# Machine Learning

## HANDS-ON FOR DEVELOPERS AND TECHNICAL PROFESSIONALS

# Machine Learning

## Hands-On for Developers and Technical Professionals

## Jason Bell

# WILEY

**Machine Learning: Hands-On for Developers and Technical Professionals**

*To Wendy and Clarissa.*

# Credits

**Executive Editor**
Carol Long

**Project Editor**
Charlotte Kughen

**Technical Editor**
Mitchell Wyle

**Production Editor**
Christine Mugnolo

**Copy Editor**
Katherine Burt

**Production Manager**
Kathleen Wisor

**Manager of Content Development
and Assembly**
Mary Beth Wakefield

**Director of Community Marketing**
David Mayhew

**Marketing Manager**
Carrie Sherrill

**Business Manager**
Amy Knies

**Professional Technology &
Strategy Director**
Barry Pruett

**Associate Publisher**
Jim Minatel

**Project Coordinator, Cover**
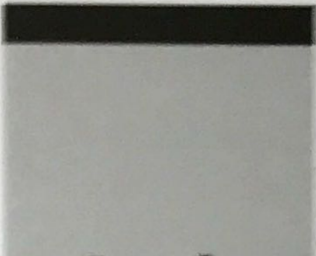Patrick Redmond

**Proofreader**
Nancy Carrasco

**Indexer**
Johnna Dinse

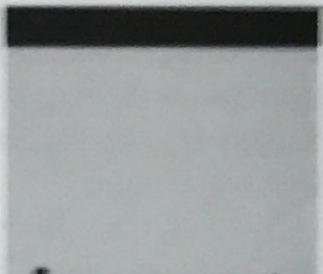**Cover Designer**
Wiley

**Cover Image**
© iStock.com/VLADGRIN

# About the Author

**Jason Bell** has been working with point-of-sale and customer-loyalty data since 2002, and he has been involved in software development for more than 25 years. He is founder of Datasentiment, a UK business that helps companies worldwide with data acquisition, processing, and insight.

# Acknowledgments

During the autumn of 2013, I was presented with some interesting options: either do a research-based PhD or co-author a book on machine learning. One would take six years and the other would take seven to eight months. Because of the speed the data industry was, and still is, progressing, the idea of the book was more appealing because I would be able to get something out while it was still fresh and relevant, and that was more important to me.

I say "co-author" because the original plan was to write a machine learning book with Aidan Rogers. Due to circumstances beyond his control he had to pull out. With Aidan's blessing, I continued under my own steam, and for that opportunity I can't thank him enough for his grace, encouragement, and support in that decision.

Many thanks goes to Wiley, especially Executive Editor, Carol Long, for letting me tweak things here and there with the original concept and bring it to a more practical level than a theoretical one; Project Editor, Charlotte Kughen, who kept me on the straight and narrow when there were times I didn't make sense; and Mitchell Wyle for reviewing the technical side of things. Also big thanks to the Wiley family as a whole for looking after me with this project.

Over the years I've met and worked with some incredible people, so in no particular order here goes: Garrett Murphy, Clare Conway, Colin Mitchell, David Crozier, Edd Dumbill, Matt Biddulph, Jim Weber, Tara Simpson, Marty Neill, John Girvin, Greg O'Hanlon, Clare Rowland, Tim Spear, Ronan Cunningham, Tom Grey, Stevie Morrow, Steve Orr, Kevin Parker, John Reid, James Blundell, Mary McKenna, Mark Nagurski, Alan Hook, Jon Brookes, Conal Loughrey, Paul Graham, Frankie Colclough, and countless others (whom I will be kicking myself that I've forgotten) for all the meetings, the chats, the ideas, and the collaborations.

# Contents