技术与应用”专科专业，近几年，全国形成了申报与建设大数据相关专业的热潮。随着专业建设的深入，大家发现一个共同的难题：没有成系列的大数据相关教材。

中南大学作为首批申报大数据专业的学校，2015 年在我校计算机科学与技术专业设立大数据方向时，信息科学与工程学院院领导便意识到系列教材缺失的严重问题，因此院领导规划由课程团队在教学的同时积累素材，形成面向大数据专业知识体系与能力体系、老师自己愿意用、同学觉得买得值、关联性强的系列教材。经过两年的准备，针对 2017 年《教育部办公厅关于推荐新工科研究与实践项目的通知》的精神，中南大学出版社组织对系列教材文稿进行相应的打磨，最终于 2018 年底出版“高等院校数据科学与大数据技术专业‘十三五’规划教材”。

该套系列教材具有如下特点：

1. 本套教材主要参照“数据科学与大数据技术”本科专业的培养方案，综合考虑专业的来源，如从计算机类专业、数学统计类专业以及经济类专业发展而来；同时适当兼顾了专科类偏向实际应用的特点。

2. 注重理论联系实际，注重能力培养。该系列教材中既有理论教材也有配套的实践教程。力图通过理论或原理教学、案例教学、课堂讨论、课程实验与实训实习等多个环节，训练学生掌握知识、运用知识分析并解决实际问题的能力，以满足学生今后就业或科研的需求；同时兼顾“全国工程教育专业认证”对学生基本能力的培养要求与复杂问题求解能力的

要求。

3. 在规范教材编写体例的同时，注重写作风格的灵活性。本套系列教材中每本书的内容都由教学目的、本章小结、思考题或练习题、实验要求等组成。每本教材都配有 PPT 电子教案及相关的电子资源，如实验要求及 DEMO、配套的实验资源管理与服务平台等。本套系列教材的文本层次分明、逻辑性强、概念清晰、图文并茂、表达准确、可读性强，同时相关配套电子资源与教材的相关性强，形成了新媒体式的立体型系列教材。

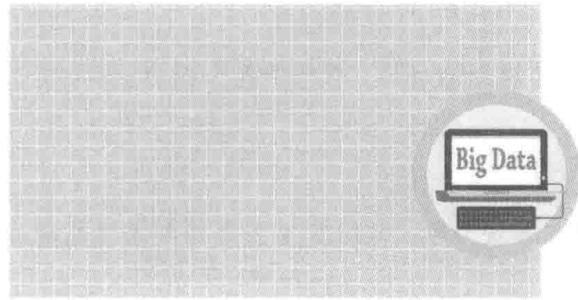
4. 响应了教育部“新工科”研究与实践项目的要求。本套教材从专业导论课开始设立相关的实验环节，作为知识主线与技术主线把相关课程串接起来，力争让学生尽早具有培养自己动手能力的意识、综合利用各种技术与平台的能力。同时为了避免新技术发展太快、教材纸质文字内容容易过时的问题，在相关技术及平台的叙述与实践中，融合了网络电子资源容易更新的特点，使新技术保持时效性。

5. 本套丛书配有丰富的多媒体教学资源，将扩展知识、习题解析思路等内容做成二维码放在书中，丰富了教材内容，增强了教学互动，增加了学生的学习积极性与主动性。

本套丛书吸纳了数据科学与大数据技术教育工作者多年教学与科研成果，凝聚了作者们的辛勤劳动，同时也得到了中南大学等院校领导和专家的大力支持。我相信本套教材的出版，对我国数据科学与大数据技术专业本科、专科教学质量的提高将有很好的促进作用。

桂卫华

2018 年 11 月



前言

Foreword

随着计算机、网络和通信等信息技术的发展，数据采集的方法越来越丰富，人类收集、存储和访问数据的能力大大增强。存储设备的容量不断提升而成本逐年下降，特别是数据库技术在各行各业的普及应用，使人类积累了海量的数据。快速增长的海量数据集已经远远超出了人类的理解能力，人类步入了大数据时代却陷入了“数据丰富、知识贫乏”的困境。人们迫切希望从所拥有的数据中获取有用的知识，以帮助其更好地进行决策。针对这一问题，数据仓库和数据挖掘技术应运而生，并且显示出强大的生命力。要将海量数据转换为有用的信息和知识，首先要有效地收集和组织数据。数据仓库是良好的数据收集和组织工具，它的任务是搜集来自各个业务系统的有用数据，存放在一个集成的储存区内。在数据仓库丰富完整的数据基础上，数据挖掘技术可以从中挖掘出有价值的知识，从而帮助决策者做出正确决策。“数据仓库与数据挖掘”已成为普通高等院校计算机、经贸管理和信息类相关专业研究生和高年级本科生的学位课程或选修科目。

按照教育部关于高等学校本科教育以培养更多应用型人才为目标的教学改革方向，以及全日制研究生以学术型和专业型两大类进行有差别培养的要求，我们迫切需要一本在教学课时限制严格的条件下，理论叙述深入浅出、实际应用具体完整、算法描述自然易懂、计算实例详略得当的数据仓库与数据挖掘方面的教材。

本书正是在这种社会需求背景和实际教学需要的情况下，在笔者总结多年教学改革与实践经验所编写的讲义基础上修改而成的。本书兼顾了应用型人才与学术型人才培养的需求，介绍了数据仓库原理、数据仓库设计和实现方法，为读者真正架起了理论与实践的桥梁。本书还以大量的计算实例来增加读者对数据挖掘原理及各种挖掘算法的理解深度。

本书主要介绍数据仓库和数据挖掘技术的基本原理和应用方法。全书共分为 13 章，主要内容包括数据仓库的概念与体系结构、数据、数据存储、OLAP 与数据立方体、数据挖掘基础、关联挖掘、聚类分析、分类、神经网络、统计分析、非结构化数据挖掘、知识图谱、大数据挖掘算法。其中，第 1 章为数据仓库的概念与体系结构，内容包括数据挖掘的兴起、数据仓库的基本概念及数据仓库的特点与组成等。第 2~4 章主要介绍数据仓库的基本原理和数据仓库系统的组建方法。第 5~12 章介绍当前流行的数据挖掘算法的主要思想和理论基础，并且给出丰富的应用实例。第 13 章为数据挖掘创新篇，内容主要为基于 MapReduce 的数据

挖掘算法及 Hadoop 的介绍等。

本书紧跟数据仓库和数据挖掘技术的发展和人才培养的目标，有以下几个特点：

(1) 可读性强，文字叙述深入浅出，易读易用。

(2) 概念清晰，条理清楚，内容取舍合理。

(3) 本书强调基础，重视实例，各章节都以经典算法为主，介绍其主要思想和基本原理，并且给出恰当和丰富的实例。

(4) 书中实例和课后习题实用、丰富，通过练习，读者可以对各个知识点从不同角度得到训练，掌握和巩固所学知识。

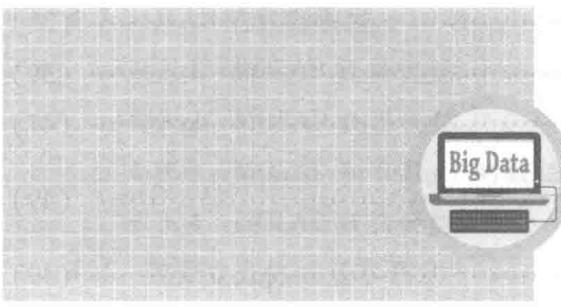
(5) 本书教学资源丰富，提供了多媒体教学课件及实验平台，方便教学。

本书各章节之间衔接自然，同时各章节又有一定的独立性，读者可按本书的自然顺序学习，也可以根据实际情况挑选需要的章节学习。

限于作者水平，加之数据仓库与数据挖掘理论技术的内容十分丰富，且发展非常迅速，疏漏和不当之处在所难免，殷切希望广大师生和读者批评指正。

编 者

2018 年 10 月



目录

Contents

第1章 数据仓库的概念与体系结构	(1)
1.1 数据仓库的兴起	(1)
1.1.1 数据管理技术的发展	(1)
1.1.2 数据仓库的萌芽	(3)
1.2 数据仓库的基本概念	(4)
1.2.1 元数据	(4)
1.2.2 数据粒度	(5)
1.2.3 数据模型	(5)
1.2.4 ETL	(6)
1.2.5 数据集市	(7)
1.3 数据仓库的特点与组成	(8)
1.3.1 数据仓库的特点	(8)
1.3.2 数据仓库的组成	(11)
1.4 数据仓库的体系结构	(15)
1.4.1 传统的数据仓库体系结构	(15)
1.4.2 传统数据仓库系统在大数据时代所面临的挑战	(16)
1.4.3 大数据时代的数据仓库	(20)
习 题	(23)
第2章 数 据	(24)
2.1 数据的概念与内容	(24)
2.2 数据属性与数据集	(28)
2.3 数据预处理	(29)
2.3.1 数据预处理概述	(30)
2.3.2 数据清洗	(31)
2.3.3 数据集成	(35)

2.3.4 数据变换	(38)
2.3.5 数据归约	(39)
习 题	(47)

第3章 数据存储 (49)

3.1 数据仓库的数据模型	(49)
3.1.1 数据仓库的概念模型	(50)
3.1.2 数据仓库的逻辑模型	(52)
3.1.3 数据仓库的物理模型	(54)
3.2 元数据存储	(55)
3.2.1 元数据的概念	(55)
3.2.2 元数据的分类方法	(55)
3.2.3 元数据的管理	(57)
3.2.4 元数据的作用	(58)
3.3 数据集市	(59)
3.3.1 数据集市的概念	(59)
3.3.2 数据集市的类型	(60)
3.3.3 数据集市的建立	(60)
3.4 大数据存储技术	(61)
3.4.1 大数据的概念	(61)
3.4.2 传统数据库的局限	(62)
3.4.3 NoSQL 数据库	(63)
3.4.4 几种主流的 NoSQL 数据库	(64)
习 题	(64)

第4章 OLAP 与数据立方体 (65)

4.1 OLAP 的概念	(65)
4.1.1 OLAP 的定义	(65)
4.1.2 OLAP 的准则	(66)
4.1.3 OLAP 的特征	(69)
4.2 多维分析的基本分析动作	(70)
4.2.1 切片	(70)
4.2.2 切块	(71)
4.2.3 钻取	(72)
4.2.4 旋转	(72)

4.3 OLAP 的数据模型	(73)
4.3.1 ROLAP 数据模型	(73)
4.3.2 MOLAP 数据模型	(75)
4.3.3 MOLAP 和 ROLAP 的数据组织与应用比较	(76)
4.3.4 HOLAP 数据模型	(77)
4.4 数据立方体的基本概念	(78)
4.4.1 数据立方体中的一些概念	(78)
4.4.2 数据立方体计算的一般策略	(79)
4.5 数据立方体的计算方法	(80)
4.5.1 多路数组策略计算完全立方体	(80)
4.5.2 从顶点方体向下计算冰山立方体	(80)
4.5.3 使用动态星树结构计算冰山立方体	(81)
4.5.4 快速高维 OLAP 预计算壳片段	(82)
习 题	(83)
第 5 章 数据挖掘基础	(84)
5.1 数据挖掘的兴起	(84)
5.1.1 数据挖掘的发展历程	(84)
5.1.2 数据挖掘的概述	(85)
5.1.3 大规模数据挖掘	(86)
5.2 数据挖掘的任务	(87)
5.2.1 关联规则	(87)
5.2.2 聚类分析	(88)
5.2.3 分类分析	(89)
5.2.4 回归分析	(90)
5.2.5 相关分析	(91)
5.2.6 异常检测	(92)
5.3 数据挖掘的流程	(92)
5.3.1 数据挖掘对象	(92)
5.3.2 数据挖掘分类	(93)
5.3.3 知识发现的过程	(94)
习 题	(96)
第 6 章 关联挖掘	(97)
6.1 关联规则的概念和分类	(97)

6.1.1	关联规则的概念	(97)
6.1.2	关联规则的分类	(99)
6.2	Apriori 算法	(100)
6.2.1	Apriori 算法概述	(100)
6.2.2	Apriori 算法的性质与步骤	(100)
6.2.3	Apriori 算法的实例	(101)
6.2.4	从频繁项集产生关联规则	(103)
6.3	FP - Growth 算法	(104)
6.3.1	FP - tree 的建立	(105)
6.3.2	FP - tree 上挖掘关联规则	(106)
6.4	挖掘算法的进阶算法	(107)
习 题		(110)

第7章 聚类分析 (112)

7.1	聚类分析概述	(112)
7.1.1	聚类分析的定义	(112)
7.1.2	聚类分析的分类	(113)
7.2	差异度的计算方法	(114)
7.2.1	聚类算法中的数据结构	(114)
7.2.2	区间标度变量的差异度计算	(115)
7.2.3	二元变量的差异度计算	(116)
7.2.4	标称型变量的差异度计算	(117)
7.2.5	序数型变量的差异度计算	(118)
7.2.6	比例标度型变量的差异度计算	(119)
7.2.7	混合类型变量的差异度计算	(119)
7.3	基于分割的聚类方法	(120)
7.3.1	分割聚类方法的描述	(120)
7.3.2	K - means 均值算法	(121)
7.3.3	PAM 算法	(122)
7.3.4	CLARA 算法和 CLARANS 算法	(125)
7.4	基于密度的聚类方法	(126)
7.4.1	基于密度的聚类方法描述	(126)
7.4.2	DBSCAN 算法	(127)
7.4.3	OPTICS 算法	(129)
7.5	谱聚类方法	(130)



7.5.1 谱聚类描述	(130)
7.5.2 谱聚类算法描述	(131)
7.5.3 谱聚类实例	(132)
7.6 ICA 聚类分析	(133)
7.6.1 ICA 的起源和目的	(133)
7.6.2 ICA 模型和应用要求	(133)
7.6.3 ICA 应用场合	(135)
习 题	(135)
第8章 分 类	(137)
8.1 分类的基本知识	(137)
8.1.1 分类的概念	(137)
8.1.2 分类的评价标准	(138)
8.1.3 分类的主要方法	(138)
8.2 决策树分类	(139)
8.2.1 决策树算法概述	(139)
8.2.2 决策树的生成	(141)
8.2.3 决策树中规则的提取	(142)
8.2.4 ID3 算法	(143)
8.2.5 C4.5 算法	(145)
8.2.6 蒙特卡洛树搜索(MCTS)算法	(146)
8.3 SVM 预测	(147)
8.3.1 线性可分的 SVM	(147)
8.3.2 线性不可分的 SVM	(150)
8.3.3 SVM 的实现——手写数字图片的识别	(153)
8.4 KNN 算法	(154)
8.4.1 KNN 算法的描述	(155)
8.4.2 KNN 算法的实现	(156)
习 题	(157)
第9章 神经网络	(159)
9.1 神经网络概述与定义	(159)
9.1.1 神经网络概述	(159)
9.1.2 神经网络的学习过程	(160)
9.2 限制玻尔兹曼机(RBM)	(161)

9.2.1 RBM 的定义	(161)
9.2.2 RBM 的能量模型与学习方法	(162)
9.3 深度信念网络	(165)
9.3.1 DBN 反向传播算法介绍与改进	(165)
9.3.2 DNN 分类与代价函数选择	(170)
9.4 卷积神经网络(CNN)	(173)
9.4.1 卷积神经网络定义与结构	(173)
9.4.2 CNN 两个特点与图形实例	(176)
9.5 循环神经网络(RNN)	(179)
9.5.1 RNN 概述	(180)
9.5.2 RNN 训练	(181)
9.5.3 LSTMs 网络与函数展示图例	(182)
习 题	(186)
第 10 章 统计分析	(188)
10.1 回归分析	(188)
10.1.1 一元线性回归	(188)
10.1.2 多元线性回归	(191)
10.1.3 非线性回归	(193)
10.2 EM 算法	(194)
10.2.1 EM 算法的引入	(194)
10.2.2 EM 算法的推导	(196)
10.2.3 EM 算法的收敛性	(197)
10.3 Bayes 分类	(199)
10.3.1 Bayes 定理	(199)
10.3.2 简单 Bayes 分类	(200)
10.3.3 Bayes 信念网络	(201)
10.3.4 Bayes 网络的应用	(203)
习 题	(203)
第 11 章 非结构化数据挖掘	(204)
11.1 文本数据挖掘	(204)
11.1.1 文本数据挖掘的概念	(204)
11.1.2 文本数据挖掘技术	(208)
11.2 Web 数据挖掘	(214)



11.2.1 Web 数据挖掘的概念	(215)
11.2.2 Web 数据挖掘的分类	(216)
11.2.3 Web 数据挖掘的应用	(220)
11.3 多媒体数据挖掘	(221)
11.3.1 多媒体数据挖掘的概念	(222)
11.3.2 多媒体数据挖掘的分类	(223)
习 题	(225)
第 12 章 知识图谱	(227)
12.1 知识图谱构建	(227)
12.1.1 知识图谱的概述	(227)
12.1.2 知识图谱的数据来源	(229)
12.1.3 多源异构数据的融合	(231)
12.1.4 知识图谱的表示	(232)
12.2 知识图谱技术	(233)
12.2.1 实体抽取	(234)
12.2.2 关系抽取	(235)
12.2.3 知识推理	(236)
12.3 知识图谱的典型应用	(238)
12.3.1 查询理解	(238)
12.3.2 自动问答	(240)
12.3.3 前景和挑战	(240)
习 题	(241)
第 13 章 大数据挖掘算法	(242)
13.1 Hadoop 介绍	(242)
13.1.1 Hadoop 的基本概念	(242)
13.1.2 Hadoop 的基本组件	(244)
13.2 基于 MapReduce 数据挖掘算法	(247)
13.2.1 基于 MapReduce 的 K-means 并行算法	(248)
13.2.2 基于 MapReduce 的分类算法	(251)
13.2.3 基于 MapReduce 的序列模式挖掘算法	(253)
习 题	(255)
参考文献	(256)



第1章 数据仓库的概念与体系结构

任何企业都希望在市场竞争中通过利用全面的数据分析能力来获取更大更持久的竞争优势。例如，银行希望知道如何有效规避信贷风险，发现欺诈和洗钱等不合法行为；电信公司希望知道如何对市场业务发展和竞争环境进行精准分析，从而为市场决策提供深入有力的分析支撑，提升营销活动的精准性；保险公司希望知道哪些理赔客户骗保的可能性更高，以及哪些客户是高价值低风险的客户群；等等。以上这些问题的解决都离不开数据的支撑和对现有数据的分析利用。

传统的信息资源管理主要依靠数据库技术，利用数据库技术进行数据的组织与存储，并使用基于数据库的信息系统进行信息资源的有效利用。但是随着计算机技术的飞速发展，一些新的需求不断提出，这些新的需求是传统数据库技术难以满足的。

传统的数据库技术以数据库为中心进行事务处理、批处理到决策分析等各种类型的数据处理工作。数据库系统作为数据管理手段，从它诞生开始，就主要用于事务处理。近年来，随着计算机应用的拓展，人们对计算机数据处理的能力提出了更高的要求，希望计算机能够更多地参与到数据分析与决策支持中。但是事务处理和分析处理有着不同的性质，直接使用事务处理环境来支持分析决策有着一定的局限性。因此，需要一种新的数据管理技术——数据仓库，来实现分析决策的目的。

本章主要介绍数据仓库的兴起(1.1节)、数据仓库的基本概念(1.2节)、数据仓库的特点与组成(1.3节)以及数据仓库的体系结构(1.4节)。

1.1 数据仓库的兴起

任何一项新技术的出现都有其深厚的背景，数据仓库技术作为一种新型的数据管理技术，它的出现也并非偶然，而是随着人们对数据处理技术新需求的不断提出，最终衍生出了数据仓库。本小节主要介绍数据仓库的兴起，并从数据管理技术的发展(1.1.1节)、数据仓库的萌芽(1.1.2节)来展开讲述。

1.1.1 数据管理技术的发展

如图1-1所示，从1946年计算机的诞生到20世纪80年代，数据管理技术经历了人工管理、文件系统以及数据库系统三个阶段，三个阶段之间彼此联系。由于人工管理的不足导致了文件系统的出现，同样由于文件系统的不足导致了数据库系统的出现。

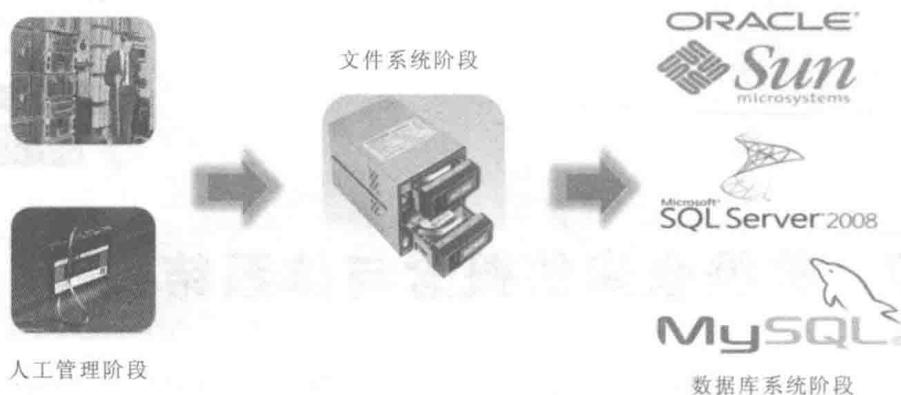


图 1-1 数据管理技术发展的三个阶段

1. 人工管理阶段

在 20 世纪 60 年代初期，创建运行于主文件上的单个应用是计算领域的主要工作。这些应用的特点表现在报表和程序上，常用的语言是 COBOL。穿孔卡是当时常用的介质，主文件存放在磁带文件上。磁带由于其廉价的特性适合于存放大量数据，但其缺点是需要顺序地访问。访问整条磁带的文件可能要花去 20~30 分钟时间，这取决于我们所需数据在磁带上的存放位置。事实上，在最糟糕的情况下，一次磁带文件的操作需要访问 100% 的记录，然而只有 1% 的记录是真正需要的。

大约在 20 世纪 60 年代中期，主文件和磁带的使用数量迅速增长，很快，处处都是主文件。随着主文件数量的增长，数据出现了大量的冗余，并由此引出了一些严重的问题：

- (1) 数据更新需要保持数据的一致性。
- (2) 程序维护的复杂性。
- (3) 开发新程序的复杂性。
- (4) 支持所有主文件所需要的硬件数量。

简言之，由于介质本身存在固有缺陷，使得主文件成为了发展的障碍。如果仍然只用磁带作为数据存储的唯一介质，很难想象现在的数据管理会是什么样子。如果除了磁带文件以外没有别的介质可以存储大量数据，那么现实生活中将永远不会有大型、快速的系统。

2. 文件系统阶段

20 世纪 60 年代中期，计算机硬件有了磁盘直接存储设备，软件也有了操作系统，于是人们将数据文件长期存储到磁盘直接存储设备上，并利用操作系统所提供的文件系统对数据进行快速的访问。磁盘存储从根本上不同于磁带存储，磁盘存储不需要经过第 1 条记录，第 2 条记录，……，第 n 条记录，才能得到第 $n+1$ 条记录。只需要知道第 $n+1$ 条记录的地址，就可以轻而易举地直接访问它。磁盘直接存储设备以及操作系统的出现，解决了磁带存储设备数据访问效率低下的问题。但是它并没有解决主文件技术所带来的问题，只是更改存储介质提升了数据访问的速度，亟需一种新的数据存储与访问技术来解决主文件技术所面临的问题。