Trevor Hastie
Robert Tibshirani
Jerome Friedman

# The Elements of Statistical Learning

## Data Mining, Inference, and Prediction

统计学习基础

Trevor Hastie
Robert Tibshirani
Jerome Friedman

# The Elements of Statistical Learning

Data Mining, Inference, and Prediction

Trevor Hastie
Department of Statistics, and Department
  of Health Research & Policy
Sequoia Hall
Stanford University
Stanford, CA 94305-5405
USA
hastie@stat.stanford.edu

Robert Tibshirani
Department of Health Research & Policy,
  and Department of Statistics
HRP Redwood Building
Stanford University
Stanford, CA 94305-5405
USA
tibs@stat.stanford.edu

Jerome Friedman
Department of Statistics
Sequoia Hall
Stanford University
Stanford, CA 94305-5405
USA
jhf@stat.stanford.edu

springer.com

*To our families:*

*Samantha, Timothy, and Lynda*

*Charlie, Ryan, Julie, and Cheryl*

*Melanie, Dora, Monika, and Ildiko*

# Preface

*We are drowning in information and starving for knowledge.*

-Rutherford D. Roger

The field of Statistics is constantly challenged by the problems that science and industry brings to its door. In the early days, these problems often came from agricultural and industrial experiments and were relatively small in scope. With the advent of computers and the information age, statistical problems have exploded both in size and complexity. Challenges in the areas of data storage, organization and searching have led to the new field of "data mining"; statistical and computational problems in biology and medicine have created "bioinformatics." Vast amounts of data are being generated in many fields, and the statistician's job is to make sense of it all: to extract important patterns and trends, and understand "what the data says." We call this *learning from data.*

The challenges in learning from data have led to a revolution in the statistical sciences. Since computation plays such a key role, it is not surprising that much of this new development has been done by researchers in other fields such as computer science and engineering.

The learning problems that we consider can be roughly categorized as either *supervised* or *unsupervised.* In supervised learning, the goal is to predict the value of an outcome measure based on a number of input measures; in unsupervised learning, there is no outcome measure, and the goal is to describe the associations and patterns among a set of input measures.

This book is our attempt to bring together many of the important new ideas in learning, and explain them in a statistical framework. While some mathematical details are needed, we emphasize the methods and their conceptual underpinnings rather than their theoretical properties. As a result, we hope that this book will appeal not just to statisticians but also to researchers and practitioners in a wide variety of fields.

Just as we have learned a great deal from researchers outside of the field of statistics, our statistical viewpoint may help others to better understand different aspects of learning:

> There is no true interpretation of anything; interpretation is a vehicle in the service of human comprehension. The value of interpretation is in enabling others to fruitfully think about an idea.

> > –Andreas Buja

> *Trevor Hastie*
> *Robert Tibshirani*
> *Jerome Friedman*
>
> Stanford, California
> May 2001

> *The quiet statisticians have changed our world; not by discovering new facts or technical developments, but by changing the ways that we reason, experiment and form our opinions ....*

> > –Ian Hacking

# Contents

# 1
# Introduction

*Statistical learning* plays a key role in many areas of science, finance and industry. Here are some examples of learning problems:

- Predict whether a patient, hospitalized due to a heart attack, will have a second heart attack. The prediction is to be based on demographic, diet and clinical measurements for that patient.

- Predict the price of a stock in 6 months from now, on the basis of company performance measures and economic data.

- Identify the numbers in a handwritten ZIP code, from a digitized image.

- Estimate the amount of glucose in the blood of a diabetic person, from the infrared absorption spectrum of that person's blood.

- Identify the risk factors for prostrate cancer, based on clinical and demographic variables.

The science of learning plays a key role in the fields of statistics, data mining and artificial intelligence, intersecting with areas of engineering and other disciplines.

This book is about learning from data. In a typical scenario, we have an outcome measurement, usually quantitative (like a stock price) or categorical (like heart attack/no heart attack), that we wish to predict based on a set of *features* (like diet and clinical measurements). We have a *training set* of data, in which we observe the outcome and feature measurements

TABLE 1.1. *Average percentage of words or characters in an email message equal to the indicated word or character. We have chosen the words and characters showing the largest difference between* spam *and* email.

|       | george | you  | your | hp   | free | hpl  | !    | our  | re   | edu  | remove |
|-------|--------|------|------|------|------|------|------|------|------|------|--------|
| spam  | 0.00   | 2.26 | 1.38 | 0.02 | 0.52 | 0.01 | 0.51 | 0.51 | 0.13 | 0.01 | 0.28   |
| email | 1.27   | 1.27 | 0.44 | 0.90 | 0.07 | 0.43 | 0.11 | 0.18 | 0.42 | 0.29 | 0.01   |

for a set of objects (such as people). Using this data we build a prediction model, or *learner*, which will enable us to predict the outcome for new unseen objects. A good learner is one that accurately predicts such an outcome.

The examples above describe what is called the *supervised learning* problem. It is called "supervised" because of the presence of the outcome variable to guide the learning process. In the *unsupervised learning problem*, we observe only the features and have no measurements of the outcome. Our task is rather to describe how the data are organized or clustered. We devote most of this book to supervised learning; the unsupervised problem is less developed in the literature, and is the focus of the last chapter.

Here are some examples of real learning problems that are discussed in this book.

## Example 1: Email Spam

The data for this example consists of information from 4601 email messages, in a study to try to predict whether the email was junk email, or "spam." The objective was to design an automatic spam detector that could filter out spam before clogging the users' mailboxes. For 3601 email messages, the true outcome (email type) email or spam is available, along with the relative frequencies of 57 of the most commonly occurring words and punctuation marks in the email message. This is a supervised learning problem, with the outcome the class variable email/spam. It is also called a *classification* problem.

Table 1.1 lists the words and characters showing the largest average difference between spam and email.

Our learning method has to decide which features to use and how: for example, we might use a rule like

$$\text{if (\%george} < 0.6) \ \& \ (\%\text{you} > 1.5) \quad \text{then spam}$$
$$\text{else email.}$$

Another form of rule would be:

$$\text{if } (0.2 \cdot \%\text{you} - 0.3 \cdot \%\text{george}) > 0 \quad \text{then spam}$$
$$\text{else email.}$$