



# 大数据分析：R语言实现

（影印版）

Big Data Analytics with R

Simon Walkowiak 著

[PACKT]  
PUBLISHING

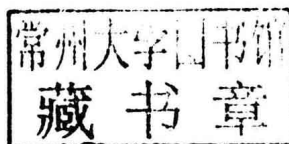


东南大学出版社  
SOUTHEAST UNIVERSITY PRESS

# 大数据分析：R语言实现(影印版)

## Big Data Analytics With R

Simon Walkowiak 著



南京 东南大学出版社

## 图书在版编目(CIP)数据

大数据分析:R 语言实现:英文/(英)西蒙·沃克威克(Simon Walkowiak)著. —影印本. —南京:东南大学出版社, 2017.10

书名原文: Big Data Analytics With R

ISBN 978-7-5641-7361-6

I. ①大… II. ①西… III. ①程序语言—程序设计—英文 IV. ①TP312

中国版本图书馆 CIP 数据核字(2017)第 192628 号

图字:10-2017-115 号

© 2016 by PACKT Publishing Ltd

Reprint of the English Edition, jointly published by PACKT Publishing Ltd and Southeast University Press, 2017.  
Authorized reprint of the original English edition, 2017 PACKT Publishing Ltd, the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

英文原版由 PACKT Publishing Ltd 出版 2016。

英文影印版由东南大学出版社出版 2017。此影印版的出版和销售得到出版权和销售权的所有者——PACKT Publishing Ltd 的许可。

版权所有,未得书面许可,本书的任何部分和全部不得以任何形式重制。

## 大数据分析:R 语言实现(影印版)

出版发行:东南大学出版社

地 址:南京四牌楼 2 号 邮编:210096

出 版 人:江建中

网 址: <http://www.seupress.com>

电子邮件: [press@seupress.com](mailto:press@seupress.com)

印 刷:常州市武进第三印刷有限公司

开 本:787 毫米×980 毫米 16 开本

印 张:31.5

字 数:617 千字

版 次:2017 年 10 月第 1 版

印 次:2017 年 10 月第 1 次印刷

书 号:ISBN 978-7-5641-7361-6

定 价:94.00 元

# Credits

**Authors**

Simon Walkowiak

**Reviewer**

Zacharias Voulgaris

Dipanjan Sarkar

**Copy Editor**

Safis Editing

**Project Coordinator**

Ulhas Kambali

**Commissioning Editor**

Akram Hussain

**Proofreader**

Safis Editing

**Acquisition Editor**

Sonali Vernekar

**Indexer**

Tejal Daruwale Soni

**Content Development Editor**

Onkar Wani

**Graphics**

Kirk D'Penha

**Technical Editor**

Sushant S Nadkar

**Production Coordinator**

Arvindkumar Gupta

# About the Author

**Simon Walkowiak** is a cognitive neuroscientist and a managing director of Mind Project Ltd – a Big Data and Predictive Analytics consultancy based in London, United Kingdom. As a former data curator at the UK Data Service (UKDS, University of Essex) – European largest socio-economic data repository, Simon has an extensive experience in processing and managing large-scale datasets such as censuses, sensor and smart meter data, telecommunication data and well-known governmental and social surveys such as the British Social Attitudes survey, Labour Force surveys, Understanding Society, National Travel survey, and many other socio-economic datasets collected and deposited by Eurostat, World Bank, Office for National Statistics, Department of Transport, NatCen and International Energy Agency, to mention just a few. Simon has delivered numerous data science and R training courses at public institutions and international companies. He has also taught a course in *Big Data Methods in R* at major UK universities and at the prestigious Big Data and Analytics Summer School organized by the Institute of Analytics and Data Science (IADS).

# Acknowledgement

The inspiration for writing this book came directly from the brilliant work and dedication of many R developers and users, whom I would like to thank first for creating a vibrant and highly-supportive community that nourishes the progress of publicly accessible data analytics and development of R language. However, this book would never be completed if I wasn't surrounded with love and unconditional support from my partner Ignacio, who always knew how to encourage and motivate me, particularly in moments of my weakness and when I lacked creativity.

I would also like to thank other members of my family, especially my father Peter, who despite not sharing my excitement of data science, always listens patiently to my stories about emerging Big Data technologies and their use cases.

Also, I dedicate this book to my friends and former colleagues from UK Data Service at the University of Essex, where I had an opportunity to work with amazing individuals and experience the best practices in robust data management and processing.

Finally, I highly appreciate the hard work, expertise and feedback offered by many people involved in the creation of this book at Packt Publishing – especially my content development editor Onkar Wani, publishers, and the reviewers, who kindly shared their knowledge with me in order to create a quality and well-received publication.

# About the Reviewers

**Dr. Zacharias Voulgaris** was born in Athens, Greece. He studied Production Engineering and Management at the Technical University of Crete, shifted to Computer Science through a Masters in Information Systems & Technology (City University, London), and then to Data Science through a PhD on Machine Learning (University of London). He has worked at Georgia Tech as a Research Fellow, at an e-marketing startup in Cyprus as an SEO manager, and as a Data Scientist in both Elavon (GA) and G2 (WA). He also was a Program Manager at Microsoft, on a data analytics pipeline for Bing.

Zacharias has authored two books and several scientific articles on Machine Learning and as well as a couple of articles on AI topics. His first book, *Data Scientist - The Definitive Guide to Becoming a Data Scientist* (Technics Publications), has been translated into Korean and Chinese, while his latest one, *Julia for Data Science* (Technics Publications) is coming out this September. He has also reviewed a number of data science books (mainly on Python and R) and has a passion for new technologies, literature, and music.

I'd like to thank the people at Packt for inviting me to review this book and for promoting Data Science and particularly Julia through their books. Also, a big thanks to all the great authors out there who choose to publish their work through the lesser-known publishers, keeping the whole process of sharing knowledge a democratic endeavor.

**Dipanjan Sarkar** is a Data Scientist at Intel, the world's largest silicon company which is on a mission to make the world more connected and productive. He primarily works on analytics, business intelligence, application development and building large scale intelligent systems. He received his Master's degree in Information Technology from the International Institute of Information Technology, Bangalore. His area of specialization includes software engineering, data science, machine learning and text analytics.

Dipanjan's interests include learning about new technology, disruptive start-ups, data science and more recently deep learning. In his spare time he loves reading, writing, gaming and watching popular sitcoms. He has authored a book on Machine Learning titled *R Machine Learning by Example*, Packt Publishing and also acted as a technical reviewer for several books on Machine Learning and Data Science from Packt Publishing.

# www.PacktPub.com

## eBooks, discount offers, and more

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](http://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 [customercare@packtpub.com](mailto:customercare@packtpub.com) for more details.

At [www.PacktPub.com](http://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.



<https://www2.packtpub.com/books/subscription/packtlib>

Do you need instant solutions to your IT questions? PacktLib is Packt's online digital book library. Here, you can search, access, and read 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 a web browser



# Table of Contents

<b>Preface</b>	1
<b>Chapter 1: The Era of Big Data</b>	7
<b>Big Data – The monster re-defined</b>	7
<b>Big Data toolbox – dealing with the giant</b>	11
Hadoop – the elephant in the room	12
Databases	15
Hadoop Spark-ed up	16
<b>R – The unsung Big Data hero</b>	17
<b>Summary</b>	24
<b>Chapter 2: Introduction to R Programming Language and Statistical Environment</b>	25
<b>Learning R</b>	25
<b>Revisiting R basics</b>	28
Getting R and RStudio ready	28
Setting the URLs to R repositories	30
<b>R data structures</b>	32
Vectors	32
Scalars	35
Matrices	35
Arrays	37
Data frames	38
Lists	41
Exporting R data objects	42
<b>Applied data science with R</b>	47
Importing data from different formats	48
Exploratory Data Analysis	50
Data aggregations and contingency tables	53
Hypothesis testing and statistical inference	56
Tests of differences	57
Independent t-test example (with power and effect size estimates)	57
ANOVA example	60
Tests of relationships	63
An example of Pearson's r correlations	63
Multiple regression example	65
Data visualization packages	70

<b>Summary</b>	71
<b>Chapter 3: Unleashing the Power of R from Within</b>	73
<b>Traditional limitations of R</b>	74
Out-of-memory data	74
Processing speed	75
<b>To the memory limits and beyond</b>	76
Data transformations and aggregations with the ff and ffbase packages	76
Generalized linear models with the ff and ffbase packages	87
Logistic regression example with ffbase and biglm	89
Expanding memory with the bigmemory package	97
<b>Parallel R</b>	106
From bigmemory to faster computations	107
An apply() example with the big.matrix object	108
A for() loop example with the ffdi object	108
Using apply() and for() loop examples on a data.frame	109
A parallel package example	110
A foreach package example	113
The future of parallel processing in R	115
Utilizing Graphics Processing Units with R	115
Multi-threading with Microsoft R Open distribution	117
Parallel machine learning with H2O and R	118
<b>Boosting R performance with the data.table package and other tools</b>	118
Fast data import and manipulation with the data.table package	118
Data import with data.table	119
Lightning-fast subsets and aggregations on data.table	120
Chaining, more complex aggregations, and pivot tables with data.table	123
Writing better R code	126
<b>Summary</b>	127
<b>Chapter 4: Hadoop and MapReduce Framework for R</b>	129
<b>Hadoop architecture</b>	130
Hadoop Distributed File System	130
MapReduce framework	131
A simple MapReduce word count example	132
Other Hadoop native tools	134
Learning Hadoop	136
<b>A single-node Hadoop in Cloud</b>	137
Deploying Hortonworks Sandbox on Azure	138
A word count example in Hadoop using Java	159
A word count example in Hadoop using the R language	169
RStudio Server on a Linux RedHat/CentOS virtual machine	169

Installing and configuring RHadoop packages	177
HDFS management and MapReduce in R – a word count example	179
<b>HDInsight – a multi-node Hadoop cluster on Azure</b>	194
Creating your first HDInsight cluster	194
Creating a new Resource Group	195
Deploying a Virtual Network	197
Creating a Network Security Group	200
Setting up and configuring an HDInsight cluster	203
Starting the cluster and exploring Ambari	211
Connecting to the HDInsight cluster and installing RStudio Server	215
Adding a new inbound security rule for port 8787	218
Editing the Virtual Network's public IP address for the head node	221
Smart energy meter readings analysis example – using R on HDInsight cluster	229
<b>Summary</b>	241
<b>Chapter 5: R with Relational Database Management Systems (RDBMSs)</b>	243
<b>Relational Database Management Systems (RDBMSs)</b>	244
A short overview of used RDBMSs	244
Structured Query Language (SQL)	245
<b>SQLite with R</b>	247
Preparing and importing data into a local SQLite database	248
Connecting to SQLite from RStudio	250
<b>MariaDB with R on a Amazon EC2 instance</b>	255
Preparing the EC2 instance and RStudio Server for use	255
Preparing MariaDB and data for use	257
Working with MariaDB from RStudio	266
<b>PostgreSQL with R on Amazon RDS</b>	281
Launching an Amazon RDS database instance	281
Preparing and uploading data to Amazon RDS	290
Remotely querying PostgreSQL on Amazon RDS from RStudio	304
<b>Summary</b>	314
<b>Chapter 6: R with Non-Relational (NoSQL) Databases</b>	315
<b>Introduction to NoSQL databases</b>	315
Review of leading non-relational databases	316
<b>MongoDB with R</b>	319
Introduction to MongoDB	319
MongoDB data models	319
Installing MongoDB with R on Amazon EC2	322

Processing Big Data using MongoDB with R	325
Importing data into MongoDB and basic MongoDB commands	326
MongoDB with R using the rmongodb package	333
MongoDB with R using the RMongo package	346
MongoDB with R using the mongolite package	350
<b>HBase with R</b>	355
Azure HDInsight with HBase and RStudio Server	355
Importing the data to HDFS and HBase	363
Reading and querying HBase using the rhbase package	367
<b>Summary</b>	372
<b>Chapter 7: Faster than Hadoop - Spark with R</b>	373
<b>Spark for Big Data analytics</b>	374
<b>Spark with R on a multi-node HDInsight cluster</b>	375
Launching HDInsight with Spark and R/RStudio	375
Reading the data into HDFS and Hive	383
Getting the data into HDFS	385
Importing data from HDFS to Hive	386
Bay Area Bike Share analysis using SparkR	393
<b>Summary</b>	411
<b>Chapter 8: Machine Learning Methods for Big Data in R</b>	413
<b>What is machine learning?</b>	414
Supervised and unsupervised machine learning methods	415
Classification and clustering algorithms	416
Machine learning methods with R	417
Big Data machine learning tools	418
<b>GLM example with Spark and R on the HDInsight cluster</b>	419
Preparing the Spark cluster and reading the data from HDFS	419
Logistic regression in Spark with R	425
<b>Naive Bayes with H2O on Hadoop with R</b>	437
Running an H2O instance on Hadoop with R	437
Reading and exploring the data in H2O	441
Naive Bayes on H2O with R	446
<b>Neural Networks with H2O on Hadoop with R</b>	458
How do Neural Networks work?	458
Running Deep Learning models on H2O	461
<b>Summary</b>	469
<b>Chapter 9: The Future of R - Big, Fast, and Smart Data</b>	471
The current state of Big Data analytics with R	471

Out-of-memory data on a single machine	471
Faster data processing with R	473
Hadoop with R	475
Spark with R	476
R with databases	477
Machine learning with R	478
<b>The future of R</b>	478
Big Data	479
Fast data	480
Smart data	481
<b>Where to go next</b>	482
<b>Summary</b>	482
<b>Index</b>	483

---

# Preface

We live in times of Internet of Things—a large, world-wide network of interconnected devices, sensors, applications, environments, and interfaces. They generate, exchange, and consume massive amounts of data on a daily basis, and the ability to harness these huge quantities of information can provide us with novel understanding of physical and social phenomena.

The recent rapid growth of various open source and proprietary big data technologies allows deep exploration of these vast amounts of data. However, many of them are limited in terms of their statistical and data analytics capabilities. Some others implement techniques and programming languages that many classically educated statisticians and data analysts are simply unfamiliar with and find them difficult to apply in real-world scenarios.

R programming language—an open source, free, extremely versatile statistical environment, has a potential to fill this gap by providing users with a large variety of highly optimized data processing methods, aggregations, statistical tests, and machine learning algorithms with a relatively user-friendly and easily customizable syntax.

This book challenges traditional preconceptions about R as a programming language that does not support big data processing and analytics. Throughout the chapters of this book, you will be exposed to a variety of core R functions and a large array of actively maintained third-party packages that enable R users to benefit from most recent cutting-edge big data technologies and frameworks, such as Hadoop, Spark, H2O, traditional SQL-based databases, such as SQLite, MariaDB, and PostgreSQL, and more flexible NoSQL databases, such as MongoDB or HBase, to mention just a few. By following the exercises and tutorials contained within this book, you will experience firsthand how all these tools can be integrated with R throughout all the stages of the Big Data Product Cycle, from data import and data management to advanced analytics and predictive modeling.

## What this book covers

Chapter 1, *The Era of "Big Data"*, gently introduces the concept of Big Data, the growing landscape of large-scale analytics tools, and the origins of R programming language and the statistical environment.

Chapter 2, *Introduction to R Programming Language and Statistical Environment*, explains the most essential data management and processing functions available to R users. This chapter also guides you through various methods of Exploratory Data Analysis and hypothesis testing in R, for instance, correlations, tests of differences, ANOVAs, and Generalized Linear Models.

Chapter 3, *Unleashing the Power of R From Within*, explores possibilities of using R language for large-scale analytics and out-of-memory data on a single machine. It presents a number of third-party packages and core R methods to address traditional limitations of Big Data processing in R.

Chapter 4, *Hadoop and MapReduce Framework for R*, explains how to create a cloud-hosted virtual machine with Hadoop and to integrate its HDFS and MapReduce frameworks with R programming language. In the second part of the chapter, you will be able to carry out a large-scale analysis of electricity meter data on a multinode Hadoop cluster directly from the R console.

Chapter 5, *R with Relational Database Management Systems (RDBMSs)*, guides you through the process of setting up and deploying traditional SQL databases, for example, SQLite, PostgreSQL and MariaDB/MySQL, which can be easily integrated with their current R-based data analytics workflows. The chapter also provides detailed information on how to build and benefit from a highly scalable Amazon Relational Database Service instance and query its records directly from R.

Chapter 6, *R with Non-Relational (NoSQL) Databases*, builds on the skills acquired in the previous chapters and allows you to connect R with two popular nonrelational databases a.) a fast and user-friendly MongoDB installed on a Linux-run virtual machine, and b.) HBase database operated on a Hadoop cluster run as part of the Azure HDInsight service.

Chapter 7, *Faster than Hadoop: Spark with R*, presents a practical example and a detailed explanation of R integration with the Apache Spark framework for faster Big Data manipulation and analysis. Additionally, the chapter shows how to use Hive database as a data source for Spark on a multinode cluster with Hadoop and Spark installed.

Chapter 8, *Machine Learning Methods for Big Data in R*, takes you on a journey through the most cutting-edge predictive analytics available in R. Firstly, you will perform fast and highly optimized Generalized Linear Models using Spark MLlib library on a multinode Spark HDInsight cluster. In the second part of the chapter, you will implement Naïve Bayes and multilayered Neural Network algorithms using R's connectivity with H2O—an award-winning, open source, big data distributed machine learning platform.

Chapter 9, *The Future of R: Big, Fast and Smart Data*, wraps up the contents of the earlier chapters by discussing potential areas of development for R language and its opportunities in the landscape of emerging Big Data tools.

*Online Chapter, Pushing R Further*, available at [https://www.packtpub.com/sites/default/files/downloads/5396\\_6457OS\\_PushingRFurther.pdf](https://www.packtpub.com/sites/default/files/downloads/5396_6457OS_PushingRFurther.pdf), enables you to configure and deploy their own scaled-up and Cloud-based virtual machine with fully operational R and RStudio Server installed and ready to use.

## What you need for this book

All the code snippets presented in the book have been tested on a Mac OS X (Yosemite) running on a personal computer equipped with 2.3 GHz Intel Core i5 processor, 1 TB Solid State hard drive, and 16 GB of RAM. It is recommended that readers run the scripts on a Mac OS X or Windows machine with at least 4 GB of RAM. In order to benefit from the instructions presented throughout the book, it is advisable that readers install most recent R and RStudio on their machines as well as at least one of the popular web browsers: Mozilla Firefox, Chrome, Safari, or Internet Explorer.

## Who this book is for

This book is intended for middle level data analysts, data engineers, statisticians, researchers, and data scientists, who consider and plan to integrate their current or future big data analytics workflows with R programming language.

It is also assumed that readers will have some previous experience in data analysis and the understanding of data management and algorithmic processing of large quantities of data. However, they may lack specific R skills related to particular open source big data tools.

## Conventions

In this book, you will find a number of text styles that distinguish between different kinds of information. Here are some examples of these styles and an explanation of their meaning.

Code words in text, database table names, folder names, filenames, file extensions, pathnames, dummy URLs, user input, and Twitter handles are shown as follows: "The `getmerge` option allows to merge all data files from a specified directory on HDFS."



Any command-line input or output is written as follows:

```
$ sudo -u hdfs hadoop fs -ls /user
```

New terms and important words are shown in **bold**. Words that you see on the screen, for example, in menus or dialog boxes, appear in the text like this: "Clicking the **Next** button moves you to the next screen."



Warnings or important notes appear in a box like this.



Tips and tricks appear like this.

## Reader feedback

Feedback from our readers is always welcome. Let us know what you think about this book—what you liked or disliked. Reader feedback is important for us as it helps us develop titles that you will really get the most out of.

To send us general feedback, simply e-mail [feedback@packtpub.com](mailto:feedback@packtpub.com), and mention the book's title in the subject of your message.

If there is a topic that you have expertise in and you are interested in either writing or contributing to a book, see our author guide at [www.packtpub.com/authors](http://www.packtpub.com/authors).

## Customer support

Now that you are the proud owner of a Packt book, we have a number of things to help you to get the most from your purchase.

## Downloading the example code

You can download the example code files for this book from your account at <http://www.packtpub.com>. If you purchased this book elsewhere, you can visit <http://www.packtpub.com/support> and register to have the files e-mailed directly to you.