



Data Science with Julia

Paul D. McNicholas
Peter A. Tait



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Data Science with Julia

By Paul D. McNicholas and Peter A. Tait



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2019 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed on acid-free paper
Version Date: 20191119

International Standard Book Number-13: 978-1-138-49998-0 (Paperback)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged, please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Library of Congress Cataloging-in-Publication Data

Names: McNicholas, Paul D., author. | Tait, Peter A., author.
Title: Data science with Julia / Paul D. McNicholas, Peter A. Tait.
Description: Boca Raton : Taylor & Francis, CRC Press, 2018. | Includes bibliographical references and index.
Identifiers: LCCN 2018025237 | ISBN 9781138499980 (pbk.)
Subjects: LCSH: Julia (Computer program language) | Data structures (Computer science)
Classification: LCC QA76.73.J85 M37 2018 | DDC 005.7/3--dc23
LC record available at <https://lccn.loc.gov/2018025237>

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

Contents

CHAPTER	1 ■ Introduction	1
1.1	DATA SCIENCE	1
1.2	BIG DATA	4
1.3	JULIA	5
1.4	JULIA AND R PACKAGES	6
1.5	DATASETS	6
1.5.1	Overview	6
1.5.2	Beer Data	6
1.5.3	Coffee Data	7
1.5.4	Leptograpsus Crabs Data	8
1.5.5	Food Preferences Data	9
1.5.6	x2 Data	9
1.5.7	Iris Data	11
1.6	OUTLINE OF THE CONTENTS OF THIS MONOGRAPH	11
CHAPTER	2 ■ Core Julia	13
2.1	VARIABLE NAMES	13
2.2	OPERATORS	14
2.3	TYPES	15
2.3.1	Numeric	15
2.3.2	Floats	17
2.3.3	Strings	19
2.3.4	Tuples	22
2.4	DATA STRUCTURES	23
2.4.1	Arrays	23

2.4.2	Dictionaries	26
2.5	CONTROL FLOW	28
2.5.1	Compound Expressions	28
2.5.2	Conditional Evaluation	29
2.5.3	Loops	30
2.5.3.1	Basics	30
2.5.3.2	Loop termination	32
2.5.3.3	Exception handling	33
2.6	FUNCTIONS	36
CHAPTER 3	Working with Data	43
3.1	DATAFRAMES	43
3.2	CATEGORICAL DATA	47
3.3	INPUT/OUTPUT	48
3.4	USEFUL DATAFRAME FUNCTIONS	54
3.5	SPLIT-APPLY-COMBINE STRATEGY	56
3.6	QUERY.JL	59
CHAPTER 4	Visualizing Data	67
4.1	GADFLY.JL	67
4.2	VISUALIZING UNIVARIATE DATA	69
4.3	DISTRIBUTIONS	72
4.4	VISUALIZING BIVARIATE DATA	83
4.5	ERROR BARS	90
4.6	FACETS	91
4.7	SAVING PLOTS	91
CHAPTER 5	Supervised Learning	93
5.1	INTRODUCTION	93
5.2	CROSS-VALIDATION	96
5.2.1	Overview	96
5.2.2	<i>K</i> -Fold Cross-Validation	97
5.3	<i>K</i> -NEAREST NEIGHBOURS CLASSIFICATION	99
5.4	CLASSIFICATION AND REGRESSION TREES	102

5.4.1	Overview	102
5.4.2	Classification Trees	103
5.4.3	Regression Trees	106
5.4.4	Comments	108
5.5	BOOTSTRAP	108
5.6	RANDOM FORESTS	111
5.7	GRADIENT BOOSTING	113
5.7.1	Overview	113
5.7.2	Beer Data	116
5.7.3	Food Data	121
5.8	COMMENTS	126
CHAPTER	6 ■ Unsupervised Learning	129
6.1	INTRODUCTION	129
6.2	PRINCIPAL COMPONENTS ANALYSIS	132
6.3	PROBABILISTIC PRINCIPAL COMPONENTS ANALYSIS	135
6.4	EM ALGORITHM FOR PPCA	137
6.4.1	Background: EM Algorithm	137
6.4.2	E-step	138
6.4.3	M-step	139
6.4.4	Woodbury Identity	140
6.4.5	Initialization	141
6.4.6	Stopping Rule	141
6.4.7	Implementing the EM Algorithm for PPCA	142
6.4.8	Comments	146
6.5	<i>K</i>-MEANS CLUSTERING	148
6.6	MIXTURE OF PROBABILISTIC PRINCIPAL COMPONENTS ANALYZERS	151
6.6.1	Model	151
6.6.2	Parameter Estimation	152
6.6.3	Illustrative Example: Coffee Data	161
6.7	COMMENTS	162

CHAPTER	7 ■ R Interoperability	165
7.1	ACCESSING R DATASETS	165
7.2	INTERACTING WITH R	166
7.3	EXAMPLE: CLUSTERING AND DATA REDUC- TION FOR THE COFFEE DATA	171
7.3.1	Coffee Data	171
7.3.2	PGMM Analysis	172
7.3.3	VSCC Analysis	175
7.4	EXAMPLE: FOOD DATA	176
7.4.1	Overview	176
7.4.2	Random Forests	176
APPENDIX	A ■ Julia and R Packages Used Herein	185
APPENDIX	B ■ Variables for Food Data	187
APPENDIX	C ■ Useful Mathematical Results	193
C.1	BRIEF OVERVIEW OF EIGENVALUES	193
C.2	SELECTED LINEAR ALGEBRA RESULTS	193
C.3	MATRIX CALCULUS RESULTS	194
APPENDIX	D ■ Performance Tips	197
D.1	FLOATING POINT NUMBERS	197
D.1.1	Do Not Test for Equality	197
D.1.2	Use Logarithms for Division	198
D.1.3	Subtracting Two Nearly Equal Numbers	198
D.2	JULIA PERFORMANCE	199
D.2.1	General Tips	199
D.2.2	Array Processing	199
D.2.3	Separate Core Computations	201
APPENDIX	E ■ Linear Algebra Functions	203
E.1	VECTOR OPERATIONS	203
E.2	MATRIX OPERATIONS	204

E.3 MATRIX DECOMPOSITIONS	205
References	208
<hr/>	
Index	217
<hr/>	



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Foreword

The 21st century will probably be the century of the data revolution. Our numerical world is creating masses of data every day and the volume of generated data is increasing more and more (the number of produced numerical data is doubling every two years according to the most recent estimates). In such a context, data science is nowadays an unavoidable field for anyone interested in exploiting data. People may be interested in either understanding a phenomenon or in predicting the future behavior of this phenomenon.

To this end, it is important to have significant knowledge of both the rationale (the theory) behind data science techniques and their practical use on real-world data. Indeed, data science is a mix of data, statistical/machine learning methods and software. Software is actually the link between data and data science techniques. It allows the practitioner to load the data and apply techniques on it for analysis. It is therefore important to master at least one of the data science languages.

The choice of the software language(s) mainly depends on your background and the expected level of analysis. **R** and Python are probably the two most popular languages for data science. On the one hand, **R** has been made by statisticians... mostly for statisticians! It is, however, an excellent tool for data science since the most recent statistical learning techniques are provided on the **R** platform (named CRAN). Using **R** is probably the best way to be directly connected to current research in statistics and data science through the packages provided by researchers. Python is, on the other hand, an actual computer science language (with all appropriate formal aspects) for which some advanced libraries for data science exist. In this context, the Julia language has the great advantage to permit users to interact with both **R** and Python (but also C, Fortran, etc.), within a software language designed for efficient and parallel numerical computing while keeping a high level of human readability.

Professor Paul McNicholas and Peter Tait propose in this book to learn both fundamental aspects of data science: theory and application. First, the book will provide you with the significant elements to understand the mathematical aspects behind the most used data science techniques. The book will also allow you to discover advanced recent techniques, such as probabilistic principal components analysis (PPCA), mixtures of PPCAs, and gradient boosting. In addition, the book will ask you to dive into the Julia language such that you directly apply the learned techniques on concrete examples. This is, in my opinion, the most efficient way to learn such an applied science. In addition, the focus made by this book on the Julia language is a great choice because of the numerous qualities of this language regarding data science practice. These include ease of learning for people familiar with **R** or Python, nice syntax, easy code debugging, the speed of the compiled language, and code reuse.

Both authors have extensive experience in data science. Professor Paul McNicholas is Canada Research Chair in Computational Statistics at McMaster University and Director of the MacDATA Institute of the same university. In his already prolific career, McNicholas has made important contributions to statistical learning. More precisely, his research is mainly focused on model-based learning with high-dimensional and skew-distributed data. He is also a researcher deeply involved in the spreading of research products through his numerous contributions to the **R** software with packages. Peter Tait is currently a Ph.D. student but, before returning to academia, he had a professional life dedicated to data science in industry. His strong knowledge of the needs of industry regarding data science problems was really an asset for the book.

This book is a great way to both start learning data science through the promising Julia language and to become an efficient data scientist.

Professor Charles Bouveyron
Professor of Statistics
INRIA Chair in Data Science
Université Côte d'Azur
Nice, France

Preface

This is a book for people who want to learn about the Julia language with a view to using it for data science. Some effort has gone into making this book suitable for someone who has familiarity with the `R` software and wants to learn about Julia. However, prior knowledge of `R` is not a requirement. While this book is not intended as a textbook for a course, some may find it a useful book to follow for a course that introduces statistics or data science students to Julia. It is our sincere hope that students, researchers and data scientists in general, who wish to learn Julia, will find this book beneficial.

More than twenty years have passed since the term data science was described by Dr. Chikio Hayashi in response to a question at a meeting of the International Federation of Classification Societies (Hayashi, 1998). Indeed, while the term data science has only gained notoriety over the past few years, much of the work it describes has been practiced for far longer. Furthermore, whatever the precise meaning of the term, there is no doubt that data science is important across virtually all areas of endeavour. This book is born out of a mixture of experiences all of which led to the conclusion that the use of Julia, as a language for data science, should be encouraged.

First, part of the motivation to write this book came from experience gained trying to teach material in data science without the benefit of a relatively easily understood base language that is effective for actually writing code. Secondly, there is the practical, and related, matter of writing efficient code while also having access to excellent code written by other researchers. This, of course, is the major advantage of `R`, where many researchers have contributed packages — sometimes based on code written in another language such as `C` or `Fortran` — for a wide variety of statistics and data science tasks. As we illustrate in this book, it is straightforward to call `R` from Julia and to thereby access whatever `R` packages are needed. Access to `R` packages and a growing selection of Julia

packages, together with an accessible, intuitive, and highly efficient base language, makes Julia a formidable platform for data science.

This book is not intended as an exhaustive introduction to data science. In fact, this book is far from an exhaustive introduction to data science. There are many very good books that one can consult to learn about different aspects of data science (e.g., Bishop, 2006; Hastie et al., 2009; Schutt, 2013; White, 2015; Efron and Hastie, 2016), but this book is primarily about Julia. Nevertheless, several important topics in data science are covered. These include data visualization, supervised learning, and unsupervised learning. When discussing supervised learning, we place some focus on gradient boosting — a machine learning technique — because we have found this approach very effective in applications. However, for unsupervised learning, we take a more statistical approach and place some focus on the use of probabilistic principal components analyzers and a mixture thereof.

This monograph is laid out to progress naturally. In [Chapter 1](#), we discuss data science and provide some historical context. Julia is also introduced as well as details of the packages and datasets used herein. [Chapters 2](#) and [3](#) cover the basics of the Julia language as well as how to work with data in Julia. After that ([Chapter 4](#)), a crucially important topic in data science is discussed: visualization. The book continues with selected techniques in supervised ([Chapter 5](#)) and unsupervised learning ([Chapter 6](#)), before concluding with details of how to call **R** functions from within Julia ([Chapter 7](#)). This last chapter also provides further examples of mixture model-based clustering as well as an example that uses random forests. Some appendices are included to provide readers with some relevant mathematics, Julia performance tips and a list of useful linear algebra functions in Julia.

There is a large volume of Julia code throughout this book, which is intended to help the reader gain familiarity with the language. We strongly encourage readers to run the code for themselves and play around with it. To make the code as easy as possible to work with, we have interlaced it with comments. As readers begin to get to grips with Julia, we encourage them to supplement or replace our comments with their own. For the reader's convenience, all of the code from this book is available on GitHub: github.com/paTait/dswj.

We are most grateful to David Grubbs of the Taylor & Francis Group for his support in this endeavour. His geniality and professionalism are always very much appreciated. Special thanks to

Professor Charles Bouveyron for kindly agreeing to lend his expertise in the form of a wonderful Foreword to this book. Thanks also to Dr. Joseph Kang and an anonymous reviewer for their very helpful comments and suggestions. McNicholas is thankful to Eamonn Mullins and Dr. Myra O'Regan for providing him with a solid foundation for data science during his time as an undergraduate student. Dr. Sharon McNicholas read a draft of this book and provided some very helpful feedback for which we are most grateful.

A final word of thanks goes to our respective families; without their patience and support, this book would not have come to fruition.

Paul D. McNicholas and Peter A. Tait
Hamilton, Ontario



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

About the Authors

Paul D. McNicholas is the Canada Research Chair in Computational Statistics at McMaster University, where he is a Professor and University Scholar in the Department of Mathematics and Statistics as well as Director of the MacDATA Institute. He has published extensively in computational statistics, with the vast majority of his work focusing on mixture model-based clustering. He is one of the leaders in this field and recently published a monograph devoted to the topic (*Mixture Model-Based Classification*; Chapman & Hall/CRC Press, 2016). He is a Senior Member of the IEEE and a Member of the College of the Royal Society of Canada.

Peter A. Tait is a Ph.D. student at the School of Computational Science and Engineering at McMaster University. His research interests span multivariate and computational statistics. Prior to returning to academia, he worked as a data scientist in the software industry, where he gained extensive practical experience.



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Introduction

DATA SCIENCE is discussed and some important connections, and contrasts, are drawn between statistics and data science. A brief discussion of big data is provided, the Julia language is briefly introduced, and all Julia packages used in this monograph are listed together with their respective version numbers. The same is done for the, albeit smaller number of, **R** packages used herein. Providing such details about the packages used helps ensure that the analyses illustrated herein can be reproduced. The datasets used in this monograph are also listed, along with some descriptive characteristics and their respective sources. Finally, the contents of this monograph are outlined.

1.1 DATA SCIENCE

What is data science? It is an interesting question and one without a widely accepted answer. Herein, we take a broad view that data science encompasses all work related to data. While this includes data analysis, it also takes in a host of other topics such as data cleaning, data curation, data ethics, research data management, etc. This monograph discusses some of those aspects of data science that are commonly handled in Julia, and similar software; hence, its title.

The place of statistics within the pantheon of data science is a topic on which much has been written. While statistics is certainly a very important part of data science, statistics should not be taken as synonymous with data science. Much has been written about the relationship between data science and statistics. On the one extreme, some might view data science — and data analysis, in particular — as a retrogression of statistics; yet, on the other

2 ■ Data Science with Julia

extreme, some may argue that data science is a manifestation of what statistics was always meant to be. In reality, it is probably an error to try to compare statistics and data science as if they were alternatives. Herein, we consider that statistics plays a crucial role in data analysis, or data analytics, which in turn is a crucial part of the data science mosaic.

Contrasting data analysis and mathematical statistics, Hayashi (1998) writes:

... mathematical statistics have been prone to be removed from reality. On the other hand, the method of data analysis has developed in the fields disregarded by mathematical statistics and has given useful results to solve complicated problems based on mathematico-statistical methods (which are not always based on statistical inference but rather are descriptive).

The views expressed by Hayashi (1998) are not altogether different from more recent observations that, insofar as analysis is concerned, data science tends to focus on prediction, while statistics has focused on modelling and inference. That is not to say that prediction is not a part of inference but rather that prediction is a part, and not the goal, of inference. We shall return to this theme, i.e., inference versus prediction, several times within this monograph.

Breiman (2001b) writes incisively about two cultures in statistical modelling, and this work is wonderfully summarized in the first few lines of its abstract:

There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown. The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems.

The viewpoint articulated here leans towards a view of data analysis as, at least partly, arising out of one culture in statistical modelling.

In a very interesting contribution, Cleveland (2001) outlines a blueprint for a university department, with knock-on implications for *curricula*. Interestingly, he casts data science as an “altered

Index

A

- Aitken acceleration, 142
- Alternating expectation-conditional maximization (AECM) algorithm, 152–153, 154–156, 159
- Arrays, 23, 45, 47, 48, 49
 - built-in functions, 24–25
 - counting, in Julia, 26
 - fast computations, 23
 - generating, 23
 - slicing, 25
- ASCII, 19

B

- Bagging, 112, 113
- Bar charts, 69, 70
- Bayesian information criterion (BIC), 146, 147, 163
- Bernoulli distribution, 168
- Big data, 4–5
- Bivariate plots, 83–84
- Boolean expression, 29–30
- Boolean values, 16, 45, 58
- Bootstrap, 93, 108–111
- Boxplots, 67, 70, 74, 75, 83
- Brier score, 95

C

- C, 5, 6
- C#, 59
- Cairo.jl package, 92
- CART. *See* classification and

- regression trees (CART)

- Categorical data, 47–49
- Char, 19
- Chi-square test, 167
- Classification and regression trees (CART), 93
 - classification trees, 103–106
 - overview, 102
 - regression trees, 106–107
- Clustering.jl, 129
- Compound expressions, 28
- Conditional evaluation, 29–30
- Constant matrix, 44, 130
- Covariance matrix, 119, 130, 131, 141
- Cross-entropy, 106
- CSV.jl library, 48–49, 59

D

- Data analysis, 3, 4
- Data matrices, 43
- Data science
 - computing, role of, 3
 - history of, 3
 - statistics, role of, 1–2, 3
- Data, big. *See* big data
- Data, labelled, 94
- Data, unlabelled, 94
- Dataframes
 - arrays (*see* arrays)
 - columns, 44
 - defining, 43
 - features, 44

- functions, 54–55
- matrix form, 43
- slicing, 45
- sorting, 58
- Datasets, 6
 - beer data, 7, 50, 60, 75, 83, 86, 116–117, 120
 - coffee data, 7, 161–162, 165, 166, 171–172, 174, 176
 - crabs data, 7, 134, 146, 147, 167
 - food preferences data, 10, 51, 165, 176
 - iris data, 10
 - wine data, 166
 - x2 data, 10
- DataStreams.jl, 60
- Deviance, 95
- Dictionaries, 26, 27
- Dot charts, 69, 70
- E**
- Eigenvalues, 131
- Eigenvectors, 131
- Empirical cumulative
 - distribution function (ECDF) plot, 80, 82, 83
- Error bars, 90
- Error() function, 35
- Exception handling, 33–35
- Expectation-maximization (EM) algorithm
 - E-step, 138–139
 - implementing, for PPCA, 142, 144
 - initialization, 141
 - M-step, 139–140
 - overview, 137
 - stopping rule, 141
- Woodbury identity, 140–141
- F**
- Facets, 91
- Factor analysis, 135
- Floating points, 61
- Floats, 16, 17–18, 44
- Fortran, 5
- Functions
 - anonymous, 38
 - defining, 36
 - inputs, 36–38
 - naming conventions, 36
 - series of, 40
 - writing, 40
- G**
- GadFly.jl, 67, 68–69, 86, 88, 90
- Gaussian mixture model, 148
- ggplot2, 67
- Gini index, 106
- Gradient boosting
 - beer data example, 116–117, 120
 - food preferences example, 121–123
 - overview, 113–115
- Grammar of Graphics (GoG), 67–68, 77, 86
- H**
- Hexbin plots, 85, 86
- Histograms, 67, 71–72, 85, 86
- I**
- IEEE 754 standard, 18
- Inter quartile range (IQR), 74
- J**
- Julia
 - interoperability, 5–6
 - syntax, 5

K

- K*-fold cross-validation, 97–98, 183
- K*-means clustering, 148–149, 159, 161
- K*-nearest neighbours (*k*NN), 99–100
- Kernel density, 71–72

L

- Language-INtegrated Query (LINQ), 59
- Loess model, 85
- Log-likelihood, 142, 144, 153
- Logistic regression model, 90
- Loops
 - continue keyword, 33
 - for loop, 30, 32
 - overview, 30
 - termination of, 32
 - while loop, 31, 32

M

- MASS package, 7
- Mean squared error (MSE), 94
- Median absolute error (MAE), 94, 121, 125, 177–178, 180, 181, 182
- Missing values, 45–46
- Mixture package, 7
- MLBase.jl, 116
- MM algorithm, 39

N

- Numeric literals, 15
- Numeric primitives, 15

O

- ODBS, 59
- Operators, 14–15

P

- Pareto distributions, 80
- Parsimonious Gaussian mixture models (PGMMs), 172–173, 174
- Perl, 21
- Plots, saving, 92
- Principal components analysis (PCA), 132–134, 135, 147
- Probabilistic principal components analysis (PPCA), 123–125, 132–134, 142, 144
 - AECM algorithm for mixture model, 159
 - mixture models, 151–152
 - parameter estimation, 152–153
- Pseudocode, 113
- Python, 5, 6, 14, 19, 37

Q

- QQ-plots, 68, 77, 80, 84
- Query.jl package
 - descending() function, 63–64
 - @group statement, 65
 - @join statement, 64
 - @let, 62
 - @orderby, 63
 - @select statement, 60
 - arrays, 60, 61
 - overview, 59, 60
 - query statement, 59
 - syntax, 60, 61

R

- R, 5, 6, 14, 19, 38, 49, 55, 56, 67, 166–171
- Random forests, 112, 113, 176–184

Random matrix, 130
 Random vector, 43, 44, 130, 131
 RCall, 166–167
 RDatasets.jl, 165, 166
 Regexes, 21
 Root mean squared error
 (RMSE), 95, 121, 125, 181

S

Scatterplots, 83–84, 85
 Scientific notation, 16, 19
 Split-apply-combine (SAC)
 strategy, 56–58
 SQLite, 59
 Statistical modelling, cultures
 within, 2–3
 Statistics, role of, in data
 science, 1–2, 4
 Strings, 19–21, 47, 49
 Supervised learning, 93–96

T

Ternary operators, 30
 Training-test paradigm, 97
 Trees, combining, 112
 Try-catch statements, 35
 Tuples, 22, 61
 Type promotion, 40

U

Unicode, 13, 19, 36
 Unsigned integers, 16
 Unsupervised learning, 129, 130, 162
 UTF-8, 13, 19

V

Validation set, 97
 Variables
 global, 40
 names, 13, 44
 symbols, 44
 Velocity, 4, 5
 Violin plots, 75, 76, 77, 83, 183
 Visualizations, data
 custom, 70
 GadFly.jl (*see* GadFly.jl)
 overview, 69–72
 VSCC technique, 175

W

Woodbury identity, 140–141, 144

X

XGBoost, 93, 115, 116, 121, 122