



Enterprise
Data
Warehouse

企业级数据仓库 (EDW)

原理、设计与实践

王彦龙 编著

大型企业级数据仓库构建者的释疑库

原理——系统介绍、通俗易懂

设计——一线经验、操作性强

实践——多年累积、详尽叙述



电子工业出版社

PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

<http://www.phei.com.cn>

内 容 简 介

本书作者曾主持建成了我国证券市场惟一涵盖全市场业务数据的大型企业级数据仓库，他将其基于日常工作实践对数据仓库的概念和理论的理解及建设经验融入本书，详细地讲述了企业级数据仓库（EDW）的基本概念、规划、设计与实现，以及解决方案，并且还详细地描述了电信业、证券业和银行业数据仓库的案例，供数据仓库建设者借鉴。本书既是一本系统介绍数据仓库技术的通俗读物，又是一本数据仓库建设的实践指南，从本身架构到技术描述，从具体内容到实际操作，都不失为一本理论基础牢固、操作性极强的数据仓库经典图书。

本书可供关注、从事数据仓库的技术人员、管理决策人员参考阅读，也适合作为大中院校研究生的参考教材。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。
版权所有，侵权必究。

图书在版编目（CIP）数据

企业级数据仓库（EDW）原理、设计与实践 / 王彦龙编著. —北京：电子工业出版社，2006.9
ISBN 7-121-03109-4

I. 企… II. 王… III. 企业管理—数据库管理系统 IV. F270.7

中国版本图书馆CIP数据核字（2006）第099545号

责任编辑：孙学瑛

印 刷：北京智力达印刷有限公司

装 订：北京中新伟业印刷有限公司

出版发行：电子工业出版社

北京市海淀区万寿路173信箱 邮编100036

开 本：787×980 1/16 印张：23.75 字数：431千字

印 次：2006年9月第1次印刷

印 数：5000册 定价：49.00元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，联系电话：（010）68279077；邮购电话：（010）88254888。

质量投诉请发邮件至 zlt@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线：（010）88258888。

序

Preface

我们处在一个信息爆炸的时代。二十多年前，我们难以掌握希望得到的信息，是因为我们拥有的数据太少。而今天，我们依然难以掌握希望得到的信息，却是因为我们拥有的数据太多。一方面我们缺乏足够的信息来支持科学的决策，而另一方面，多年积累的丰富数据却没有发挥作用。

我们在日常工作中希望将拥有的大量业务数据用于统计和分析，以帮助我们做出正确的决策。这个想法虽然简单而自然，但是在实际操作中我们却发现事情往往并不那么容易。丰富的业务数据往往被存放于分散的异构环境中，很难进行统一的查询访问；而且这些业务数据是按照日常业务处理流程来组织的，并不方便根据决策的主题进行汇总和分析。

数据仓库技术正是为了解决这一问题而产生的，它提供了一种进行分析处理的数据环境，从而使一个企业的信息化建设从支持日常业务操作上升到支持管理层的分析决策。在国外已有不少数据仓库建设的成功案例，它们使企业获得了可喜的效益。

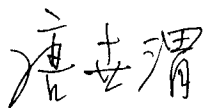
目前，国内数据仓库的建设方兴未艾。随着我国加入WTO，银行、电信、保险、证券等传统的垄断行业不得不对国内外市场的激烈竞争。各个企业为了提高自己的市场竞争力，需要对自身企业的经营状况及整个市场相关行业的态势进行深入分析，进而做出有利的决策。这就需要建设自己的数据仓库系统。在这种背景下，《企业级数据仓库（EDW）原理、设计与实践》一书的出版适逢其时。

《企业级数据仓库（EDW）原理、设计与实践》一书的作者长期参与证券行业技术系统的一线建设工作和管理工作，直接参与了我国最早的证券交易联机事务处理系统的建设工作，主持建成了我国证券市场惟一涵盖全市场业务数据的大型企业级数据仓库，

在数据仓库系统的研发方面具有丰富的实践经验。

作者在书中将其基于日常工作实践对数据仓库的概念和理论的理解进行了系统的总结,对数据仓库涉及的关键技术作了深入浅出的介绍,对数据仓库建设中可能遇到的实际问题进行了深入的分析,并给出了切实可行的解决方案。同时,作者给出了多个行业的数据仓库典型案例,供建设者借鉴。

可以说,《企业级数据仓库(EDW)原理、设计与实践》一书既是一本系统介绍数据仓库技术的通俗读物,又是一本数据仓库建设的实践指南。此书的出版,对我国企业的大型数据仓库建设将不无裨益。



2006年8月

唐世渭教授简介

唐世渭教授,1964年毕业于北京大学数学力学系计算数学专业,毕业后留校任教至今。现为北京大学信息科学技术学院教授,博士生导师,中国计算机学会数据库专业委员会副主任。曾任北京大学信息科学中心主任、视觉与听觉信息处理国家重点实验室主任。多年来承担并完成“973”、“863”、国家科技攻关、国家自然科学基金等多项国家重点科研项目。曾获国家科技进步二等奖、三等奖各一项,省部级科技进步奖多项。在国内外杂志及会议上发表论文百余篇,著译作多部。目前主要研究方向为数据库系统、数据仓库和数据挖掘、Web环境下的信息集成与共享、典型应用领域的信息系统等。

前 言

Introduction

在过去的几十年里，我们生成、收集和存储数据的能力不断提高。人们被淹没在数据的海洋中，却很难从繁杂的数据中获得决策的依据。我们需要一种有效的技术，帮助我们对数据库中的数据进行再加工，形成一个综合的、面向分析的环境，从而更好地支持决策制定。这种有效的技术就是数据仓库技术。

数据仓库是分析型数据库，是与操作系统相分离的、基于标准企业模型集成的、带有时间属性的、面向主题及不可更新的数据集合。它作为分析处理服务的基础，为制定决策提供所需的信息。

传统的数据库系统主要以面向事务处理为主的联机事务处理应用，无法满足决策制定时的分析处理要求。事务处理和分析处理在本质上存在很大差异，这种差异也导致了它们对数据有着不同的要求。数据仓库之父——W.H. Inmon 在其“Building the Data Warehouse”一书中，指出数据仓库中的数据应具备以下四个基本特征：

- 数据仓库的数据是面向主题的 (subject-oriented)
- 数据仓库的数据是集成的 (integrated)
- 数据仓库的数据是非易失的 (nonvolatile)
- 数据仓库的数据是随时间不断变化的 (time-variant)

数据仓库技术近几年得到了迅猛发展。随着市场竞争的加剧，尤其是我国加入 WTO 以后国外竞争者的加入，各个企业为了提高自己的市场竞争力，需要对自身企业的经营状况及整个市场相关行业的态势进行深入分析，进而做出有利的决策。在此过程中，很多大型企业纷纷开始建设自己的数据仓库系统。

作者长期参与证券行业技术系统的一线建设工作和管理工作，直接参与了我国最早的证券交易联机事务处理系统的建设工作，主持建成了我国证券市场唯一涵盖全市场业务数据的大型企业级数据仓库。在此，作者很高兴能有机会把在数据仓库系统研发和建设方面积累的实践经验进行总结并与大家分享，以期对数据仓库的建设者们有所帮助。

本书系统地总结了作者基于日常工作实践对数据仓库概念和理论的理解，对数据仓库涉及的关键技术做了深入浅出的介绍，深入地分析数据仓库建设中可能遇到的实际问题，并给出了切实可行的解决方案，同时，书中还给出了多个行业的数据仓库典型案例供大家参考。

本书是针对数据仓库的开发者、管理者、设计者以及数据仓库建设的其他相关人员而写的，作者真诚地希望读者在数据仓库实践中遇到的大多数问题都能在本书中找到答案。由于本书对数据仓库的概念和理论都做了系统的介绍，因此对于数据仓库的初学者而言，也是一本很好的入门读物。

作者简介



王彦龙 工学硕士、经济学博士、全国金融标准化技术委员会证券分技术委员会副主任委员。曾在深圳证券交易所、中国证监会工作多年，历任深圳证券交易所电脑工程部经理、深圳巨潮信息公司总经理、中国证监会市场监管部交易监管处处长，现任中国证券登记结算公司总工程师。作者对中国证券交易及登记结算技术系统建设具有实际的领导经验，曾主持并参与深圳证券交易所交易系统建设、中国证券登记结算领域多项技术系统的建设，主持搭建了国内首家证券类互联网信息系统，并主持搭建了涵盖沪、深证券市场的统一的TB级数据仓库。

目 录

Contents

第 1 章	数据仓库的基本概念	1
1.1	什么是数据仓库	1
1.1.1	数据仓库是面向主题的	3
1.1.2	数据仓库是集成的	3
1.1.3	数据仓库是非易失的	4
1.1.4	数据仓库是随时间不断变化的	4
1.2	数据仓库和 OLTP 数据库、数据集市的区别	5
1.2.1	数据仓库和 OLTP 数据库	5
1.2.2	数据仓库和数据集市	7
1.3	数据仓库技术的发展	11
1.3.1	数据仓库的起步阶段	11
1.3.2	企业级信息集成	12
1.3.3	企业级数据仓库	13
1.3.4	数据集市	14
1.3.5	争吵与混乱	14
1.3.6	合并	15
1.4	数据仓库的投资回报	15
第 2 章	数据仓库方法论	19
2.1	数据仓库规划	21

2.1.1	业务探索	21
2.1.2	信息调研	22
2.1.3	逻辑数据建模	22
2.1.4	数据仓库解决方案准备	23
2.2	数据仓库的设计与实现	23
2.2.1	系统体系结构设计	24
2.2.2	物理数据库和物理数据模型设计	24
2.2.3	数据转换	25
2.2.4	应用开发	26
2.2.5	数据挖掘	26
2.2.6	数据仓库管理	27
2.2.7	元数据管理	28
2.2.8	数据仓库评估	29
2.3	数据仓库的支持与增强	30
2.3.1	系统维护和支持	30
2.3.2	逻辑数据模型回顾	31
2.3.3	物理数据模型回顾	32
2.3.4	性能调整	32
2.3.5	容量规划	33
第 3 章	数据仓库解决方案	35
3.1	Teradata 数据仓库解决方案	38
3.1.1	产品简介与特点介绍	38
3.1.2	数据装载	44
3.1.3	数据仓库管理工具	47
3.1.4	数据挖掘工具	50
3.2	IBM 数据仓库解决方案	56
3.2.1	产品简介与特点介绍	56
3.2.2	ETL 工具介绍	59

3.2.3	数据仓库工具介绍	60
3.2.4	联机分析工具介绍	63
3.2.5	前端图形工具介绍	65
3.2.6	数据挖掘工具介绍	66
3.3	Oracle 数据仓库解决方案	67
3.3.1	产品简介与特点介绍	67
3.3.2	数据仓库工具介绍	69
3.3.3	联机分析工具介绍	71
3.3.4	数据挖掘工具介绍	74
3.3.5	展现工具介绍	76
第 4 章	实施规划	81
4.1	实施规划阶段的任务	81
4.2	业务探索	83
4.2.1	业务探索阶段的任务	84
4.2.2	业务探索阶段的产出	86
4.3	信息探索	89
4.3.1	信息探索阶段的任务	89
4.3.2	信息探索阶段的产出	93
第 5 章	系统设计	97
5.1	系统体系结构设计	97
5.1.1	设计原则	98
5.1.2	主要需求	99
5.1.3	层次架构	101
5.1.4	组件的详细设计	103
5.2	逻辑数据模型设计	110
5.2.1	设计方法	111

5.2.2	设计目标	114
5.2.3	设计过程	116
5.2.4	设计中的关键问题	121
5.3	物理数据模型设计	125
5.3.1	设计目标	126
5.3.2	技术手段	127
第 6 章	数据的抽取转换加载	137
6.1	数据接口	138
6.1.1	接口流程及要求	138
6.1.2	接口文件说明及格式	143
6.2	数据映射	145
6.2.1	数据映射阶段的任务	146
6.2.2	数据映射阶段的产出	148
6.3	ETL 设计及流程管理	152
6.3.1	ETL 阶段的任务	152
6.3.2	ETL 系统的设计	156
第 7 章	数据汇总	163
7.1	数据汇总的概念	163
7.2	数据汇总的类型	165
7.2.1	实体化视图	166
7.2.2	中间汇总层	166
7.2.3	两种方式的优缺点	167
7.3	中间汇总层的设计原则	167
7.4	中间表的设计模板	169
7.4.1	源表与目标表的对应关系	169

7.4.2	抽取过程说明	170
7.5	数据挖掘基础数据集的设计与开发	170
7.5.1	水平结构的挖掘数据集	172
7.5.2	垂直结构的挖掘数据集	172
7.5.3	两种组织形式的比较	173
7.5.4	基础数据集的开发	173
7.6	数据汇总的典型案例分析	174
7.6.1	数据量	174
7.6.2	基础表结构	174
7.6.3	应用需求	176
7.6.4	设计中间表	177
第 8 章	关键绩效指标 (KPI) 分析	181
8.1	KPI 概述	181
8.2	KPI 设计方法	184
8.2.1	基本方法	184
8.2.2	结合平衡计分卡设计 KPI	186
8.3	基于数据仓库的 KPI 应用	188
8.3.1	KPI 设计	188
8.3.2	KPI 应用系统	190
第 9 章	报表与即席查询	195
9.1	报表	196
9.1.1	固定报表	196
9.1.2	自定义报表	200
9.2	即席查询	204
9.2.1	查询方法	204
9.2.2	工具与技术	208

10.1	OLAP 的概念	211
10.1.1	E.F.Codd 的定义	211
10.1.2	OLAP 委员会的定义	216
10.1.3	FASMI 测试	216
10.2	OLAP 相关术语	218
10.2.1	维 (Dimension)	218
10.2.2	度量值 (Measure)	218
10.2.3	维层次 (Dimension Hierarchy)	219
10.2.4	维成员 (Dimension Member)	219
10.2.5	多维模型 (Multi-Dimensional Model)	220
10.2.6	数据立方体 (Cube)	220
10.2.7	数据单元格 (Cell)	221
10.3	OLAP 操作	221
10.3.1	切片 (Slice)	224
10.3.2	切块 (Dice)	225
10.3.3	下钻 (Drill Down)	225
10.3.4	上卷 (Roll up)	226
10.3.5	旋转 (Rotate)	226
10.4	OLAP 主题的选择	228
10.4.1	自顶向下——业务探索	229
10.4.2	自底向上——信息探索	231
10.4.3	技术实现	232
10.4.4	主题确定——自顶向下与自底向上相结合	233
10.5	构造数据立方体	234
10.5.1	定义维度和度量信息	235
10.5.2	定义数据抽取和转换规则	235
10.5.3	Cube 的存储	236
10.5.4	定义 Cube 的刷新方式	241

10.6	OLAP 分析的方法	241
10.6.1	趋势分析	241
10.6.2	排名分析	243
10.6.3	构成分析	243
10.6.4	意外分析	244
10.6.5	比较分析	244

第 11 章 数据挖掘 247

11.1	数据挖掘的定义	248
11.1.1	技术上的定义	248
11.1.2	商业上的定义	249
11.1.3	数据挖掘和传统分析方法的区别	250
11.1.4	数据挖掘和数据仓库	250
11.2	数据挖掘方法论	251
11.2.1	阶段 1: 定义业务问题范围	251
11.2.2	阶段 2: 选择和抽样	252
11.2.3	阶段 3: 探索型数据分析	252
11.2.4	阶段 4: 建模	253
11.2.5	阶段 5: 实施	253
11.3	数据挖掘实施步骤	254
11.3.1	步骤 1: 准备数据	255
11.3.2	步骤 2: 抽样	255
11.3.3	步骤 3, 5: 建立模型	255
11.3.4	步骤 4: 验证模型	255
11.3.5	步骤 6: 模型评分	256
11.3.6	步骤 7, 8: 执行	256
11.3.7	步骤 9: 模型监测	256
11.4	数据挖掘案例	257
11.4.1	定义业务问题范围	257

11.4.2	数据准备	257
11.4.3	探索型数据分析	260
11.4.4	建模	262
11.4.5	模型评估	264
11.4.6	模型发布	267
第 12 章	数据质量	269
12.1	数据质量的定义	269
12.2	数据质量问题产生的影响	271
12.3	数据质量问题来源	271
12.4	数据质量检查	273
12.4.1	典型问题	273
12.4.2	检查原则	274
12.4.3	管理流程	275
第 13 章	元数据管理	281
13.1	元数据的概念及分类	281
13.1.1	按用途分类	281
13.1.2	按作用分类	283
13.2	元数据的作用	284
13.3	元数据管理标准化	286
13.3.1	OIM 和 CWM 标准	286
13.3.2	CWM 标准	287
13.4	元数据管理系统的设计原则	291
13.5	元数据管理系统举例	292
13.5.1	整体结构	293
13.5.2	元模型	294

13.5.3	元数据采集	294
13.5.4	元数据应用	295
第 14 章	性能调优	297
14.1	获取高性能的关键因素	297
14.1.1	应用需求	298
14.1.2	数据量	300
14.1.3	平台	302
14.2	性能调优的方法	304
14.3	应用级性能调优	305
14.3.1	索引技术	305
14.3.2	实体化视图	308
14.3.3	连接索引	308
14.3.4	数据库压缩	310
14.3.5	抽样近似	311
14.4	产品级性能调整	312
14.4.1	内存调整	312
14.4.2	I/O 调整	313
14.4.3	并行度的调整	314
14.4.4	收集统计信息	315
第 15 章	数据集市	317
15.1	数据集市结构的发展历程	317
15.2	数据集市的概念	319
15.3	数据集市的几种架构	321
15.3.1	独立数据集市	321
15.3.2	从属数据集市	324
15.3.3	逻辑数据集市	326

16.1 电信业数据仓库案例	329
16.1.1 市场背景	329
16.1.2 项目背景	330
16.1.3 数据仓库选型	331
16.1.4 解决方案	331
16.1.5 实施效果	336
16.2 证券业数据仓库案例	337
16.2.1 市场背景	337
16.2.2 项目背景	339
16.2.3 数据仓库选型	342
16.2.4 解决方案	345
16.3 银行业数据仓库案例	347
16.3.1 市场背景	347
16.3.2 项目背景	348
16.3.3 数据仓库选型	349
16.3.4 解决方案	350
16.3.5 实施效果	358

第 1 章 数据仓库的基本概念

本章向读者介绍关于数据仓库的基本概念。首先介绍什么是数据仓库，然后将数据仓库和 OLTP 数据库、数据集市进行比较，阐述它们之间的不同。随着数据仓库技术的不断发展，许多新技术涌现了出来，在数据仓库技术的发展一节将对此进行详细描述。最后，向读者描绘数据仓库的投资回报，IDC 公司的调查结果表明，数据仓库的确为企业提供了巨大的收益。

1.1 什么是数据仓库

社会的需求极大地推动了技术的发展。数据库技术最初产生于 20 世纪 60 年代中期，从最初的网状、层次数据库，到日益成熟并广泛应用的关系数据库，大量的数据已经在数据库系统中积累起来。然而，数据的丰富却未能很好地解决知识贫乏的问题，人们不断尝试对数据库中的数据进行再加工，形成一个综合的、面向分析的环境，从而更好地支持决策制定。这些尝试促使数据仓库的思想逐渐形成。但对于什么是数据仓库，许多人提出了不同的定义。

“数据仓库是分析型数据库，作为分析处理服务的基础，用来存放海量的只读数据，为制定决策提供所需的信息。”

“数据仓库是与操作型系统相分离的、基于标准企业模型集成的、带有时间属性的（即与企业定义的时间区段相关）、面向主题及不可更新的数据集合。”