# Practical Statistics for Data Scientists

## 50+ Essential Concepts Using R and Python

Peter Bruce, Andrew Bruce & Peter Gedeck

# Practical Statistics for Data Scientists

Statistical methods are a key part of data science, yet few data scientists have formal statistical training. Courses and books on basic statistics rarely cover the topic from a data science perspective. The second edition of this popular guide adds comprehensive examples in Python, provides practical guidance on applying statistical methods to data science, tells you how to avoid their misuse, and gives you advice on what's important and what's not.

Many data science resources incorporate statistical methods but lack a deeper statistical perspective. If you're familiar with the R or Python programming languages and have some exposure to statistics, this quick reference bridges the gap in an accessible, readable format.

With this book, you'll learn:

- Why exploratory data analysis is a key preliminary step in data science
- How random sampling can reduce bias and yield a higher-quality dataset, even with big data
- How the principles of experimental design yield definitive answers to questions
- How to use regression to estimate outcomes and detect anomalies
- Key classification techniques for predicting which categories a record belongs to
- Statistical machine learning methods that "learn" from data
- Unsupervised learning methods for extracting meaning from unlabeled data

"This book is not another statistics textbook nor a machine learning manual. It's much better: it makes the connection between useful statistical terms and principles and today's data mining lingo and practices, with clear explanations and plenty of examples. This is a terrific reference for data science beginners and old timers."

**—Galit Shmueli**
lead author of the best-selling series *Data Mining for Business Analytics* and Distinguished Professor, National Tsing Hua University, Taiwan

**Peter Bruce is** the founder of the Institute for Statistics Education at Statistics.com.

**Andrew Bruce** is a principal research scientist at Amazon and has over 30 years of experience in statistics and data science.

**Peter Gedeck** is a senior data scientist at Collaborative Drug Discovery, developing machine learning algorithms to predict properties of drug candidates.

---

DATA SCIENCE | STATISTICS

US $69.99          CAN $92.99
ISBN: 978-1-492-07294-2

9 781492 072942

56999

# Practical Statistics for Data Scientists

## 50+ Essential Concepts Using R and Python

*Peter Bruce, Andrew Bruce, and Peter Gedeck*

**Practical Statistics for Data Scientists**

by Peter Bruce, Andrew Bruce, and Peter Gedeck

Printed in the United States of America.

*Peter Bruce and Andrew Bruce would like to dedicate this book to the memories of our parents, Victor G. Bruce and Nancy C. Bruce, who cultivated a passion for math and science; and to our early mentors John W. Tukey and Julian Simon and our lifelong friend Geoff Watson, who helped inspire us to pursue a career in statistics.*

*Peter Gedeck would like to dedicate this book to Tim Clark and Christian Kramer, with deep thanks for their scientific collaboration and friendship.*

This page intentionally left blank

# Table of Contents