

Financial Analytics with R

Building a Laptop Laboratory for Data Science

Are you innately curious about dynamically inter-operating financial markets? Since the crisis of 2008, there is a need for professionals with more understanding about statistics and data analysis, who can discuss the various risk metrics, particularly those involving extreme events.

By providing a resource for training students and professionals in basic and sophisticated analytics, this book meets that need. It offers both the intuition and basic vocabulary as a step toward the financial, statistical, and algorithmic knowledge required to resolve the industry problems, and it depicts a systematic way of developing analytical programs for finance in the statistical language R. Build a hands-on laboratory and run many simulations. Explore the analytical fringes of investments and risk management.

Bennett and Hugen help profit-seeking investors and data science students sharpen their skills in many areas, including time-series, forecasting, portfolio selection, covariance clustering, prediction, and derivative securities.

Mark J. Bennett is a senior data scientist with a major investment bank and a lecturer in the University of Chicago's Master's program in Analytics. He has held software positions at Argonne National Laboratory, Unisys Corporation, AT&T Bell Laboratories, Northrop Grumman, and XR Trading Securities.

Dirk L. Hugen is a graduate student in the Department of Statistics and Actuarial Science at the University of Iowa. He previously worked as a signal processing engineer.

Financial Analytics with R

Building a Laptop Laboratory for Data Science

MARK J. BENNETT

University of Chicago

DIRK L. HUGEN

University of Iowa



CAMBRIDGE
UNIVERSITY PRESS

CAMBRIDGE
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781107150751

© Mark J. Bennett and Dirk L. Hugen 2016

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2016

Printed in the United Kingdom by Clays, St Ives plc

A catalogue record for this publication is available from the British Library.

Library of Congress Cataloging-in-Publication Data

Names: Bennett, Mark J. (Mark Joseph), 1959– author. | Hugen, Dirk L., author.

Title: Financial analytics with R : building a laptop laboratory for data science / Mark J. Bennett, University of Chicago, Dirk L. Hugen, University of Iowa.

Description: Cambridge, UK : Cambridge University Press, 2016.

Identifiers: LCCN 2016026635 | ISBN 9781107150751

Subjects: LCSH: Finance—Mathematical models—Data processing. | Finance—Databases. | R (Computer program language)

Classification: LCC HG104 .B46 2016 | DDC 332.0285/513--dc23

LC record available at <https://lccn.loc.gov/2016026635>

ISBN 978-1-107-15075-1 Hardback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party Internet Web sites referred to in this publication and does not guarantee that any content on such Web sites is, or will remain, accurate or appropriate.

**To our parents:
Mary and Herb and Patricia and Bernard
and family:
Rachel, Austin, and Cheryl
for all their kindness, love, and support**

Contents

	<i>Preface</i>	page xiii
	<i>Acknowledgments</i>	xvii
1	Analytical Thinking	1
	1.1 What Is Financial Analytics?	2
	1.2 What Is the Laptop Laboratory for Data Science?	3
	1.3 What Is R and How Can It Be Used in the Professional Analytics World?	5
	1.4 Exercises	6
2	The R Language for Statistical Computing	7
	2.1 Getting Started with R	7
	2.2 Language Features: Functions, Assignment, Arguments, and Types	10
	2.3 Language Features: Binding and Arrays	13
	2.4 Error Handling	17
	2.5 Numeric, Statistical, and Character Functions	18
	2.6 Data Frames and Input–Output	19
	2.7 Lists	20
	2.8 Exercises	22
3	Financial Statistics	23
	3.1 Probability	23
	3.2 Combinatorics	24
	3.3 Mathematical Expectation	31
	3.4 Sample Mean, Standard Deviation, and Variance	35
	3.5 Sample Skewness and Kurtosis	36
	3.6 Sample Covariance and Correlation	36
	3.7 Financial Returns	39
	3.8 Capital Asset Pricing Model	40
	3.9 Exercises	42
4	Financial Securities	44
	4.1 Bond Investments	45
	4.2 Stock Investments	48

4.3	The Housing Crisis	49
4.4	The Euro Crisis	50
4.5	Securities Datasets and Visualization	52
4.6	Adjusting for Stock Splits	55
4.7	Adjusting for Mergers	61
4.8	Plotting Multiple Series	62
4.9	Securities Data Importing	64
4.10	Securities Data Cleansing	71
4.11	Securities Quoting	74
4.12	Exercises	75
5	Dataset Analytics and Risk Measurement	77
5.1	Generating Prices from Log Returns	77
5.2	Normal Mixture Models of Price Movements	80
5.3	Sudden Currency Price Movement in 2015	86
5.4	Exercises	90
6	Time Series Analysis	92
6.1	Examining Time Series	92
6.2	Stationary Time Series	97
6.3	Auto-Regressive Moving Average Processes	98
6.4	Power Transformations	98
6.5	The TSA Package	99
6.6	Auto-Regressive Integrated Moving Average Processes	109
6.7	Case Study: Earnings of Johnson & Johnson	110
6.8	Case Study: Monthly Airline Passengers	114
6.9	Case Study: Electricity Production	117
6.10	Generalized Auto-Regressive Conditional Heteroskedasticity	120
6.11	Case Study: Volatility of Google Stock Returns	121
6.12	Exercises	128
7	The Sharpe Ratio	130
7.1	Sharpe Ratio Formula	131
7.2	Time Periods and Annualizing	131
7.3	Ranking Investment Candidates	132
7.4	The Quantmod Package	136
7.5	Measuring Income Statement Growth	141
7.6	Sharpe Ratios for Income Statement Growth	144
7.7	Exercises	155
8	Markowitz Mean-Variance Optimization	157
8.1	Optimal Portfolio of Two Risky Assets	157
8.2	Quadratic Programming	160
8.3	Data Mining with Portfolio Optimization	162

8.4	Constraints, Penalization, and the Lasso	165
8.5	Extending to High Dimensions	171
8.6	Case Study: Surviving Stocks of the S&P 500 Index from 2003 to 2008	179
8.7	Case Study: Thousands of Candidate Stocks from 2008 to 2014	182
8.8	Case Study: Exchange-Traded Funds	186
8.9	Exercises	195
9	Cluster Analysis	197
9.1	K-Means Clustering	197
9.2	Dissecting the K-Means Algorithm	204
9.3	Sparsity and Connectedness of Undirected Graphs	208
9.4	Covariance and Precision Matrices	211
9.5	Visualizing Covariance	215
9.6	The Wishart Distribution	221
9.7	Glasso: Penalization for Undirected Graphs	225
9.8	Running the Glasso Algorithm	225
9.9	Tracking a Value Stock through the Years	226
9.10	Regression on Yearly Sparsity	231
9.11	Regression on Quarterly Sparsity	235
9.12	Regression on Monthly Sparsity	236
9.13	Architecture and Extension	238
9.14	Exercises	239
10	Gauging the Market Sentiment	240
10.1	Markov Regime Switching Model	241
10.2	Reading the Market Data	244
10.3	Bayesian Reasoning	247
10.4	The Beta Distribution	250
10.5	Prior and Posterior Distributions	250
10.6	Examining Log Returns for Correlation	253
10.7	Momentum Graphs	255
10.8	Exercises	259
11	Simulating Trading Strategies	261
11.1	Foreign Exchange Markets	261
11.2	Chart Analytics	263
11.3	Initialization and Finalization	264
11.4	Momentum Indicators	265
11.5	Bayesian Reasoning within Positions	266
11.6	Entries	268
11.7	Exits	269
11.8	Profitability	270
11.9	Short-Term Volatility	270
11.10	The State Machine	271

x	Contents	
	11.11 Simulation Summary	278
	11.12 Exercises	280
12	Data Exploration Using Fundamentals	281
	12.1 The RSQLite Package	281
	12.2 Finding Market-to-Book Ratios	283
	12.3 The Reshape2 Package	285
	12.4 Case Study: Google	288
	12.5 Case Study: Walmart	289
	12.6 Value Investing	290
	12.7 Lab: Trying to Beat the Market	294
	12.8 Lab: Financial Strength	295
	12.9 Exercises	296
13	Prediction Using Fundamentals	297
	13.1 Best Income Statement Portfolio	298
	13.2 Reformatting Income Statement Growth Figures	298
	13.3 Obtaining Price Statistics	300
	13.4 Combining the Income Statement with Price Statistics	306
	13.5 Prediction Using Classification Trees and Recursive Partitioning	308
	13.6 Comparing Prediction Rates among Classifiers	314
	13.7 Exercises	316
14	Binomial Model for Options	318
	14.1 Applying Computational Finance	318
	14.2 Risk-Neutral Pricing and No Arbitrage	322
	14.3 High Risk-Free Rate Environment	322
	14.4 Convergence of Binomial Model for Option Data	324
	14.5 Put–Call Parity	327
	14.6 From Binomial to Log-Normal	328
	14.7 Exercises	330
15	Black–Scholes Model and Option-Implied Volatility	331
	15.1 Geometric Brownian Motion	332
	15.2 Monte Carlo Simulation of Geometric Brownian Motion	333
	15.3 Black–Scholes Derivation	335
	15.4 Algorithm for Implied Volatility	338
	15.5 Implementation of Implied Volatility	339
	15.6 The Rcpp Package	345
	15.7 Exercises	348
Appendix	Probability Distributions and Statistical Analysis	350
	A.1 Distributions	350
	A.2 Bernoulli Distribution	350

A.3	Binomial Distribution	351
A.4	Geometric Distribution	352
A.5	Poisson Distribution	352
A.6	Functions for Continuous Distributions	354
A.7	The Uniform Distribution	356
A.8	Exponential Distribution	357
A.9	Normal Distribution	359
A.10	Log-Normal Distribution	359
A.11	The t_v Distribution	360
A.12	Multivariate Normal Distribution	361
A.13	Gamma Distribution	361
A.14	Estimation via Maximum Likelihood	362
A.15	Central Limit Theorem	364
A.16	Confidence Intervals	366
A.17	Hypothesis Testing	366
A.18	Regression	367
A.19	Model Selection Criteria	369
A.20	Required Packages	370
	<i>References</i>	372
	<i>Index</i>	376

Preface

In 1994 the Channel Tunnel opened between England and France, allowing high-speed Eurostar trains to whisk passengers from the continent to the United Kingdom and back on a grand scale. What an amazing engineering feat it was for the time (beyond many people's earlier imaginations), yet we take it for granted today. In 1994, Grumman Aerospace Corporation, the chief contractor on the Apollo Lunar Module, was acquired by Northrop Corporation to form the new aerospace giant, Northrop Grumman. It was the prime contractor of the newly deployed advanced technology B-2 Stealth Bomber. On a much more mundane and personal scale, also in 1994, in a townhouse just outside the City of Chicago, I was performing a tedious daily exercise: looking up daily closing prices each evening in a stack of *Investor's Business Daily* newspapers for the two stock investments that were about to be purchased. This was not only to find out their running rate of return but also to find out their historical volatility relative to other stocks before entering into the positions. Doing this manual calculation was slow and tedious. The World Wide Web was introduced in the form of the Mosaic browser the next year. It was not long before Yahoo! was posting stock quotes and historical price charts, as well as technical indicators on the charts, available on demand for free in just a few seconds via the new web browsers.

The advent of spreadsheet software took analysts to a new level of analytical thinking. No longer were live, human-operated calculations limited to a single dimension. Each row or column could present a time dimension, a production category, a business scenario. And the automated dependency feature made revisions quite easy. Now spreadsheets can be used for a prototype for a more sophisticated and permanent analytical product: the large-scale, analytical computer program.

With modern programming languages like R and Python[®], a skilled analyst can now design their analytic logic with significantly less effort than before, using resources such as Yahoo! or other free services for historical quotes. It has been said that Python's terse syntax allows for programs with the same functionality as their Java equivalents, yet four times smaller, and we suspect that R is similar. A small financial laboratory can be built on a laptop costing less than \$200 in a matter of weeks, simulating multiple market variables as required. Or, by obtaining a higher-end laptop model with more drive space, the entire market can be loaded with 10 to 20 years' worth of historical data on a scale never before possible.

Once that laboratory is built, one can start to gain insights. Knowledge discovery was once a term for a human process. Now we're talking about computer automation.

Knowledge discovery seems like a bold term, a little too ambitious for anything a computer program could create. For example, the computer science professional society, the Association for Computing Machinery (ACM), has a special interest group called Knowledge Discovery and Data Mining (KDD). Hardly anyone would challenge the “data mining” part of that. After all, statisticians and computer scientists having at it with data is what they do. But discovering knowledge with a machine? Really? Automatically? Now that seems a little too exaggerated to be true. Then again, experiencing firsthand the algorithms that will be described in this book, we soon realized that the programs, using data science techniques, can not only automate very tedious calculations but then very positively yield insights into the human thinking level: insights that would otherwise not be found.

Perhaps one can view the experience with a sports analogy. In many sports we have defense to protect our current position and prevent the opponents from scoring further points. Offense is the ability to piece together athletic feats sequentially to put more points on the scoreboard. The data mining portion of KDD can be thought of as defense: the more disciplined, regimented side of the sport. Single achievements can be effective: thrusting up one’s hand to block a pass, throwing a curve ball to prevent a hitter from connecting on the pitch. On the other hand, knowledge discovery is the offensive skill set, going beyond the required and expected data analysis and into the creative side. On offense, an entire series of events needs to be successful to yield progress: a full-field soccer scoring drive or a series of three successive baseball base hits before three outs to score a run. The likelihood of a success on offense is less.

So in the KDD model and sports analogy, data mining is the defense and knowledge discovery is the offense. Achieving knowledge discovery is a rare event with amazing impact. The discovery can be as powerful as human-made ideas, and can certainly enhance them. For example, we may discover that there is a publicly traded stock with uniquely desirable properties. The KDD domain touches the limits of what these machines can do with all the advancements in computer science.

In 1968, a Hollywood movie and novel by author Arthur C. Clarke, *2001: A Space Odyssey*, predicted automated reasoning, natural language speech recognition, video calls, and face recognition. The HAL9000 computer controls the flight to Jupiter while conversing and playing chess with astronaut Dr. Frank Poole and monitoring life conditions of over 300 astronauts. Since then, computer science, specifically simulation, has greatly impacted the research and discovery process in many fields and effectively achieved many of these science fiction goals now. Among others, there are fields of computational biology, computational cosmology, and computational linguistics. Images from these fields are shown in Figure 1.

Throughout this book we are concerned with computer simulation. Computer simulation has become so successful that it is now widely accepted that, after theory and physical experimentation, it is a third scientific method. As the subtitle says, this book can be used to build a simulation laboratory for finance. The book was developed as study material for the graduate Financial Analytics course in the Graham School at the University of Chicago Master of Science in Analytics program and from the undergraduate Investments course in the Department of Finance in the Tippie College of Business

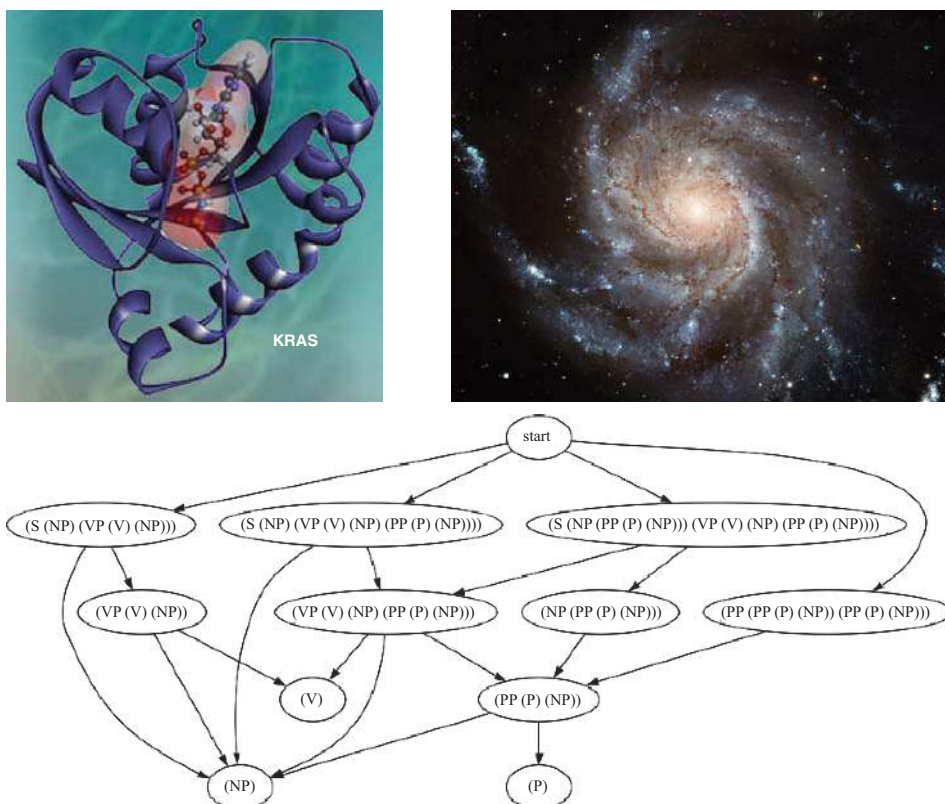


Figure 1 Sample images from computational biology, computational cosmology, and computational linguistics.

at the University of Iowa. It is recommended as a graduate textbook when used at a college or university. With the proper mathematical and computer science background, it could also be used at the advanced undergraduate level.

It is best to have taken a course in statistical analysis, probability and statistics, or, ideally, mathematical statistics for the material in the book, but much of the required material is introduced within the main text and the Appendix. It is best to have an undergraduate-level background in calculus, linear algebra, and enough computer science to be familiar with array manipulation with one or more scientifically based programming languages such as C, C++, Java, C#, Python, or Matlab[®]. A finance background is not necessary. Any experience with R is, of course, useful.

Financial computer simulation in the R language can be more intricate and challenging than building a spreadsheet. A quantitative optimizer can be better controlled and tailored when its logic is immediately apparent from the surrounding program code. More computer science knowledge is required by our reader to build more robust and sophisticated platforms, and more goes into the compiler and run-time system behind the scenes. But as the pieces are completed, the builder, or operator, or student of financial analytics begins to realize the benefits of simulation performed in a language designed

for statistical simulation. The insights that can be gained from building simulators and from observing the simulations will help deepen understanding for upcoming professional venues. Just as it is now for machines, for people it has always been about learning.

Regarding the exercises at the end of each chapter, data science involves the study of statistical and computational models. In this book, that means that we are unlocking the economic value which exists in the financial markets. Data engineering is the process of implementing models on computers as applied to large datasets, using files, program logic, testing, and continuous improvement. With these exercises, we take advantage of the data science principles of the prior chapters to build and engineer our financial laboratory.

As the reader performs these exercises, they may need to install various R packages from time to time. Various pages found by internet searching will steer the reader to proper instructions for loading R packages and troubleshooting any failed attempts. There are too many packages, conditions, and cases to repeat those instructions here.

The exercises focus individually on the various components so that we obtain an understanding of the logic and data. Each new component builds upon prior components in order to provide the level of sophistication required to answer our financial analytics inquiries.

Acknowledgments

Through our careers we are influenced greatly by people who are special to us. We can never pay them back for the gifts they have bestowed upon us: conversation, insight, opportunity, knowledge, and companionship.

Mark would like to thank Ron Krupp and Barry Finkel at Argonne; Phil Ridinger, Pat Baldwin, Howard Seckler, David Rouse, Bill Neidfelt, and Tom Bishop at Bell Labs; Christopher Marlin and Adrienne Critcher at the University of Iowa; Per Brinch-Hansen and Orna Berry at the University of Southern California; Dave Martin, Milos Ercegovac, Stott Parker, and Dan Berry at the University of California, Los Angeles; Bob MacGregor at System Development; Mark Christensen, Paul Schmitz, Tim Bancroft, Clare Morgan, Margaret Lakins, and Joe Dvorak at Northrop; Shelly Reis, Greg Brim, Brian Ostrow, Ron Netzel, David Joffe, Mack Amin, Thayer Allison, Dilip Nair, Gregg Berger, Haider Sajjad, Yuri Salnikov, Jim Bohmbach, Dante Lomibao, Paul Lee, Li Chen, Steve Zhu, Mario Konrad, Brian Philpott, Nancy Goldberg, Laura Lang, Raja Afandi, Ashish Batra, Krishna Bhamidi, and Chris Leakeas at Nationsbanc-CRT; Harry Georgakopoulos and YeeMan Bergstrom at XR Trading; and Adam Ginensky, Yuri Balasanov, and Sema Barlas at the University of Chicago and Marc Tempkin at the Chicago ACM for being tremendous mentors and colleagues over the years. Inspiration for this effort was provided at these wonderful locales: the bohemian Wolverine Farm book store in Fort Collins, the stately Pentacrest at the University of Iowa campus, and the Crerar Library of the gothic University of Chicago campus.

Dirk would like to thank his advisor Joseph Lang for all his help as well as Luke Tierney, Kate Cowles, Joyee Ghosh, Dale Zimmerman, Rhonda DeCook, Kung-Sik Chan, Erik Lie, David Bates, Ashish Tiwari, Wei Li, and Paul Weller for the excellent statistics and finance coursework at the University of Iowa. Thanks also to Kung-Sik Chan, Brian Ripley, Dirk Eddelbuettel, Hadley Wickham, David A. James, Seth Falcon, Winston Chang, Romain Francois, J.J. Allaire, Kevin Ushey, Qiang Kou, Douglas Bates, Jeffrey A. Ryan, Joshua M. Ulrich, Wouter Thielen, and John Chambers for development of the R packages used in the book.

The authors would particularly like to thank managing editor David Tranah for insightful comments based upon his vast experience, Clare Dennison for the timely content decisions, and Austin Bennett for the creative data scientist caricature on the cover.