

Ральф Кимбалл Марджи Росс

ИНСТРУМЕНТАРИЙ ХРАНЕНИЯ И АНАЛИЗА ДАННЫХ



ПОЛНОЕ РУКОВОДСТВО
ПО РАЗМЕРНОМУ МОДЕЛИРОВАНИЮ

УДК 004.6
ББК 32.973.2-018.2
К40

The Data Warehouse Toolkit: The Definitive Guide to
Dimensional Modeling, 3rd edition

Ralph Kimball and Margy Ross

Copyright © 2013 by Ralph Kimball and Margy Ross

All rights reserved. This translation published under license with the original
publisher John Wiley and Sons, Inc. via Igor Korzhenevskiy
of Alexander Korzhenevskiy Agency.

Wiley and the Wiley logo are trademarks or registered trademarks of John
Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries,
and may not be used without written permission.

All other trademarks are the property of their respective owners. John Wiley & Sons,
Inc. is not associated with any product or vendor mentioned in this book.

Кимбалл, Ральф.

К40 Инструментарий хранения и анализа данных : полное руко-
водство по размерному моделированию / Ральф Кимбалл, Марджи
Росс ; [перевод с английского М. А. Райтмана]. — Москва : Эксмо,
2024. — 656 с. — (Data Science. Лучшие книги о науке о данных).

ISBN 978-5-04-108040-2

Сегодня тысячи компаний собирают и сохраняют большие данные о по-
ведении своих клиентов, ассортименте, производственном процессе и других
немаловажных для бизнеса вещах. Однако, чтобы принимать обоснованные
решения на основе этих данных, недостаточно их просто собрать — нужно
правильно их обработать и провести грамотный анализ. Благодаря этой книге
вы освоите все необходимые инструменты для хранения и анализа большого
количества данных, научитесь правильно управлять ими и извлекать полезную
информацию для развития бизнеса.

УДК 004.6
ББК 32.973.2-018.2

ISBN 978-5-04-108040-2

© Райтман М.А., перевод на русский язык, 2024
© Оформление. ООО «Издательство «Эксмо», 2024

Оглавление

Благодарности	23
Введение	24
Для кого эта книга	25
Структура книги	26
Глава 1. Хранение данных, анализ данных и основы размерного моделирования	27
Глава 2. Обзор методов размерного моделирования Кимбалла	27
Глава 3. Розничные продажи	27
Глава 4. Склад	27
Глава 5. Закупки	28
Глава 6. Управление заказами	28
Глава 7. Бухгалтерский учет	28
Глава 8. Управление взаимоотношениями с клиентами	28
Глава 9. Управление персоналом	29
Глава 10. Финансовые услуги	29
Глава 11. Телекоммуникации	29
Глава 12. Транспортировка	29
Глава 13. Образование	29
Глава 14. Здравоохранение	30
Глава 15. Электронная коммерция	30
Глава 16. Страхование	30
Глава 17. Обзор жизненного цикла Кимбалла	30
Глава 18. Процессы и задачи размерного моделирования	30
Глава 19. Подсистемы и методы ETL	30
Глава 20. Задачи и процессы разработки и проектирования систем ETL	31
Глава 21. Аналитика больших данных	31
Веб-сайты	31
Выводы	32

1	Хранение данных, анализ данных и основы размерного моделирования	33
	Различные миры сбора и анализа данных	34
	Цели хранения и анализа данных	35
	Сравнение обязанностей менеджеров DW/BI с издательским бизнесом	37
	Введение в размерное моделирование	40
	Схема «звезда» против кубов OLAP	41
	Таблицы фактов для измерений	44
	Таблицы измерений для описательного контекста	47
	Факты и измерения, объединенные в схему «звезда»	50
	DW/BI-архитектура Кимбалла	53
	Операционные исходные системы	53
	Система извлечения, преобразования и загрузки	54
	Область представления для поддержки анализа данных	56
	Приложения по анализу данных	58
	Архитектура Кимбалла и метафора ресторана	58
	Альтернативные архитектуры DW/BI	62
	Независимая архитектура Data Mart («Витрина данных»)	62
	Веерная архитектура корпоративной информационной фабрики Инмона	64
	Гибридная веерная и кимбалловская архитектура	66
	Мифы о размерном моделировании	67
	Миф 1: размерные модели только для сводных данных	67
	Миф 2: размерные модели для отделов, а не для предприятий в целом	67
	Миф 3: размерные модели не масштабируемы	68
	Миф 4: размерные модели только для заранее определенного использования	68
	Миф 5: размерные модели не интегрируемы	69
	Еще больше причин мыслить многомерно	69
	Соглашения Agile	70
	Выводы	72
2	Обзор методов размерного моделирования Кимбалла	73
	Основные концепции	73
	Соберите бизнес-требования и реалии данных	74
	Совместные рабочие сессии по пространственному моделированию	74

Четырехэтапный процесс размерного проектирования	75
Бизнес-процессы	75
Зернистость	76
Измерения для описательного контекста	76
Факты для измерений	77
Схема «звезда» и кубы OLAP	78
Стабильные расширения размерных моделей	78
Основные методы работы с таблицами фактов	79
Структура таблиц фактов	79
Аддитивные, полуаддитивные и неаддитивные факты	79
Пустые значения (Null) в таблицах фактов	80
Согласованные факты	80
Таблицы фактов транзакций	81
Таблицы фактов периодических моментальных снимков	81
Накопительные таблицы фактов моментальных снимков	82
Таблицы фактов без показателей	83
Агрегированные таблицы фактов, или кубы OLAP	83
Консолидированные таблицы фактов	84
Основные методы работы с таблицами измерений	84
Структура таблицы измерений	84
Суррогатные ключи измерений	85
Натуральные, стойкие и сверхнатуральные ключи	85
Детализация	86
Вырожденные измерения	86
Денормализованные плоские измерения	87
Несколько иерархий в измерениях	87
Флаги и индикаторы как текстовые атрибуты	87
Пустые атрибуты в измерениях	88
Измерение «Календарная дата»	88
Важные ролевые изменения	89
Мусорные измерения	89
Измерения в виде «снежинки»	90
Измерения с внешней опорой	90
Интеграция через согласованные измерения	90
Согласованные измерения	91
Сжатые измерения	91
Копаем вширь	92
Цепочка значений	92
Архитектура шины корпоративного хранилища данных	92

Матрица шины корпоративного хранилища данных	93
Подробная матрица шины реализации	93
Матрица возможностей/заинтересованных сторон	94
Работа с атрибутами медленно изменяющегося измерения	94
Тип 0: сохранение оригинала	94
Тип 1: перезапись	95
Тип 2: добавление новой строки	95
Тип 3: добавление нового атрибута	96
Тип 4: добавление мини-измерения	96
Тип 5: добавление мини-измерения и внешней опоры типа 1 ...	96
Тип 6: добавление атрибута типа 1 к измерению типа 2	97
Тип 7: двойные измерения типа 1 и типа 2	97
Работа с иерархиями измерений	98
Позиционные иерархии с фиксированной глубиной	98
Иерархии с пропущенными уровнями / иерархии переменной глубины	98
Неровные иерархии/иерархии переменной глубины с соединительными таблицами иерархии	99
Рваные иерархии/иерархии переменной глубины с атрибутами пути	99
Продвинутые методы работы с таблицами фактов	100
Суррогатные ключи таблицы фактов	100
Таблицы-«сороконожки» с фактами	100
Числовые значения как атрибуты или факты	101
Факты о задержке/продолжительности	101
Заголовок/строка в таблице фактов	102
Выделенные факты	102
Таблицы фактов прибылей и убытков с выделением фактов ...	102
Факты разных валют	103
Факты с множественными единицами измерения	103
Факты текущего года (Year-to-date)	104
Многопроходный SQL, чтобы избежать объединения таблиц «факт — факт»	104
Отслеживание промежутка времени в таблицах фактов	104
Факты, появляющиеся с опозданием	105
Расширенные методы работы с измерениями	105
Соединения таблиц «измерение — измерение»	105
Многозначные измерения и соединительные таблицы	106
Многозначные соединительные таблицы, меняющиеся во времени	106

Временной ряд тега поведения	107
Исследовательские группы изучения поведения	107
Агрегированные факты как атрибуты измерения	107
Динамические диапазоны значений	108
Измерение «Текстовые комментарии»	108
Несколько часовых поясов	109
Измерения типа «Показатель»	109
Измерения «Шаг»	109
Измерения с возможностью горячей замены	110
Абстрактные общие измерения	110
Измерения «Аудит»	110
Измерения, прибывающие с опозданием	111
Схемы специального назначения	112
Схемы супертипа и подтипа для гетерогенных продуктов	112
Таблицы фактов в реальном времени	112
Схемы событий ошибок	113
3 Розничные продажи	115
Четырехэтапный процесс размерного проектирования	116
Шаг 1: выбор бизнес-процесса	116
Шаг 2: объявление зернистости	117
Шаг 3: определение измерений	118
Шаг 4: определение фактов	118
Пример использования в розничной торговле	119
Шаг 1: выбор бизнес-процесса	121
Шаг 2: объявление зернистости	121
Шаг 3: определение измерений	123
Шаг 4: определение фактов	123
Подробная информация о таблице измерений	127
Измерение «Дата»	127
Измерение «Продукт»	132
Измерение «Магазин»	136
Измерение «Промоакция»	138
Прочие измерения розничных продаж	142
Вырожденные измерения для номеров транзакций	143
Розничная схема в действии	144
Расширяемость схемы розничных продаж	145
Таблицы фактов без метрик	147

Ключи таблиц измерений и фактов	148
Суррогатные ключи таблицы измерений	148
Натуральные, стойкие и сверхнатуральные ключи	151
Суррогатные ключи вырожденного измерения	152
Умные ключи измерения «Дата»	152
Суррогатные ключи таблицы фактов	153
Спротивление стремлению к нормализации	155
Схемы «снежинки» с нормализованными измерениями	155
Внешняя опора	158
Таблицы фактов «сороконожка» со «слишком большим количеством измерений»	159
Выводы	161
4 Склад	163
Введение в цепочку ценности	163
Модели инвентаризации	165
Периодический моментальный снимок инвентаризации	165
Полуаддитивные факты	167
Расширенные сведения о запасах	168
Операции с запасами	169
Накопительный моментальный снимок запасов	171
Типы таблиц фактов	172
Таблицы фактов транзакций	173
Таблицы фактов периодических снимков	173
Накопительные таблицы фактов моментальных снимков	174
Задержки между этапами и количество этапов	175
Накопление обновлений снимков и кубов OLAP	175
Дополнительные типы таблиц фактов	175
Интеграция цепочки ценности	176
Архитектура шины хранилища корпоративных данных	177
Понимание архитектуры шины	177
Матрица шины корпоративного хранилища данных	179
Согласованные измерения	185
Горизонтальный анализ таблицы фактов	185
Идентичные согласованные измерения	186
Сжатие согласованного измерения с подмножеством атрибутов	187
Сжатие согласованного измерения с подмножеством строк	187
Сжатые согласованные размеры на матрице шины	189

Ограниченная согласованность	190
Важность управления данными и ответственности за данные	191
Согласованные измерения и движение к Agile	193
Согласованные факты	194
Выводы	195
5 Закупки	197
Закупки: практический пример	197
Закупочные операции и матрица шины	198
Одна или много таблиц фактов транзакций	199
Дополнительный моментальный снимок закупок	203
Основные сведения о медленно меняющихся измерениях	204
Тип 0: сохранение оригинала	205
Тип 1: перезапись	206
Тип 2: добавление новой строки	208
Тип 3: добавление нового атрибута	211
Тип 4: добавление мини-измерения	214
Гибридные методы медленно изменяющихся измерений	217
Тип 5: добавление мини-измерения и внешней опоры типа 1	218
Тип 6: добавление атрибута типа 1 к измерению типа 2	219
Тип 7: Двойные измерения типа 1 и типа 2	220
Тип 7 для незапланированных отчетов «По состоянию на»	222
Обобщение медленно меняющихся измерений	223
Выводы	224
6 Управление заказами	225
Матрица шины управления заказами	226
Транзакции по заказам	226
Нормализация фактов	227
Ролевые измерения	228
Еще раз об измерении «Продукт»	230
Измерение «Клиент»	233
Измерение «Сделка»	237
Вырожденное измерение для номера заказа	238
Мусорные измерения	239
Паттерн заголовков/строк, которого следует избегать	241
Несколько валют	243
Факты о транзакциях с разной зернистостью	245

Еще один паттерн заголовка/строк, которого следует избегать	247
Операции по выставлению счета	248
Показатели уровня обслуживания в виде фактов, измерений или того и другого	250
Факты о прибылях и убытках	251
Измерение «Аудит»	254
Накопление снимков для конвейера выполнения заказов	256
Расчет задержек	259
Несколько единиц измерения	259
За пределами зеркала заднего вида	261
Выводы	262
7 Бухгалтерский учет	263
Тематическое исследование по бухгалтерскому учету и матрица шин	264
Данные Главной бухгалтерской книги	265
Периодический моментальный снимок Главной бухгалтерской книги	266
План счетов	266
Закрытие периода	267
Факты типа «С начала года и до сегодняшнего дня» (year-to-date/YTD)	269
Пересмотр нескольких валют	269
Транзакции журнала Главной бухгалтерской книги	270
Несколько календарей финансового учета	271
Детализация по многоуровневой иерархии	272
Финансовые отчеты	273
Процесс составления бюджета	274
Иерархии атрибутов измерений	278
Позиционные иерархии с фиксированной глубиной	278
Прерывающиеся иерархии переменной глубины	279
Прерывающиеся иерархии переменной глубины	280
Совместный доступ при прерывающейся иерархии	284
Неравномерная иерархия, изменяющаяся во времени	285
Изменение прерывающихся иерархий	285
Альтернативные подходы к моделированию прерывающейся иерархии	287
Преимущества подхода с соединительной таблицей для неравномерных иерархий	289

Консолидированные таблицы фактов	290
Роль OLAP и комплексных аналитических решений	292
Выводы	293
8 Управление взаимоотношениями с клиентами	295
Обзор CRM-системы	296
Операционная и аналитическая CRM	298
Атрибуты измерения «Клиент»	300
Синтаксический анализ имени и адреса	300
Некоторые соображения по поводу интернационализации имен и адресов	303
Даты, ориентированные на клиента	306
Агрегированные факты как атрибуты измерений	307
Взаимосвязь между интеллектуальным анализом данных и системой DW/BI	310
Различные счетчики в измерениях типа 2	311
Выносное внешнее измерения для атрибутов с низкой кардинальностью	312
Соображения об иерархии клиентов	313
Соединительные таблицы для многозначных измерений	314
Соединительные таблицы для разреженных атрибутов	316
Соединительная таблица для нескольких контактов с клиентами	317
Сложное поведение клиента	318
Группы по изучению поведения для когорт	318
Измерение «Шаг» для последовательного поведения	320
Таблицы фактов временного интервала	321
Пометка таблиц фактов показателями удовлетворенности	324
Пометка таблиц фактов индикаторами ненормальных сценариев	325
Подходы к интеграции клиентских данных	326
Управление основными данными при создании единого измерения «Клиент»	326
Частичная согласованность нескольких измерений «Клиент»	328
Избегание соединений таблицы фактов с таблицами фактов	329
Проверка реальности с низкими задержками	331
Выводы	332

9	Управление персоналом	333
	Отслеживание профиля сотрудников	333
	Точное время вступления в силу и истечения срока действия ...	336
	Отслеживание причин изменения параметров	336
	Изменения профиля как атрибуты типа 2 или события факта ...	337
	Периодический снимок численности персонала	338
	Матрица шин для HR-процессов	339
	Комплексные аналитические решения и модели данных	341
	Рекурсивные иерархии сотрудников	342
	Отслеживание изменений на встроенном ключе менеджера ...	344
	Детализация иерархий управления: вверх и вниз	344
	Многозначные атрибуты ключевых навыков сотрудников	346
	Ключевые слова для навыков сотрудников	347
	Текстовая строка для ключевого слова навыка	348
	Данные анкеты-опросника	349
	Текстовые комментарии	350
	Выводы	351
10	Финансовые услуги	353
	Тематическое исследование банковского дела и матрица шин ...	354
	Рассмотрение измерений для исключения «недостаточного количества измерений»	355
	Измерение «Домохозяйство»	359
	Многозначные измерения и весовые коэффициенты	360
	Пересмотр мини-измерений	362
	Добавление мини-измерений к соединительным таблицам ...	364
	Динамическая группировка значений в таблицах фактов	365
	Схемы супертипов и подтипов для разнородных продуктов	366
	Супертипы и подтипы банковских продуктов с общими фактами	369
	Измерения с возможностью быстрой замены	370
	Выводы	370
11	Телекоммуникации	371
	Тематическое исследование телекоммуникаций и матрица шин ...	371
	Общие соображения по рассмотрению и оценке проекта	374
	Сбалансируйте бизнес-требования и исходные реалии	374
	Сосредоточьтесь на бизнес-процессах	374

Зернистость	375
Единая зернистость фактов	376
Зернистость измерений и иерархии	376
Измерение «Дата»	377
Вырожденные измерения	378
Суррогатные ключи	378
Расшифровки и описания в измерениях	379
Приверженность согласованности	379
Рекомендации по рассмотрению проекта	380
Обсуждение эскизного проекта	382
Изменение существующих структур данных	385
Измерение «Географическое положение»	386
Выводы	387
12 Транспортировка	389
Тематическое исследование авиакомпаний	
и матрица шин для них	390
Зернистость нескольких таблиц фактов	391
Объединение сегментов в поездки	394
Таблицы связанных фактов	395
Расширения для других отраслей промышленности	396
Грузоотправитель	396
Туристические услуги	397
Объединение коррелированных измерений	398
Класс обслуживания	398
Пункты отправки и назначения	399
Дополнительные соображения о дате и времени	401
Календари для конкретных стран в качестве внешних	
выносных измерений	401
Дата и время в нескольких часовых поясах	403
Краткое описание локализации	404
Выводы	404
13 Образование	405
Тематическое исследование университета и матрица шин	405
Таблицы фактов накопительных моментальных снимков	406
Конвейер кандидатов	408
Конвейер предложений по исследовательским грантам	410

Таблицы фактов без фактов	410
События приема	410
Регистрация на курсы	411
Использование объекта	415
Посещаемость студентов	416
Увеличение образовательных аналитических возможностей	417
Выводы	418
14 Здоровоохранение	419
Тематическое исследование здравоохранения и матрица шин	419
Выставление счетов и платежей по претензиям	423
Ролевое измерение «Дата»	426
Многозначные диагнозы	426
Супертипы и подтипы для запросов оплаты	429
Электронные медицинские записи	430
Измерение «Тип измерения» для разреженных фактов	431
Текстовые комментарии произвольной формы	432
Изображения	433
Использование инвентаря помещений, оборудования	433
Работа с ретроактивными изменениями	434
Выводы	434
15 Электронная коммерция	437
Источники данных для потока кликов	437
Проблемы с данными из потоков кликов	438
Размерные модели потока кликов	442
Измерение «Событие»	444
Измерение «Сеанс»	444
Измерение «Направление»	445
Таблица фактов сеанса потока кликов	446
Таблица фактов события страницы потока кликов	449
Измерение «Шаг»	452
Агрегированные таблицы фактов потока кликов	452
Google Analytics	453
Интеграция потока кликов в матрицу шин интернет-магазина	454
Прибыльность по всем каналам, включая веб	456
Выводы	460

16 Страхование	461
Изучение предметной области «Страхование»	462
Цепочка создания стоимости страхования	463
Проект матрицы шины	465
Транзакции по страховому полису	465
Ролевые измерения	466
Медленно меняющиеся измерения	467
Мини-измерения для больших или быстро меняющихся атрибутов	468
Многозначные атрибуты измерений	469
Числовые атрибуты как факты или измерения	469
Вырожденные измерения	470
Таблицы измерений с низкой кардинальностью	470
Измерение «Аудит»	470
Таблица фактов транзакций по полису	470
Гетерогенные продукты: супертипы и подтипы	471
Дополнительная стратегия, накапливающая моментальный снимок	472
Премиальный периодический страховых взносов	473
Согласованные измерения	473
Согласованные факты	474
Факты о предоплате	474
Пересмотр гетерогенных супертипов и подтипов	475
Пересмотр многозначных измерений	476
Более подробная информация об изучении страховых случаев ...	476
Обновленная матрица страховой шины	477
Подробная матрица шины реализации	478
Операции с претензиями	478
Транзакция в сравнении с мусорным измерением «Профиль претензии»	481
Накопительный моментальный снимок для претензий	481
Накопление моментальных снимков для сложных рабочих процессов	482
Накопительный моментальный снимок во времени	483
Моментальный снимок вместо периодического	484
Консолидированный периодический снимок полисов/претензий ..	484
События без фактов, связанные с несчастными случаями	485
Типичные ошибки размерного моделирования, которых следует избегать	486

Ошибка 10: размещение текстовых атрибутов в таблице фактов	487
Ошибка 9: ограничивать подробные дескрипторы для экономии места	487
Ошибка 8: разделение иерархий на несколько измерений	487
Ошибка 7: игнорировать необходимость отслеживать изменения измерений	488
Ошибка 6: решение всех проблем с производительностью с помощью большего количества оборудования	488
Ошибка 5: использование натуральных ключей для соединения измерений и фактов	489
Ошибка 4: пренебрежение декларированием и соблюдением зернистости таблиц фактов.	489
Ошибка 3: использовать отчет для разработки размерной модели	489
Ошибка 2: ожидать, что пользователи будут запрашивать нормализованные атомарные данные	490
Ошибка 1: терпеть неудачу при попытке построить согласованные измерения и факты	490
Выводы	491

17 Обзор жизненного цикла хранилища данных по Кимбаллу	493
Дорожная карта жизненного цикла	494
Дорожная карта и верстовые столбы	495
Мероприятия по запуску жизненного цикла	496
Планирование и управление программами/проектами	496
Определение бизнес-требований	501
Отслеживание технологии жизненного цикла	508
Технический архитектурный дизайн	508
Выбор и установка продукта	511
Отслеживание данных жизненного цикла	513
Размерное моделирование	513
Физический дизайн	513
Проектирование и разработка ETL	515
Отслеживание приложений BI жизненного цикла	516
Спецификация приложения BI	516
Разработка приложений BI	517
Мероприятия по завершении жизненного цикла	517
Развертывание	518

Поддержка и рост	518
Распространенные подводные камни, которых следует избегать	520
Выводы	521
18 Процессы и задачи размерного моделирования	523
Обзор процесса моделирования	523
Подготовка к процессу размерного моделирования	525
Определение участников, особенно представителей бизнеса	525
Ознакомление с бизнес-требованиями	526
Использование инструмента моделирования	527
Использование инструмента профилирования данных	527
Использование или определение соглашения об именовании	527
Координация календарей и помещений	528
Разработка размерной модели	529
Достижение консенсуса по пузырьковой диаграмме высокого уровня	530
Разработка детальной размерной модели	531
Просмотр и подтверждение модели	535
Доработка проектной документации	537
Выводы	537
19 Подсистемы и методы ETL	539
Обобщение требований	540
Потребности бизнеса	540
Соответствие	541
Качество данных	541
Безопасность	542
Интеграция данных	542
Задержка передачи данных	543
Архивирование и происхождение	544
Интерфейсы доставки BI	544
Доступные навыки	545
Существующие на предприятии лицензии	546
34 подсистемы ETL	546
Извлечение: получение данных в хранилище данных	547
Подсистема 1: профилирование данных	547
Подсистема 2: изменение системы сбора данных	548

Подсистема 3: система извлечения	551
Очистка и согласование данных	553
Улучшение культуры и процессов качества данных	553
Подсистема 4: система очистки данных	555
Подсистема 5: схема событий ошибок	557
Подсистема 6: ассемблер измерений «Аудит»	559
Подсистема 7: система дедупликации	559
Подсистема 8: согласующая система	560
Доставка данных: подготовка к презентации	562
Подсистема 9: управление медленно меняющимися измерениями	563
Подсистема 10: генератор суррогатных ключей	569
Подсистема 11: менеджер иерархии	570
Подсистема 12: менеджер специальных измерений	570
Подсистема 13: составители таблиц фактов	573
Подсистема 14: конвейер суррогатных ключей	576
Подсистема 15: конструктор соединительных многозначных параметров	578
Подсистема 16: обработчик данных с задержкой	579
Подсистема 17: система управления измерениями	580
Подсистема 18: система предоставления фактов	581
Подсистема 19: агрегатный конструктор	582
Подсистема 20: конструктор кубов OLAP	583
Подсистема 21: менеджер распространения данных	584
Управление средой ETL	584
Подсистема 22: планировщик заданий	585
Подсистема 23: система резервного копирования	587
Подсистема 24: восстановление и перезапуск системы	589
Подсистема 25: система контроля версий	591
Подсистема 26: система миграции версий	591
Подсистема 27: монитор рабочего процесса	592
Подсистема 28: система сортировки	593
Подсистема 29: анализатор происхождения и зависимостей	594
Подсистема 30: система эскалации проблем	594
Подсистема 31: система распараллеливания/конвейеризации	596
Подсистема 32: система безопасности	596
Подсистема 33: менеджер по соблюдению требований	597
Подсистема 34: менеджер хранилища метаданных	600
Выводы	600

20	Задачи и процессы разработки и проектирования систем ETL	601
	Обзор процесса ETL	601
	Разработка плана ETL	602
	Шаг 1: разработка плана высокого уровня	602
	Шаг 2: выбор инструмента ETL	603
	Шаг 3: разработка стратегий по умолчанию	604
	Шаг 4: детализация по целевой таблице	605
	Разработка системы разовой загрузки исторических данных	608
	Шаг 5: заполнение таблицы измерений историческими данными	608
	Шаг 6: загрузка истории таблицы фактов	614
	Инкрементная обработка ETL	619
	Шаг 7: инкрементная обработка таблицы измерений	619
	Шаг 8: инкрементная обработка таблицы фактов	622
	Шаг 9: сводная таблица и загрузка OLAP	626
	Шаг 10: эксплуатация и автоматизация системы ETL	627
	Последствия обработки в реальном времени	628
	Рассмотрение в реальном времени	628
	Компромиссы в архитектуре реального времени	630
	Разделы данных в режиме реального времени на сервере презентаций	632
	Резюме	634
21	Аналитика Big data	637
	Обзор больших данных	637
	Расширенная архитектура СУБД	639
	Архитектура MapReduce/Hadoop	640
	Сравнение архитектур Big data	641
	Рекомендуемые лучшие практики для Big data	641
	Лучшие практики менеджмента Big data	642
	Лучшие практики архитектуры Big data	644
	Лучшие практики моделирования Big data	650
	Лучшие практики управления Big data	654
	Резюме	655

Введение

Отрасль хранения данных (DW, Data Warehousing) и анализа данных (BI, Business Intelligence) стала развиваться с момента публикации первого издания «Инструментарий хранилищ данных» Ральфа Кимбалла в 1996 г. Хотя ранние крупные корпоративные последователи проложили путь, концепции DW/BI с тех пор были охвачены организациями всех размеров. Отрасль создала тысячи систем DW/BI. Объем данных продолжает расти, поскольку хранилища наполнены растущим количеством атомарных данных и обновляются со все большей частотой. На протяжении нашей карьеры мы видели, что базы данных растут с мегабайтов и гигабайтов до терабайтов и петабайтов, но все же основная задача систем DW/BI остается на удивление постоянной.

Наша работа состоит в том, чтобы передать их бизнес-клиентам для помощи в принятии решений. Мы пришли к этой цели коллективно, бизнес-профессионалы всегда и везде принимают лучшие решения и окупают свои инвестиции в системы DW/BI.

С момента, когда было опубликовано первое издание книги «Инструментарий хранилищ данных», размерное моделирование было принято в качестве основной техники для представления DW/BI. Практики и эксперты единогласно признали, что представление данных должно быть основано на простоте, если вы хотите иметь какие-либо шансы на успех. Простота — фундаментальный ключ, который позволяет пользователям легко понимать базы данных и программное обеспечение для эффективного использования баз данных. Последовательно возвращаясь к перспективе, диктуемой бизнесом, и не компрометируя скорость выполнения запросов и легкость понимания структур пользователями, вы устанавливаете последовательный дизайн, удовлетворяющий аналитическим нуждам организации. Этот многомерно смоделированный фреймворк становится платформой BI. Основываясь на нашем опыте и ошеломительной обратной связи многочисленных практиков из разных компаний, таких как вы, мы верим, что размерное моделирование критически важно для успешной поддержки DW/BI.

Размерное моделирование также выступает в качестве ведущей архитектуры для создания интегрированных систем DW/BI. Когда вы используете согласованные измерения и факты набора размерных моделей, вы имеете практичный и предсказуемый фреймворк для пошагового создания комплексных систем DW/BI, которые изначально распределены.

Учитывая, как много изменений произошло в нашей отрасли, основные методы размерного моделирования, опубликованные Ральфом Кимбаллом 17 лет назад, претерпели проверку временем. Концепции, такие как согласованные измерения, медленно меняющиеся измерения, разнородные продукты, таблицы фактов без метрик, а также матрица шины предприятия продолжают обсуждаться на симпозиумах по всему земному шару. Оригинальные понятия были детализированы и расширены новыми дополнительными методами. Мы решили опубликовать это третье издание оригинальной работы Кимбалла, так как чувствовали, что будет полезно суммировать наш коллективный опыт размерного моделирования под общей обложкой. Каждый из нас концентрировался исключительно на поддержке принятия решений, хранилищах и анализе данных в течение более чем трех десятилетий. Мы хотим поделиться паттернами размерного моделирования, которые неоднократно появлялись в ходе нашей карьеры. Книга нагружена специфическими практическими рекомендациями по проектированию, основанными на реальных сценариях.

Цель этой книги состоит в том, чтобы предоставить универсальный магазин для методов размерного моделирования. Соответственно своему названию книга представляет собой набор инструментов, принципов и методов размерного проектирования. Мы стремимся удовлетворить потребности тех, кто только начинает знакомство с DW/BI, и описываем расширенные понятия для тех из вас, кто уже давно знаком с DW/BI. Мы убеждены, что эта книга выделяется глубиной освещения темы размерного моделирования. Это исчерпывающее руководство.

Для кого эта книга

Книга предназначена для разработчиков в сфере хранилищ и анализа данных, инженеров и менеджеров. Кроме того, книгу найдут полезной аналитики данных и люди, управляющие данными, активно взаимодействующие с DW/BI.

Даже если вы прямо не связаны с размерной моделью, мы считаем важным, чтобы все члены команды проекта были знакомы с понятиями размерного моделирования. Размерная модель оказывает влияние на большинство аспектов реализации DW/BI, начиная от перевода бизнес-требований через процессы извлечения, преобразования и загрузки (extract, transformation, load — ETL) и заканчивая открытием хранилища данных через приложения аналитики данных. Из-за широкого спектра применения вы должны быть сведущими в размерном моделировании независимо от того, за что именно вы ответственны — за управление проектом, анализ данных, архитектуру данных, проектирование баз данных, ETL, приложения BI или обучение и поддержку. Мы написали эту книгу таким образом, что она доступна для широкой аудитории.

Те из вас, кто читал ранние издания этой книги, найдут некоторые уже знакомые тематические исследования, однако они были значительно обновлены и обогащены содержанием, включая примеры матриц шин предприятий почти для каждого тематического исследования. Мы разработали краткие описания для новых предметных областей, включая анализ больших данных.

Содержание этой книги техническое. Впервые мы обсудим размерное моделирование в контексте реляционных баз данных с нюансами для кубов аналитической обработки в реальном времени (OLAP — online analytical processing), что отмечено в соответствующих случаях. Мы предполагаем, что вы имеете базовые понятия о реляционных базах данных, такие как таблицы, строки, ключи и объединения таблиц (join). Поскольку мы будем обсуждать размерные модели в нетрадиционном ключе, мы не станем углубляться в конкретные рекомендации по физическому проектированию и настройке любых описываемых СУБД.

Структура книги

Книга организована вокруг серий кратких описаний или примеров. Мы считаем, что разработка техник проектирования с помощью примеров — весьма эффективный способ, поскольку позволяет делиться реальным руководством и извлекать выгоду из опыта реального мира. Хотя эти примеры не предназначены для того, чтобы стать полномасштабными приложениями или решениями для отрасли, они служат основой для обсуждения паттернов, появляющихся в размерном моделировании. По нашему опыту, зачастую легче уловить главные элементы метода разработки, если отступить от слишком знакомых сложностей собственного бизнеса. Читатели предыдущих изданий благоприятно отзывались о таком подходе.

Предупреждаем вас, что мы отклоняемся от подхода учебных примеров в главе 2. Учитывая широкое промышленное принятие методов размерного моделирования, разработанных Kimball Group, мы составили официальный список наших методов наряду с краткими описаниями и указателями на более подробное освещение и иллюстрации этих методов в последующих главах. Хотя эту главу необязательно читать от начала и до конца, как другие, мы считаем, что эта центральная техническая глава может быть полезной ссылкой и может даже послужить профессиональным контрольным списком для разработчиков DW/BI.

За исключением главы 2 остальные части этой книги опираются друг на друга. Мы начинаем с базовых понятий и представляем более расширенное содержание по мере чтения. Главы следует читать по порядку каждому читателю. Например, может быть сложно осмыслить главу 16, пока вы не прочтете предыдущие главы о розничных продажах, закупках, управлении персоналом и управлении отношениями с клиентами.

Тот, кто читал предыдущее издание, может пропустить первые несколько глав. Хотя некоторые ранние факты и изучение измерений могут быть вам знакомы, мы бы не хотели, чтобы вы забегали слишком далеко вперед. Вы упустите новое в фундаментальных понятиях, если пролистаете много.

ПРИМЕЧАНИЕ Книга наполнена подсказками (такими, как это примечание), списками ключевых понятий и указателями на главы, которые будут более полезны и на которые будет легче ссылаться в будущем.

Глава 1. Хранение данных, анализ данных и основы размерного моделирования

Книга начинается с основ хранения данных, анализа данных и размерного моделирования. Мы исследуем компоненты всей архитектуры DW/BI и устанавливаем основной словарь, который будет использоваться в книге. Будут развеяны некоторые мифы и заблуждения о размерном моделировании.

Глава 2. Обзор методов размерного моделирования Кимбалла

Эта глава описывает более 75 методов размерного моделирования и паттернов. Этот официальный список методов Кимбалла включает указатели на последующие главы, где методы воплощены на примерах.

Глава 3. Розничные продажи

Розничные продажи — это классический пример, используемый для иллюстрирования размерного моделирования. Мы начинаем с классики, поскольку это то, что все понимают. Надеемся, вам не понадобится глубоко задумываться об отрасли, поскольку мы хотим, чтобы вы сфокусировались на основных понятиях размерного моделирования. Мы начнем с обсуждения четырехэтапного процесса разработки размерных моделей. Мы подробно изучим таблицы измерений, включая измерение «Дата», которое будет регулярно повторно использоваться в книге. Мы также обсудим вырожденные измерения, понятие «снежинка» и суррогатные ключи. Даже если вы не разрабатываете проект в сфере продаж, эта глава важна, поскольку наполнена основными принципами.

Глава 4. Склад

Мы останемся в отрасли розничных продаж для изучения следующего примера, но обратим ваше внимание к другому бизнес-процессу. Эта глава представляет архитектуру шины предприятия по хранению данных и матрицу шины

с согласованными измерениями. Эти понятия очень важны для любого, кто стремится к созданию интегрированной и расширяемой архитектуры DW/BI. Мы также сравним три основных типа таблиц фактов: транзакционные таблицы, таблицы фактов периодических моментальных снимков и таблицы фактов кумулятивных моментальных снимков.

Глава 5. Закупки

Эта глава укрепляет важность рассмотрения цепочки ценностей организации в процессе подготовки среды DW/BI. Мы также исследуем серии базовых и расширенных методов обработки признаков медленно меняющихся измерений и будем основываться на изменениях типа 1 (перезапись), типа 2 (добавить строку) и типа 3 (добавить столбец), когда будем представлять читателям тип 0 и типы с 4 по 7.

Глава 6. Управление заказами

В этом примере мы изучим бизнес-процессы, которые, как правило, реализуются первыми в системах DW/BI, поскольку они предоставляют основные метрики исполнения бизнеса — что мы продаем каждому клиенту и по какой цене? Мы обсудим измерения, играющие множество ролей в этой схеме. Мы также исследуем наиболее частые проблемы, с которыми сталкиваются создатели моделей при работе с информацией об управлении заказами, такими как соглашения об именовании заголовков/строк, множественные значения атрибутов или единицы измерения и мусорные измерения с разными операционными индикаторами.

Глава 7. Бухгалтерский учет

Мы обсудим моделирование главной бухгалтерской книги для хранилищ данных в этой главе. Мы опишем соответствующую обработку фактов YTD (Year-to-date) и нескольких фискальных календарей, а также объединенные таблицы данных, которые объединяют данные различных бизнес-процессов. Мы также предоставим детальное руководство по иерархии признаков измерений, начиная с простых денормализованных фиксированных глубоких иерархий и заканчивая сопоставительными таблицами для навигации по более сложным рваным переменным глубоким иерархий.

Глава 8. Управление взаимоотношениями с клиентами

Многие системы DW/BI были построены на предпосылке, что вы должны лучше понимать, чтобы обслужить своих клиентов. В этой главе рассматриваются измерение «Клиент», включая стандартизацию адреса и сопоставительные таблицы

для многозначных признаков измерений. Мы также описываем паттерны моделирования сложного поведения потребителей, такие как консолидация данных клиентов из нескольких источников.

Глава 9. Управление персоналом

В этой главе рассматривается несколько уникальных аспектов моделей измерений «Управление персоналом», включая ситуацию, в которой таблица измерений начинает вести себя как таблица фактов. Мы обсуждаем комплексные аналитические решения, обработку рекурсивных иерархий управления и опросные анкеты. Сравниваются несколько методов обработки многозначных атрибутов.

Глава 10. Финансовые услуги

В банковском примере рассматривается концепция схем супертипа и подтипа для разнородных продуктов, в которых каждая сфера деятельности имеет уникальные описательные признаки и показатели эффективности. Очевидно, что необходимость обрабатывать разнородные продукты не уникальна для финансовых услуг. Мы также обсуждаем сложные отношения между счетами, клиентами и домашними хозяйствами.

Глава 11. Телекоммуникации

Эта глава структурирована несколько иначе, чтобы побудить вас критически мыслить при выполнении анализа архитектуры размерной модели. Начнем с размерного дизайна, который на первый взгляд выглядит правдоподобно. Сможете ли вы найти проблемы? Кроме того, мы исследуем особенности географических измерений местоположения.

Глава 12. Транспортировка

В этом примере мы рассмотрим связанные таблицы фактов на разных уровнях детализации, указывая на уникальные характеристики таблиц фактов, описывающих сегменты в путешествии или сети. Мы более детально рассмотрим измерения даты и времени, охватывая специфичные для страны календари и синхронизацию между несколькими часовыми поясами.

Глава 13. Образование

В этой главе мы рассмотрим несколько таблиц фактов без метрик. Кроме того, мы исследуем таблицы фактов кумулятивных моментальных снимков для обработки заявок студентов и заявок на гранты на исследования. В этой главе вы оцените разнообразие бизнес-процессов в образовательном учреждении.

Глава 14. Здравоохранение

Некоторые из самых сложных моделей, с которыми мы когда-либо работали, относятся к сфере здравоохранения. Эта глава иллюстрирует обработку таких трудных случаев, включая использование соединительных таблиц для моделирования множества диагнозов и поставщиков медицинских услуг, которые осуществляют лечение пациентов.

Глава 15. Электронная коммерция

Эта глава фокусируется на нюансах веб-данных о потоке кликов, включая их уникальные измерения. Мы также представляем измерение «Шаг», которое используется для лучшего понимания любого процесса, состоящего из последовательных шагов.

Глава 16. Страхование

Последний пример объединяет многие паттерны, которые мы обсуждали ранее в книге, в единый набор взаимосвязанных схем. Эту главу можно рассматривать как главу, которая сводит все воедино, поскольку методы моделирования накладываются друг на друга.

Глава 17. Обзор жизненного цикла Кимбалла

Теперь, когда вам удобно разрабатывать размерные модели, мы предлагаем общий обзор действий, с которыми сталкиваются в течение срока службы типичного проекта DW/BI. Эта глава представляет собой краткий обзор второго издания книги «Инструментарий хранилища данных», над которым мы работали в соавторстве с Бобом Беккером, Джой Манди и Уорреном Торнтуэйтом.

Глава 18. Процессы и задачи размерного моделирования

В этой главе изложены конкретные рекомендации по решению задач размерного моделирования в рамках жизненного цикла Кимбалла. Первые 16 глав этой книги охватывают методы размерного моделирования и паттерны проектирования, здесь же описываются обязанности, инструкции и результаты деятельности по проектированию размерного моделирования.

Глава 19. Подсистемы и методы ETL

Система извлечения, преобразования и загрузки потребляет непропорционально большую долю времени и усилий, необходимых для создания среды DW/BI. Тщательное рассмотрение передового опыта выявило 34 подсистемы, найденные в основе почти каждого хранилища данных. Эта глава начинается с требований

и ограничений, которые необходимо учитывать перед проектированием системы ETL, а затем она описывает 34 подсистемы извлечения, очистки, согласования, доставки и управления.

Глава 20. Задачи и процессы разработки и проектирования систем ETL

Эта глава посвящена конкретным тактическим действиям и антипаттернам, связанным с проектированием и разработкой ETL. Она обязательна к прочтению всем, кто отвечает за обязанности ETL.

Глава 21. Аналитика больших данных

В заключительной главе мы сосредоточимся на популярной теме больших данных (Big Data). Мы считаем, что большие данные — естественное продолжение ваших обязанностей в области DW/BI. Мы начнем с обзора нескольких архитектурных альтернатив, в том числе MapReduce и Hadoop, и опишем, как эти альтернативы могут сосуществовать с вашей текущей архитектурой DW/BI. Затем мы рассмотрим лучшие практики управления, архитектуры, моделирования и управления большими данными.

Веб-сайты

Сайт Kimball Group наполнен дополнительными материалами и ресурсами для моделирования:

1. Зарегистрируйтесь для получения советов Кимбалла по дизайну, чтобы получить практическое руководство по размерному моделированию и темам DW/BI.
2. Получите доступ к архиву из более чем 300 советов и статей по архитектуре.
3. Узнайте о публичных и выездных занятиях в Университете Кимбалла по качественному, независимому от поставщиков обучению в соответствии с нашим опытом и публикациями.
4. Узнайте о консалтинговых услугах Kimball Group, чтобы использовать наш многолетний опыт DW/BI.
5. Задайте вопросы другим осведомленным в измерении участникам на форуме Kimball.

Выводы

Цель этой книги — рассказать о принятых в индустрии и проверенных годами методах размерного моделирования и разработки, основанных на более чем 60-летнем опыте авторов и извлеченных уроках в реальной бизнес-среде. DW/BI-системы должны основываться на потребностях бизнес-клиентов и поэтому разрабатываться и представляться с простой размерной точки зрения. Мы уверены, что вы окажетесь на один гигантский шаг ближе к успеху в области DW/BI, если купите это издание.

Теперь, когда вы знаете, куда направляетесь, пришло время погрузиться в детали. Мы начнем с учебника по DW/BI и размерному моделированию в главе 1, чтобы убедиться, что все находятся на одной волне в отношении ключевых терминов и архитектурных концепций.

1

Хранение данных, анализ данных и основы размерного моделирования

Эта глава закладывает основу для следующих. Мы начнем изучение систем хранения и анализа данных (data warehousing and business intelligence, DW/BI) с перспективы высокого уровня. Вы можете быть разочарованы, узнав, что мы не начинаем с технологий и инструментов, но прежде всего система DW/BI должна учитывать потребности бизнеса. Сформулировав четкие потребности бизнеса, мы двинемся в обратном направлении через логический, а затем физический дизайн наряду с решениями по поводу технологии и инструментов.

В этой главе мы закладываем основы понимания в отношении целей хранения и анализа данных, наблюдая при этом странное сходство между обязанностями менеджера DW/BI и издателя.

С этой высокоуровневой перспективы мы исследуем основные концепции размерного моделирования и установим основной словарь. В этой главе мы обсудим основные компоненты DW/BI-архитектуры Кимбалла, а также сравним альтернативные архитектурные подходы. К счастью, независимо от архитектуры вашей системы, в ней всегда найдется место для размерного моделирования. Наконец, мы рассмотрим общие мифы о моделировании. К концу этой главы вы поймете, что вам нужно быть наполовину администратором базы данных (DBA) и наполовину бизнес-аналитиком (MBA) при работе с проектом DW/BI.

В главе 1 обсуждаются следующие понятия:

- бизнес-ориентированные цели DW/BI;
- основные понятия и словарь размерного моделирования, включая таблицы фактов и измерений;
- компоненты и принципы DW/BI архитектуры Кимбалла;

- сравнение альтернативных архитектур DW/BI и роли размерного моделирования в каждой из них;
- заблуждения о размерном моделировании.

Различные миры сбора и анализа данных

Один из важнейших активов любой организации — ее данные. Этот актив почти всегда используется для двух целей: ведение оперативного учета и принятие аналитических решений. Проще говоря, операционные системы (OLTP) — это то место, куда вы вводите данные, а система DW/BI — это то место, откуда вы извлекаете данные.

Пользователи операционной системы крутят колеса организации. Они принимают заказы, регистрируют новых клиентов, отслеживают состояние оперативной деятельности и регистрируют жалобы. Операционные системы оптимизированы для быстрой обработки транзакций. Эти системы почти всегда имеют дело с записью одной транзакции за раз. Они предсказуемо выполняют одни и те же операционные задачи снова и снова, выполняя бизнес-процессы организации. С учетом этого факта операционные системы обычно не поддерживают историю изменения данных, а обновляют их, чтобы отразить наиболее актуальное состояние.

Пользователи DW/BI-систем, с другой стороны, наблюдают за показателями организации, чтобы оценить производительность. Они подсчитывают новые заказы, сравнивают их с заказами прошлой недели и спрашивают, почему новые клиенты зарегистрировались и на что они жаловались. Они беспокоятся, правильно ли выполняются рабочие процессы. Хотя им нужны подробные данные для поддержки их постоянно меняющихся вопросов, пользователи DW/BI почти никогда не имеют дела с одной транзакцией за раз. Эти системы оптимизированы для высокопроизводительных запросов, поскольку вопросы пользователей часто требуют поиска и сжатия сотен или сотен тысяч транзакций в набор ответов. Чтобы еще больше усложнить ситуацию, пользователи DW/BI-систем обычно требуют сохранения исторического контекста для точной оценки эффективности работы организации с течением времени.

В первом издании «Инструментарий хранилища данных» Ральф Кимбалл посвятил целую главу описанию дихотомии между мирами оперативной обработки и хранения данных. В настоящее время широко признано, что система DW/BI имеет совершенно другие потребности, клиентов, структуры и ритмы, нежели операционные системы записи информации. К сожалению, мы все еще сталкиваемся с так называемыми системами DW/BI, которые скорее просто копии операционных систем записи, хранящихся на отдельной аппаратной платформе. Хотя эти среды могут учитывать необходимость изолировать

операционную и аналитическую среды по соображениям производительности, они ничего не делают для устранения других внутренних различий между двумя типами систем. Бизнес-клиенты не в восторге от удобства использования и производительности, предоставляемых этими псевдохранилищами данных. Такие самозванцы оказывают плохую услугу DW/BI, потому что они не признают, что у их клиентов потребности совершенно не такие, как у пользователей операционной системы.

Цели хранения и анализа данных

Прежде чем мы углубимся в детали размерного моделирования, полезно сосредоточиться на основных целях хранения и анализа данных. Цели легко можно сформулировать, прогулявшись по холлу любой организации и прислушавшись к менеджерам. Эти повторяющиеся темы существуют уже более трех десятилетий:

- «Мы собираем тонны данных, но не можем получить к ним доступ»;
- «Нам нужно сгруппировать и отфильтровать данные различными способами»;
- «Люди, принимающие важные для бизнеса решения (бизнес-клиенты), должны иметь доступ к необходимой им информации»;
- «Просто покажи мне, что важно»;
- «Мы проводим целые собрания, где спорим, у кого верные цифры, а не принимаем решения»;
- «Мы хотим, чтобы люди использовали информацию для поддержки принятия решений на основе фактов».

Исходя из нашего опыта скажем: эти проблемы все еще настолько универсальны, что определяют основные требования к системе DW/BI. Теперь превратим эти предложения по управлению бизнесом в требования.

- **Система DW/BI должна сделать информацию доступной.**
Содержание системы DW/BI должно быть понятным. Данные должны быть интуитивно понятными и очевидными для бизнес-клиента, а не только для разработчика. Структуры и метки данных должны имитировать мыслительные процессы и словарь бизнес-клиентов. Бизнес-клиенты хотят разделять и объединять аналитические данные в бесконечных комбинациях. Инструменты и приложения аналитики данных, обеспечивающие доступ к данным, должны быть простыми и удобными в использовании. Они также должны возвращать результаты запроса пользователю с минимальным

временем ожидания. Мы можем обобщить это требование, сказав «просто и быстро».

- **Система DW/BI должна представлять информацию последовательно.**
Данные в системе DW/BI должны быть достоверными, тщательно собранными из различных источников, очищенными, проверенными на качество и выпущенными только тогда, когда они пригодны для использования пользователем. Согласованность также подразумевает использование общих меток и определений для содержимого системы DW/BI, которые используются в разных источниках данных. Если два показателя производительности имеют одинаковое имя, они должны означать одно и то же. И наоборот, если две меры не означают одно и то же, они должны быть помечены по-разному.
- **Система DW/BI должна адаптироваться к изменениям.**
Потребности пользователей, условия бизнеса, данные и технологии могут меняться. Система DW/BI должна быть спроектирована так, чтобы корректно обрабатывать это неизбежное изменение и чтобы не превращать в недействительные существующие данные или приложения. Существующие данные и приложения не должны изменяться или разрушаться, когда бизнес-сообщество задает новые вопросы или когда новые данные добавляются в хранилище. Наконец, если необходимо изменить описательные данные в системе DW/BI, вы должны соответствующим образом учесть эти изменения и сделать их прозрачными для пользователей.
- **Система DW/BI должна предоставлять информацию своевременно.**
Поскольку система DW/BI более интенсивно используется для принятия оперативных решений, сырые данные, возможно, придется преобразовать в полезную информацию за часы, минуты или даже секунды. Команда DW/BI и бизнес-клиенты должны иметь реалистичные ожидания в отношении того, что означает предоставление данных, когда у них мало времени для их очистки или проверки.
- **Система DW/BI должна быть бастионом, защищающим информационные активы.**
Информационные «жемчужины» организации помещаются в хранилище данных. Как минимум это хранилище содержит информацию о том, что и кому вы продаете, а также по какой цене — это потенциально опасные детали в руках не тех людей. Система DW/BI должна эффективно контролировать доступ к конфиденциальной информации организации.

- **Система DW/BI должна служить авторитетной и заслуживающей доверия основой для более эффективного принятия решений.**

Хранилище данных должно иметь правильные данные для поддержки принятия решений. Наиболее важные результаты работы системы DW/BI — решения, которые принимаются на основе представленных аналитических данных. Эти решения обеспечивают влияние на бизнес и ценность, придаваемую системе DW/BI. Первоначальное название, предшествующее названию DW/BI, — по-прежнему лучшее описание того, что вы разрабатываете: система поддержки принятия решений.

- **Деловое сообщество должно принять систему DW/BI, чтобы считать ее успешной.**

Неважно, что вы создали элегантное решение, используя лучшие в своем роде продукты и платформы. Если бизнес-сообщество не принимает эту DW/BI-среду и не использует ее активно, вы не прошли вступительный тест. В отличие от реализации операционной системы, где у бизнес-клиентов нет выбора, кроме как использовать новую систему, использование DW/BI иногда необязательно. Бизнес-клиенты будут использовать систему DW/BI, если это простой и быстрый источник полезной информации.

Хотя каждое требование в этом списке значимо, последние два наиболее важны и, к сожалению, часто игнорируются. Успешное хранение и анализ данных требуют большего, чем быть звездным архитектором, техником, разработчиком моделей или администратором базы данных. Благодаря инициативе DW/BI вы одной ногой стоите в своей зоне комфорта информационных технологий (ИТ), в то время как ваша другая нога находится в незнакомой среде деловых пользователей. Вы должны управлять обеими, модифицируя некоторые проверенные навыки, чтобы приспособиться к уникальным требованиям DW/BI. Очевидно, что вам нужно иметь недюжинный талант, чтобы вести себя так, как будто вы гибрид DBA и MBA.

Сравнение обязанностей менеджеров DW/BI с издательским бизнесом

Исходя из целей DW/BI давайте сравним обязанности менеджеров DW/BI с обязанностями главного редактора. Как редактор высококачественного журнала, вы будете иметь широкие возможности по управлению содержанием, стилем и доставкой журнала. Любой, кто занимает эту должность, занимается следующими видами деятельности:

- Понимать читателей:
 - определить их демографические характеристики;
 - найти, что читатели хотят видеть в этом журнале;
 - определить «лучших» читателей, которые обновят свои подписки и купят товары у рекламодателей журнала;
 - найти потенциальных новых читателей и ознакомить их с журналом.
- Убедиться, что журнал обращается к читателям:
 - выбирать интересный и востребованный контент;
 - принимать решения о макете и визуализации, чтобы максимизировать удовольствие читателей;
 - поддерживать высококачественные стандарты письма и редактирования, придерживаясь единого стиля представления;
 - постоянно контролировать точность статей и претензий рекламодателей;
 - адаптироваться к изменяющимся профилям читателей и обеспечить доступность новых материалов от писателей и участников.
- Поддерживать публикации:
 - привлекать рекламодателей и вести журнал с прибылью;
 - публиковаться на регулярной основе;
 - поддерживать доверие читателей;
 - удовлетворять требования владельцев бизнеса.

Вы также можете определить пункты, которые должны быть нецелевыми для главного редактора. Это может быть создание журнала на основе конкретной технологии печати или использование энергии руководства исключительно для повышения эффективности работы. Сюда отнесем навязывание технического стиля письма, который читатели не могут легко понять, или создание запутанного, перегруженного макета, который трудно читать.

Если издательский бизнес построен на основе эффективного обслуживания читателей, то журнал, вероятно, будет успешным. И наоборот: просмотрите список и представьте, что произойдет, если вы пропустите какой-либо элемент — в конце концов у журнала начнутся серьезные проблемы.

Можно провести яркие параллели между обычным издателем и менеджером DW/BI. Руководствуясь потребностями бизнеса, менеджеры DW/BI должны публиковать данные, которые были собраны из различных источников и отредактированы для обеспечения качества и согласованности. Основная обязанность заключается в обслуживании читателей, иначе известных как бизнес-клиенты. Образ публикации подчеркивает необходимость сосредоточиться на ваших клиентах, а не просто фокусироваться на продуктах и процессах. Хотя вы используете технологию для доставки системы DW/BI, эта технология — в лучшем случае

средство для достижения цели. Таким образом, технологии и методы, использованные для разработки системы, не должны прямо отражаться в ваших главных должностных обязанностях.

Теперь измените обязанности издателя журнала на обязанности менеджера DW/BI:

- Понимать бизнес-клиентов:
 - понимать свои должностные обязанности, цели и задачи;
 - определять решения, которые бизнес-клиенты хотят принимать с помощью системы DW/BI;
 - определять «лучших» пользователей, которые принимают эффективные решения;
 - найти потенциальных новых пользователей и ознакомить их с возможностями системы DW/BI.
- Предоставлять бизнес-клиентам высококачественную, актуальную и доступную информацию и аналитику:
 - выбирать наиболее надежные и действенные данные для представления в системе DW/BI, тщательно отобранные из множества возможных источников данных в вашей организации;
 - сделать пользовательские интерфейсы и приложения простыми и управляемыми паттернами, явно согласованными с профилями когнитивной обработки пользователей;
 - убедиться, что данные точны и им можно доверять, последовательно маркируя их по всему предприятию;
 - постоянно контролировать точность данных и анализа;
 - адаптироваться к изменяющимся профилям пользователей, требованиям и бизнес-приоритетам, а также обеспечивать доступность новых источников данных.
- Поддерживать среду DW/BI:
 - использовать часть кредита для бизнес-решений, принятых с использованием системы DW/BI, и использовать эти успехи, чтобы оправдать расходы на персонал и текущие расходы;
 - регулярно обновлять систему DW/BI;
 - поддерживать доверие бизнес-клиентов;
 - удовлетворять требования бизнес-клиентов, исполнительных спонсоров и ИТ-менеджмента.

Если вы хорошо справитесь со всеми этими обязанностями, вы станете отличным менеджером DW/BI! И наоборот: просмотрите список и представьте, что

произойдет, если вы пропустите какой-либо пункт — в итоге у среды начнутся серьезные проблемы. Теперь сопоставьте этот взгляд на работу менеджера DW/BI с собственным описанием работы. Скорее всего, предыдущий список больше ориентирован на пользователей и бизнес-проблемы и может даже не походить на работу в сфере ИТ. По нашему мнению, именно это делает хранение и анализ данных интересными.

Введение в размерное моделирование

Теперь, когда вы понимаете цели системы DW/BI, давайте рассмотрим основы размерного моделирования. Размерное моделирование широко признано в качестве предпочтительного метода представления аналитических данных, поскольку оно одновременно учитывает два требования:

- предоставлять данные, понятные бизнес-клиентам;
- обеспечивать высокую производительность выполнения запросов.

Размерное моделирование — это давно разработанный метод, упрощающий работу с базами данных. В каждом случае в течение более пяти десятилетий ИТ-организации, консультанты и бизнес-клиенты естественным образом стремились к простой размерной структуре, чтобы соответствовать фундаментальной потребности человека в простоте. Простота имеет решающее значение, поскольку она гарантирует, что пользователи смогут легко понять данные, а также позволяет программному обеспечению быстро и эффективно осуществлять навигацию и возвращать результаты.

Представьте себе руководителя, который описывает свой бизнес следующим образом: «Мы продаем продукты на различных рынках и измеряем наши результаты с течением времени». Разработчики моделей внимательно прислушаются к акценту на продукт, рынок и время. Большинство людей интуитивно представляют такой бизнес как куб данных, грани которого обозначены как «Продукт», «Рынок» и «Время». Представьте себе нарезку вдоль каждого из этих измерений. Точки внутри куба — это места, где хранятся измерения, такие как объем продаж или прибыль для этой комбинации продукта, рынка и времени. Способность визуализировать что-то настолько абстрактное, как набор данных, конкретным и осязаемым способом — секрет понятности. Если эта перспектива кажется слишком простой, хорошо! Модель данных, которая начинается с простого, имеет шанс остаться простой в конце проекта. Модель, которая начинается сложнее, несомненно, будет слишком сложной в конце, что приведет к снижению производительности выполнения запросов и сокращению количества бизнес-клиентов. Альберт Эйнштейн сформировал основную философию,

управляющую размерным дизайном, когда сказал: «Все следует упрощать до тех пор, пока это возможно, но не более того».

Хотя размерные модели часто создаются в системах управления реляционными базами данных, они довольно сильно отличаются от моделей третьей нормальной формы (3NF), которые стремятся устранить избыточность данных. Нормализованные структуры 3NF делят данные на множество отдельных объектов, каждый из которых становится реляционной таблицей. База данных заказов на продажу может начинаться с записи для каждой строки заказа, но превращаться в сложную паутину в виде модели 3NF, возможно, состоящей из сотен нормализованных таблиц.

Промышленность иногда называет 3NF-модели моделями отношения сущностей (ER). Диаграммы отношений сущностей (ER-диаграммы или ERD) — это чертежи, которые сообщают о взаимосвязях между таблицами. И 3NF-модели, и размерные модели могут быть представлены в виде ER-диаграмм, поскольку состоят из объединенных реляционных таблиц. Ключевое различие между 3NF-моделями и размерными моделями заключается в степени нормализации. Поскольку оба типа моделей могут быть представлены как ERD, мы не будем называть 3NF-модели ER-моделями. Вместо этого мы называем их нормализованными моделями, чтобы минимизировать путаницу.

Нормализованные структуры 3NF чрезвычайно полезны для оперативной обработки, поскольку транзакция обновления или вставки затрагивает базу данных только в одном месте. Нормализованные модели, однако, слишком сложны для запросов анализа данных. Пользователи не могут понимать, перемещаться или помнить нормализованные модели, которые напоминают карту системы автострад Лос-Анджелеса. Аналогичным образом большинство реляционных СУБД не могут эффективно запрашивать нормализованную модель. Сложность непредсказуемых запросов пользователей перегружает оптимизаторы базы данных, что приводит к катастрофическим результатам.

ПРИМЕЧАНИЕ Размерная модель содержит ту же информацию, что и нормализованная модель, но упаковывает данные в формате, который обеспечивает понятность для пользователя, производительность запросов и устойчивость к изменениям.

Схема «звезда» против кубов OLAP

Размерные модели, реализованные в реляционных СУБД, называются схемами типа «звезда» из-за их сходства со звездообразной структурой. Размерные модели, реализованные в размерных средах баз данных, называются кубами

оперативной аналитической обработки (OLAP — online analytical processing), как показано на рис. 1-1.

Если ваша среда DW/BI включает в себя схемы типа «звезда» или кубы OLAP, она управляет концепциями измерений. И звезды, и кубы имеют общий логический дизайн с узнаваемыми измерениями, однако физическая реализация отличается.

Когда данные загружаются в куб OLAP, они сохраняются и индексируются с использованием форматов и методов, разработанных для размерных данных. Агрегирование или предварительно рассчитанные сводные таблицы часто создаются и управляются механизмом куба OLAP. Следовательно, кубы обеспечивают превосходную производительность запросов благодаря предварительным расчетам, стратегиям индексации и другим способам оптимизации. Бизнес-клиенты могут производить свертку и развертку, добавляя или удаляя атрибуты из своих анализов с превосходной производительностью, не выполняя новых запросов. Кубы OLAP также предоставляют более аналитически устойчивые функции, которые превосходят функции, доступные в SQL. Недостатком в том, что вы платите производительностью за эти возможности, особенно с большими наборами данных.

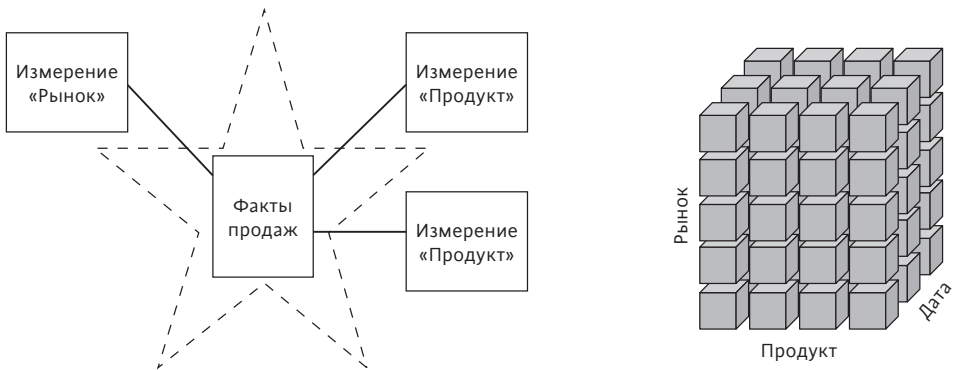


Рис. 1-1. Схема «звезда» в сравнении с кубом OLAP

К счастью, большинство рекомендаций в этой книге не зависят от платформы размерных реляционных баз данных. Хотя возможности технологии OLAP постоянно совершенствуются, мы обычно рекомендуем загружать подробную атомарную информацию в схему типа «звезда», дополнительные кубы OLAP затем заполняются из схемы типа «звезда». По этой причине большинство методов размерного моделирования в этой книге сформулированы в терминах реляционной схемы «звезда».

Соглашение об именовании при развертке OLAP

При развертывании данных в кубах OLAP необходимо учитывать некоторые моменты:

1. Схема типа «звезда», размещенная в реляционной базе данных, — это хорошая физическая основа для создания куба OLAP и обычно рассматривается как более стабильная основа для резервного копирования и восстановления.
2. Кубы OLAP традиционно отличались исключительными преимуществами в производительности по сравнению с реляционными СУБД, но это различие стало менее важным благодаря достижениям в области компьютерного оборудования, такого как устройства и базы данных в памяти, а также в области программного обеспечения реляционных СУБД, такого как столбчатые базы данных.
3. Структуры данных куба OLAP сильнее отличаются у разных поставщиков, чем реляционные СУБД, поэтому подробные сведения о развертывании часто зависят от того, какой поставщик OLAP выбран. Как правило, переносить приложения BI между различными инструментами OLAP сложнее, чем переносить приложения BI на разные реляционные базы данных.
4. Кубы OLAP обычно предлагают более сложные параметры безопасности, чем СУБД, например ограничивают доступ к подробным данным, но предоставляют более открытый доступ к сводным данным.
5. Кубы OLAP предлагают значительно более широкие возможности анализа, чем реляционные СУБД, которые обременены ограничениями SQL. Это может быть основным оправданием использования продукта OLAP.
6. Кубы OLAP исправно поддерживают медленно меняющиеся измерения типа изменения 2 (которые обсуждаются в главе 5), но кубы часто необходимо обрабатывать частично или полностью с использованием альтернативных техник медленно меняющихся измерений, когда данные перезаписываются.
7. Кубы OLAP стабильно поддерживают транзакционные таблицы фактов и таблицы фактов периодических моментальных снимков, но не обрабатывают таблицы фактов кумулятивных моментальных снимков из-за ограничений на перезапись данных, описанных в предыдущем пункте.
8. Кубы OLAP обычно поддерживают сложные рваные иерархии неопределенной глубины, такие как организационные диаграммы или ведомости материалов, с использованием собственного синтаксиса запросов, превосходящего по функциональным возможностям подходы, применяемые в РСУБД.

9. Кубы OLAP могут накладывать детальные ограничения на структуру ключей измерений, которые реализуют более детальные иерархии по сравнению с реляционными базами данных.
10. Некоторые продукты OLAP не поддерживают размерные роли или псевдонимы, поэтому необходимо определять отдельные физические измерения.

Мы вернемся в мир размерного моделирования на реляционной платформе, рассматривая два ключевых компонента схемы типа «звезда».

Таблицы фактов для измерений

Таблица фактов в размерной модели хранит измерения производительности, полученные в результате событий бизнес-процессов организации. Вы должны стараться хранить данные самых детальных измерений, возникающие в результате бизнес-процесса, в одной размерной модели. Поскольку эти данные в подавляющем большинстве — самые большие наборы данных, их не следует реплицировать в нескольких местах для нескольких организационных функций по всему предприятию. Предоставление бизнес-клиентам из нескольких организаций доступа к единому централизованному репозиторию для каждого набора данных измерений обеспечивает использование согласованных данных по всему предприятию.

Термин «факт» представляет собой бизнес-измерение. Представьте себе, что вы стоите на рынке, наблюдаете, как продаются продукты, и записываете количество единиц товара и объем продаж в долларах для каждого продукта в каждой транзакции. Эти измерения регистрируются при сканировании продуктов в регистре, как показано на рис. 1-2.

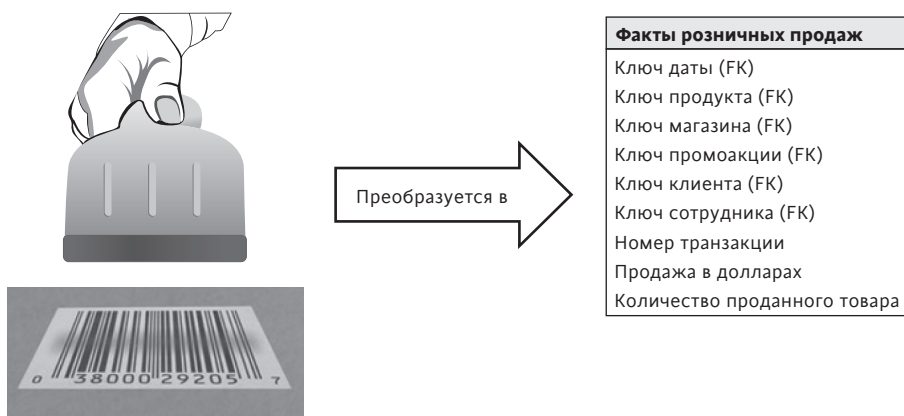


Рис. 1-2. События измерения бизнес-процессов переводятся в таблицы фактов