# Pentaho for Big Data Analytics

Enhance your knowledge of Big Data and leverage the power of Pentaho to extract its treasures

Manoj R Patil        Feris Thia

# Pentaho for Big Data Analytics

Enhance your knowledge of Big Data and leverage the power of Pentaho to extract its treasures

**Manoj R Patil**

**Feris Thia**

# Pentaho for Big Data Analytics

# Credits

**Authors**
Manoj R Patil

Feris Thia

**Reviewers**
Rio Bastian

Paritosh H. Chandorkar

Vikram Takkar

**Acquisition Editors**
Kartikey Pandey

Rebecca Youe

**Commissioning Editor**
Mohammed Fahad

**Technical Editors**
Manan Badani

Pankaj Kadam

**Copy Editors**
Alisha Aranha

Sarang Chari

Brandt D'Mello

Tanvi Gaitonde

Dipti Kapadia

Laxmi Subramanian

**Project Coordinator**
Sageer Parkar

**Proofreaders**
Ameesha Green

Maria Gould

**Indexer**
Rekha Nair

**Graphics**
Sheetal Aute

Ronak Dhruv

Disha Haria

Abhinash Sahu

**Production Coordinator**
Arvindkumar Gupta

**Cover Work**
Arvindkumar Gupta

# About the Authors

**Manoj R Patil** is the Chief Architect in Big Data at Compassites Software Solutions Pvt. Ltd. where he overlooks the overall platform architecture related to Big Data solutions, and he also has a hands-on contribution to some assignments. He has been working in the IT industry for the last 15 years. He started as a programmer and, on the way, acquired skills in architecting and designing solutions, managing projects keeping each stakeholder's interest in mind, and deploying and maintaining the solution on a cloud infrastructure. He has been working on the Pentaho-related stack for the last 5 years, providing solutions while working with employers and as a freelancer as well.

Manoj has extensive experience in JavaEE, MySQL, various frameworks, and Business Intelligence, and is keen to pursue his interest in predictive analysis.

He was also associated with TalentBeat, Inc. and Persistent Systems, and implemented interesting solutions in logistics, data masking, and data-intensive life sciences.

**Feris Thia** is a founder of PHI-Integration, a Jakarta-based IT consulting company that focuses on data management, data warehousing and Business Intelligence solutions. As a technical consultant, he has spent the last seven years delivering solutions with Pentaho and the Microsoft Business Intelligence platform across various industries, including retail, trading, finance/banking, and telecommunication.

He is also a member and maintainer of two very active local Indonesian discussion groups related to Pentaho (`pentaho-id@googlegroups.com`) and Microsoft Excel (the BelajarExcel.info Facebook group).

His current activities include research and building software based on Big Data and the data mining platform, that is, Apache Hadoop, R, and Mahout.

He would like to work on a book with a topic on analyzing customer behavior using the Apache Mahout platform.

# About the Reviewers

**Rio Bastian** is a happy software developer already working on several IT projects. He is interested in Data Integration, and tuning SQL and Java code. He has also been a Pentaho Business Intelligence trainer for several companies in Indonesia and Malaysia. Rio is currently working as a software developer in PT. Aero Systems Indonesia, a company that focuses on the development of airline customer loyalty programs. It's an IT consultant company specializing in the airline industry. 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 (`rio.bastian`) or e-mail him at `altanovela@gmail.com`.

**Paritosh H. Chandorkar** is a young and dynamic IT professional with more than 11 years of information technology management experience in diverse domains, such as telecom and banking.

He has both strong technical (in Java/JEE) and project management skills. He has expertise in handling large customer engagements. Furthermore, he has expertise in the design and development of very critical projects for clients such as BNP Paribas, Zon TVCabo, and Novell. He is an impressive communicator with strong leadership, coordination, relationship management, analytical, and team management skills. He is comfortable interacting with people across hierarchical levels for ensuring smooth project execution as per client specifications. He is always eager to invest in improving knowledge and skills.

He is currently studying at Manipal University for a full-time M.S. in Software Design and Engineering.

His last designation was Technology Architect at Infosys Ltd.

I would like to thank Manoj R Patil for giving me the opportunity to review this book.

**Vikram Takkar** is a freelance Business Intelligence and Data Integration professional with nine years of rich, hands-on experience in multiple BI and ETL tools. He has strong expertise in tools such as Talend, Jaspersoft, Pentaho, Big Data-MongoDB, Oracle, and MySQL. He has managed and successfully executed multiple projects in data warehousing and data migration developed for both UNIX and Windows environments.

Apart from this, he is a blogger and publishes articles and videos on open source BI and ETL tools along with supporting technologies. You can visit his blog at `www.vikramtakkar.com`.

His YouTube channel is `www.youtube.com/vtakkar`. His Twitter handle is `@VikTakkar`. You can also follow him on his blog at `www.vikramtakkar.com`.

# 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