

Making Everything Easier!™

Big Data

FOR
DUMMIES[®]
A Wiley Brand

Learn to:

- Leverage big data tools and architectures
- Explore how big data can transform your business
- Integrate structured and unstructured data into your big data environment
- Use predictive analytics to make better decisions

Judith Hurwitz
Alan Nugent
Dr. Fern Halper
Marcia Kaufman



Get More and Do More at Dummies.com®



Start with **FREE** Cheat Sheets

Cheat Sheets include

- Checklists
- Charts
- Common Instructions
- And Other Good Stuff!

To access the Cheat Sheet created specifically for this book, go to
www.dummies.com/cheatsheet/bigdata

Get Smart at Dummies.com

Dummies.com makes your life easier with 1,000s of answers on everything from removing wallpaper to using the latest version of Windows.

Check out our

- Videos
- Illustrated Articles
- Step-by-Step Instructions

Plus, each month you can win valuable prizes by entering our Dummies.com sweepstakes. *

Want a weekly dose of Dummies? Sign up for Newsletters on

- Digital Photography
- Microsoft Windows & Office
- Personal Finance & Investing
- Health & Wellness
- Computing, iPods & Cell Phones
- eBay
- Internet
- Food, Home & Garden

Find out "HOW" at Dummies.com

*Sweepstakes not currently available in all countries; visit Dummies.com for official rules.





Big Data

FOR

DUMMIES[®]

A Wiley Brand



**by Judith Hurwitz, Alan Nugent, Dr. Fern Halper,
and Marcia Kaufman**



Big Data For Dummies®

Published by
John Wiley & Sons, Inc.
111 River Street
Hoboken, NJ 07030-5774

www.wiley.com

Copyright © 2013 by John Wiley & Sons, Inc., Hoboken, New Jersey

Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

Trademarks: Wiley, the Wiley logo, For Dummies, the Dummies Man logo, A Reference for the Rest of Us!, The Dummies Way, Dummies Daily, The Fun and Easy Way, Dummies.com, Making Everything Easier, and related trade dress 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.

LIMIT OF LIABILITY/DISCLAIMER OF WARRANTY: THE PUBLISHER AND THE AUTHOR MAKE NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE ACCURACY OR COMPLETENESS OF THE CONTENTS OF THIS WORK AND SPECIFICALLY DISCLAIM ALL WARRANTIES, INCLUDING WITHOUT LIMITATION WARRANTIES OF FITNESS FOR A PARTICULAR PURPOSE. NO WARRANTY MAY BE CREATED OR EXTENDED BY SALES OR PROMOTIONAL MATERIALS. THE ADVICE AND STRATEGIES CONTAINED HEREIN MAY NOT BE SUITABLE FOR EVERY SITUATION. THIS WORK IS SOLD WITH THE UNDERSTANDING THAT THE PUBLISHER IS NOT ENGAGED IN RENDERING LEGAL, ACCOUNTING, OR OTHER PROFESSIONAL SERVICES. IF PROFESSIONAL ASSISTANCE IS REQUIRED, THE SERVICES OF A COMPETENT PROFESSIONAL PERSON SHOULD BE SOUGHT. NEITHER THE PUBLISHER NOR THE AUTHOR SHALL BE LIABLE FOR DAMAGES ARISING HEREFROM. THE FACT THAT AN ORGANIZATION OR WEBSITE IS REFERRED TO IN THIS WORK AS A CITATION AND/OR A POTENTIAL SOURCE OF FURTHER INFORMATION DOES NOT MEAN THAT THE AUTHOR OR THE PUBLISHER ENDORSES THE INFORMATION THE ORGANIZATION OR WEBSITE MAY PROVIDE OR RECOMMENDATIONS IT MAY MAKE. FURTHER, READERS SHOULD BE AWARE THAT INTERNET WEBSITES LISTED IN THIS WORK MAY HAVE CHANGED OR DISAPPEARED BETWEEN WHEN THIS WORK WAS WRITTEN AND WHEN IT IS READ.

For general information on our other products and services, please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993, or fax 317-572-4002.

For technical support, please visit www.wiley.com/techsupport.

Wiley publishes in a variety of print and electronic formats and by print-on-demand. Some material included with standard print versions of this book may not be included in e-books or in print-on-demand. If this book refers to media such as a CD or DVD that is not included in the version you purchased, you may download this material at <http://booksupport.wiley.com>. For more information about Wiley products, visit www.wiley.com.

Library of Congress Control Number: 2013933950

ISBN: 978-1-118-50422-2 (pbk); ISBN 978-1-118-64417-1 (ebk); ISBN 978-1-118-64396-9 (ebk); ISBN 978-1-118-64401-0 (ebk)

Manufactured in the United States of America

10 9 8 7 6 5 4 3 2 1

About the Authors

Judith S. Hurwitz is President and CEO of Hurwitz & Associates, a research and consulting firm focused on emerging technology, including cloud computing, big data, analytics, software development, service management, and security and governance. She is a technology strategist, thought leader, and author. A pioneer in anticipating technology innovation and adoption, she has served as a trusted advisor to many industry leaders over the years. Judith has helped these companies make the transition to a new business model focused on the business value of emerging platforms. She was the founder of Hurwitz Group. She has worked in various corporations, including Apollo Computer and John Hancock. She has written extensively about all aspects of distributed software. In 2011 she authored *Smart or Lucky? How Technology Leaders Turn Chance into Success* (Jossey Bass, 2011). Judith is a co-author on five retail *For Dummies* titles including *Hybrid Cloud For Dummies* (John Wiley & Sons, Inc., 2012), *Cloud Computing For Dummies* (John Wiley & Sons, Inc., 2010), *Service Management For Dummies*, and *Service Oriented Architecture For Dummies, 2nd Edition* (both John Wiley & Sons, Inc., 2009). She is also a co-author on many custom published *For Dummies* titles including *Platform as a Service For Dummies, CloudBees Special Edition* (John Wiley & Sons, Inc., 2012), *Cloud For Dummies, IBM Midsize Company Limited Edition* (John Wiley & Sons, Inc., 2011), *Private Cloud For Dummies, IBM Limited Edition* (2011), and *Information on Demand For Dummies, IBM Limited Edition* (2008) (both John Wiley & Sons, Inc.).

Judith holds BS and MS degrees from Boston University, serves on several advisory boards of emerging companies, and was named a distinguished alumnus of Boston University's College of Arts & Sciences in 2005. She serves on Boston University's Alumni Council. She is also a recipient of the 2005 Massachusetts Technology Leadership Council award.

Alan F. Nugent is a Principal Consultant with Hurwitz & Associates. Al is an experienced technology leader and industry veteran of more than three decades. Most recently, he was the Chief Executive and Chief Technology Officer at Mzinga, Inc., a leader in the development and delivery of cloud-based solutions for big data, real-time analytics, social intelligence, and community management. Prior to Mzinga, he was executive vice president and Chief Technology Officer at CA, Inc. where he was responsible for setting the strategic technology direction for the company. He joined CA as senior vice president and general manager of CA's Enterprise Systems Management (ESM) business unit and managed the product portfolio for infrastructure and data management. Prior to joining CA in April of 2005, Al was senior vice president and CTO of Novell, where he was the innovator behind the company's moves into open source and identity-driven solutions. As consulting CTO for BellSouth he led the corporate initiative to consolidate and transform all of BellSouth's disparate customer and operational data into a single data instance.

Al is the independent member of the Board of Directors of Adaptive Computing in Provo, UT, chairman of the advisory board of SpaceCurve in Seattle, WA, and a member of the advisory board of N-of-one in Waltham, MA. He is a frequent writer on business and technology topics and has shared his thoughts and expertise at many industry events throughout the years.

He is an instrument rated private pilot and has played professional poker for the past three decades. In his sparse spare time he enjoys rebuilding older American muscle cars and motorcycles, collecting antiquarian books, epicurean cooking, and has passion for cellaring American and Italian wines.

Fern Halper, PhD, is a Fellow with Hurwitz & Associates and Director of TDWI Research for Advanced Analytics. She has more than 20 years of experience in data analysis, business analysis, and strategy development. Fern has published numerous articles on data analysis and advanced analytics. She has done extensive research, writing, and speaking on the topic of predictive analytics and text analytics. Fern publishes a regular technology blog. She has held key positions at AT&T Bell Laboratories and Lucent Technologies, where she was responsible for developing innovative data analysis systems as well as developing strategy and product-line plans for Internet businesses. Fern has taught courses in information technology at several universities. She received her BA from Colgate University and her PhD from Texas A&M University.

Fern is a co-author on four retail *For Dummies* titles including *Hybrid Cloud For Dummies* (John Wiley & Sons, Inc., 2012), *Cloud Computing For Dummies* (John Wiley & Sons, Inc., 2010), *Service Oriented Architecture For Dummies*, 2nd Edition, and *Service Management For Dummies* (both John Wiley & Sons, Inc., 2009). She is also a co-author on many custom published *For Dummies* titles including *Cloud For Dummies*, IBM Midsize Company Limited Edition (John Wiley & Sons, Inc., 2011), *Platform as a Service For Dummies*, CloudBees Special Edition (John Wiley & Sons, Inc., 2012), and *Information on Demand For Dummies*, IBM Limited Edition (John Wiley & Sons, Inc., 2008).

Marcia A. Kaufman is a founding Partner and COO of Hurwitz & Associates, a research and consulting firm focused on emerging technology, including cloud computing, big data, analytics, software development, service management, and security and governance. She has written extensively on the business value of virtualization and cloud computing, with an emphasis on evolving cloud infrastructure and business models, data-encryption and end-point security, and online transaction processing in cloud environments. Marcia has more than 20 years of experience in business strategy, industry research, distributed software, software quality, information management, and analytics. Marcia has worked within the financial services, manufacturing, and services industries. During her tenure at Data Resources, Inc. (DRI), she developed sophisticated industry models and forecasts. She holds an AB from Connecticut College in mathematics and economics and an MBA from Boston University.

Marcia is a co-author on five retail *For Dummies* titles including *Hybrid Cloud For Dummies* (John Wiley & Sons, Inc., 2012), *Cloud Computing For Dummies* (John Wiley & Sons, Inc., 2010), *Service Oriented Architecture For Dummies*, 2nd Edition, and *Service Management For Dummies* (both John Wiley & Sons, Inc., 2009). She is also a co-author on many custom published *For Dummies* titles including *Platform as a Service For Dummies*, CloudBees Special Edition (John Wiley & Sons, Inc., 2012), *Cloud For Dummies*, IBM Midsize Company Limited Edition (John Wiley & Sons, Inc., 2011), *Private Cloud For Dummies*, IBM Limited Edition (2011), and *Information on Demand For Dummies* (2008) (both John Wiley & Sons, Inc.).

Dedication

Judith dedicates this book to her husband, Warren, her children, Sara and David, and her mother, Elaine. She also dedicates this book in memory of her father, David.

Alan dedicates this book to his wife Jane for all her love and support; his three children Chris, Jeff, and Greg; and the memory of his parents who started him on this journey.

Fern dedicates this book to her husband, Clay, daughters, Katie and Lindsay, and her sister Adrienne.

Marcia dedicates this book to her husband, Matthew, her children, Sara and Emily, and her parents, Gloria and Larry.

Authors' Acknowledgments

We heartily thank our friends at Wiley, most especially our editor, Nicole Sholly. In addition, we would like to thank our technical editor, Brenda Michelson, for her insightful contributions.

The authors would like to acknowledge the contribution of the following technology industry thought leaders who graciously offered their time to share their technical and business knowledge on a wide range of issues related to hybrid cloud. Their assistance was provided in many ways, including technology briefings, sharing of research, case study examples, and reviewing content. We thank the following people and their organizations for their valuable assistance:

Context Relevant: Forrest Carman

Dell: Matt Walken

Epsilon: Bob Zurek

IBM: Rick Clements, David Corrigan, Phil Francisco, Stephen Gold, Glen Hintze, Jeff Jones, Nancy Kop, Dave Lindquist, Angel Luis Diaz, Bill Mathews, Kim Minor, Tracey Mustacchio, Bob Palmer, Craig Rhinehart, Jan Shauer, Brian Vile, Glen Zimmerman

Kognitio: Michael Hiskey, Steve Millard

Opera Solutions: Jacob Spoelstra

RainStor: Ramon Chen, Deidre Mahon

SAS Institute: Malcom Alexander, Michael Ames

VMware: Chris Keene

Xtremedata: Michael Lamble

Publisher's Acknowledgments

We're proud of this book; please send us your comments at <http://dummies.custhelp.com>. For other comments, please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993, or fax 317-572-4002.

Some of the people who helped bring this book to market include the following:

Acquisitions, Editorial

Senior Project Editor: Nicole Sholly

Project Editor: Dean Miller

Acquisitions Editor: Constance Santisteban

Copy Editor: John Edwards

Technical Editor: Brenda Michelson

Editorial Manager: Kevin Kirschner

Editorial Assistant: Anne Sullivan

Sr. Editorial Assistant: Cherie Case

Cover Photo: © Baris Simsek / iStockphoto

Composition Services

Project Coordinator: Sheree Montgomery

Layout and Graphics: Jennifer Creasey,
Joyce Haughey

Proofreaders: Debbye Butler, Lauren
Mandelbaum

Indexer: Valerie Haynes Perry

Publishing and Editorial for Technology Dummies

Richard Swadley, Vice President and Executive Group Publisher

Andy Cummings, Vice President and Publisher

Mary Bednarek, Executive Acquisitions Director

Mary C. Corder, Editorial Director

Publishing for Consumer Dummies

Kathleen Nebenhaus, Vice President and Executive Publisher

Composition Services

Debbie Stailey, Director of Composition Services

Contents at a Glance

| | |
|---|------------|
| <i>Introduction</i> | 1 |
| <i>Part I: Getting Started with Big Data</i> | 7 |
| Chapter 1: Grasping the Fundamentals of Big Data..... | 9 |
| Chapter 2: Examining Big Data Types | 25 |
| Chapter 3: Old Meets New: Distributed Computing..... | 37 |
| <i>Part II: Technology Foundations for Big Data</i> | 45 |
| Chapter 4: Digging into Big Data Technology Components | 47 |
| Chapter 5: Virtualization and How It Supports Distributed Computing..... | 61 |
| Chapter 6: Examining the Cloud and Big Data | 71 |
| <i>Part III: Big Data Management</i> | 83 |
| Chapter 7: Operational Databases..... | 85 |
| Chapter 8: MapReduce Fundamentals | 101 |
| Chapter 9: Exploring the World of Hadoop | 111 |
| Chapter 10: The Hadoop Foundation and Ecosystem..... | 121 |
| Chapter 11: Appliances and Big Data Warehouses | 129 |
| <i>Part IV: Analytics and Big Data</i> | 139 |
| Chapter 12: Defining Big Data Analytics | 141 |
| Chapter 13: Understanding Text Analytics and Big Data..... | 153 |
| Chapter 14: Customized Approaches for Analysis of Big Data..... | 167 |
| <i>Part V: Big Data Implementation</i> | 179 |
| Chapter 15: Integrating Data Sources..... | 181 |
| Chapter 16: Dealing with Real-Time Data Streams and Complex Event Processing | 193 |
| Chapter 17: Operationalizing Big Data..... | 201 |
| Chapter 18: Applying Big Data within Your Organization | 211 |
| Chapter 19: Security and Governance for Big Data Environments | 225 |

| | |
|--|------------|
| <i>Part VI: Big Data Solutions in the Real World</i> | 235 |
| Chapter 20: The Importance of Big Data to Business | 237 |
| Chapter 21: Analyzing Data in Motion: A Real-World View | 245 |
| Chapter 22: Improving Business Processes with Big Data Analytics: A Real-World View | 255 |
| <i>Part VII: The Part of Tens</i> | 263 |
| Chapter 23: Ten Big Data Best Practices | 265 |
| Chapter 24: Ten Great Big Data Resources | 271 |
| Chapter 25: Ten Big Data Do's and Don'ts..... | 275 |
| <i>Glossary</i> | 279 |
| <i>Index</i> | 295 |

Table of Contents

***Introduction*..... 1**

| | |
|--|---|
| About This Book..... | 2 |
| Foolish Assumptions..... | 2 |
| How This Book Is Organized..... | 3 |
| Part I: Getting Started with Big Data..... | 3 |
| Part II: Technology Foundations for Big Data..... | 3 |
| Part III: Big Data Management..... | 3 |
| Part IV: Analytics and Big Data..... | 4 |
| Part V: Big Data Implementation..... | 4 |
| Part VI: Big Data Solutions in the Real World..... | 4 |
| Part VII: The Part of Tens..... | 4 |
| Glossary..... | 4 |
| Icons Used in This Book..... | 5 |
| Where to Go from Here..... | 5 |

***Part 1: Getting Started with Big Data*..... 7**

Chapter 1: Grasping the Fundamentals of Big Data..... 9

| | |
|--|----|
| The Evolution of Data Management..... | 10 |
| Understanding the Waves of Managing Data..... | 11 |
| Wave 1: Creating manageable data structures..... | 11 |
| Wave 2: Web and content management..... | 13 |
| Wave 3: Managing big data..... | 14 |
| Defining Big Data..... | 15 |
| Building a Successful Big Data Management Architecture..... | 16 |
| Beginning with capture, organize, integrate, analyze, and act..... | 16 |
| Setting the architectural foundation..... | 17 |
| Performance matters..... | 20 |
| Traditional and advanced analytics..... | 22 |
| The Big Data Journey..... | 23 |

Chapter 2: Examining Big Data Types..... 25

| | |
|---|----|
| Defining Structured Data..... | 26 |
| Exploring sources of big structured data..... | 26 |
| Understanding the role of relational databases in big data..... | 27 |
| Defining Unstructured Data..... | 29 |
| Exploring sources of unstructured data..... | 29 |
| Understanding the role of a CMS in big data management..... | 31 |

| | |
|---|----|
| Looking at Real-Time and Non-Real-Time Requirements | 32 |
| Putting Big Data Together | 33 |
| Managing different data types..... | 33 |
| Integrating data types into a big data environment | 34 |

Chapter 3: Old Meets New: Distributed Computing 37

| | |
|---|----|
| A Brief History of Distributed Computing | 37 |
| Giving thanks to DARPA..... | 38 |
| The value of a consistent model | 39 |
| Understanding the Basics of Distributed Computing | 40 |
| Why we need distributed computing for big data..... | 40 |
| The changing economics of computing..... | 40 |
| The problem with latency..... | 41 |
| Demand meets solutions..... | 41 |
| Getting Performance Right | 42 |

***Part II: Technology Foundations for Big Data* 45**

Chapter 4: Digging into Big Data Technology Components 47

| | |
|---|----|
| Exploring the Big Data Stack | 48 |
| Layer 0: Redundant Physical Infrastructure | 49 |
| Physical redundant networks | 51 |
| Managing hardware: Storage and servers | 51 |
| Infrastructure operations | 51 |
| Layer 1: Security Infrastructure..... | 52 |
| Interfaces and Feeds to and from Applications and the Internet..... | 53 |
| Layer 2: Operational Databases..... | 54 |
| Layer 3: Organizing Data Services and Tools..... | 56 |
| Layer 4: Analytical Data Warehouses | 56 |
| Big Data Analytics..... | 58 |
| Big Data Applications..... | 58 |

Chapter 5: Virtualization and How It Supports Distributed Computing 61

| | |
|---|----|
| Understanding the Basics of Virtualization | 61 |
| The importance of virtualization to big data | 63 |
| Server virtualization | 64 |
| Application virtualization | 65 |
| Network virtualization..... | 66 |
| Processor and memory virtualization..... | 66 |
| Data and storage virtualization..... | 67 |
| Managing Virtualization with the Hypervisor | 68 |
| Abstraction and Virtualization | 69 |
| Implementing Virtualization to Work with Big Data | 69 |

Chapter 6: Examining the Cloud and Big Data 71

| | |
|--|----|
| Defining the Cloud in the Context of Big Data | 71 |
| Understanding Cloud Deployment and Delivery Models | 72 |
| Cloud deployment models..... | 73 |
| Cloud delivery models | 74 |
| The Cloud as an Imperative for Big Data..... | 75 |
| Making Use of the Cloud for Big Data | 77 |
| Providers in the Big Data Cloud Market | 78 |
| Amazon's Public Elastic Compute Cloud..... | 78 |
| Google big data services | 79 |
| Microsoft Azure..... | 80 |
| OpenStack..... | 80 |
| Where to be careful when using cloud services | 81 |

Part III: Big Data Management 83**Chapter 7: Operational Databases 85**

| | |
|---|----|
| RDBMSs Are Important in a Big Data Environment..... | 87 |
| PostgreSQL relational database..... | 87 |
| Nonrelational Databases..... | 88 |
| Key-Value Pair Databases | 89 |
| Riak key-value database..... | 90 |
| Document Databases | 91 |
| MongoDB..... | 92 |
| CouchDB | 93 |
| Columnar Databases | 94 |
| HBase columnar database | 94 |
| Graph Databases..... | 95 |
| Neo4J graph database | 96 |
| Spatial Databases..... | 97 |
| PostGIS/OpenGEO Suite..... | 98 |
| Polyglot Persistence..... | 99 |

Chapter 8: MapReduce Fundamentals 101

| | |
|---------------------------------------|-----|
| Tracing the Origins of MapReduce..... | 101 |
| Understanding the map Function..... | 103 |
| Adding the reduce Function..... | 104 |
| Putting map and reduce Together | 105 |
| Optimizing MapReduce Tasks | 108 |
| Hardware/network topology | 108 |
| Synchronization | 108 |
| File system | 108 |

| | |
|--|------------|
| Chapter 9: Exploring the World of Hadoop | 111 |
| Explaining Hadoop | 111 |
| Understanding the Hadoop Distributed File System (HDFS) | 112 |
| NameNodes | 113 |
| Data nodes | 114 |
| Under the covers of HDFS | 115 |
| Hadoop MapReduce | 116 |
| Getting the data ready | 117 |
| Let the mapping begin | 118 |
| Reduce and combine | 118 |
| Chapter 10: The Hadoop Foundation and Ecosystem | 121 |
| Building a Big Data Foundation with the Hadoop Ecosystem | 121 |
| Managing Resources and Applications with Hadoop YARN | 122 |
| Storing Big Data with HBase | 123 |
| Mining Big Data with Hive | 124 |
| Interacting with the Hadoop Ecosystem | 125 |
| Pig and Pig Latin | 125 |
| Sqoop | 126 |
| Zookeeper | 127 |
| Chapter 11: Appliances and Big Data Warehouses | 129 |
| Integrating Big Data with the Traditional Data Warehouse | 129 |
| Optimizing the data warehouse | 130 |
| Differentiating big data structures from data warehouse data ... | 130 |
| Examining a hybrid process case study | 131 |
| Big Data Analysis and the Data Warehouse | 133 |
| The integration lynchpin | 134 |
| Rethinking extraction, transformation, and loading | 134 |
| Changing the Role of the Data Warehouse | 135 |
| Changing Deployment Models in the Big Data Era | 136 |
| The appliance model | 136 |
| The cloud model | 137 |
| Examining the Future of Data Warehouses | 137 |
| Part IV: Analytics and Big Data | 139 |
| Chapter 12: Defining Big Data Analytics | 141 |
| Using Big Data to Get Results | 142 |
| Basic analytics | 142 |
| Advanced analytics | 143 |
| Operationalized analytics | 146 |
| Monetizing analytics | 146 |

Modifying Business Intelligence Products to Handle Big Data..... 147
 Data..... 147
 Analytical algorithms 148
 Infrastructure support 148
 Studying Big Data Analytics Examples..... 149
 Orbitz..... 149
 Nokia..... 150
 NASA..... 150
 Big Data Analytics Solutions 151

Chapter 13: Understanding Text Analytics and Big Data 153

Exploring Unstructured Data 154
 Understanding Text Analytics 155
 The difference between text analytics and search..... 156
 Analysis and Extraction Techniques..... 157
 Understanding the extracted information..... 159
 Taxonomies 160
 Putting Your Results Together with Structured Data 160
 Putting Big Data to Use 161
 Voice of the customer 161
 Social media analytics 162
 Text Analytics Tools for Big Data..... 164
 Attensity..... 164
 Clarabridge 165
 IBM..... 165
 OpenText 165
 SAS 166

Chapter 14: Customized Approaches for Analysis of Big Data 167

Building New Models and Approaches to Support Big Data..... 168
 Characteristics of big data analysis 168
 Understanding Different Approaches to Big Data Analysis 170
 Custom applications for big data analysis 171
 Semi-custom applications for big data analysis..... 173
 Characteristics of a Big Data Analysis Framework 174
 Big to Small: A Big Data Paradox 177

Part V: Big Data Implementation..... 179

Chapter 15: Integrating Data Sources. 181

Identifying the Data You Need 181
 Exploratory stage..... 182
 Codifying stage..... 184
 Integration and incorporation stage 184

| | |
|--|-----|
| Understanding the Fundamentals of Big Data Integration | 186 |
| Defining Traditional ETL..... | 187 |
| Data transformation | 188 |
| Understanding ELT — Extract, Load, and Transform..... | 189 |
| Prioritizing Big Data Quality..... | 189 |
| Using Hadoop as ETL..... | 191 |
| Best Practices for Data Integration in a Big Data World..... | 191 |

Chapter 16: Dealing with Real-Time Data Streams and Complex Event Processing 193

| | |
|---|-----|
| Explaining Streaming Data and Complex Event Processing..... | 194 |
| Using Streaming Data | 194 |
| Data streaming | 195 |
| The need for metadata in streams..... | 196 |
| Using Complex Event Processing | 198 |
| Differentiating CEP from Streams | 199 |
| Understanding the Impact of Streaming Data and CEP on Business | 200 |

Chapter 17: Operationalizing Big Data 201

| | |
|--|-----|
| Making Big Data a Part of Your Operational Process | 201 |
| Integrating big data..... | 202 |
| Incorporating big data into the diagnosis of diseases | 203 |
| Understanding Big Data Workflows | 205 |
| Workload in context to the business problem..... | 206 |
| Ensuring the Validity, Veracity, and Volatility of Big Data..... | 207 |
| Data validity..... | 207 |
| Data volatility | 208 |

Chapter 18: Applying Big Data within Your Organization. 211

| | |
|--|-----|
| Figuring the Economics of Big Data | 212 |
| Identification of data types and sources..... | 212 |
| Business process modifications or new process creation | 215 |
| The technology impact of big data workflows..... | 215 |
| Finding the talent to support big data projects | 216 |
| Calculating the return on investment (ROI) from big data investments..... | 216 |
| Enterprise Data Management and Big Data..... | 217 |
| Defining Enterprise Data Management..... | 217 |
| Creating a Big Data Implementation Road Map..... | 218 |
| Understanding business urgency | 218 |
| Projecting the right amount of capacity | 219 |
| Selecting the right software development methodology..... | 219 |
| Balancing budgets and skill sets | 219 |
| Determining your appetite for risk..... | 220 |
| Starting Your Big Data Road Map | 220 |

Chapter 19: Security and Governance for Big Data Environments . . . 225

| | |
|--|-----|
| Security in Context with Big Data..... | 225 |
| Assessing the risk for the business | 226 |
| Risks lurking inside big data..... | 226 |
| Understanding Data Protection Options | 227 |
| The Data Governance Challenge | 228 |
| Auditing your big data process..... | 230 |
| Identifying the key stakeholders..... | 231 |
| Putting the Right Organizational Structure in Place | 231 |
| Preparing for stewardship and management of risk..... | 232 |
| Setting the right governance and quality policies | 232 |
| Developing a Well-Governed and Secure Big Data Environment | 233 |

Part VI: Big Data Solutions in the Real World..... 235**Chapter 20: The Importance of Big Data to Business 237**

| | |
|---|-----|
| Big Data as a Business Planning Tool | 238 |
| Stage 1: Planning with data..... | 238 |
| Stage 2: Doing the analysis | 239 |
| Stage 3: Checking the results..... | 239 |
| Stage 4: Acting on the plan | 240 |
| Adding New Dimensions to the Planning Cycle..... | 240 |
| Stage 5: Monitoring in real time | 240 |
| Stage 6: Adjusting the impact..... | 241 |
| Stage 7: Enabling experimentation | 241 |
| Keeping Data Analytics in Perspective | 241 |
| Getting Started with the Right Foundation | 242 |
| Getting your big data strategy started | 242 |
| Planning for Big Data..... | 243 |
| Transforming Business Processes with Big Data | 244 |

Chapter 21: Analyzing Data in Motion: A Real-World View. 245

| | |
|---|-----|
| Understanding Companies' Needs for Data in Motion | 246 |
| The value of streaming data..... | 247 |
| Streaming Data with an Environmental Impact | 247 |
| Using sensors to provide real-time information about rivers and oceans | 248 |
| The benefits of real-time data | 249 |
| Streaming Data with a Public Policy Impact | 249 |
| Streaming Data in the Healthcare Industry | 251 |
| Capturing the data stream..... | 251 |

| | |
|--|-----|
| Streaming Data in the Energy Industry | 252 |
| Using streaming data to increase energy efficiency | 252 |
| Using streaming data to advance the production of alternative sources of energy | 252 |
| Connecting Streaming Data to Historical and Other Real-Time Data Sources | 253 |

Chapter 22: Improving Business Processes with Big Data Analytics: A Real-World View 255

| | |
|--|-----|
| Understanding Companies' Needs for Big Data Analytics | 256 |
| Improving the Customer Experience with Text Analytics | 256 |
| The business value to the big data analytics implementation | 257 |
| Using Big Data Analytics to Determine Next Best Action | 257 |
| Preventing Fraud with Big Data Analytics | 260 |
| The Business Benefit of Integrating New Sources of Data | 262 |

***Part VII: The Part of Tens*** 263

Chapter 23: Ten Big Data Best Practices 265

| | |
|--|-----|
| Understand Your Goals | 265 |
| Establish a Road Map | 266 |
| Discover Your Data | 266 |
| Figure Out What Data You Don't Have | 267 |
| Understand the Technology Options | 267 |
| Plan for Security in Context with Big Data | 268 |
| Plan a Data Governance Strategy | 268 |
| Plan for Data Stewardship | 268 |
| Continually Test Your Assumptions | 269 |
| Study Best Practices and Leverage Patterns | 269 |

Chapter 24: Ten Great Big Data Resources 271

| | |
|--|-----|
| Hurwitz & Associates | 271 |
| Standards Organizations | 271 |
| The Open Data Foundation | 272 |
| The Cloud Security Alliance | 272 |
| National Institute of Standards and Technology | 272 |
| Apache Software Foundation | 273 |
| OASIS | 273 |
| Vendor Sites | 273 |
| Online Collaborative Sites | 274 |
| Big Data Conferences | 274 |

Chapter 25: Ten Big Data Do's and Don'ts 275

Do Involve All Business Units in Your Big Data Strategy 275
Do Evaluate All Delivery Models for Big Data 276
Do Think about Your Traditional Data Sources as Part of
Your Big Data Strategy 276
Do Plan for Consistent Metadata 276
Do Distribute Your Data 277
Don't Rely on a Single Approach to Big Data Analytics 277
Don't Go Big Before You Are Ready 277
Don't Overlook the Need to Integrate Data 277
Don't Forget to Manage Data Securely 278
Don't Overlook the Need to Manage the Performance of Your Data 278

***Glossary* 279**

***Index* 295**

Introduction

Welcome to *Big Data For Dummies*. Big data is becoming one of the most important technology trends that has the potential for dramatically changing the way organizations use information to enhance the customer experience and transform their business models. How does a company go about using data to the best advantage? What does it mean to transform massive amounts of data into knowledge? In this book, we provide you with insights into how technology transitions in software, hardware, and delivery models are changing the way that data can be used in new ways.

Big data is not a single market. Rather, it is a combination of data-management technologies that have evolved over time. Big data enables organizations to store, manage, and manipulate vast amounts of data at the right speed and at the right time to gain the right insights. The key to understanding big data is that data has to be managed so that it can meet the business requirement a given solution is designed to support. Most companies are at an early stage with their big data journey. Many companies are experimenting with techniques that allow them to collect massive amounts of data to determine whether hidden patterns exist within that data that might be an early indication of an important change. Some data may indicate that customer buying patterns are changing or that new elements are in the business that need to be addressed before it is too late.

As companies begin to evaluate new types of big data solutions, many new opportunities will unfold. For example, manufacturing companies may be able to monitor data coming from machine sensors to determine how processes need to be modified before a catastrophic event happens. It will be possible for retailers to monitor data in real time to upsell customers related products as they are executing a transaction. Big data solutions can be used in healthcare to determine the cause of an illness and provide a physician with guidance on treatment options.

Big data is not an isolated solution, however. Implementing a big data solution requires that the infrastructure be in place to support the scalability, distribution, and management of that data. Therefore, it is important to put both a business and technical strategy in place to make use of this important technology trend.

For many important reasons, we think that it is important for you to understand big data technologies and know the ways that companies are using emerging technologies such as Hadoop, MapReduce, and new database

engines to transform the value of their data. We wrote this book to provide a perspective on what big data is and how it's changing the way that organizations can leverage more data than was possible in the past. We think that this book will give you the context to make informed decisions.

About This Book

Big data is new to many people, so it requires some investigation and understanding of both the technical and business requirements. Many different people need knowledge about big data. Some of you want to delve into the technical details, while others want to understand the economic implications of making use of big data technologies. Other executives need to know enough to be able to understand how big data can affect business decisions. Implementing a big data environment requires both an architectural and a business approach — and lots of planning.

No matter what your goal is in reading this book, we address the following issues to help you understand big data and the impact it can have on your business:

- ✔ What is the architecture for big data? How can you manage huge volumes of data without causing major disruptions in your data center?
- ✔ When should you integrate the outcome of your big data analysis with your data warehouse?
- ✔ What are the implications of security and governance on the use of big data? How can you keep your company safe?
- ✔ What is the value of different data technologies, and when should you consider them as part of your big data strategy?
- ✔ What types of data sources can you take advantage of with big data analytics? How can you apply different types of analytics to business problems?

Foolish Assumptions

Try as we might to be all things to all people, when it came to writing this book, we had to pick who we thought would be most interested in *Big Data For Dummies*. Here's who we think you are:

- ✔ **You're smart.** You're no dummy, yet the topic of big data gives you an uneasy feeling. You can't quite get your head around it, and if you're pressed for a definition, you might try to change the subject.

- ✔ **You're a businessperson who wants little or nothing to do with technology.** But you live in the 21st century, so you can't escape it. People are saying, "It's all about big data," so you think that you better find out what they're talking about.
- ✔ **You're an IT person who knows a heck of a lot about technology.** The thing is, you're new to big data. Everybody says it's something different. Once and for all, you want the whole picture.

Whoever you are, welcome. We're here to help.

How This Book Is Organized

We divided our book into seven parts for easy reading. Feel free to skip about.

Part I: Getting Started with Big Data

In this part, we explain the basic concepts you need for a full understanding of big data, from both a technical and a business perspective. We also introduce you to the major concepts and components so that you can hold your own in any meaningful conversation about big data.

Part II: Technology Foundations for Big Data

Part II is for both technical and business professionals who need to understand the different types of big data components and the underlying technology concepts that support big data. In this section, we give you an understanding about the type of infrastructure that will make big data more practical.

Part III: Big Data Management

Part III is for both technical and business professionals, but it gets into a lot more of the details of different database options and emerging technologies such as MapReduce and Hadoop. Understanding these underlying technologies can help you understand what is behind this important trend.

Part IV: Analytics and Big Data

How do you analyze the massive amounts of data that become part of your big data infrastructure? In this part of the book, we go deeper into the different types of analytics that are helpful in getting real meaning from your data. This part helps you think about ways that you can turn big data into action for your business.

Part V: Big Data Implementation

This part gets to the details of what it means to actually manage data, including issues such as operationalizing your data and protecting the security and privacy of that data. This section gives you plenty to think about in this critical area.

Part VI: Big Data Solutions in the Real World

In this section, you get an understanding of how companies are beginning to use big data to transform their business operations. If you want to get a peek into the future at what you might be able to do with data, this section is for you.

Part VII: The Part of Tens

If you're new to the *For Dummies* treasure-trove, you're no doubt unfamiliar with The Part of Tens. In this section, Wiley editors torture *For Dummies* authors into creating useful bits of information that are easily accessible in lists containing ten (or so) elucidating elements. We started these chapters kicking and screaming but are ultimately very glad that they're here. After you read through the big data best practices, and the do's and don'ts we provide in The Part of Tens, we think you'll be glad, too.

Glossary

We include a glossary of terms frequently used when people discuss big data. Although we strive to define terms as we introduce them in this book, we think you'll find the glossary a useful resource.

Icons Used in This Book



Pay attention. The bother you save may be your own.



You may be sorry if this little tidbit slips your mind.



With this icon, we mark particularly useful points to pay attention to.



Here you find tidbits for the more technically inclined.

Where to Go from Here

We've created an overview of big data and introduced you to all its significant components. We recommend that you read the first four chapters to give you the context for what big data is about and what technologies are in place to make implementations a reality. The next two chapters introduce you to some of the underlying infrastructure issues that are important to understand. The following eight chapters get into a lot more detail about the different types of data structures that are foundational to big data.

You can read the book from cover to cover, but if you're not that kind of person, we've tried to adhere to the *For Dummies* style of keeping chapters self-contained so that you can go straight to the topics that interest you most. Wherever you start, we wish you well.

Many of these chapters could be expanded into full-length books of their own. Big data and the emerging technology landscape are a big focus for us at Hurwitz & Associates, and we invite you to visit our website and read our blogs and insights at www.hurwitz.com.

Occasionally, John Wiley & Sons, Inc., has updates to its technology books. If this book has technical updates, they will be posted at www.dummies.com/go/bigdatafdupdates.

