



C o m m u n i t y E x p e r i e n c e D i s t i l l e d

Pentaho Analytics for MongoDB

Combine Pentaho Analytics and MongoDB to create powerful
analysis and reporting solutions

Bo Borland

[PACKT] open source*
PUBLISHING community experience distilled

Pentaho Analytics for MongoDB

Combine Pentaho Analytics and MongoDB to create powerful analysis and reporting solutions

Bo Borland



BIRMINGHAM - MUMBAI

Pentaho Analytics for MongoDB

Copyright © 2014 Packt Publishing

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author, nor Packt Publishing, and its dealers and distributors will be held liable for any damages caused or alleged to be caused directly or indirectly by this book.

Packt Publishing has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, Packt Publishing cannot guarantee the accuracy of this information.

First published: February 2014

Production Reference: 1180214

Published by Packt Publishing Ltd.
Livery Place
35 Livery Street
Birmingham B3 2PB, UK.

ISBN 978-1-78216-835-5

www.packtpub.com

Cover Image by Andrey Bayda (<http://shutterstock.com>)

Credits

Author

Bo Borland

Reviewers

Rio Bastian

Pooya Esfandiar

Gretchen Moran

Khaled Tannir

Acquisition Editors

James Jones

Meeta Rajani

Content Development Editor

Govindan K

Technical Editors

Dennis John

Gaurav Thingalaya

Project Coordinator

Aboli Ambardekar

Proofreaders

Simran Bhogal

Maria Gould

Ameesha Green

Indexers

Priya Subramani

Hemangini Bari

Monica Ajmera Mehta

Rekha Nair

Graphics

Abhinash Sahu

Yuvraj Mannari

Production Coordinators

Conidon Miranda

Kyle Albuquerque

Cover Work

Kyle Albuquerque

About the Author

Bo Borland is the vice president of field technical sales at Pentaho, a leading Big Data analytics software provider. He has a passion for building teams and helping companies improve performance with data analytics. Prior to joining Pentaho, Bo worked as a management consultant at Deloitte and as a solution architect at Cognos, an IBM company. He founded a successful analytics consulting company, Management Signals, which merged with Pentaho in 2012. His 14 years' experience in professional analytics includes roles in management, sales, consulting, and sales engineering.

Pentaho Corporation is a leading Big Data analytics company headquartered in Orlando, FL, with offices in San Francisco, London, and Portugal. Pentaho tightly couples data integration with full business analytics capabilities into a single, integrated software platform. The company's subscription-based software offers SMEs and global enterprises the ability to reduce the time taken to design, develop, and deploy Big Data analytics solutions.

I wish to thank the wonderful team at Pentaho for their support and encouragement of this project, especially Rebecca Shomair and Monie TenBroeck for proactively reaching out to offer a helping hand, and Gretchen Moran for her valuable assistance as a technical reviewer for this book.

I also want to thank my wife, Alison, and daughters, Lin and Greta, for graciously accepting the fact that many nights and weekends were spent writing this book.

About the Reviewers

Rio Bastian is a happy software developer already working on several IT projects. His interests include business intelligence, data integration, tuning SQL, and tuning Java code. He has also been a Pentaho Business Intelligence trainer for a server company in Indonesia and Malaysia. He is currently focused on the development of an airline customer loyalty program in PT. Aero Systems Indonesia, an IT consultant specializing in airline industries. In his spare time, he tries to share his experience in developing software through his personal blog altanovela.wordpress.com. You can reach him on Skype via `rio.bastian` or via e-mail at altanovela@gmail.com.

Pooya Esfandiar is a software engineer and data analyst. He received an MS in Computer Science from the University of British Columbia in 2010 while working at the data management and mining labs. His research interests include simulation, data mining, and machine learning in social networks. He has worked with a few companies and organizations in Iran and Canada to solve business problems in information systems, general software engineering, and data analysis. He once designed an in-house analytics solution using Pentaho and MongoDB and is now working as a software engineer at Amazon.com, Inc.

Gretchen Moran is working as an independent Pentaho consultant on a variety of business analytics and Big Data projects. She has 15 years' experience in the business intelligence industry, developing software and providing services for a number of companies, including Hyperion Solutions and Pentaho Corporation.

Gretchen continues to contribute to Pentaho Corporation's latest software initiatives while managing the daily adventures of her two children, Isabella and Jack, with her husband Doug.

Khaled Tannir has been working with computers since 1980. He began programming with the legendary Sinclair Zx81 and afterwards with all Commodore home computer products (VIC-20, Commodore 64, Commodore 128D, and Amiga 500).

He has a Bachelor's degree in Electronics, a Master's degree in System Information Architectures in which he graduated with a professional thesis, and he completed his education with a Master of Research degree.

He is a Microsoft Certified Solution Developer (MCSD) and has more than 20 years' technical experience leading the development and implementation of software solutions and giving technical presentations. He works as an independent IT consultant and has worked as an infrastructure engineer, senior developer, and enterprise/solution architect for many companies in France and Canada.

With a significant experience in Microsoft .NET/Servers and Oracle Java technologies, he has extensive skills in online/offline application design, system conversions, and multilanguage applications on both the Internet and desktops.

He spends his time researching on new technologies, learning about them, and looking for new adventures in France, North America, and the Middle East. He owns an IT and electronics laboratory with many servers, monitors, open electronics boards (such as Arduino, Netduino, Raspberry Pi, and .NET Gadgeteer), and some smartphone devices based on Windows Phone, Android, and iOS operating systems.

In 2012, he contributed to the EGC 2012 (International Complex Data Mining forum at Bordeaux University, France) and presented his work on *how to optimize data distribution in a cloud computing environment* in a workshop session. This work aims to define an approach to optimize the use of data mining algorithms such as k-means and Apriori in a cloud computing environment.

He is the author of *RavenDB 2.x Beginner's Guide*, Packt Publishing and *Optimizing Hadoop MapReduce*, Packt Publishing.

He aims to get a PhD degree in Cloud Computing and Big Data and wants to learn more and more about these technologies.

He enjoys taking landscape and night photos, travelling, playing video games, creating funny electronics gadgets with Arduino / .NET Gadgeteer and of course, spending time with his wife and family. You can reach him at contact@khaledtannir.net.

www.PacktPub.com

Support files, eBooks, discount offers and more

You might want to visit www.PacktPub.com for support files and downloads related to your book.

Did you know that Packt offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at www.PacktPub.com and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at service@packtpub.com for more details.

At www.PacktPub.com, you can also read a collection of free technical articles, sign up for a range of free newsletters and receive exclusive discounts and offers on Packt books and eBooks.



<http://PacktLib.PacktPub.com>

Do you need instant solutions to your IT questions? PacktLib is Packt's online digital book library. Here, you can access, read and search across Packt's entire library of books.

Why Subscribe?

- Fully searchable across every book published by Packt
- Copy and paste, print and bookmark content
- On demand and accessible via web browser

Free Access for Packt account holders

If you have an account with Packt at www.PacktPub.com, you can use this to access PacktLib today and view nine entirely free books. Simply use your login credentials for immediate access.

Table of Contents

Preface	1
Chapter 1: Getting Started with Pentaho and MongoDB	5
MongoDB technology overview	6
Pentaho technology overview	9
Installing MongoDB	10
Installing MongoDB as a Windows service	12
Restoring the sample clickstream MongoDB database	13
Installing Pentaho	14
Summary	15
Chapter 2: MongoDB Database Fundamentals	17
MongoDB database objects	17
Sample clickstream database objects	19
MongoDB data modeling	21
Normalized models	21
Denormalized models	22
MongoDB query methods	24
Query exercise 1	24
Read operations	25
Query exercise 2	26
Query operators	26
Querying arrays	27
Summary	29
Chapter 3: Using Pentaho Instaview	31
Accessing and connecting Instaview to MongoDB	31
Parsing and profiling a MongoDB collection	33
Adding a MongoDB query expression	35
Creating and saving an analysis view and Instaview	38
Summary	42

Chapter 4: Modifying and Enhancing Instaview Transformations	43
Opening an existing Instaview	44
Data integration	45
Adding a new data source	45
CSV file input	47
Stream lookup	48
Creating a new analysis view from blended data	50
Summary	52
Chapter 5: Modifying and Enhancing Instaview Metadata	53
Model design with dimensions and measures	53
Open an existing Instaview	54
Modifying measures and dimensions	55
Session duration measure	56
Session count measure	57
Event count measure	57
Referring URL dimension	57
Other dimension changes	58
Creating a new analysis view	60
Summary	61
Chapter 6: Pentaho Report Designer Fundamentals	63
Pentaho Report Designer features	63
Data sources	64
Report elements	64
Aggregations and calculations	65
Formatting and output	65
Navigating through Pentaho Report Designer	65
Report workspace	66
The Structure tab	68
The Data tab	68
The Style and Attributes tabs	68
The palette	69
The main menu and toolbar	71
The tab toolbar	71
Interface reference	72
Creating a MongoDB connection and query	73
Adding a MongoDB data source	74
Adding and formatting report elements	76
Adding a message field to your report	76
Adding number-fields to your report	78

Adding calculated values to your report	79
Summary	81
Chapter 7: Pentaho Report Designer Prompting and Charting	83
Adding additional MongoDB queries	83
Adding a bar chart query	84
Adding a pie chart query	85
Visualizing your data with charts	86
JFreeChart chart types	87
Subreports	87
Chart data collectors and properties	88
Creating a bar chart	89
Modifying bar chart properties	91
Creating a pie chart	92
Creating a report prompt	95
Creating a new parameter	95
Adding parameters to existing report queries	97
Creating subreport import parameters	98
Summary	100
Chapter 8: Deploying Pentaho Analytics to the Web	101
Publishing a Report Designer report to the Web	102
Publishing the clickstream report	102
An introduction to the Pentaho User Console	104
Running and scheduling the clickstream report	106
Enabling your Instaview output for the Web	108
Copying and modifying the Instaview transformation	109
Using the Data Source Wizard to model your data	112
Creating a JDBC connection and default metadata model	113
Customizing the metadata model	114
Creating Analyzer Views and Dashboard Designer dashboards	118
Creating a map view in Analyzer	118
Creating a heat grid in Analyzer	120
Creating a dashboard using Dashboard Designer	121
Summary	123
Index	125

Preface

MongoDB and Pentaho go together like yin and yang. They are emerging as a powerful combination for scalable data storage, processing, and analytics. Leading companies are pairing these complementary technologies together in development labs and production to deliver innovative analytics. These innovations are creating worldwide demand for developers with skills in both Pentaho and MongoDB.

You want to make an impact by creating innovative data storage capabilities or eye-catching data visualizations. Wouldn't it be great if you could quickly ramp up on both technologies to develop a turn-key solution for your organization? However, as with any new and emerging technology combination, the availability of organized knowledge on the combined topic is scarce.

Pentaho Analytics for MongoDB will show you how to develop an analytic solution that you can demonstrate to your colleagues. It is a practical guide to get you started with both Pentaho and MongoDB, beginning with basic MongoDB data modeling and querying and then advancing to data integration, analysis, and reporting with Pentaho. Each chapter guides you through using different components of the Pentaho platform to create analytic models and reports using a sample MongoDB database.

What this book covers

Chapter 1, Getting Started with Pentaho and MongoDB, introduces you to the powerful combination of MongoDB and Pentaho and provides step-by-step guidance on how to install and configure both technologies and restore the sample MongoDB data provided with this book.

Chapter 2, MongoDB Database Fundamentals, expands on the topic of data modeling and explains MongoDB database concepts essential to querying MongoDB data with Pentaho.

Chapter 3, Using Pentaho Instaview, shows you how to visualize data by connecting Pentaho to MongoDB. You use Instaview with the sample MongoDB database to analyze and visualize the website clickstream data.

Chapter 4, Modifying and Enhancing Instaview Transformations, introduces Pentaho Data Integration (PDI) – the ETL tool used by Instaview to extract, load, and transform data from various data sources.

Chapter 5, Modifying and Enhancing Instaview Metadata, explores metadata by explaining dimensional modeling concepts and how to model metadata to better reflect business requirements.

Chapter 6, Pentaho Report Designer Fundamentals, teaches you the basics of Pentaho Report Designer (PRD) to build pixel-perfect reports sourced directly from MongoDB databases.

Chapter 7, Pentaho Report Designer Prompting and Charting, expands on the previous chapter by teaching you additional advanced PRD features. You can enhance your report with new queries, charts, and a prompt designed to make the report more interactive.

Chapter 8, Deploying Pentaho Analytics to the Web, is all about web-enabling your MongoDB data using Pentaho methods and web interfaces for connecting to, modeling, and analyzing our sample clickstream data in a web browser.

What you need for this book

We need the following software for this book:

- Pentaho Business Analytics v5.0.2 (64-bit for Windows)
- MongoDB v2.2.3 (64-bit for Windows)

This book provides two data sources for use throughout the book, a MongoDB database of sample web clickstream data, and an associated comma-separated (CSV) file containing geographic data. Both files are available as a free download from: <http://www.packtpub.com/support>.

Who this book is for

This book is intended for business analysts, data architects, and developers new to either Pentaho or MongoDB, who want to be able to deliver a complete solution for storing, processing, and visualizing data. It's assumed that you already have experience in defining the data requirements needed to support business processes and exposure to database modeling, SQL query, and reporting techniques.