

STATISTICAL LEARNING - INTRODUCTION

统计学习

课程介绍

王涛 neowangtao@sjtu.edu.cn

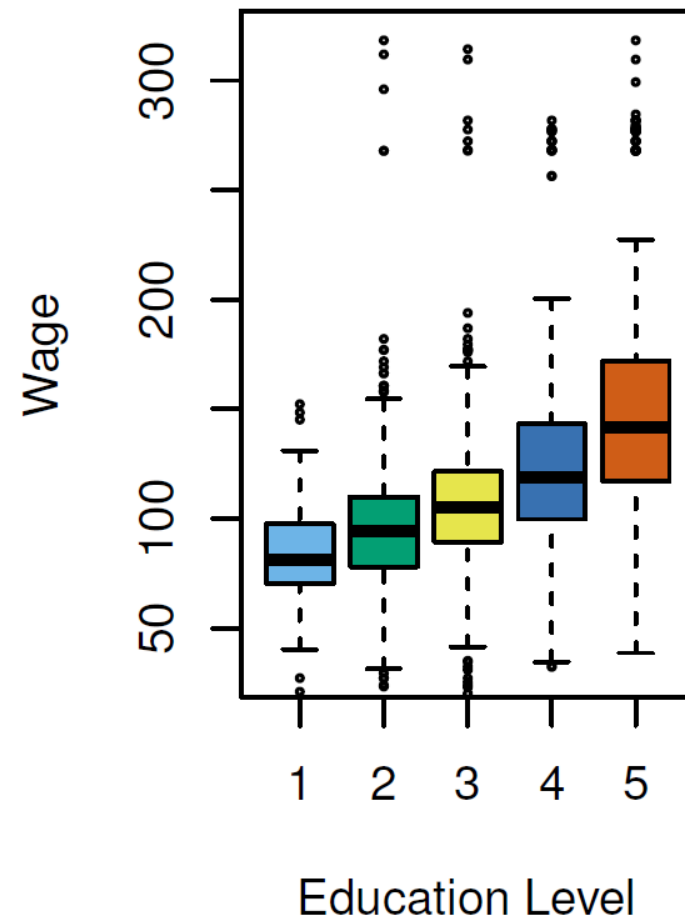
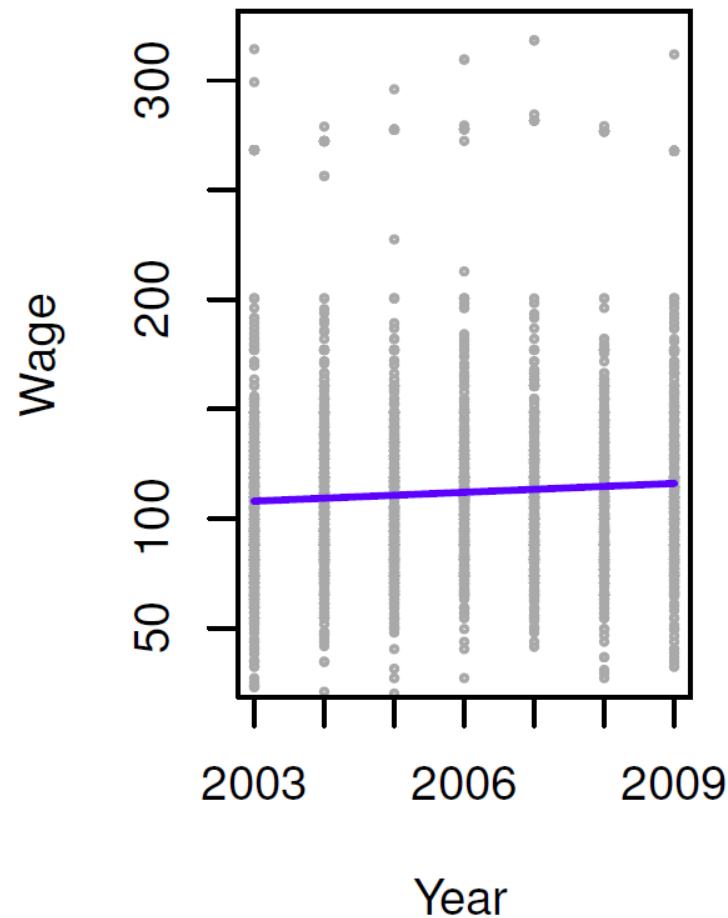
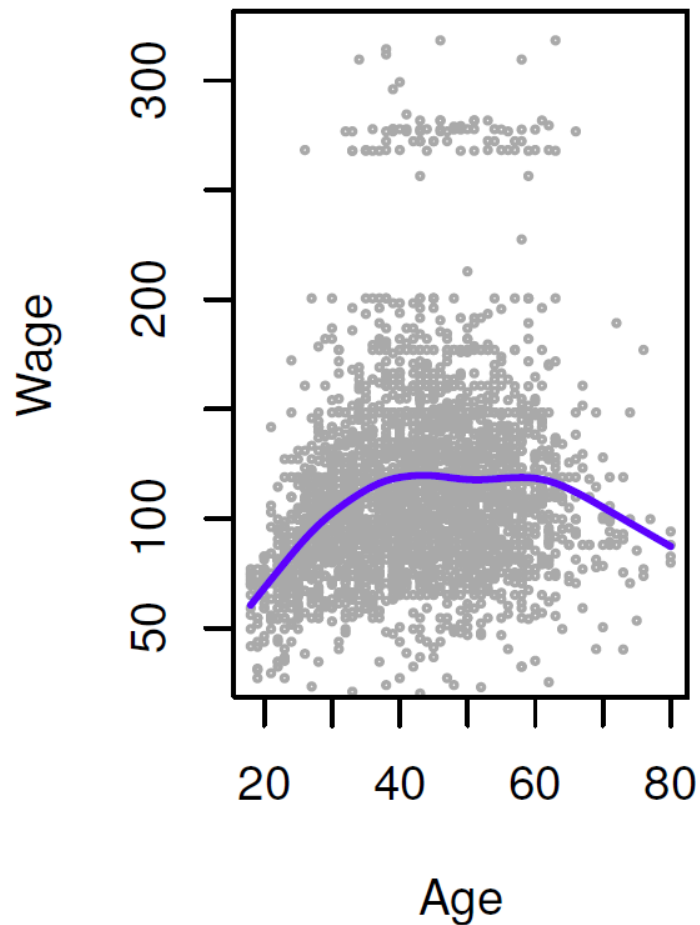
Outline

- An overview of statistical learning
- Schedules (教学安排)
- Grading policy (考核方法)
- Textbooks & references
- A brief history of statistical learning

An overview of statistical learning

Wage data

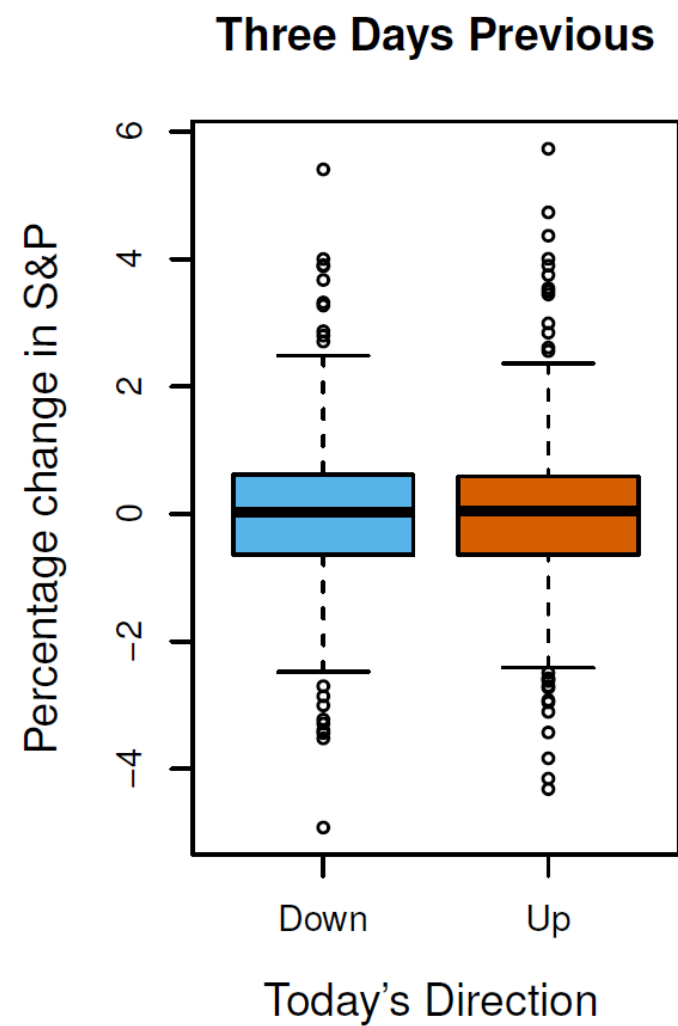
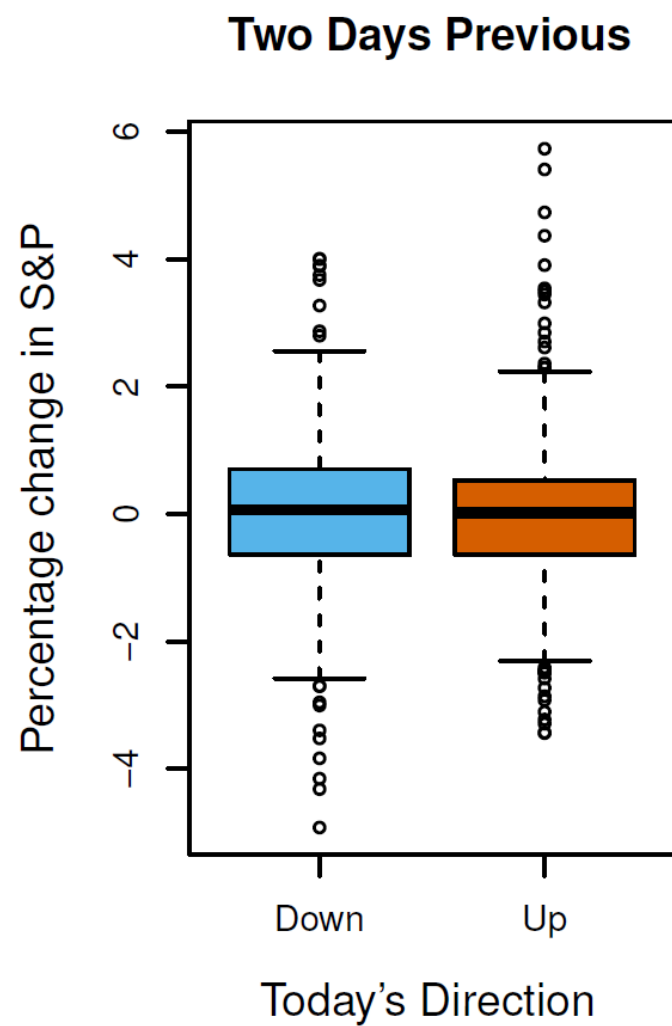
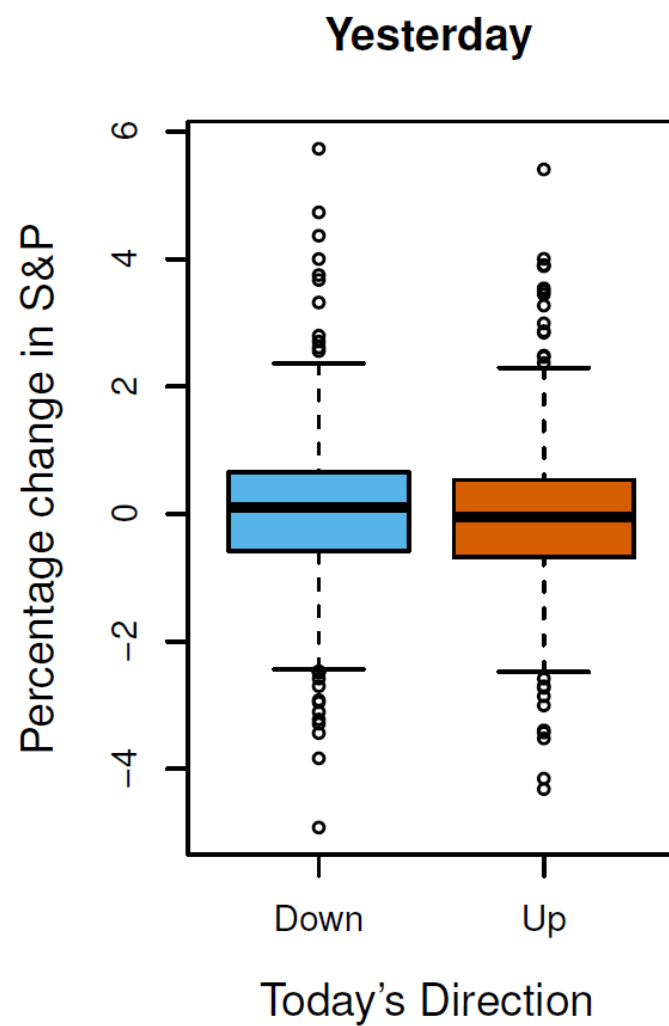
- Income survey data for males in central Atlantic region of USA (**Wage**)
- Understand the association between an employee's **wage** and a number of factors, such as **age**, **education**, and the calendar **year**



Some of the figures and tables in this presentation are taken from "*An Introduction to Statistical Learning, with Applications in R*" (Springer) with permission from the authors: G. James, D. Witten, T. Hastie, and R. Tibshirani

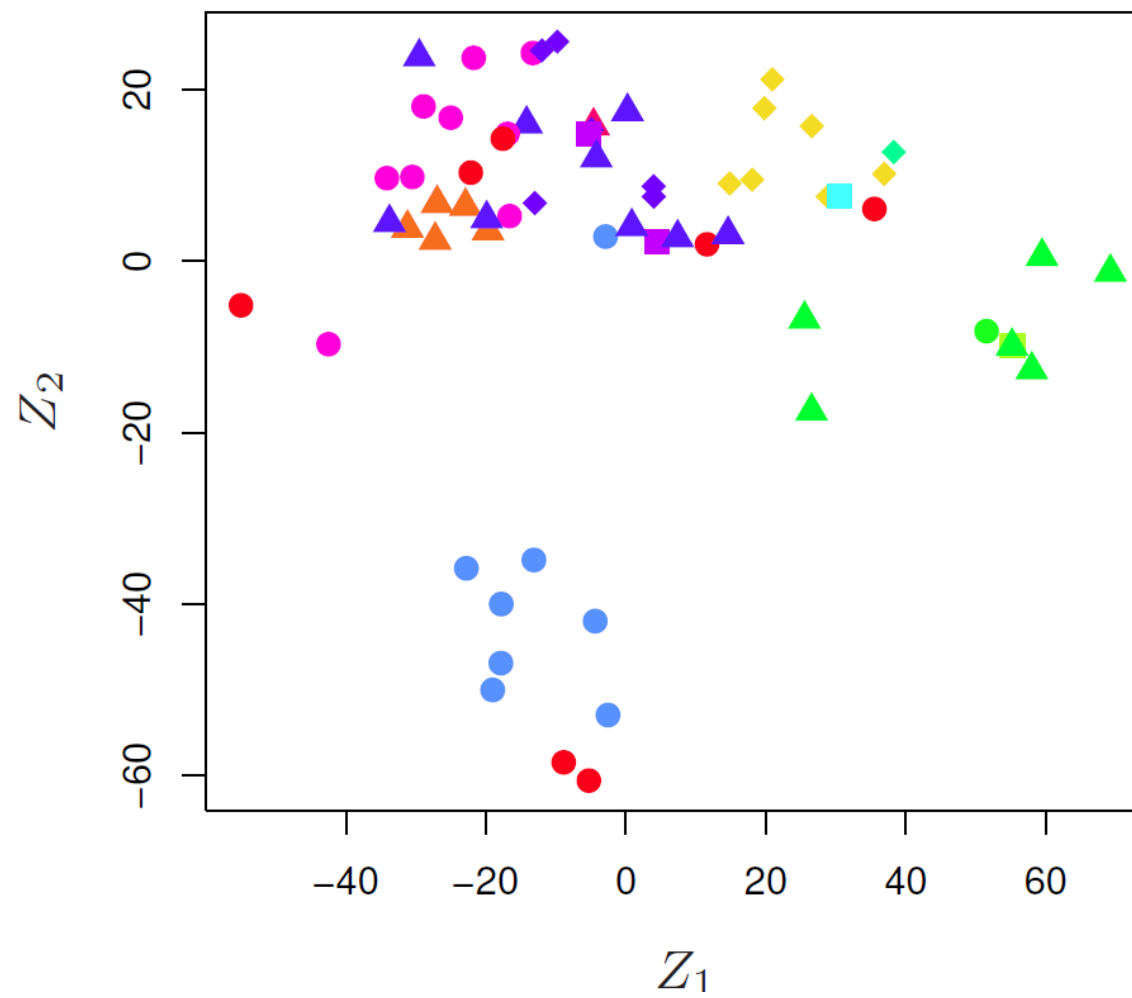
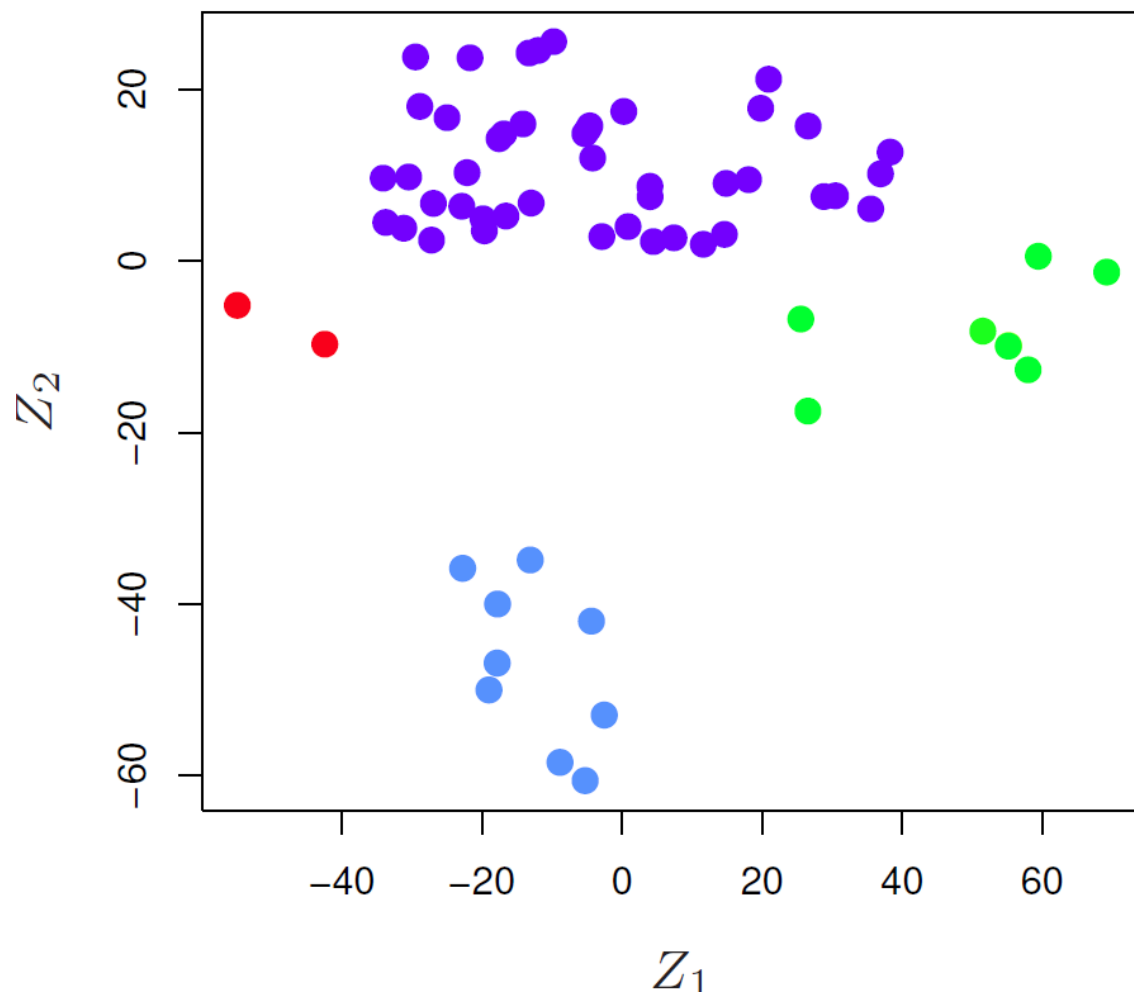
Stock market data

- Daily percentage returns for S&P 500 stock index over a 5-year period (**Smarket**)
- Predict whether the index will increase or decrease on a given day using the past 5 day's percentage changes in the index

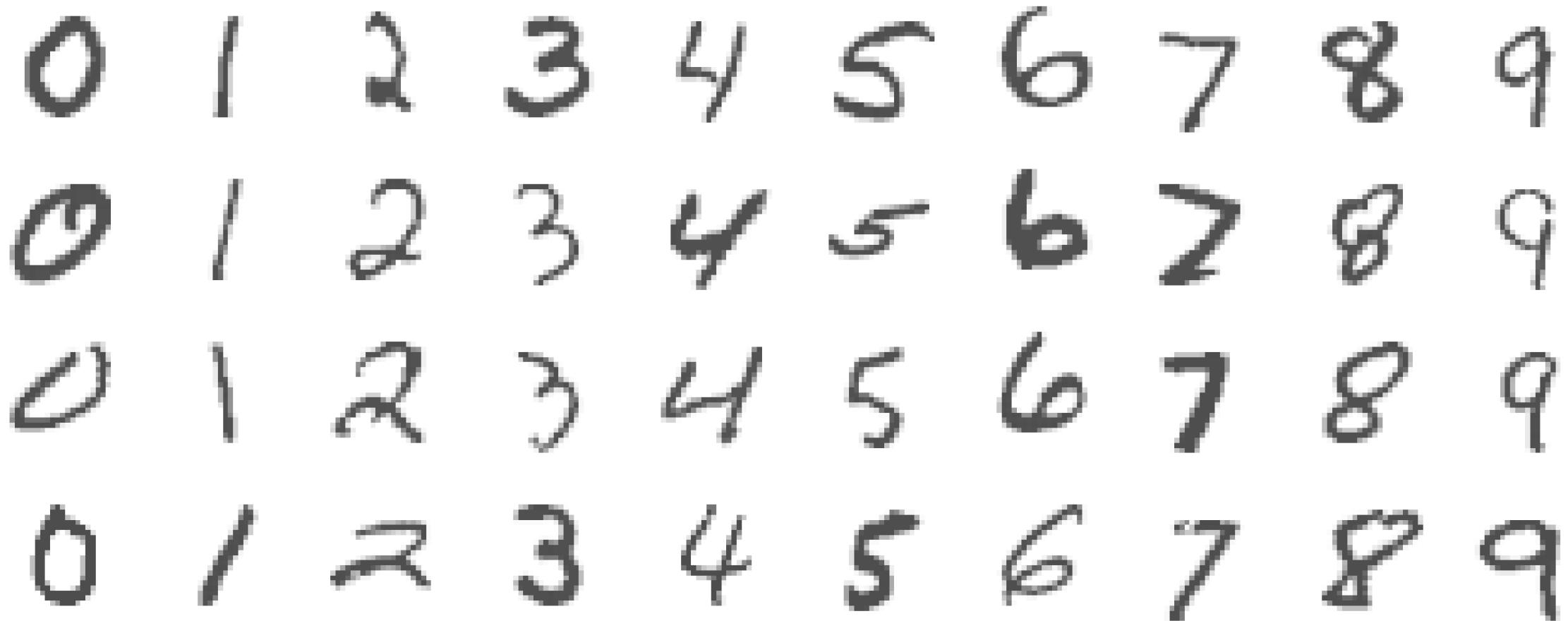


Gene expression data

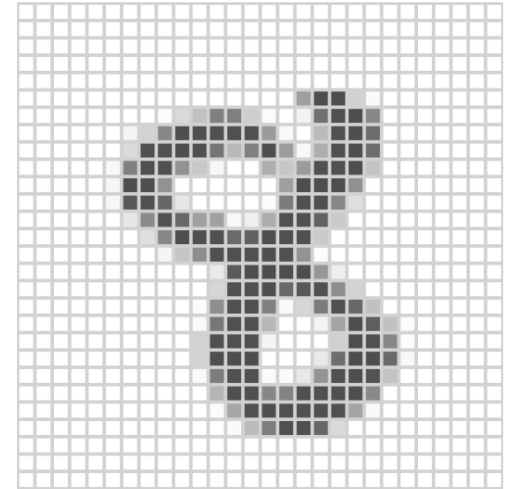
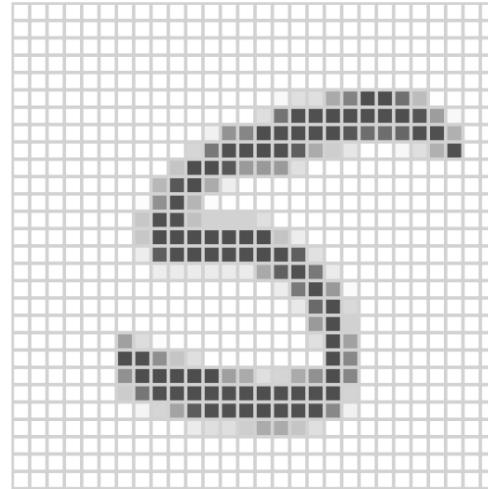
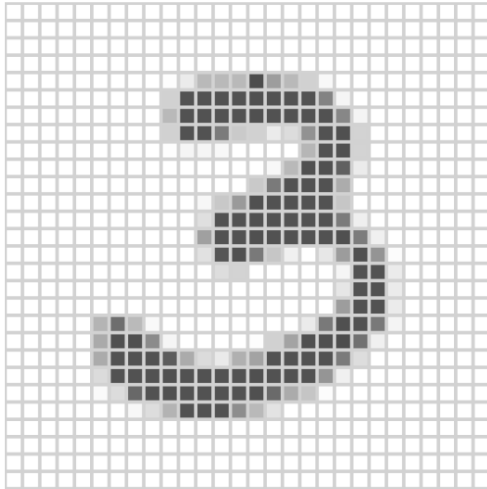
- Gene expression measurements for 64 cancer cell lines (NCI60)
- Determine whether there are subgroups among 64 cancer cell lines based on 6830 gene expression measurements



Handwritten digits from the MNIST corpus



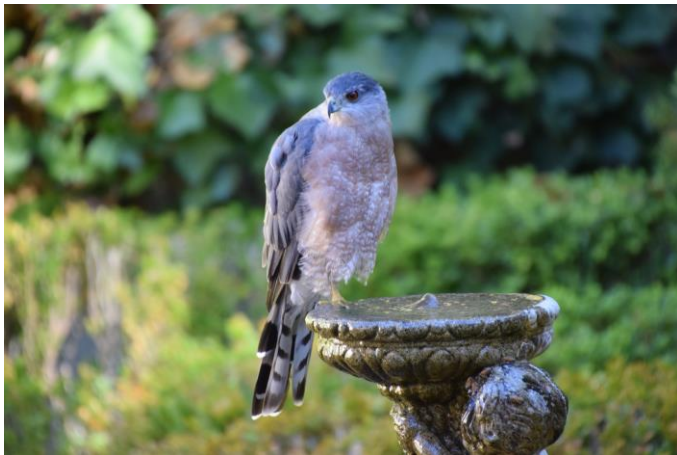
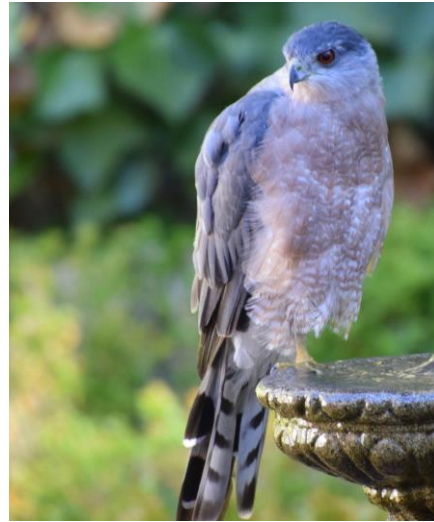
- Modified National Institute of Standards and Technology
- 28×28 grayscale images
- Features: the 784 pixel grayscale values (0-255)
- Labels: the digit classes (0-9)
- 60K training and 10K test images



Images from the CIFAR100 database



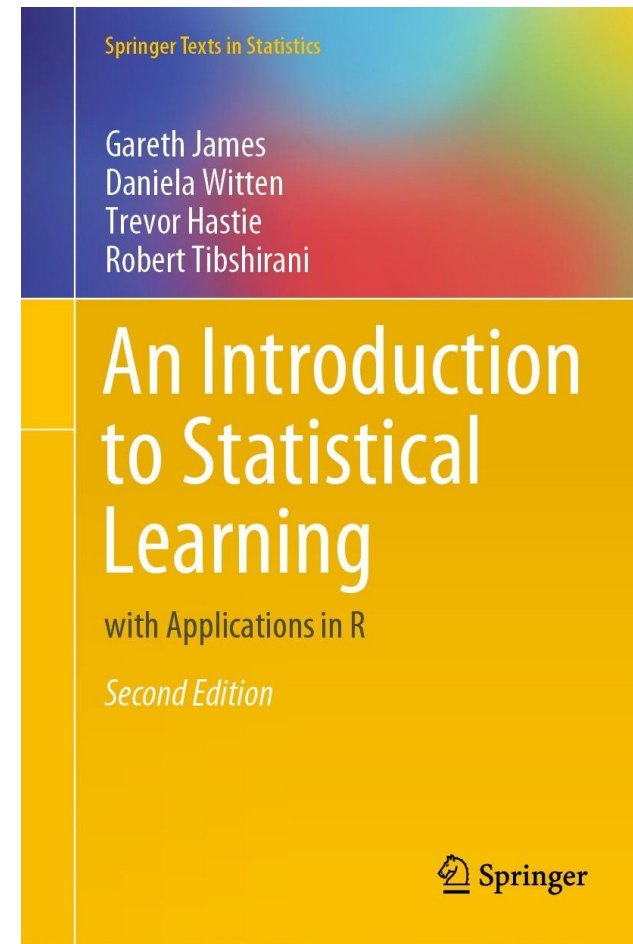
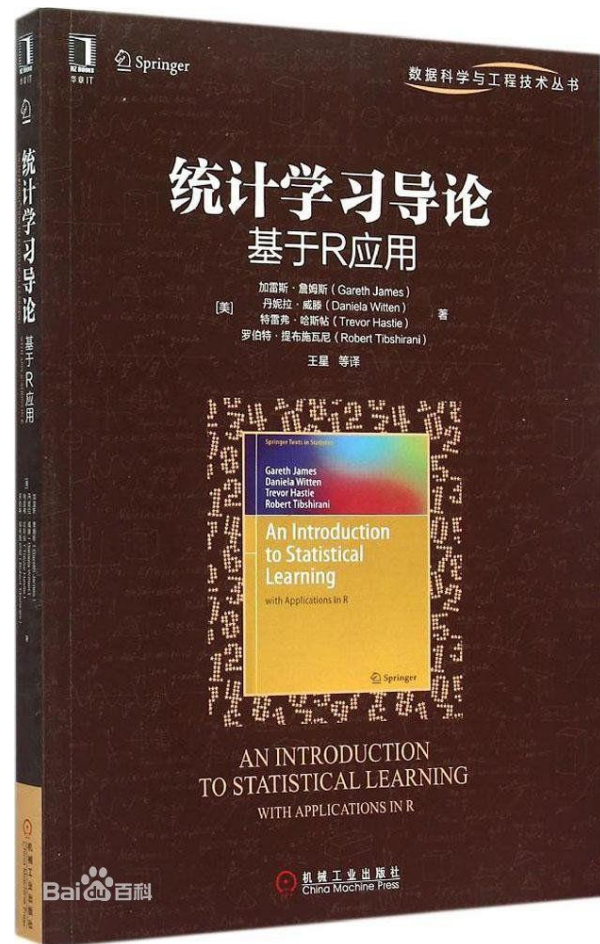
- Canadian Institute for Advanced Research
- 32×32 color natural images, with 100 classes
- 50K training and 10K test images
- Each image is a $32 \times 32 \times 3$ array of 8-bit numbers (0-255)
- The last dimension represents the color channel (R, G, B)



Schedules (教学安排)

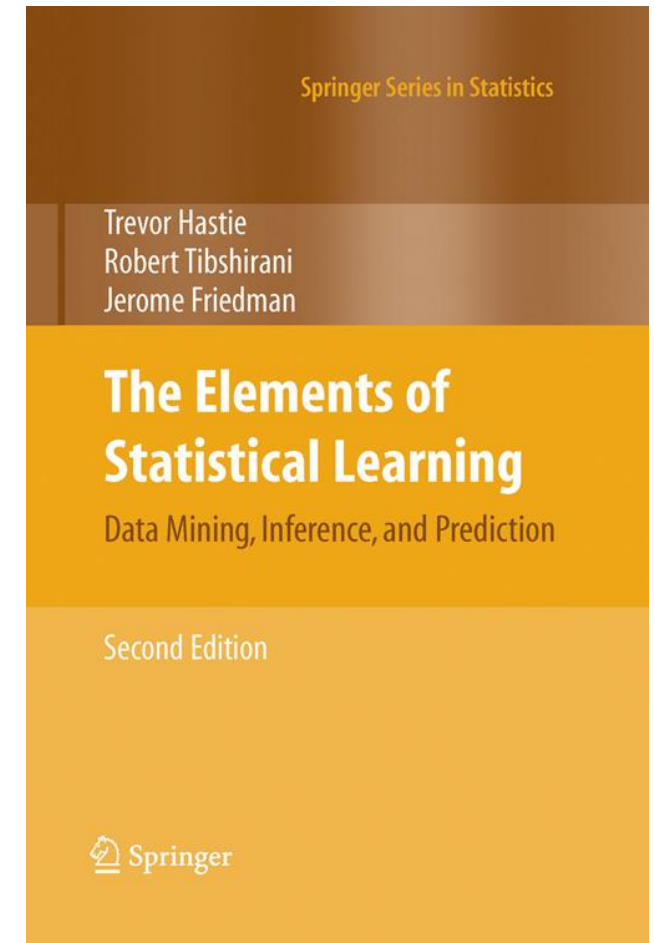
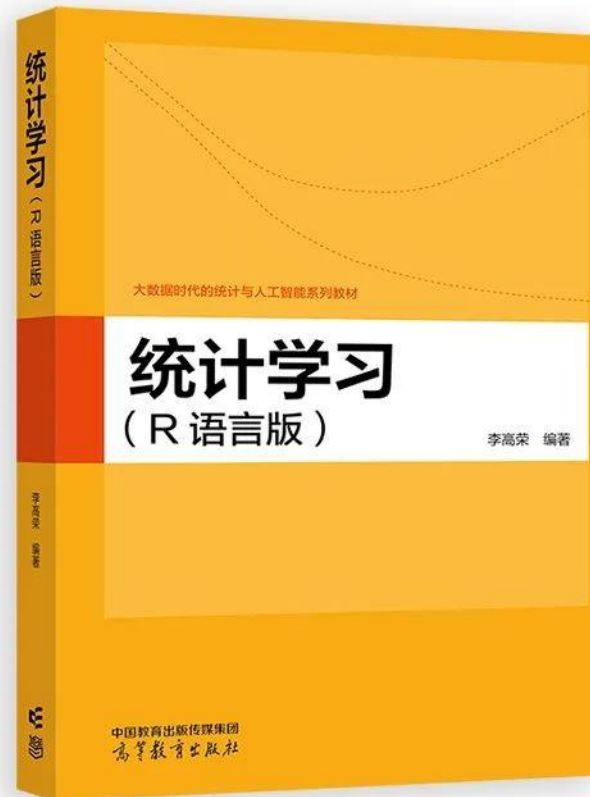
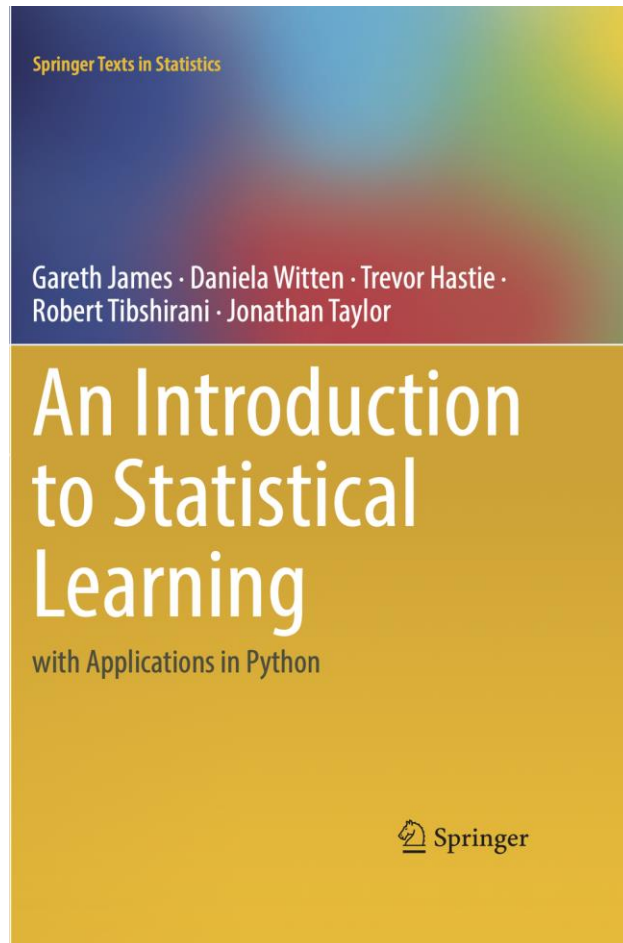
- Basic terminology and concepts
- Linear regression and classification
- Resampling methods
- Model selection and regularization
- Nonlinear methods
- Tree-based methods
- Support vector machines
- Neural networks and deep learning
- Unsupervised learning
- Oral presentation and quizzes

Textbooks & references



An Introduction of Statistical Learning

- by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani
- Book website: <https://www.statlearning.com/>
 - Data sets, .R/Rmarkdown files
 - Slides and video lectures
 - Book PDF and errata



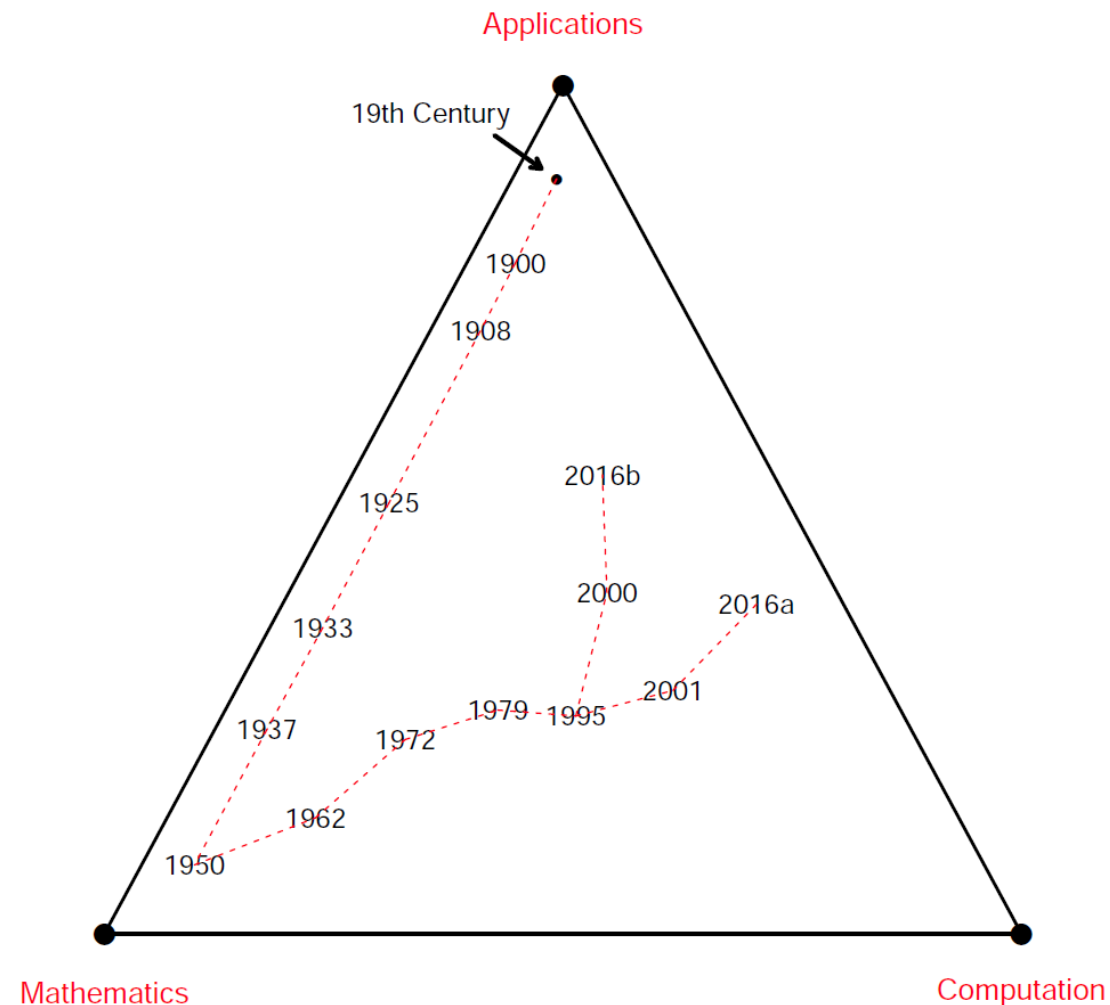
Grading policy (考核方法)

- Homework (30%)
 - 不按时交作业1次扣3分
- Group project (40%)
 - 大作业按贡献打分
- In-class quizz (30%)

Group project

- A proposal (May 08)
 - Questions/problems, data sets, methods
- Oral presentations (June 05)
- A report (June 08)

A brief history of statistical learning



Development of the statistics discipline since the end of the 19th century. This figure is taken from "*Computer Age Statistical Inference*" (Cambridge)

