

The Data Webhouse Toolkit

Building the Web-Enabled
Data Warehouse

Ralph Kimball
Richard Merz

WILEY COMPUTER PUBLISHING



John Wiley & Sons, Inc.

New York • Chichester • Weinheim • Brisbane • Singapore • Toronto

Contents

Introduction	1
Building the Infrastructure for the Revolution	1
<i>The Data Webhouse</i>	3
What This Book Is About	4
<i>Who the Book Is For</i>	7
<i>What You Need to Know</i>	8
<i>How to Use This Book</i>	9
<i>The Purpose of Each Chapter</i>	10
<i>Goals of a Data Webhouse</i>	18
<i>Goals of the Book</i>	19
PART ONE BRINGING THE WEB TO THE WAREHOUSE	21
Chapter 1 Why Bring the Web to the Warehouse?	23
Why the Chckstream Is Not Just Another Data Source	24
Analyzing Behavior	26
Ensuring Privacy	29
The Webhouse Architecture	30
<i>The User and the ISP</i>	33
<i>The Public Web Server and Business Transactions</i>	33
<i>The Hot Response Cache</i>	34
<i>The Data Webhouse System</i>	37
Summary	38

Chapter 2	Tracking Website User Actions	41
	A Brief Catalog of User Actions	44
	Steps in Product Purchase	46
	<i>Recognition of Need</i>	46
	<i>Trying to Find What's Needed</i>	46
	<i>Searching for Information about Alternatives</i>	47
	<i>Selection</i>	47
	<i>Cross-Selling and Up-Selling</i>	47
	<i>Checkout</i>	50
	<i>Post-Order Processing</i>	51
	Steps in Software or Content Purchase	51
	<i>Trials and Demos</i>	52
	Elements of Tracking	52
	<i>User Origin</i>	53
	<i>Session Identification</i>	54
	<i>User Identification</i>	56
	Behavioral Analysis	59
	<i>Entry Point</i>	60
	<i>Dwell</i>	60
	<i>Querying</i>	61
	<i>Intra-Site Navigation</i>	61
	<i>Exit Point</i>	63
	Associating Diverse Actions	63
	The Requirements of Personalization	64
	<i>Recognition of Re-visits</i>	64
	<i>User Interface and Content Personalization</i>	65
	<i>Collateral and Impulse Sales</i>	65
	<i>Active Collaborative Filtering</i>	66
	<i>Calendar and Lifestyle Events</i>	66
	<i>Localization</i>	66
	Summary	68
Chapter 3	Using the Clickstream to Make Decisions	69
	Decisions About Identifying and Recognizing Customers	71
	<i>Customizing Marketing Activities by Identifying</i>	
	<i>Your Customers</i>	71
	<i>Targeting Marketing Activities by Clustering Your</i>	
	<i>Customers</i>	73
	<i>Deciding Whether to Encourage or Support a</i>	
	<i>Referring Cross-Link</i>	75
	<i>Deciding Whether a Customer Is About to Leave Us</i>	76

Decisions About Communicating	77
<i>Deciding Whether a Particular Web Ad Is Working</i>	11
<i>Deciding If Custom Greetings Are Working</i>	78
<i>Deciding If a Promotion Is Profitable</i>	80
<i>Responding to a Customer's Life Change</i>	81
<i>Improving the Effectiveness of Your Website</i>	82
<i>Fostering a Sense of Community</i>	83
Fundamental Decisions About Your Web Business	83
<i>Deciding which Products and Services We Provide over the Web</i>	84
<i>Providing Real Time Status Tracking of Our Operations</i>	85
<i>Determining If Our Web Business Is Profitable</i>	87
Summary	89

Chapter 4 Understanding the Clickstream as a Data Source 91

Web Client/Server Interactions—A Brief Tutorial	92
<i>Basic Client I Server Interaction</i>	92
<i>Advertisements</i>	94
<i>The Referrer</i>	94
<i>The Profiler</i>	94
<i>Composite Sites</i>	95
Proxy Servers and Browser Caches	95
<i>Browser Caches</i>	97
Web Server Logs	97
<i>Host</i>	99
<i>Ident</i>	101
<i>Authuser</i>	101
<i>Time</i>	101
<i>Request</i>	101
<i>Status</i>	102
<i>Bytes</i>	102
<i>Referrer</i>	102
<i>User-Agent</i>	102
<i>Filename</i>	104
<i>Time-to-Serve</i>	104
<i>IP Address</i>	104
<i>Server Port</i>	104
<i>Process ID</i>	105
<i>URL</i>	105
Cookies	105
<i>Cookie Contents</i>	107
<i>Cookie Tutorial—Examining Your Own Cookie File</i>	108

Universal System Identifiers	110
Query Strings	110
<i>Templates</i>	111
Summary	112
Chapter 5 Designing the Website to Support Warehousing	113
Monolithic vs. Distributed Web Servers	114
Synchronize Your Servers	115
<i>Time Synchronization Tools and Techniques</i>	116
Content Labels for Pages	119
<i>Content Indexes for Static HTML</i>	120
<i>Content Indexes for Dynamic HTML</i>	121
<i>A Simple Content Index Application</i>	121
Consistent Cookies	122
Null Logging Server	123
Personal Data Repository	126
Building Trust	127
<i>Special Issues in Collecting Information from Children</i>	127
Summary	128
Chapter 6 Building Ciickstream Data Marts	129
A Lightning Tour of Dimensional Modeling	129
<i>Stringing Stars Together</i>	134
Ciickstream Dimensions	139
<i>Calendar Date Dimension</i>	139
<i>Time of Day Dimension</i>	142
<i>Customer Dimension</i>	143
<i>Page Dimension</i>	148
<i>Event Dimension</i>	149
<i>Session Dimension</i>	150
<i>Referral Dimension</i>	151
<i>Product (or Service) Dimension</i>	152
<i>Causal Dimension</i>	154
<i>Business Entity Dimension</i>	155
<i>Ciickstream Tracking Keys</i>	157
The Ciickstream Data Mart	158
<i>A Ciickstream Fact Table to Analyze Complete Sessions</i>	159

	<i>A Ciickstream Fact Table to Analyze Individual Page Aggregate Ciickstream Fact Tables Summary</i>	163 167 168
Chapter 7	Assembling Ciickstream Value Chains	169
	The Sales Transaction Data Mart	170
	The Customer Communication Data Mart	171
	The Web Profitability Data Mart	172
	A Supply Chain for a Web Retailer	176
	A Policies and Claims Chain for Insurance	178
	A Sales Pipeline Chain	180
	A Health Care Value Circle	182
	Summary	185
Chapter 8	Implementing the Ciickstream Post-Processor	187
	Post-Processor Architecture	189
	<i>The Page Event Extractor</i>	191
	<i>The Content Resolver</i>	192
	<i>The Session Identifier</i>	192
	<i>Computing Dwell Time</i>	193
	<i>Host and Referrer Resolver</i>	195
	Summary	197
PART TWO	BRINGING THE WAREHOUSE TO THE WEB	199
Chapter 9	Why Bring the Warehouse to the Web?	201
	The Web Pulls the Data Warehouse	202
	The Web Pushes the Data Warehouse	204
	<i>Tightening the User Interface Feedback Loop</i>	205
	<i>Mixing Query and Update</i>	206
	<i>Speed Is Nonnegotiable</i>	206
	<i>The Sun Never Sets on the Data Webhouse</i>	207
	<i>Multimedia Merges into Communication</i>	208
	<i>The Web Is Mass Customization</i>	209
	<i>The Webhouse Is Profoundly Distributed</i>	210
	<i>We Must Face Security and Its Cousin, Privacy</i>	211
	Summary	212

Chapter 10 Designing the User Experience	215
How the Second Revolution Differs from the First	215
Second Generation User Interface Guidelines	217
<i>Ensure Near-Instantaneous Performance</i>	217
<i>Meet User Expectations</i>	226
<i>Make Each Page a Pleasant Experience</i>	234
<i>Streamline Processes</i>	237
<i>Reassure Users</i>	239
<i>Provide a Means for Resolving Problems</i>	241
<i>Build Trust</i>	243
<i>Provide Communication Hooks</i> "*	246
<i>Support International Transparency</i>	247
Summary	248
Chapter 11 Driving Data Mining from the Webhouse	251
The Roots of Data Mining	252
The Activities of Data Mining	253
Preparing for Data Mining	255
<i>Data Transformations for Webhouses in General</i>	255
<i>Data Transformations for all Forms of Data Mining</i>	256
<i>Special Data Transformations Depending on the</i>	
<i>Data Mining Tool</i>	259
Handing the Data to the Data Miner	261
OLAP, Data Mining, and the Webhouse	265
Summary	266
Chapter 12 Creating an International Data Webhouse	269
The Evolving International Web	270
<i>UNICODE</i>	271
<i>Parallel Hypertext and Machine Translation</i>	273
<i>Multilingual Search</i>	275
<i>Time Zone Converter Services</i>	276
<i>Holiday Lookup Services</i>	277
International Webhouse Techniques	278
<i>Synchronize Multiple Time Zones and Time Formats</i>	278
<i>Support Multiple National Calendars and Date</i>	
<i>Formats</i>	280
<i>Collect Revenue in Multiple Currencies</i>	281
<i>Handle International Names and Addresses</i>	284
<i>Support Variable Number Formats</i>	290

<i>Support International Telephone Numbers</i>	290
<i>Handle Multinational Queries, Reports, and Collating Sequences</i>	290
<i>Apply Localization in the Data Webhouse</i>	292
Summary	292
Chapter 13 Data Webhouse Security	295
Recommended Security Techniques	297
<i>Provide Two-Factor Authentication</i>	297
<i>Secure the Connection</i>	300
<i>Connect the Authenticated User to a Role</i>	302
<i>Access All Webhouse Objects Through the Roles</i>	304
Manage a Security Process, Not a Solution	305
Summary	306
Chapter 14 Scaling the Webhouse	307
The Webhouse Is Not the Web Server	308
Explosive Changes in Clickstream Activity	309
<i>Web-Enabled Population Growth</i>	310
<i>Increasing Click Rates</i>	311
<i>User-Level Auto-Search</i>	312
<i>Deeper Economic Penetration</i>	312
<i>Sudden Fame</i>	312
<i>IP as a Universal Transport Protocol</i>	313
<i>XML—Universal Transfer</i>	313
Explosive Changes in Demand for Data Warehouse Services	314
Critical Bottlenecks in Hardware and Software	314
<i>Avoiding the Single Bottleneck</i>	315
<i>Avoiding Process Duplication</i>	317
<i>Physical Considerations: Co-Location</i>	317
<i>Operating Systems</i>	318
<i>Programming Languages</i>	319
<i>Databases</i>	319
<i>Query and Reporting Software</i>	320
<i>Balance the Use of E-Mail and Links</i>	321
<i>Hardware Characteristics</i>	321
The Granularity Tradeoff	322
Summary	323

Chapter 15 Managing the Webhouse Project	325
Define the Project	326
Identify the Roles	328
<i>Front Office: Sponsors and Drivers</i>	328
<i>Coaches: Project Managers and Leads</i>	330
<i>Regular Lineup: Core Project Team</i>	331
Gather Business Requirements and Audit Data	337
Plan and Manage the Implementation	339
Launch the System	340
Loop Back and Do It Again	341
Summary	341
Chapter 16 The Future of Webhousing	343
CRM Will Continue to Drive Data Webhousing	344
Describing Behavior Better	345
We Will Finally Need Data Mining	346
ISPs Own a Gold Mine	348
Wanted: Better Search Engines	349
Is Data Winning the War over Storage and Speed?	350
Full Inversion of Databases	351
Website Application Logs	351
Everything Is a Module	352
Summary	353
Glossary of Abbreviations and Terms	355
Bibliography	387
Index	391