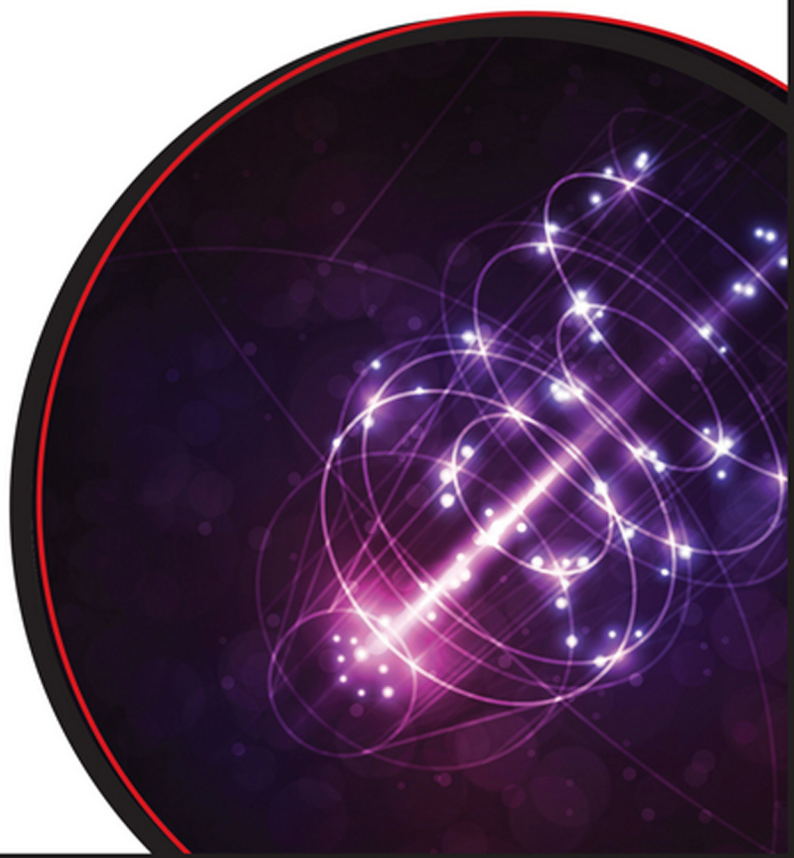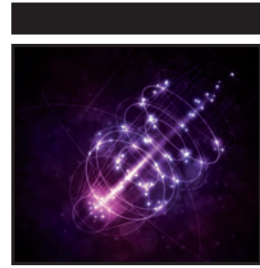# Data Analysis Using SQL and Excel®

## Second Edition

Gordon S. Linoff

# Data Analysis Using SQL and Excel®

# Data Analysis Using SQL and Excel®

## Second Edition

Gordon S. Linoff

WILEY

**Data Analysis Using SQL and Excel®, Second Edition**

*To Giuseppe—for twenty five years, five books, and counting . . .*

# About the Author

Gordon S. Linoff has been working with databases, big data, and data mining for almost longer than he can remember. With decades of experience on the practice of using data effectively, he is a recognized expert in the field of data mining.

Gordon started using spreadsheets while a student at MIT, on the original Compaq Portable, the world's first luggable computer. Not very many years later, he managed a development group at the now-defunct Thinking Machines Corporation, tasked with building a massively parallel relational database for decision support.

After Thinking Machines' demise, he founded Data Miners in 1998 with his friend and former colleague Michael J. A. Berry (who left in 2012). Since then, he has worked on a wide diversity of projects across many different companies. He has taught hundreds of classes around the world on data mining and survival analysis through SAS Institute, a leader in statistical and business analytics software. He is also an avid contributor to Stack Overflow, particularly on questions related to databases, having the highest score in 2014.

Together with Michael Berry, Gordon has written several influential books on data mining, including *Data Mining Techniques for Marketing, Sales, and Customer Support,* the first book on data mining to achieve a third edition.

Gordon lives in New York with Giuseppe Scalia, his partner of 25 years.

# Credits

**Project Editor**
John Sleeva

**Technical Editor**
Michael Berry

**Production Editor**
Dassi Zeidel

**Copy Editor**
Mike La Bonne

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

**Marketing Director**
David Mayhew

**Marketing Manager**
Carrie Sherrill

**Professional Technology & Strategy Director**
Barry Pruett

**Business Manager**
Amy Knies

**Associate Publisher**
Jim Minatel

**Project Coordinator, Cover**
Brent Savage

**Proofreader**
Sara Wilson

**Indexer**
Johnna VanHoose Dinse

**Cover Designer**
Wiley

**Cover Image**
©iStock.com/Nobi_Prizue

# Acknowledgments

Although this book has only one name on the cover, many people have helped me both specifically on this book and more generally in understanding data, analysis, and presentation.

I first met Michael Berry in 1990. We later founded Data Miners together, and he has been helpful on all fronts. He reviewed the chapters, tested the SQL code in the examples, and helped anonymize the data. His insights have been helpful and his debugging skills have made the examples much more accurate. His wife, Stephanie Jack, also deserves special praise for her patience and willingness to share Michael's time.

The original idea for the book came from Nick Drake, who then worked at Datran Media. A statistician by training, Nick was looking for a book that would help him use databases for data analysis. Bob Elliott, at the time my editor at Wiley, liked the idea.

Throughout the chapters, the understanding of data processing is based on dataflows, which Craig Stanfill of Ab Initio Corporation first introduced me to long ago when we worked together at Thinking Machines Corporation.

Along the way, I have learned a lot from many people. Anne Milley of SAS Institute first suggested that I learn survival analysis. Will Potts, now working at CapitalOne, then taught me much of what I know about the subject. Brij Masand helped extend the ideas to practical forecasting applications. Chi Kong Ho and his team at the *New York Times* provided valuable feedback for applying survival analysis to customer value calculations.

Stuart Ward from the *New York Times* and Zaiying Huang spent countless hours explaining and discussing statistical concepts. Harrison Sohmer, also of the *New York Times*, taught me many Excel tricks, some of which I've been able to include in the book.

# Contents at a Glance

# Contents