

数据仓库

原理与实践

林宇 等 编著

人民邮电出版社
POSTS & TELECOMMUNICATIONS PRESS

数据仓库

原理与实践

林宇 等 编著

人民邮电出版社

图书在版编目 (CIP) 数据

数据仓库原理与实践 / 林宇等编著. —北京: 人民邮电出版社, 2003.1
ISBN 7-115-10044-6

I. 数... II. 林... III. 数据库系统IV. TP311.13

中国版本图书馆 CIP 数据核字 (2002) 第 093538 号

内容提要

本书比较全面系统地介绍了数据仓库 (Data Warehouse)、联机分析处理 (OLAP)、数据挖掘 (Data Mining) 等 3 个层次的基本概念、原理和应用技术。全书分成 4 篇, 基本原理篇和设计建模篇的内容主要包括: 数据仓库的基本概念、体系结构、创建过程、建模设计、项目规划, OLAP 的基本概念、ROLAP 和 MOLAP 的实现原理、OLAP 模型设计, 数据挖掘的基本概念、基本过程、常见模型和算法。产品介绍篇介绍了现有数据仓库厂商产品工具的基本情况, 并对产品选择进行了一些分析。应用实践篇结合电信领域的实例, 介绍了数据仓库项目在设计 and 实施中的关键问题。

本书的编写以理论联系实际为原则, 内容系统全面, 对于从事数据仓库研究、设计、开发等工作的人员具有宝贵的参考价值, 对于需要了解数据仓库技术的系统集成人员、系统分析师、系统设计师也具有一定的参考价值。

数据仓库原理与实践

- ◆ 编 著 林 宇 等
责任编辑 张立科
- ◆ 人民邮电出版社出版发行 北京市崇文区夕照寺街 14 号
邮编 100061 电子函件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
读者热线 010-67132692
北京汉魂图文设计有限公司制作
北京密云春雷印刷厂印刷
新华书店总店北京发行所经销
- ◆ 开本: 787×1092 1/16
印张: 23.25
字数: 563 千字
印数: 1-4 000 册
- 2003 年 1 月第 1 版
2003 年 1 月北京第 1 次印刷

ISBN 7-115-10044-6/TP · 2729

定价: 42.00 元

本书如有印装质量问题, 请与本社联系 电话: (010)67129223

前 言

随着计算机应用和网络计算的发展，“计算”正向两个不同的方向拓展：一是广度计算，二是深度计算。广度计算是把计算机的应用范围尽量扩大，同时实现广泛的数据交流。深度计算是人们对以往计算机的简单数据操作。目前对深度操作已提出了更高的要求，希望计算机能够更多地参与数据分析与制定决策的工作。传统的数据库技术是单一的数据库资源，它适合操作型事务处理，但分析型处理（或信息型处理）能力较弱。数据仓库的出现，将操作型环境和分析型环境进行了分离，划清了数据处理的分析型环境与操作型环境之间的界限，从而由原来的以单一数据库为中心的数据环境发展为以数据仓库为中心的一种新的体系化环境。

数据仓库技术以改进后的数据库技术作为存储数据和管理资源的基本手段，以统计分析技术作为分析数据和提取信息的有效方法，通过人工智能、神经网络、知识推理等数据挖掘方法来发现数据背后隐藏的规律，从而实现从“数据→信息→知识”的过程，为企业的管理阶层提供各种层次的决策支持。

本书从数据仓库、联机分析（OLAP）、数据挖掘等 3 个层次对数据仓库系统的关键技术进行了深入细致的介绍和分析。

全书共分成 4 篇，分别是基本原理篇、设计建模篇、产品工具篇、应用实践篇。

基本原理篇从第 1 章至第 4 章。第 1 章“数据仓库概述”介绍从数据库到数据仓库的演变过程，着重说明了要解决“蜘蛛网”问题，必须在体系结构上加以变革，将操作型环境和分析型环境分离。第 2 章“数据仓库基本原理”主要介绍与数据仓库相关的基本概念和数据处理的流程，着重说明了数据仓库（Data Warehouse）中数据清洗、数据转化、元数据、数据集市等一系列重要的概念，比较了数据仓库和数据库的特点和关键技术，详细介绍了数据仓库面向主题、数据集成、数据稳定、数据随时间变化的根本特征。第 3 章“OLAP 的基本原理”主要介绍 OLAP 技术的概念和展现方式。第 4 章“数据挖掘的基础”主要介绍与数据挖掘相关的概念，本章首先给出数据挖掘的形式化定义，然后介绍数据挖掘的描述型挖掘算法和预测性挖掘算法，并列举了几个数据挖掘在网络中的实际应用。

设计建模篇从第 5 章至第 10 章，它在基本原理篇的基础上进一步介绍了数据仓库模型设计、OLAP 模型设计、数据挖掘模型、数据仓库的规划以及数据仓库的维护等问题。第 5 章“企业模型设计”和第 6 章“数据仓库的模型设计”用于介绍数据仓库的建模知识。第 7 章“OLAP 建模方法”对 OLAP 建模方法进行了比较细致的介绍。第 8 章“数据仓库规划和开发方法”介绍了数据仓库的投资分析和两种软件开发方法学。第 9 章“数据挖掘的实施过程”介绍了数据挖掘的基本方法。第 10 章“数据仓库的建立和维护”从业务系统、数据仓库、OLAP、数据挖掘、数据展现等方面完整地说明了数据仓库建立的全过程，并介绍了在数据仓库维护阶段设计人员和维护人员需要了解的技术问题。

产品介绍篇包括了第 11 章和第 12 章。在第 11 章中，简单介绍了现有的著名的数据仓库厂商及其主要产品，并选择了有代表性的产品进行说明。在第 12 章中，从技术和需求两方面对产品选择问题作了一些分析。

应用实践篇从第 13 章至第 16 章。本篇结合电信领域的特点，通过一个电信领域的应用

实例来说明数据仓库的具体实施方案。第 13 章“项目的需求和目标分析”主要说明数据仓库在电信领域的应用规划。第 14 章“系统结构和模型设计”介绍了系统结构设计、数据仓库模型设计、OLAP 模型设计三部分。第 15 章“数据装载设计和数据挖掘”介绍了数据载入、建立多维数据库、数据挖掘相关的设计和实现问题。第 16 章“界面设计和项目总结”说明了界面设计的相关问题，并对项目的总体工作流程、数据流程、模块构成进行了总结。

本书的写作过程也是笔者们学习、探讨、实践的过程。在此过程中，笔者阅读了大量的国内外著作、论文，考察了数据仓库各大厂商的产品特征和性能，对大量的资料进行了归纳和整理。

在本书的编写过程中，得到了程时端、卢美莲、金跃辉、龚向阳、阙喜戎等教授的息心指导。参加本书研讨和写作的有林宇、郭凌云、王令成、伦勇、高波、翟朝阳、胡颖峰、李劲松、徐自军、王重钢、乌海涛、王茜、白刚、范锐、彭泳、柴平宣、乐辉华等。另外，高波为本书的筹备和组织做了大量工作。在此对他们的辛勤工作表示深深的谢意！

由于数据仓库技术是一个新的技术，加之笔者水平有限，书中错误在所难免，希望各位专家和广大读者给予批评指正，多提宝贵意见。

编者

2002.11

目 录

第一篇 基本原理篇

第 1 章 数据仓库概述.....	2
1.1 数据库到数据仓库的演变.....	2
1.1.1 蜘蛛网问题.....	2
1.1.2 操作型系统和分析型系统的分离.....	5
1.2 商业智能系统的功能和构成.....	7
1.2.1 商业智能系统的数据处理循环.....	7
1.2.2 决策支持系统的总体功能.....	7
1.3 仓库的应用前景.....	9
第 2 章 数据仓库的基本原理.....	12
2.1 数据仓库的体系结构.....	12
2.1.1 数据仓库的体系结构.....	12
2.1.2 数据仓库中的关键名词.....	13
2.2 数据仓库的特点.....	16
2.2.1 数据库的功能和特征.....	16
2.2.2 数据仓库的功能和特征.....	17
2.2.3 面向主题.....	17
2.2.4 数据的集成性.....	19
2.2.5 数据的稳定性.....	21
2.2.6 数据随时间变化的特点.....	22
2.3 数据仓库的数据组织.....	23
2.3.1 数据仓库的数据组织结构.....	23
2.3.2 数据颗粒度.....	25
2.3.3 数据的分割.....	29
2.3.4 数据仓库的数据组织形式.....	30
2.3.5 数据仓库的数据追加技术.....	32
2.3.6 清理数据仓库的数据.....	36
2.4 数据仓库建设的两条技术路线.....	36
2.5 操作数据存储 ODS.....	38
2.5.1 ODS 的概念.....	38
2.5.2 ODS 的应用.....	39
2.5.3 DB—ODS—DW 的 3 层体系结构.....	41
2.5.4 ODS/DW、ODS/DB 间的比较.....	43
2.6 外部数据和非结构数据.....	44
2.6.1 外部数据的特征.....	44

2.6.2	为什么将外部数据放在数据仓库	45
2.6.3	对外部数据进行管理的元数据	46
2.6.4	外部数据/非结构化数据的存储	47
2.6.5	外部数据的使用	49
第 3 章	OLAP 的基本原理	50
3.1	OLAP 的基本概念	50
3.1.1	OLAP 的基本概念	50
3.1.2	OLAP 的基本分析动作	53
3.1.3	OLAP 的展现方式	57
3.1.4	OLAP 和 OLTP	59
3.1.5	OLAP 的体系结构和分类	60
3.2	以多维数据库为基础的 OLAP 服务器	61
3.2.1	多维数据库 (Multi-Dimensional Database)	61
3.2.2	MDDDB 产品实例	64
3.2.3	MOLAP 产品的结构	66
3.3	基于关系型数据库的 OLAP (ROLAP)	66
3.3.1	维表	67
3.3.2	事实表	67
3.3.3	星型结构	69
3.3.4	ROLAP 和 MOLAP 的比较	73
3.3.5	HOLAP (Hybrid OLAP)	76
3.4	OLAP 的特征和衡量	76
3.4.1	OLAP 的 12 准则	76
3.4.2	OLAP 的简洁准则	79
3.5	OLAP 的前端展现方式	80
3.5.1	OLAP 的 C/S 方式	80
3.5.2	OLAP 的 Web	80
3.5.3	瘦客户机方式	81
3.5.4	OLAP 的局限性	82
第 4 章	数据挖掘基础	83
4.1	数据挖掘的概念	83
4.1.1	什么是数据挖掘	83
4.1.2	数据挖掘的形式化定义	84
4.1.3	数据挖掘的基本步骤	85
4.1.4	数据挖掘方法的分类	85
4.2	描述性挖掘分析	87
4.2.1	关联规则	87
4.2.2	序列模型分析	92

4.2.3 聚类分析 (Clustering)	93
4.3 预测类的挖掘算法	98
4.3.1 分类问题	99
4.3.2 回归问题	102
4.3.3 时间序列问题	102
4.3.4 神经网络	105
4.3.5 决策树分析	110
4.4 数据挖掘的体系结构	112
4.5 数据挖掘系统的应用实例	113
4.5.1 从用户的行为模式来自动地生成动态链接	113
4.5.2 用数据挖掘的方法来解决网络拥挤的问题	119
4.6 数据挖掘和相关系统的联系	120
4.6.1 数据挖掘和人工智能、统计学的关系	120
4.6.2 数据挖掘和数据仓库的关系	121
4.6.3 数据挖掘和 OLAP 的区别和联系	122
4.6.4 数据挖掘不是万能的	124
 第二篇 设计建模篇 	
第 5 章 企业模型设计	126
5.1 数据仓库设计和数据库设计的区别	126
5.2 企业模型	129
5.2.1 什么是企业模型	129
5.2.2 面向对象的分析方法	133
5.2.3 使用面向对象的方法建立企业模型	139
5.2.4 面向对象方法和 ER 模型的对比	144
5.3 企业模型到数据库模型的映射	145
5.3.1 限定集成的范围	145
5.3.2 映射到关系模型	146
5.3.3 对比映射结果和现有系统	148
5.4 将企业模型映射到数据仓库概念模型	149
第 6 章 数据仓库的模型设计	152
6.1 数据仓库的逻辑模型设计	152
6.1.1 系统数据量估算	152
6.1.2 数据颗粒度的选择	153
6.1.3 表的分割	157
6.1.4 增加时间字段	157
6.1.5 去除纯操作型数据	157

6.1.6	合理的表划分	158
6.1.7	定义关系模式	161
6.1.8	增加导出字段	161
6.1.9	记录系统的定义	162
6.2	数据仓库的物理模型设计	163
6.2.1	确定数据的存储结构	163
6.2.2	索引策略	166
6.2.3	数据存储策略	173
6.2.4	存储分配优化	176
6.3	数据装载接口设计	176
6.4	并行优化	177
6.4.1	数据仓库中并行优化的必要性和可能性	177
6.4.2	并行体系结构	179
第 7 章	OLAP 建模方法	183
7.1	维表	183
7.1.1	维表的特征	183
7.1.2	维的变化	184
7.1.3	维表的共享	187
7.1.4	雪花型结构处理多对多关系	189
7.1.5	层次信息和分类信息的位置	191
7.1.6	非分析数据的分离	194
7.1.7	典型的维层次	195
7.2	事实表	198
7.2.1	事实	198
7.2.2	事实表的特征	199
7.2.3	数据的粒度	199
7.2.4	聚合操作和聚合表	200
7.2.5	没有度量变量的事实表	201
7.2.6	通用数据和专用数据事实表	202
7.3	OLAP 的并行优化	204
7.3.1	B-TREE 索引、位图索引和 Bit-Wise 索引	204
7.3.2	星型查询优化	204
7.3.3	预连接技术	207
7.4	多维数据库	207
第 8 章	数据仓库规划和开发方法	210
8.1	数据仓库的投资分析	210
8.1.1	数据仓库的应用目标	210
8.1.2	建设数据仓库的必要性	211

8.1.3	数据仓库的投资回报分析	211
8.1.4	技术选择分析	212
8.1.5	IDC 的统计结果	213
8.2	数据仓库的开发方法	213
8.2.1	瀑布式开发	213
8.2.2	螺旋式开发	214
8.3	数据仓库主题的选择和阶段规划	216
8.3.1	阶段规划的原则	216
8.3.2	维护阶段	216
8.4	现有数据库系统的改造问题	216
8.5	数据仓库和数据库系统的相互作用	217
8.5.1	相互促进的过程	217
8.5.2	解决“蜘蛛网”问题	218
8.5.3	数据仓库的“间接使用”	218
8.6	分布式数据仓库	219
8.6.1	采用分布式数据仓库的原因	219
8.6.2	分布式下的模型建立和数据划分	221
8.6.3	分布式数据仓库的建设策略	224
8.6.4	分布式数据仓库技术的缺点	227
8.7	需要避免的错误	228
第 9 章	数据挖掘的实施过程	231
9.1	数据挖掘过程模型 5A	231
9.2	数据挖掘过程模型 CRISP-DM	233
9.3	数据挖掘过程中的相关问题	235
9.3.1	定义商业问题	235
9.3.2	建立数据挖掘库	236
9.3.3	分析数据 / 选择变量	241
9.3.4	模型训练方法	245
9.3.5	数据挖掘模型的评价方法	248
9.3.6	数据仓库的实施和维护	250
9.3.7	模型实例	250
第 10 章	数据仓库的建立和维护	252
10.1	数据仓库建立的过程	252
10.1.1	建立企业模型	252
10.1.2	阶段规划/主题选取	253
10.1.3	技术准备工作	253
10.1.4	逻辑设计	254
10.1.5	物理设计	255

10.1.6	数据载入接口设计	255
10.1.7	装载一个主题的数据和数据校验	256
10.1.8	OLAP 模型设计和应用开发	256
10.1.9	数据准备程序设计	257
10.1.10	数据挖掘模型设计	257
10.1.11	界面系统设计	258
10.1.12	装载其他主题数据	258
10.1.13	同客户交流	259
10.1.14	重新开始循环	260
10.2	数据仓库的维护工作	260
10.2.1	数据周期	260
10.2.2	参照完整性	261
10.2.3	数据环境信息	262

第三篇 产品介绍篇

第 11 章	数据仓库产品的介绍	266
11.1	数据仓库工具	266
11.2	INFORMIX 数据仓库产品简介	268
11.2.1	INFORMIX 数据仓库解决方案	268
11.2.2	数据抽取、转换和装载	269
11.2.3	数据存储	270
11.2.4	数据访问/呈现	274
11.3	SAS 产品简介	278
第 12 章	数据仓库产品的选择	283
12.1	数据仓库产品应具备的关键技术	283
12.2	各数据仓库厂商产品的比较	285
12.3	数据仓库工具的选择	286
12.4	数据仓库工具的互通问题	287

第四篇 应用实践篇

第 13 章	项目的需求和目标分析	292
13.1	电信领域建立数据仓库的常见主题	292
13.2	电信领域常见的分析问题	295
13.2.1	客户群体划分	295
13.2.2	客户流失分析	296
13.2.3	客户欺诈分析	296
13.2.4	网络规划优化	297

13.2.5 网管中的分析问题	298
13.3 项目规划	301
13.4 需求分析的形成	301
13.4.1 任务说明书	301
13.4.2 需求说明书	302
第 14 章 系统结构和模型设计	307
14.1 系统结构设计	307
14.1.1 数据量的估算	307
14.1.2 系统硬件结构/软件结构选择	307
14.2 数据仓库模型的设计	310
14.2.1 可利用的数据	310
14.2.2 粒度的确定	311
14.3 OLAP 模型设计	314
14.3.1 项目涉及的维度分析	315
14.3.2 各个主题的维度设计	317
第 15 章 系统装载、数据挖掘和界面设计	324
15.1 数据装载/数据综合模块设计	324
15.2 OLAP 模型生成程序	329
15.3 数据挖掘宽表设计和生成	335
15.3.1 确定同目标变量相关的数据	335
15.3.2 创建新变量	336
15.3.3 准备训练集合与验证集合	342
15.3.4 确定分析的次序	343
15.3.5 变量选择	343
15.3.6 模型的维护和完善	344
15.4 创建多维数据库模块设计	345
第 16 章 界面设计和项目总结	348
16.1 界面展现设计	348
16.1.1 三层体系结构	348
16.1.2 按照内容对界面进行规划	349
16.2 系统的工作流程总结	351
16.3 系统的数据流程总结	353
16.4 系统的模块组成	355
附录 常用名词表	356
参考文献	360

第一篇 基本原理篇

本篇主要对数据仓库、OLAP 分析、数据挖掘的基本原理进行了详细介绍，以求读者能够对基本概念有比较明确的认识。

本篇共 4 章。第 1 章“数据仓库概述”介绍了数据库到数据仓库的演变过程，重点说明了“蜘蛛网”问题的产生原因以及随之而来的种种问题。要解决这些问题就必须在体系结构上加以变革，将操作型环境和分析环境分离，使企业由原先以数据库为中心的生产环境过渡到以数据仓库为中心的生产环境。

第 2 章“数据仓库基本原理”主要介绍数据仓库相关的基本概念和数据处理的流程。重点说明了数据仓库 (Data Warehouse) 中数据清洗、数据转化、元数据、数据市场等一系列重要的概念。我们还比较数据仓库和数据库的特点和关键技术，详细阐明数据仓库面向主题、数据集成、数据稳定、数据随时间变化的根本特征。接着，介绍数据仓库的数据组织结构和组织方式，如何进行数据追加，以及数据颗粒度和数据分割的概念。

第 3 章“OLAP 的基本原理”主要介绍 OLAP 技术的概念和展现方式。首先介绍了 OLAP 的基本概念：变量、维、维层次和维分类、事实、多维数据立方体。然后说明了 OLAP 分析的基本动作：数据钻取、数据聚合、数据切片、数据切块、数据旋转等。接着介绍了 OLAP 丰富多彩的展现方式，比较了 OLAP 和 OLTP 处理的差异。最后介绍了以多维数据库为基础的 MOLAP 的实现原理和以关系型数据库为基础的 ROLAP 的实现原理，详细比较了 ROLAP 和 MOLAP 技术的优缺点。

第 4 章“数据挖掘的基础”主要介绍了与数据挖掘相关的概念。本章首先给出了数据挖掘的形式化定义，并根据形式化定义简要说明了数据挖掘的过程。然后介绍了数据挖掘的算法分类，说明了描述型挖掘算法，它包括关联分析、序列分析、分类分析、聚类分析。接着介绍了预测性挖掘算法，它包括决策树算法、神经网络算法等。最后讲述了数据挖掘的体系结构，并列举了几个数据挖掘在网络中的实际应用。

第 1 章 数据仓库概述

本章介绍了从数据库到数据仓库的演变过程，着重说明了“蜘蛛网”问题的产生原因以及随之而来的种种问题。要解决这些问题就必须在体系结构上加以变革，将操作型环境和分析环境分离，使企业由原先以数据库为中心的生产环境过渡到以数据仓库为中心的生产环境。本章还介绍了以数据仓库为基础的商业智能系统的功能和构成以及商业智能系统中的核心技术。最后还简要介绍了数据仓库技术的应用前景。

1.1 数据库到数据仓库的演变

1.1.1 蜘蛛网问题

随着数据库技术的广泛运用，企业的运营环境逐渐转化成以数据库为中心的运营环境。企业对数据的需求是多方面的，除了在企业中建立企业级的数据库外，常常还要建立部门级数据库。比如，市场部人员通常只关心企业的销售、市场策划方面的信息，而不注重企业的研发、生产等其他环节。因此，将销售、市场策划方面的信息抽取出来单独建立部门级的数据库很有必要，这样可以提高数据的访问效率。在部门级数据的基础上还要建立个人级的数据库，比如，专门负责制作公司财务报表的数据人员，常常需要从财务部门的数据库系统中抽取数据。又如，部门经理可能经常抽取常用的数据到本地，有针对性的建立个人级数据库就显得尤为重要。

随着数据的逐层抽取，很可能会形成如图 1-1 所示的“蜘蛛网”现象，使数据的抽取和访问显得错综复杂。一个大型的公司每天进行上万次的数据抽取很普通。这种演变不是人为制造的，而是自然演变的结果，如果不再体系结构上进行调整，“蜘蛛网”问题将越来越严重。

错综复杂的抽取与访问将产生很多的问题，诸如数据分析的结果缺乏可靠性、数据处理效率很低、难于将数据转化成信息。

1. 数据分析的结果缺乏可靠性

图 1-2 中展示了某电信公司的市场部和计划部对业务 A 是否具有市场前景的分析过程和结果。市场部认为“业务 A 的市场前景很好”，而计划部却得到截然相反的结果——“业务 A 没有市场前景”。作为企业的最终决策者，将如何根据这样的结论进行决策呢？

两个分析的数据都来自于企业数据库，但是结论却不同，下面通过分析两个过程的差异来寻求原因。

首先，市场部门和计划部门从企业数据库中抽取的数据可能不同，比如，市场部抽取的是在大客户中对业务 A 的使用情况，而计划部抽取的是在普通客户中对业务 A 的使用情况。两者分析数据的内容存在差异。

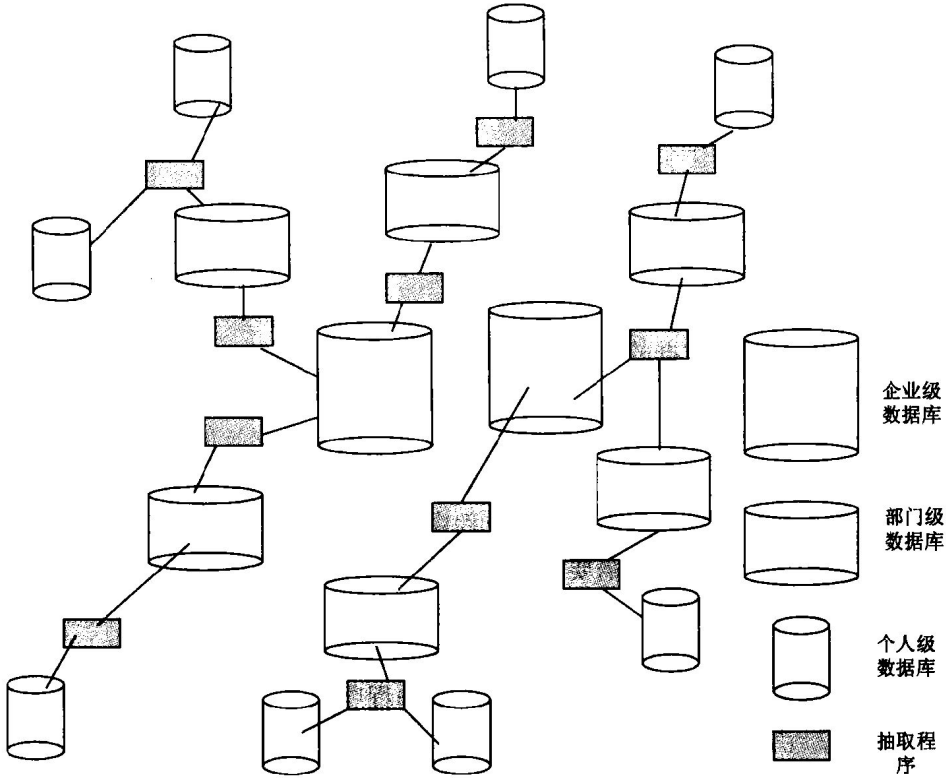


图 1-1 企业中存在的“蜘蛛网”现象

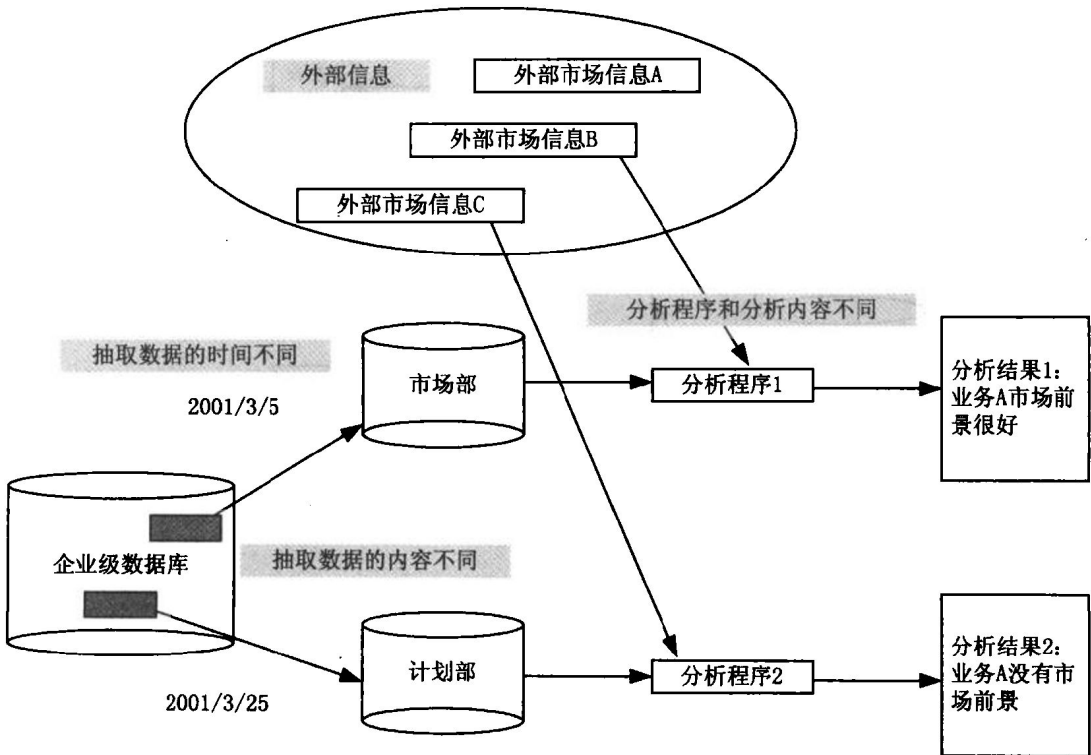


图 1-2 两个分析过程的差异

其次，市场部分分析的时间是 2001 年 3 月 5 日，而计划部分分析的时间是 2001 年 3 月 25 日，两个分析过程从企业数据库中抽取数据的时间不同，企业数据库中的内容已经发生了变化，这很可能导致分析的结果不同。

再次，分析业务的发展趋势常常需要引入企业外部的信息，比如客户的满意程度、国家的政策等。市场部门引用的外部信息来源可能与计划部不同，而外部信息自然是仁者见仁，智者见智，这也可能是导致最终分析结果不同的原因。

最后，市场部门使用的分析程序可能与计划部不同，分析的内容和指标也可能不同。通过上面的分析，我们可以看出导致两个分析过程出现截然相反的结论的根本原因是数据的来源不一致，对于不同来源的数据的分析结果显然是不一致的。

2. 数据处理的效率很低

在错综复杂的体系结构中，不同级别的数据库可能使用不同类型的数据库系统，对于拥有巨型数据量的企业级数据库可能使用 IBM DB2，对于中小型数据库可能使用 SQL Server。各种数据库的开发工具和开发环境不同，抽取程序应用的技术不同，因而难以集成。

如果一个大型企业的决策领导需要一份关于公司整体运营情况的报表，通常需要动用大量的人力和物力才能达到。首先，需要确定报表涉及的内容分布在哪个数据库的哪个位置，然后调动各个部门的程序员/分析员对应用进行分析、设计和编码。

由于数据分散在各个数据库中，因此需要编写的程序很多。由于在企业中使用的数据库类型很多，因此可能需要使用多种技术来实现。程序的重用性很差，因为决策者明天想看的内容很可能与今天不同。可见，动用大量的人力、物力和时间才能完成的报表不仅时效性很差，数据处理的效率也很低。报表形成的过程如图 1-3 所示。

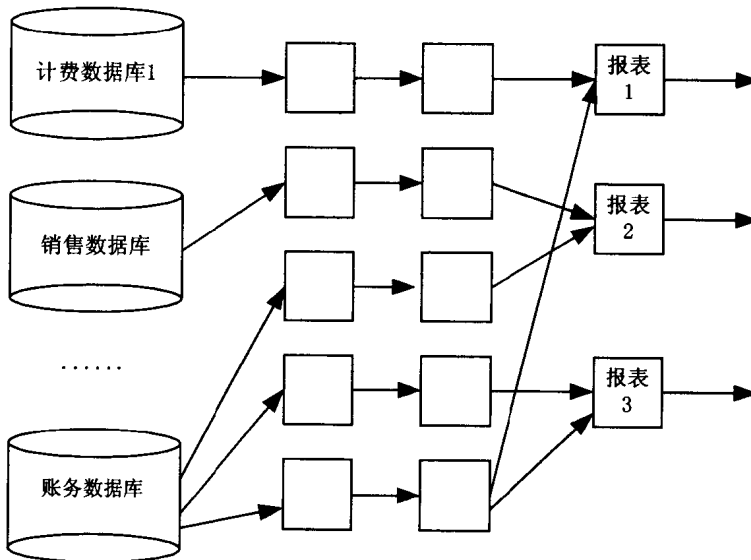


图 1-3 报表形成过程示意图

3. 难以将数据转化成信息

除了数据处理效率和数据可信度的问题之外，“蜘蛛网”式的结构还难以将数据转化为信息。比如，某电信公司想分析某个大客户今年的情况和过去 3 年有什么不同？大客户的情况可能包括客户的呼叫行为、话费情况、交费情况、咨询问题等。因此要想比较完整地回答

这个问题，实际上需要将客户多方面的数据综合成信息。

在实际的数据库系统中，记录客户呼叫行为的数据库通常只保留客户最近3个月的呼叫话单，账务数据库可能保留客户今年的交费情况，客户咨询数据库可能只保留客户2年内的咨询信息，如图1-4所示。每个数据库由于其数据量和业务处理的需求不同，对历史数据的存储时间也不同，因此以现有的数据库系统难以提供完整的历史数据。鉴于这样的原因，用户根本不可能从这些数据中提取出完整的信息。

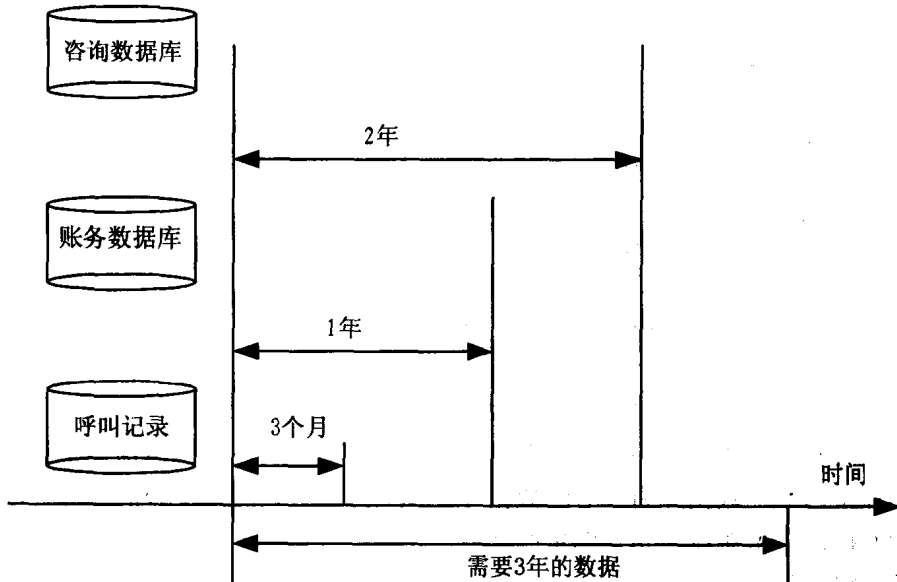


图 1-4 现有数据库系统难以提供完整的历史数据

1.1.2 操作型系统和分析型系统的分离

随着人们应用需求的发展，要解决“蜘蛛网”问题，必须将用于事务处理的数据环境和用于数据分析的数据环境分离开。

这样，数据处理被分为操作型处理和分析型处理（或信息型处理）两大类。操作型处理以传统的数据库为中心进行企业的日常业务处理。比如电信部门的计费数据库用于记录客户的通信消费情况，银行的数据库用于记录客户的账号、密码、存入和支出等一系列业务行为。

分析型处理以数据仓库为中心分析数据背后的关联和规律，为企业的决策提供可靠有效的依据。比如，通过对超市近期数据进行分析可以发现近期畅销的产品，从而为公司的采购部门提供指导信息。又如，对于一个大型的连锁超市，如果能够将这些各个营业点不同时期的营业情况以非常直观的方式展现给管理人员，则管理人员可以根据这些分析结果决定是否需要撤销营业情况极差的营业点，而在客户流量特别大的超市附近增设营业网点。

操作型系统的使用人员通常是企业的具体操作人员，处理的数据通常是企业业务的细节信息，其目标是实现企业的业务运营；而分析型系统的使用人员通常是企业的中高层的管理者，或者是从事数据分析的工程师。分析型系统包含的信息往往是企业的宏观信息而非具体的细节，其目的是为企业的决策者提供支持信息。操作型系统和分析型系统的划分如图1-5所示。