# Learning Data Mining with Python

The next step in the information age is to gain insights from the deluge of data coming our way. Data mining provides a way of finding this insight, and Python is one of the most popular languages for data mining, providing both power and flexibility in analysis.

This book teaches you to design and develop data mining applications using a variety of datasets, starting with basic classification and affinity analysis. Next, we move on to more complex data types including text, images, and graphs. In every chapter, we create models that solve real-world problems.
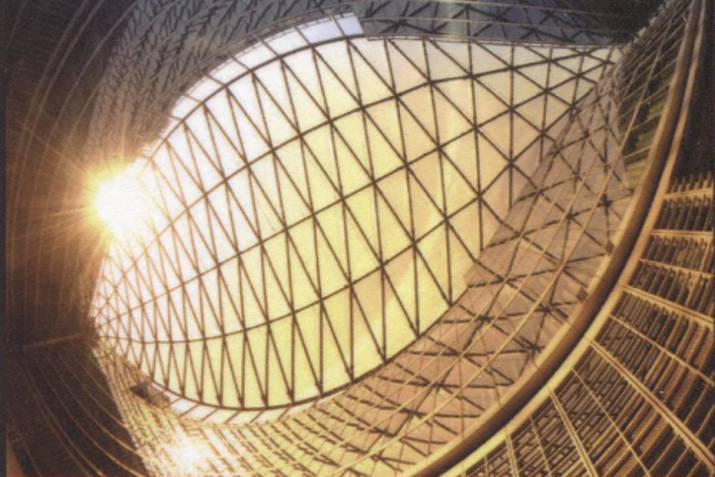
There is a rich and varied set of libraries available in Python for data mining. This book covers a large number, including the IPython Notebook, pandas, scikit-learn and NLTK.

Each chapter of this book introduces you to new algorithms and techniques. By the end of the book, you will gain a large insight into using Python for data mining, with a good knowledge and understanding of the algorithms and implementations.

## Who this book is written for

If you are a programmer who wants to get started with data mining, then this book is for you.
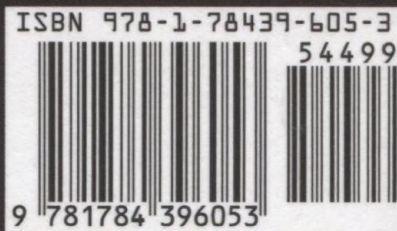
## What you will learn from this book

- Apply data mining concepts to real-world problems

- Predict the outcome of sports matches based on past results

- Determine the author of a document based on their writing style

- Use APIs to download datasets from social media and other online services

- Find and extract good features from difficult datasets

- Create models that solve real-world problems

- Design and develop data mining applications using a variety of datasets

- Set up reproducible experiments and generate robust results

- Recommend movies, online celebrities, and news articles based on personal preferences

- Compute on big data, including real-time data from the Internet

$ 44.99 US
£ 28.99 UK

Prices do not include local sales tax or VAT where applicable

[PACKT] open source*
PUBLISHING   community experience distilled

Visit **www.PacktPub.com** for books, eBooks, code, downloads, and PacktLib.