

Ralph Kimball 著
Richard Merz 译
张丽萍 等 译



Web数据仓库 构建指南

Z



清华大学出版社

Web 数据仓库构建指南

Ralph Kimball 著
Richard Merz
张丽萍 等译

清华大学出版社
北京

Ralph Kimball & Richard Mervin
The Data Webhouse Toolkit
EISBN: 0-471-37680-9

Copyright © 2000 by Ralph Kimball, Richard Merz.

All Rights Reserved. Authorized translation from the English language edition published by John Wiley & Sons, Inc.

Simplified Chinese translation edition is published and distributed exclusively by Tsinghua University Press under the authorization by John Wiley & Sons, Inc., within the territory of the People's Republic of China only (excluding Hong Kong, Macao SAR and Taiwan). Unauthorized export of this edition is a violation of the Copyright Act. Violation of this Law is subject to Civil and Criminal Penalties.

本书中文简体字翻译版由美国 John Wiley & Sons, Inc. 公司授权清华大学出版社在中华人民共和国境内(不包括中国香港、澳门特别行政区和中国台湾地区)独家出版发行。未经许可之出口视为违反著作权法,将受法律之制裁。未经出版者预先书面许可,不得以任何方式复制或抄袭本书的任何部分。

北京市版权局著作权合同登记号 图字: 01-2004-1234

版权所有,翻印必究。举报电话: 010-62782989 13501256678 13501256678

本书防伪标签采用清华大学核研院专有核径迹膜防伪技术，用户可通过在图案表面涂抹清水，图案消失，水干后图案复现；或将表面膜揭下，放在白纸上用彩笔涂抹，图案在白纸上再现的方法识别真伪。

图书在版编目(CIP)数据

、Web 数据仓库构建指南/金贝尔(Kimball, R.), 莫兹(Merz, R.)著; 张丽萍等译.—北京: 清华大学出版社, 2005.4

书名原文：The Data Webhouse Toolkit

ISBN 7-302-10502-2

I .W… II .①金… ②莫… ③张… III .计算机网 - 数据库 IV .TP311.13

中国版本图书馆 CIP 数据核字 (2005) 第 011983 号

出版者：清华大学出版社

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

社总机：010-62770175

地 址：北京清华大学学研大厦

邮 编：100084

客户服务：010-62776969

责任编辑：常晓波

封面设计：立日新

印 刷 者：北京密云胶印厂

装订者：北京鑫海金澳胶印有限公司

发 行 者：新华书店总店北京发行所

开 本：185×230 印张：16.5 字数：365千字

版 次：2005年4月第1版 2005年4月第1次印刷

书 号：ISBN 7-302-10502-2/TP·7128

印 数：1~3000

定 价：29.80 元

Teh Data Webhouse Toolkit : Building the Web-Enabled Data Warehouse

所赢得的赞誉：

“在当今充斥着商业 Web 服务器日志分析工具的世界中, Kimball 的新作给人们带来了新鲜的活力, 就如同他的处女作一样。 Kimball 针对点击流数据提出的维度模型证明了标准的 RDBMS 加上超越标准的大脑才是全球电子商务站点的需要。”

——ArsDigita.com 公司 CEO Philip Greenspun
“Ralph 最新的力作”

“Ralph 的最新力作引发了 Internet 的第二次浪潮。在 Ralph 看来, 公司不仅可以把 Internet 作为一种销售手段, 更可将其用做一种可以增进客户关系的企业战略工具。最关键的是, 此书提供了一种崭新的视角, 可以帮助公司综合利用基于 Internet 的商业智能和由 Internet 带来的海量客户资料, 将‘吸引眼球’转变为利润。”

——IBM NUMA - Q 市场营销部全球解决方案方管
William Schmarzo

“作为一种社会现象, 我们开始接受这样的现实: Internet 真的改变了一切。伟大的公司开始乐于接受这个现实, 并大胆地投资于基于 Web 的信息资产。此书提出了这种设想, 也给出了使之付诸实现的实践建议。”

——Oracle 全球服务部商业智能和仓库副总裁
David Fenton

感 谢

本书两位作者都为能够作为专家生活在历史中这样一个激动人心的时期而感到激动。我们在 IBM 大型机还只有不到 8kB 内存的时候就已经开始自己的职业生涯了！由 Web 引起的改变正在以惊人的速度发生，而且很明显我们还只是处于 Internet 时代的初期。因此我们首先要承认 Web 本身作为一种有趣的现象创造了如此多的机会。

本书的编写过程是一件真正的乐事，部分原因是我们在 John Wiley & Sons 的编辑 Bob Elliott 那里得到了极好支持。Bob 和他的职员是这个行业中最优秀的。我们享受着在 Wiley 办公室与 Emilie Herman、Susan McCarty 和 Brian Snapp 共同工作的乐趣。

《智能企业》(*Intelligent Enterprise*)杂志的主编 Justin Kestelyn 特别慷慨地允许我们从 Ralph Kimball 的一些文章中摘录材料以便在本书中使用。

我们还要感谢 Verio 公司的 Neal Kite 和 Kent Millington 对这项工作的鼓励和支持，此外还要特别感谢 Verio 所属 WebCom 子公司的网站管理员 James Mitchell 为第 10 章绘制的大量插图。

与 Ralph Kimball 的前两本工具书类似，本书在将原稿提交给 Wiley 之前经过了详尽地审查。尽管这样的审查步骤会给作者带来更多的工作，但是这些努力使得本书精益求精。此次我们有许多非常出色的审查人员，这里也要感谢他们，他们的名字(按照字母顺序)为：Caroline Ahdab、Maurice Frank、Mike Heathman、Joy Mundy、Margy Ross、Susan Roston、Dick Smallwood、Mark Sweiger、Jerry Tattershall 和 Warren Thornthwaite。

我们还要感谢 Julie Kimball、Ralph Kimball 的生意伙伴和妻子，她在本书的整个写作过程中起着不可或缺的作用。Julie 谈好了商业合同，安排了审查人员，并激励我们前进。

最后，Ralph 和 Richard 还想要感谢我们的妻子 Julie Kimball 和 Toni Merz 在本书写作过程中给予的坚定支持和热情，感谢我们已长大成人的孩子 Sara Hayden Smith、Brian Kimball、Alexey Merz 和 Chris Merz 能够容忍他们父亲的脾气。

目 录

引言	1
----------	---

第 1 部分 将 Web 引入数据仓库

第 1 章 为什么要把 Web 引入数据仓库	14
1.1 为什么说点击流不只是另一种数据资源	15
1.2 行为分析	16
1.3 隐私保障	18
1.4 Web 数据仓库的结构	19
1.4.1 用户和 ISP	20
1.4.2 公共 Web 服务器和商业事务	20
1.4.3 热响应缓存	22
1.4.4 Web 数据仓库系统	24
1.5 小结	25
第 2 章 跟踪网站用户的操作	26
2.1 用户操作的主要类别	29
2.2 购买产品的步骤	29
2.2.1 察觉需要	30
2.2.2 试着找到所需商品	30
2.2.3 搜索关于可替换产品的信息	30
2.2.4 选择	31
2.2.5 相关商品销售和升级商品销售	31
2.2.6 结账	31
2.2.7 订单之后的处理过程	33
2.3 购买软件或内容的步骤	34
2.4 跟踪的要素	34
2.4.1 用户来源	35
2.4.2 会话识别	35
2.4.3 用户识别	37
2.5 行为分析	40

2.5.1 入口点	40
2.5.2 驻留时间	40
2.5.3 查询	41
2.5.4 站点导航	41
2.5.5 出口点	41
2.6 关联各种操作	43
2.7 个性化的要求	43
2.7.1 重识别	44
2.7.2 用户界面和内容的个性化	44
2.7.3 相关销售和刺激性销售	44
2.7.4 有效协作过滤	44
2.7.5 日程表和有关生活方式的事件	45
2.7.6 本地化	46
2.8 小结	46
第 3 章 利用点击流来进行决策	47
3.1 关于辨认客户的决策	48
3.1.1 通过识别客户来定制营销活动	48
3.1.2 通过对客户进行集群分析来锁定营销活动的目标	49
3.1.3 决定是否鼓励或者支持引荐相关链接	51
3.1.4 判定客户是否将要离开	52
3.2 关于通信的决策	52
3.2.1 判定一个特定的 Web 广告是否有效	52
3.2.2 判定客户问候是否有效	53
3.2.3 判定促销是否有利可图	54
3.2.4 对客户的生活变化做出响应	55
3.2.5 提高网站的有效性	56
3.2.6 培育社区感觉	57
3.3 关于网络业务的基本决策	57
3.3.1 决定通过网络提供哪些产品和服务	57
3.3.2 提供对操作的实时状态跟踪	58
3.3.3 判断我们的 Web 业务是否有利可图	59
3.4 小结	61
第 4 章 把点击流理解为数据源	62
4.1 Web 客户端/服务器互动——概要指南	62

4.1.1 基本的客户端/服务器互动	63
4.1.2 广告	64
4.1.3 引用站点	64
4.1.4 特征文件	65
4.1.5 复合站点	65
4.2 代理服务器和浏览器缓冲	65
4.3 Web 服务器日志	67
4.3.1 主机	68
4.3.2 识别	70
4.3.3 审核	70
4.3.4 时间	70
4.3.5 请求	70
4.3.6 状态	71
4.3.7 字节数	72
4.3.8 访问者	72
4.3.9 用户代理	72
4.3.10 文件名	73
4.3.11 服务时间	73
4.3.12 IP 地址	73
4.3.13 服务器端口	73
4.3.14 进程 ID	73
4.3.15 URL	73
4.4 Cookie	74
4.4.1 Cookie 内容	75
4.4.2 Cookie 手册——检查自己的 cookie 文件	76
4.5 统一系统标识符	77
4.6 查询字符串	78
4.7 小结	79
第 5 章 支持数据仓库的网站设计	80
5.1 单块集成电路与分布式 Web 服务器的比较	80
5.2 使你的服务器同步	81
5.3 页面内容标签	84
5.3.1 用于静态 HTML 的内容索引	85
5.3.2 用于动态 HTML 的内容索引	85

5.3.3 一个简单的内容索引应用程序	86
5.4 一致的 Cookie	87
5.5 空日志记录服务器	88
5.6 个人数据仓库	90
5.7 建立信任	90
5.8 小结	92
第6章 创建点击流数据集市	93
6.1 多维建模快速浏览	93
6.2 点击流维	99
6.2.1 日历日期维	100
6.2.2 一日时间维	102
6.2.3 客户维	103
6.2.4 页面维	106
6.2.5 事件维	107
6.2.6 会话维	108
6.2.7 引荐维	108
6.2.8 产品(或服务)维	109
6.2.9 因果关系维	111
6.2.10 商业实体维	111
6.2.11 点击流追踪关键字	113
6.3 点击流数据集市	113
6.3.1 用于分析完整会话的点击流事实表	114
6.3.2 用于分析个体页面使用情况的点击流事实表	117
6.3.3 聚集点击流事实表	119
6.4 小结	120
第7章 装配点击流价值链	122
7.1 销售交易数据集市	122
7.2 客户通信数据集市	123
7.3 网络收益率数据集市	124
7.4 针对网络零售商的供应链	127
7.5 保险业中的保险单和索赔链	129
7.6 销售流水线链	130
7.7 卫生保健价值链	132
7.8 小结	133

第 8 章 实现点击流后处理机	135
8.1 后处理机体系结构	136
8.1.1 页面事件提取器	137
8.1.2 内容解析器	138
8.1.3 会话识别器	138
8.1.4 计算驻留时间	139
8.1.5 主机和引用站点解析器	140
8.2 小结	142

第 2 部分 把数据仓库引入 Web

第 9 章 为什么要把数据仓库引入 Web	143
9.1 Web 拉动数据仓库	144
9.2 Web 推动数据仓库	145
9.2.1 紧化用户界面反馈回路	145
9.2.2 查询与更新的整合	146
9.2.3 速度是不可商量的	146
9.2.4 Sun 从不设置 Web 数据仓库	147
9.2.5 多媒体合并到通信中	148
9.2.6 Web 是大规模定制的	149
9.2.7 网络仓库是深度分布式的	149
9.2.8 我们必须面对安全及其隐私问题	150
9.3 小结	151
第 10 章 设计用户体验	152
10.1 两次革命之间的差异	152
10.2 第二代用户界面指导方针	153
10.2.1 确保准瞬态性能	154
10.2.2 满足用户要求	160
10.2.3 让每个页面都成为愉快的体验	165
10.2.4 简单化过程	168
10.2.5 安定用户	170
10.2.6 提供分解问题的手段	171
10.2.7 建立信任	173
10.2.8 提供通信钩子(Communication Hook)	175
10.2.9 支持国际化透明	175

10.3 小结	176
第 11 章 通过网库驱动的数据挖掘	177
11.1 数据挖掘的起源	177
11.2 数据挖掘的行为	178
11.3 数据挖掘的准备工作	179
11.3.1 一般的网库数据转换	180
11.3.2 适用于所有形式的数据挖掘的数据转换	181
11.3.3 依赖于数据挖掘工具的特定的数据转换	182
11.4 将数据提交给数据挖掘工具	184
11.5 OLAP、数据挖掘和网库	187
11.6 小结	188
第 12 章 创建国际 Web 数据仓库	189
12.1 发展中的国际 Web	190
12.1.1 UNICODE	190
12.1.2 并行的超文本和机器翻译	192
12.1.3 多语言搜索	193
12.1.4 时区转换服务	193
12.1.5 节日查找服务	195
12.2 国际网库技术	195
12.2.1 在多个时区和时间格式之间实现同步	195
12.2.2 支持多国日历和日期格式	196
12.2.3 集中多种货币形式的收入	197
12.2.4 处理国际名称和地址	199
12.2.5 支持不同的数字格式	204
12.2.6 支持国际电话号码	204
12.2.7 处理跨国查询、报表和对照序列	204
12.2.8 本地化在 Web 数据仓库中的应用	205
12.3 小结	206
第 13 章 Web 数据仓库安全	207
13.1 推荐的安全技术	208
13.1.1 提供双因素认证	208
13.1.2 保护连接	210
13.1.3 将通过认证的用户与角色联系在一起	212

13.1.4 通过角色访问所有的网库对象	213
13.2 管理安全过程,而不是解决方案	214
13.3 小结	215
第 14 章 网库的缩放	216
14.1 网库不是 Web 服务器	216
14.2 点击流活动突变	217
14.2.1 上网人数增长	218
14.2.2 越来越多的点击率	219
14.2.3 用户级自动搜索	219
14.2.4 更深的经济渗透	220
14.2.5 一夜成名	220
14.2.6 IP 成为一种通用传输协议	220
14.2.7 XML——通用传输	221
14.3 对数据仓库服务需求的剧变	221
14.4 软硬件的严重瓶颈	221
14.4.1 避免单一瓶颈	222
14.4.2 避免进程重复	224
14.4.3 物理上的考虑:托管	224
14.4.4 操作系统	224
14.4.5 编程语言	225
14.4.6 数据库	225
14.4.7 查询和报告软件	226
14.4.8 平衡电子邮件和链接的使用	226
14.4.9 硬件特性	227
14.5 粒度权衡	227
14.6 小结	228
第 15 章 管理网库项目	229
15.1 定义项目	229
15.2 确定角色	230
15.2.1 全体决策人员:主管与监督人	232
15.2.2 教练:项目经理和领导	232
15.2.3 常规阵容:核心项目团队	233
15.3 搜集业务需求和审计数据	236
15.4 计划并管理实现	237

15.5 启动系统	238
15.6 回过头来再做一遍	239
15.7 小结	239
第 16 章 网库的未来	240
16.1 CRM 将继续推动 Web 数据仓库	240
16.2 更好地描述行为	241
16.3 我们最终将需要数据挖掘	242
16.4 ISP 拥有一座金矿	243
16.5 寻求更好的搜索引擎	244
16.6 数据能否战胜存储和速度	245
16.7 数据库的完全转置	246
16.8 网站应用程序日志	246
16.9 每件东西都是一个模块	247
16.10 小结	248

引　　言

就在刚刚过去的一年里,由 Web 产生的影响把 IT 的任务从支持应用继承转变成了发送内容、信息和事务处理能力,而这一切都是通过浏览器界面实现的。或者看起来是这个样子。IT 组织在此之前从来没有受到这么大的压力,要接受多种新的思维方式并重组如此多的界面。

在应对更大的挑战——响应 Web 需求的过程中,千年虫问题只是一个短暂的停顿。Web 不仅仅是用来连接分布式处理设备的一项技术,它还是一种更新更廉价的通信形式。

为变革构建基础

在我们的历史中,有过多次使通信费用骤然降低的变革。在所有这些事件中,当花费明显下降时,通信量就会突然增加,人们通过新的媒介手段变得更具文化性,社会也由此不断改变。通过下面的列表以及图 I.1 展示的内容我们可以知道:

- 16 世纪的印刷术把阅读和出版引入了社会的方方面面
- 19 世纪 40 年代出现的便士邮递(Penny Post)使人们仅花 1 便士就能把一封信寄到英格兰境内的任何地方
- 19 世纪 40、50 年代出现在欧洲和美国的电报是第一种即时长距离通信的形式
- 20 世纪早期出现在美国的乡村免费邮政派送(Rural free postal delivery)把新闻、商务和个人通信带给了遍布在整个国家的许多人
- 20 世纪 20 年代出现的电话使得任何市民都可以与其他市民进行实时交谈
- 20 世纪 30 年代,无线电广播可以在瞬间跨国度传播新闻、文化、语言和言论
- 20 世纪 50 年代,电视延续了由无线电广播所产生的影响,区别在于使用了更引人注目的媒体
- 20 世纪 90 年代的 Web 则是对邮局、电话、无线电广播和电视的一种动力强劲通往高带宽方向的继承发展

20 世纪 90 年代的 Web 革命与早期革命的最大不同可能在于此前者发展的惊人速度。在六年里,这个世界的重要组成部分已经改变了它的通信方式、交易方式,以及使用信息的方式。

这一次,轮到作为 IT 专业人员的我们来为变革构建基础了。我们要使旧应用程序适用于 Web,这样每个人都可以在任何地方使用 Web 进行输入和输出。我们要能够为我们的客

户、生意伙伴以及内部雇员提供适用于 Web 的界面。我们还要能从档案中得到 Web 要求的各种格式的多媒体数据(包括数据库、电子数据表、未组织的文本、图像、地图,甚至可能是音频和视频等等)。与此同时,我们还应该处理由于保护和适当发布所有这些信息而引起的安全机密性问题。

构建这个新的 Web 基础结构既令人兴奋同时也让人毫无头绪。我们应该从哪儿开始? 我们如何才能把问题缩小到恰当的范围,使之易于处理并且使得我们可以针对这一挑战应用已有的系统构建技巧?

Web 数据仓库

过去的 10 年中,在 Web 革命形成气候之前,IT 组织已经学习了如何向内部分析人员和管理人员发布组织的数据资源。发布就是数据仓库的中心任务。

数据仓库是 IT 的核心功能之一。通过多种方式,数据仓库在基于 OLTP(联机事务处理)系统“采集进数据”后,实现“读出数据”。根据在数据仓库方面十多年的经验,我们已经对什么是数据仓库及 IT 如何引入技术以有效发布所有这些数据有了相对成熟的理解。

Web 革命并没有取代人们对数据仓库的需求。事实上,Web 革命把每个人对通过 Web 浏览器界面无缝发布各种信息的期望大大提高了。数据仓库数据的拥护者从内部管理人员扩展到了周边客户、合作伙伴以及更多的内部员工。Web 关注“顾客经验”(customer experience),这使得许多组织进一步意识到要了解顾客并给予顾客有用信息。

由于在许多情况下控制或分析 Web 经验的引擎必定是数据仓库,因此 Web 革命把数据仓库推到了主要舞台上。为了实现这一强化功能,必须对数据仓库进行调整。数据仓库的本质要求和它过去 10 年中的样子有所不同。图 I.2 为网站顾客和网库之间的关系。

我们将把数据仓库的重生称为 Web 数据仓库。

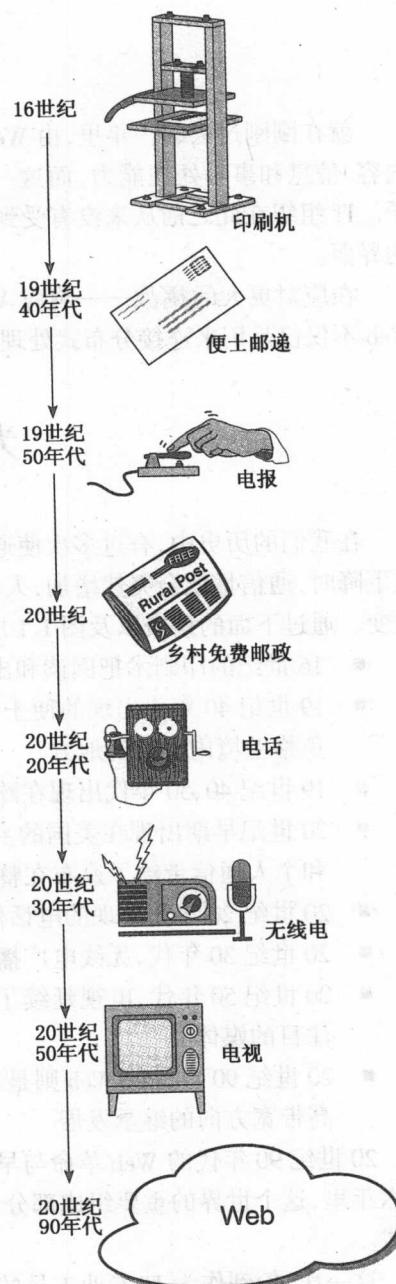


图 I.1 对社会产生深刻影响的媒体

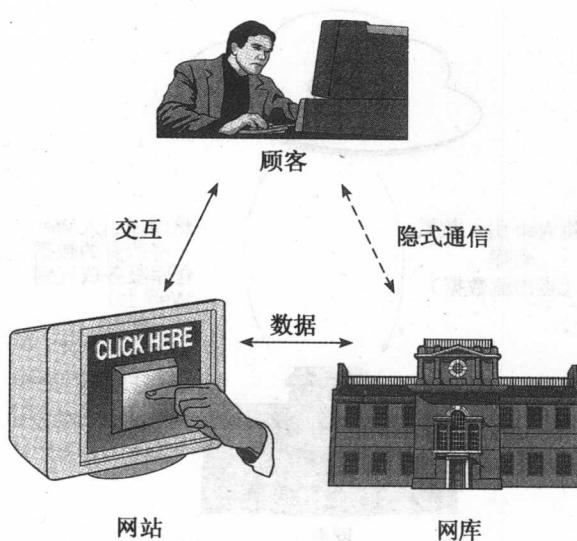


图 I.2 顾客、网站和网库

本书内容

本书是关于 Web 数据仓库的。我们将定义 Web 数据仓库并说明它与传统数据仓库的区别。最重要的是，我们会详细说明如何建立 Web 数据仓库，或者是从零开始，或者更可能的是改写并扩展现有的数据仓库使之成为完全意义上的 Web 数据仓库。

我们将看到 Web 数据仓库有两种特性，这些将体现在本书的两个部分中，请参见图I.3。前半部分为将 Web 引入 Web 数据仓库。当个体通过浏览器与远程网站进行交互时，Web 本身就是一种极好的行为数据源。尽管这种点击流(dickstream)数据在许多情况下是原始且未经修饰的，但是它具有提供使用 Web 媒体的任何人所做每个动作详细细节的潜能。像亚马逊那样的点击流数据形成的原始数据源已经成为现有已知的、最大的文本与数字数据库，这使得从电信到保险这样著名的大型数据库都显得相形见绌。当我们把 Web 引入 Web 数据仓库时，我们把这个巨大的、未经提炼的数据源导入 Web 数据仓库使之自我分析，除此之外，也可以使之与现有的较为传统的数据源保持一致并互相结合。

本书的整个前半部分是关于了解和掌握点击流资源，以便把这些资源存储在 Web 数据仓库中并予以有效地利用。

Web 数据仓库的第二个特性将在本书的第二部分中进行介绍，该特性把现有的数据仓库引入 Web。描述这个过程的最简洁的方式是我们已经不再处于客户端/服务器环境下，我们是处于适用于 Web 的环境中，并且体系结构中有比过去更多的层。

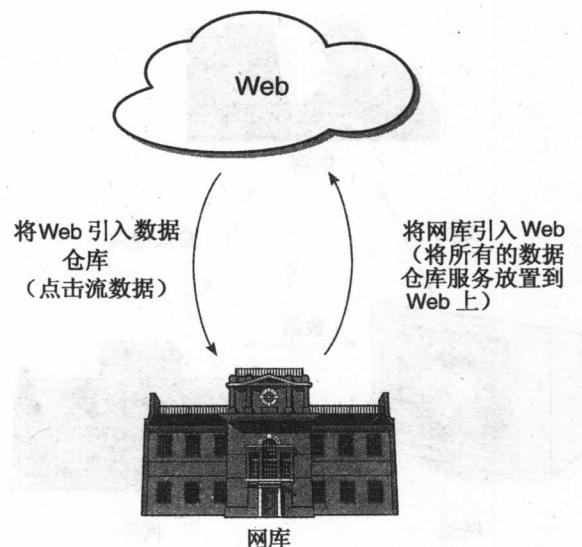


图 I.3 网库的两种特性

把数据仓库引入 Web 意味着要使所有的数据仓库接口通过 Web 浏览器就能使用。如一来,整个过程可分成数据记录、简单特殊查询、复杂报告、正式应用程序开发,以及最后的数据库和系统管理。

把数据仓库引入 Web 也意味着要一次并彻底解决完全分布式环境的问题。Web 数据仓库是比完全集中式数据仓库方法更有意义的替代产品。Web 数据仓库与 Internet 本身一样,也难以实现集中化。但是不管怎样,我们必须设法接受这样一种设计理念:要允许分布在 Web 之上的分散岛式数据仓库以一种有效的方式彼此可见并通信。本书的目的之一就是要说明这种理念的含义。我们将学习如何建立并控制分布式 Web 数据仓库。

本书不对 intranet 网库和 Internet 网库进行更多的区分,这是因为我们认为二者在设计问题、技术和开放程度方面实质上是一样的。我们确信在实际操作中我们无法限制一个组织只使用纯粹的 intranet 方法。几乎从一开始,人们就想通过 Web 连接骨架远程使用他们的浏览器。而这时,人们就会遇到 Web 配置的所有问题。所以我们从一开始就可能需要处理这类较大的问题。

本书适用对象

本书适用于那些像之前介绍过的在工厂组织中担当各种 Web 和数据仓库职责、打算建立部分或整个 Web 数据仓库的设计人员和项目经理。如果你正是这样的专业人员,那么你或你的员工可能需要: