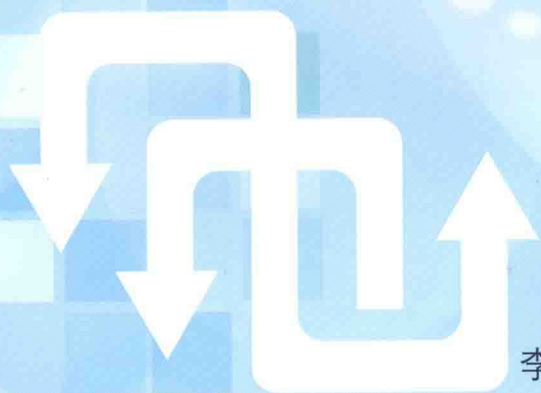


21世纪高等学校电子商务专业规划教材



李春葆 蒋林 陈良臣 喻丹丹 曾平 编著

数据仓库与数据挖掘 应用教程



清华大学出版社

21世纪高等学校电子商务专业规划教材

李春葆 蒋林 陈良臣 喻丹丹 曾平 编著

数据仓库与数据挖掘 应用教程

清华大学出版社
北京

内 容 简 介

本书以 SQL Server 分析服务为环境介绍数据仓库和数据挖掘应用技术,包括数据仓库和数据挖掘概述、OLAP 和多维数据模型、数据仓库设计和 SQL Server 数据仓库开发实例、关联分析算法、决策树分类算法、贝叶斯分类算法、神经网络算法、回归分析算法、时间序列分析和聚类算法。

本书内容翔实,循序渐进地介绍各个知识点,并提供全面而丰富的教学资源,可作为各类高等院校计算机及相关专业“数据仓库和数据挖掘应用技术”和“SQL Server 高级应用”课程的教学用书,也适合计算机应用人员和计算机爱好者参考。

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

数据仓库与数据挖掘应用教程/李春葆等编著.--北京:清华大学出版社,2016
21世纪高等学校电子商务专业规划教材
ISBN 978-7-302-43077-3

I. ①数… II. ①李… III. ①数据库系统—教材 ②数据采集—教材 IV. ①TP311.13 ②TP274
中国版本图书馆 CIP 数据核字(2016)第 034105 号

责任编辑:魏江江 薛 阳

封面设计:常雪影

责任校对:时翠兰

责任印制:沈 露

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

课 件 下 载: <http://www.tup.com.cn>, 010-62795954

印 刷 者:北京季蜂印刷有限公司

装 订 者:三河市溧源装订厂

经 销:全国新华书店

开 本:185mm×260mm

印 张:20

字 数:509千字

版 次:2016年10月第1版

印 次:2016年10月第1次印刷

印 数:1~2000

定 价:39.50元

产品编号:067892-01

出版说明

电子商务是以信息技术为手段,以商品交换为中心的商务活动,是“互联网+”的杰作之一。特别是在 2015 年年初的政府工作报告中,李克强总理首次提出“制定‘互联网+’行动计划”,极大地推进了电子商务在我国的蓬勃发展,改造和影响众多传统行业。电子商务系统是保证以电子商务为基础的网上交易实现的体系。

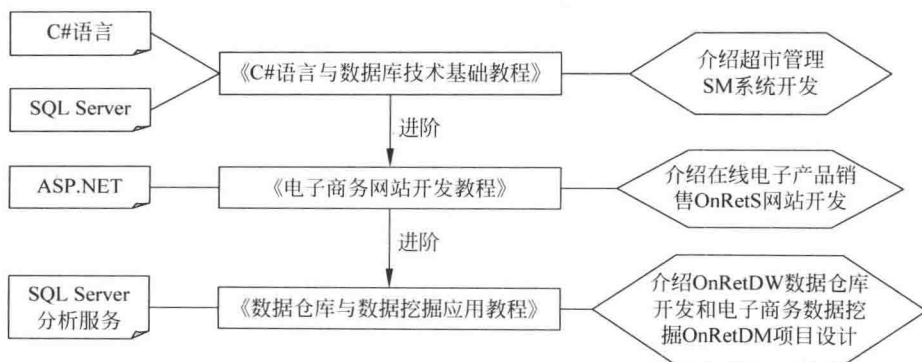
电子商务应用的快速发展,需要大量的专业技术人员,据专家测算,未来 10 年我国电子商务人才缺口将达到 200 万。为了加快电子商务系统人才培养,我们以国家卓越工程师计划为契机规划并出版本系列教材。

电子商务网站是电子商务系统的核心,电子商务网站开发涉及多方面的技术。本系列教材以 Visual Studio 为开发环境,以案例为向导,全面介绍电子商务网站的开发技术,涵盖的教程如下:

- C# 语言与数据库技术基础教程
- 电子商务网站开发教程
- 数据仓库与数据挖掘应用教程

本系列教材具有专业培养定位清晰、可操作性强的特点。《C# 语言与数据库技术基础》从零基础开始,循序渐进地介绍 C# 语言的基本语法、面向对象编程、Windows 窗体应用程序设计、SQL Server 数据库操作、C# 访问数据库方法以及 Windows 界面的电子商务系统开发技术。《电子商务网站开发教程》以 ASP.NET 为背景,介绍动态网站的开发技术。《数据仓库与数据挖掘应用教程》介绍数据仓库数据和电子商务数据分析技术。

教程中涉及的相关案例如下。



本系列教材是武汉大学计算机学院和解放军理工大学在探索电子商务人才培养并结合国家卓越工程师计划的教学实践,总结和提炼的教学成果。教学改革是教育工作永恒不变的主题,也是需要不断探索的课题,需要不断地努力实践和完善。本系列教材虽然经过细致的编写与校订,仍然难免有疏漏和不足之处,需要不断地补充、修订和完善,我们热情欢迎使用本系列教材的教师、学生和读者朋友提出宝贵意见和建议,使之更臻成熟。

“21 世纪高等学校电子商务专业规划教材”编委会

前 言

数据仓库是企业决策支持系统和联机分析处理(OLAP)的结构化数据环境,具有面向主题、集成性、稳定性和随时间变化的(时变性)的特征。数据挖掘(Data Mining)是从大量的、有噪声的、不完全的数据中提取隐含的、人们事先未知的有用知识和信息的过程。数据仓库和数据挖掘是电子商务数据分析的有效手段。本书讨论数据仓库和数据挖掘应用的相关技术,其内容组织如下。

第1章为数据仓库和数据挖掘概述,介绍数据仓库的特征、数据仓库系统及开发工具、商业智能和数据仓库的关系、数据挖掘的定义和数据挖掘过程。

第2章为 OLAP 和多维数据模型,介绍 OLAP 定义和特性、多维数据模型和数据仓库的维度建模。

第3章为数据仓库设计,介绍数据仓库规划与需求分析、数据仓库建模、数据仓库物理模型设计和数据仓库部署与维护。

第4章为 SQL Server 数据仓库开发实例,介绍一个基于在线电子产品销售数据的 OnRetDW 数据仓库的设计过程,包括需求分析、建模、数据抽取工具设计等。

第5章为关联分析算法,介绍关联分析的相关概念、Apriori 算法、SQL Server 挖掘关联规则方法和电子商务数据的关联规则挖掘过程。

第6章为决策树分类算法,介绍基本分类步骤、决策树分类、SQL Server 决策树分类方法和电子商务数据的决策树分类过程。

第7章为贝叶斯分类算法,介绍贝叶斯公式、朴素贝叶斯分类原理、SQL Server 朴素贝叶斯分类方法和电子商务数据的贝叶斯分类过程。

第8章为神经网络算法,介绍人工神经网络相关概念、用于分类的前馈神经网络、SQL Server 神经网络分类方法和电子商务数据的神经网络分类过程。

第9章为回归分析算法,介绍回归分析相关概念、线性回归分析、非线性回归分析、逻辑回归分析方法和电子商务数据的逻辑回归分析过程。

第10章为时间序列分析,介绍时间序列分析相关概念、确定性时间序列分析、随机时间序列模型、SQL Server 时间序列分析方法和电子商务数据的时间序列分析过程。

第11章为聚类算法,介绍聚类相关概念、 k -均值算法及其应用、EM 算法及其应用、电子商务数据的聚类分析过程以及 Microsoft 顺序分析和聚类分析算法。

书中提供了大量的练习题和上机实验题供读者选用,附录 A 给出了部分练习题参考答案,附录 B 给出了所有上机实验题参考答案,附录 C 给出了书中数据库和包含的数据表。其中带“*”部分为选修内容。

本书紧扣数据仓库和数据挖掘开发所需要的知识、技能和素质要求,以技术应用能力培养为主线构建教材内容,具有以下特色:

- ☑ 内容全面、知识点翔实:在内容讲授上力求翔实和全面,细致解析每个知识点和各知识点的联系。
- ☑ 条理清晰、讲解透彻:从介绍数据仓库和数据挖掘的基本概念出发,由简单到复杂,循序渐进介绍数据仓库和数据挖掘系统的开发过程。
- ☑ 精选实例、实用性强:列举了大量的应用示例,读者通过上机模仿可以大大提高使用应用系统开发能力。
- ☑ 配套教学资源丰富:提供了教学 PPT、书中所有示例代码、相关数据库文件和 ETL 源程序,便于读者打开和调试。配套的教学资源可以从清华大学出版社网站下载。

本教材的编写工作得到武汉大学教务部教改项目的资助,解放军理工大学和清华大学出版社给予了大力支持,连续多届选课的同学提出了许多宝贵的建议,编者在此表示衷心感谢。

编 者

2016 年 4 月

目 录

第 1 章 数据仓库和数据挖掘概述	1
1.1 数据仓库概述	2
1.1.1 数据仓库的定义.....	2
1.1.2 数据仓库与操作型数据库的关系.....	4
1.1.3 数据仓库的应用.....	6
1.2 数据仓库系统及开发工具	7
1.2.1 数据仓库系统的组成.....	7
1.2.2 数据仓库系统开发工具	10
1.3 商业智能和数据仓库.....	12
1.3.1 什么是商业智能	12
1.3.2 商业智能和数据仓库的关系	13
1.4 数据挖掘概述.....	14
1.4.1 数据挖掘的定义	14
1.4.2 数据挖掘的主要任务	15
1.4.3 数据挖掘的对象	15
1.4.4 数据挖掘的知识表示	16
1.4.5 数据挖掘与数据仓库及 OLAP 的关系	17
1.4.6 数据挖掘的应用	17
1.5 数据挖掘过程.....	19
1.5.1 数据挖掘步骤	19
1.5.2 数据清理	19
1.5.3 数据集成	21
1.5.4 数据变换	22
1.5.5 数据归约	23
1.5.6 离散化和概念分层生成	23
1.5.7 数据挖掘的算法	25
练习题	27
第 2 章 OLAP 和多维数据模型	29
2.1 OLAP 概述	30
2.1.1 什么是 OLAP	30

2.1.2	OLAP 和 OLTP 的区别	30
2.1.3	数据仓库与 OLAP 的关系	31
2.2	多维数据模型	31
2.2.1	多维数据模型的相关概念	32
2.2.2	OLAP 的基本分析操作	34
2.2.3	多维数据模型的实现途径	38
2.3	数据仓库的维度建模	40
2.3.1	数据仓库建模概述	40
2.3.2	星形模型	40
2.3.3	雪花模型	41
2.3.4	事实星座模型	43
	练习题	44
第 3 章	数据仓库设计	46
3.1	数据仓库设计概述	47
3.1.1	数据仓库设计原则	47
3.1.2	建立数据仓库系统的两种模式	47
3.1.3	数据仓库设计过程	48
3.2	数据仓库规划与需求分析	48
3.2.1	数据仓库规划	49
3.2.2	数据仓库需求分析	49
3.3	数据仓库建模	50
3.3.1	数据仓库建模的主要工作	50
3.3.2	维表设计	53
3.3.3	事实表设计	54
3.4	数据仓库物理模型设计	55
3.4.1	确定数据的存储结构	56
3.4.2	确定索引策略	56
3.4.3	确定存储分配	57
3.5	数据仓库部署与维护	57
3.5.1	数据仓库的部署	57
3.5.2	数据仓库的维护	58
	练习题	58
第 4 章	SQL Server 数据仓库开发实例	60
4.1	OnRetDW 系统需求分析	61
4.1.1	OnRetDW 系统的主题	61
4.1.2	OnRetDW 系统的功能	62
4.2	OnRetDW 的建模	62
4.2.1	维表设计	62
4.2.2	事实表设计	66
4.3	数据抽取工具设计	67

4.4	基于 SQL Server 2012 设计 OnRetDW	75
4.4.1	创建数据仓库 OnRetDW 项目	75
4.4.2	创建数据源	77
4.4.3	创建数据源视图	78
4.4.4	创建维表	80
4.4.5	创建多维数据集	84
4.4.6	部署 SDWS	85
4.4.7	浏览已部署的多维数据集	85
4.5	MDX 简介*	90
4.5.1	MDX 语言概述	90
4.5.2	执行 MDX 查询	91
4.5.3	多维数据查询	92
	练习题	95
	上机实验题	96
第 5 章	关联分析算法	97
5.1	关联分析概述	98
5.1.1	什么是关联分析	98
5.1.2	事务数据库	98
5.1.3	关联规则及其度量	99
5.1.4	频繁项集	101
5.1.5	挖掘关联规则的基本过程	101
5.2	Apriori 算法	102
5.2.1	Apriori 性质	102
5.2.2	Apriori 算法求频繁项集	103
5.2.3	由频繁项集产生强关联规则	108
5.3	SQL Server 挖掘关联规则	109
5.3.1	创建 DMK 数据库	109
5.3.2	建立关联挖掘项目	110
5.3.3	部署关联挖掘项目并浏览结果	116
5.4	电子商务数据的关联规则挖掘	119
5.4.1	创建 OnRetDMK 数据库	119
5.4.2	数据加载功能设计	120
5.4.3	建立关联挖掘项目	121
5.4.4	部署关联挖掘项目并浏览结果	121
	练习题	124
	上机实验题	126
第 6 章	决策树分类算法	127
6.1	分类过程	128
6.1.1	分类概述	128
6.1.2	分类过程的学习阶段	128

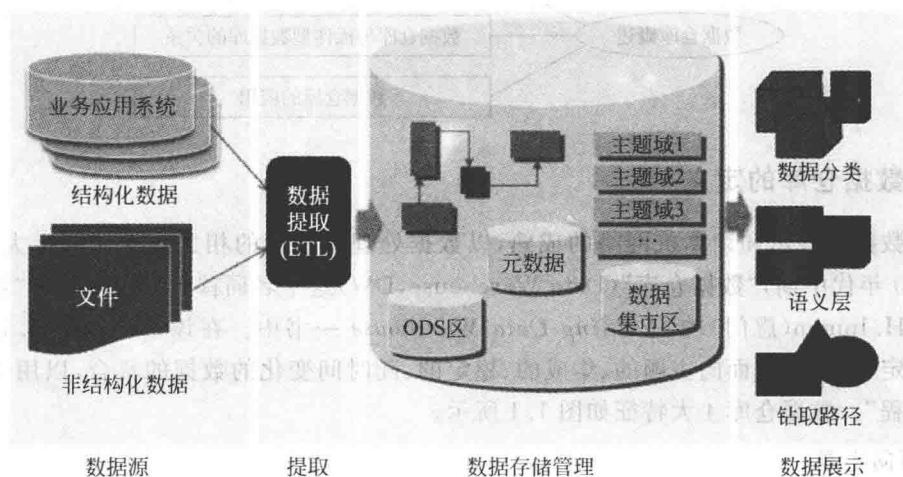
6.1.3 分类过程的分类阶段	130
6.2 决策树分类	130
6.2.1 决策树	130
6.2.2 建立决策树的 ID3 算法	131
6.3 SQL Server 决策树分类	139
6.3.1 建立数据表	139
6.3.2 建立决策树分类挖掘模型	140
6.3.3 浏览决策树模型和分类预测	143
6.4 电子商务数据的决策树分类	146
6.4.1 创建 OnRetDMK.DST 数据表	146
6.4.2 数据加载功能设计	146
6.4.3 建立决策树分类模型	148
6.4.4 浏览决策树	150
练习题	150
上机实验题	151
第 7 章 贝叶斯分类算法	153
7.1 贝叶斯分类概述	154
7.1.1 贝叶斯定理	154
7.1.2 贝叶斯信念网络	155
7.2 朴素贝叶斯分类	156
7.2.1 朴素贝叶斯分类原理	157
7.2.2 朴素贝叶斯分类算法	159
7.3 SQL Server 朴素贝叶斯分类	161
7.3.1 建立朴素贝叶斯分类挖掘模型	161
7.3.2 浏览朴素贝叶斯分类模型和分类预测	164
7.4 电子商务数据的贝叶斯分类	168
7.4.1 建立朴素贝叶斯分类挖掘模型	168
7.4.2 浏览挖掘结果及分析	169
练习题	171
上机实验题	172
第 8 章 神经网络算法	173
8.1 人工神经网络概述	174
8.1.1 人工神经元	174
8.1.2 人工神经网络	176
8.1.3 神经网络应用	177
8.2 用于分类的前馈神经网络	177
8.2.1 前馈神经网络的学习过程	177
8.2.2 前馈神经网络用于分类的算法	180
8.3 SQL Server 神经网络分类	184
8.3.1 建立神经网络分类挖掘模型	184

8.3.2 浏览神经网络分类模型和分类预测	186
8.4 电子商务数据的神经网络分类	189
8.4.1 建立神经网络分类挖掘模型	189
8.4.2 浏览挖掘结果及分析	189
练习题	191
上机实验题	192
第9章 回归分析算法	194
9.1 回归分析概述	195
9.2 线性回归分析	196
9.2.1 一元线性回归分析	196
9.2.2 多元线性回归分析	197
9.2.3 SQL Server 线性回归分析	199
9.3 非线性回归分析	206
9.3.1 非线性回归分析的处理方法	206
9.3.2 可转换成线性回归的非线性回归	206
9.3.3 不可转换成线性回归的非线性回归分析*	208
9.4 逻辑回归分析	209
9.4.1 逻辑回归原理	209
9.4.2 逻辑回归模型	210
9.4.3 SQL Server 逻辑回归分析	211
9.5 电子商务数据的逻辑回归分析	218
9.5.1 建立逻辑回归挖掘模型	218
9.5.2 浏览挖掘结果及分析	219
练习题	220
上机实验题	221
第10章 时间序列分析	222
10.1 时间序列分析概述	223
10.1.1 什么是时间序列和时间序列分析	223
10.1.2 时间序列的分类和平稳性判断	224
10.1.3 时间序列建模的两种基本假设	225
10.1.4 回归分析与时间序列分析	226
10.2 确定性时间序列分析	226
10.2.1 移动平均模型	226
10.2.2 指数平滑模型	228
10.3 随机时间序列模型*	230
10.3.1 随机时间序列模型概述	230
10.3.2 自回归模型 $AR(p)$	231
10.4 SQL Server 时间序列分析	232
10.4.1 建立数据表	232
10.4.2 建立时间序列分析模型	233

10.4.3	浏览时间序列分析模型	236
10.5	电子商务数据的时间序列分析	238
10.5.1	创建 OnRetDMK.TS 数据表	238
10.5.2	数据加载功能设计	238
10.5.3	建立时间序列分析模型	239
10.5.4	浏览时间序列分析模型	241
	练习题	242
	上机实验题	242
第 11 章	聚类算法	243
11.1	聚类概述	244
11.1.1	什么是聚类	244
11.1.2	相似性度量	245
11.1.3	聚类过程	245
11.1.4	常见的聚类算法	246
11.1.5	聚类分析的应用	246
11.2	k -均值算法及其应用	247
11.2.1	k -均值算法	247
11.2.2	SQL Server 的 k -均值算法应用	250
11.3	EM 算法及其应用	256
11.3.1	EM 算法	256
11.3.2	SQL Server 中 EM 算法	260
11.4	电子商务数据的聚类分析	264
11.4.1	建立聚类挖掘模型	264
11.4.2	两种算法结果的比较	266
11.5	Microsoft 顺序分析和聚类分析算法*	269
11.5.1	Microsoft 顺序分析和聚类分析算法概述	269
11.5.2	Microsoft 顺序分析和聚类分析算法的应用	270
	练习题	276
	上机实验题	278
附录 A	部分练习题参考答案	279
第 1 章	279
第 2 章	279
第 3 章	281
第 4 章	281
第 5 章	281
第 6 章	283
第 7 章	284
第 8 章	285
第 9 章	286
第 10 章	286

第 11 章	286
附录 B 上机实验题参考答案	288
第 4 章	288
第 5 章	288
第 6 章	290
第 7 章	291
第 8 章	292
第 9 章	292
第 10 章	293
第 11 章	295
附录 C 书中数据库和包含的数据表	298
1. OnRet 数据库	298
2. SDW 数据库	300
3. OnRetDMK 数据库	301
4. DMK 数据库	301
参考文献	303

第1章 数据仓库和数据挖掘概述



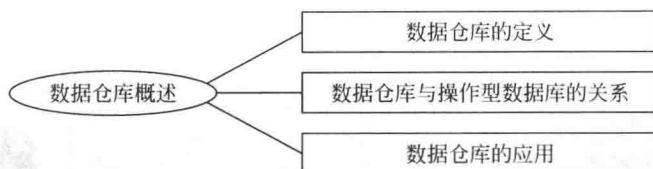
数据仓库体系结构

本章指南

- 数据仓库概述
- 数据仓库系统及开发工具
- 商业智能和数据仓库
- 数据挖掘概述
- 数据挖掘过程

1.1 数据仓库概述

知识梳理



1.1.1 数据仓库的定义

随着数据库技术和计算机网络的成熟,以数据处理为基础的相关技术得到巨大的发展。20世纪80年代中期,“数据仓库”(Data Warehouse, DW)这个名词首次出现在号称“数据仓库之父”W. H. Inmon(恩门)的 *Building Data Warehouse* 一书中。在该书中, W. H. Inmon 把数据仓库定义为“一个面向主题的、集成的、稳定的、随时间变化的数据的集合,以用于支持管理决策过程”。数据仓库4大特征如图1.1所示。

1. 面向主题

主题是指用户使用数据仓库进行决策时所关心的重点领域,也就是在一个较高的管理层次上对信息系统的数据按照某一具体的管理对象进行综合、归类所形成的分析对象。例如,某保险公司有人寿保险和财产保险两类业务,构建有人寿保险和财产保险两个管理信息系统,如果要对所有顾客进行分析,需要构建面向顾客主题的数据仓库;如果要对所有保单进行分析,需要构建面向保单主题的数据仓库;如果要对所有保费进行分析,需要构建面向保费主题的数据仓库,如图1.2所示。

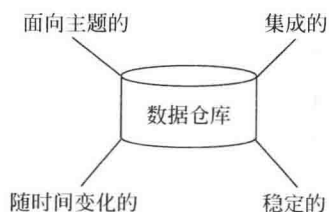


图 1.1 数据仓库的特征

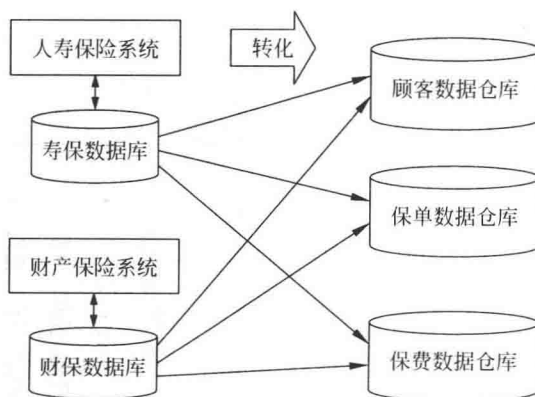


图 1.2 面向主题的示例

从数据组织的角度看,主题是一些数据集合,这些数据集合对分析对象做了比较完整的、一致的描述,这种描述不仅涉及数据自身,而且涉及数据之间的关系。面向主题的数据组织方式,就是在较高层次上对分析对象的数据的一个完整、一致的描述,能完整、统一地刻画各个分

析对象所涉及的企业的各项数据,以及数据之间的联系。

操作型数据库(如人寿保险数据管理系统)中的数据针对事务处理任务(如处理某顾客的人寿保险),各个业务系统之间各自分离,而数据仓库中的数据是按照一定的主题进行组织的。

面向主题组织的数据具有以下特点:

- (1) 各个主题有完整、一致的内容以便在此基础上做分析处理。
- (2) 主题之间有重叠的内容,反映主题间的联系。重叠是逻辑上的,不是物理上的。
- (3) 各主题的综合方式存在不同。
- (4) 主题域应该具有独立性(数据是否属于该主题有明确的界限)和完备性(对该主题进行分析所涉及的内容均要在主题域内)。

2. 集成性

数据仓库中存储的数据一般从企业原来已建立的数据库系统中提取出来,但并不是原有数据的简单拷贝,而是经过了抽取、筛选、清理、转换、综合等工作得到的数据。例如,某顾客数据仓库中的数据是从应用A、B、C中集成的,则需要将性别数据统一转换成m、f,如图1.3所示。

原有数据库系统记录的是每一项业务处理的流水账,这些数据不适合于分析处理。在进入数据仓库之前必须经过综合、计算,同时抛弃一些分析处理不需要的数据项,必要时还要增加一些可能涉及的外部数据。

数据仓库每一个主题所对应的源数据在源分散数据库中有许多重复或不一致之处,必须将这些数据转换成全局统一的定义,消除不一致和错误之处,以保证数据的质量。显然,对不准确,甚至不正确的数据分析得出的结果将不能用于指导企业做出科学的决策。

源数据加载到数据仓库后,还要根据决策分析的需要对这些数据进行概括、聚集处理。

3. 稳定性

数据仓库在某个时间段内来看是保持不变的。

操作型数据库系统中一般只存储短期数据,因此其数据是不稳定的,它记录的是系统中数据变化的瞬态。但对于决策分析而言,历史数据是相当重要的,许多分析方法必须以大量的历史数据为依托。没有大量历史数据的支持是难以进行企业的决策分析的,因此数据仓库中的数据大多表示过去某一时刻的数据,主要用于查询、分析,不像业务系统中的数据库那样,要经常进行修改、添加,除非数据仓库中的数据是错误的。

例如,操作型应用数据库中的数据可以随时被插入、更新、删除和访问(查询),可以从中抽取10年的数据构建数据仓库,用于对这10年的数据进行分析,一旦数据仓库构建完成,它主要用于访问,一般不会被修改,具有相对的稳定性,如图1.4所示。

4. 随时间而变化

数据仓库大多关注的是历史数据,其中数据是批量载入的,即定期从操作型应用系统中接收新的数据内容,这使得数据仓库中的数据总是拥有时间维度。从这个角度看,数据仓库实际是记录了系统的各个瞬态(快照),并通过将各个瞬态连接起来形成动画(即数据仓库的快照集

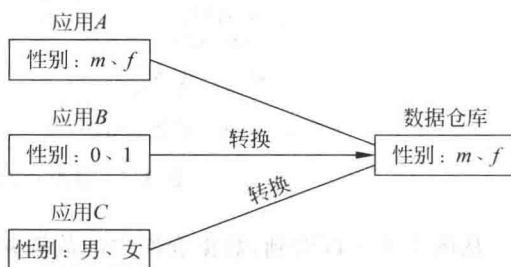


图 1.3 性别的集成