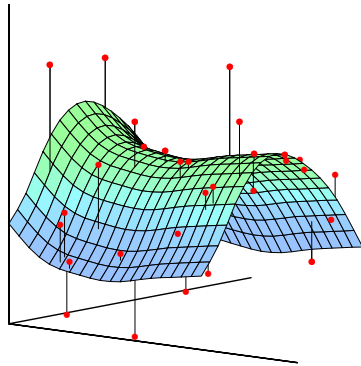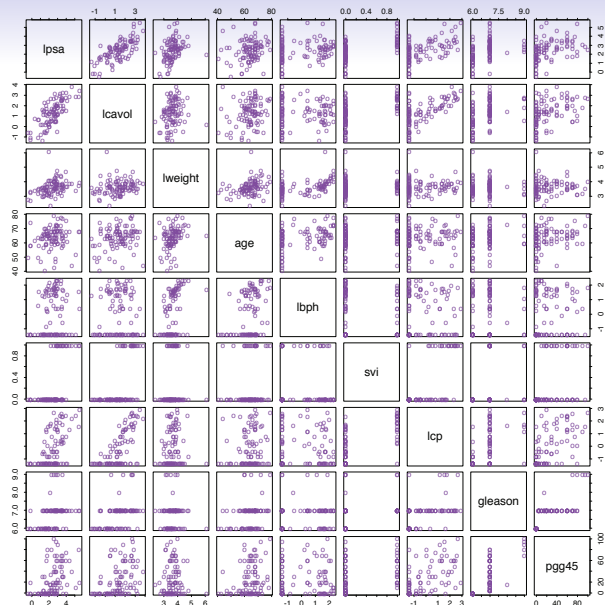# Statistical Learning



*Trevor Hastie and Robert Tibshirani*
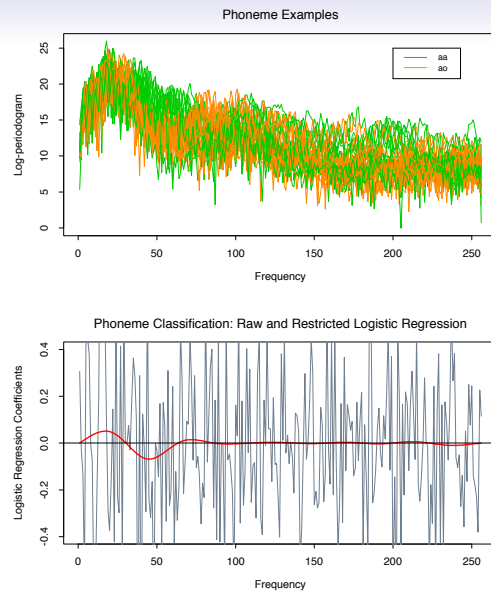
## Statistical Learning Problems

• Identify the risk factors for prostate cancer.

• Classify a recorded phoneme based on a log-periodogram.

• Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.

• Customize an email spam detection system.

• Identify the numbers in a handwritten zip code.

• Classify a tissue sample into one of several cancer classes, based on a gene expression profile.

• Establish the relationship between salary and demographic variables in population survey data.

• Classify the pixels in a LANDSAT image, by usage.

## Statistical Learning Problems

• Identify the risk factors for prostate cancer.

• Classify a recorded phoneme based on a log-periodogram.

• Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.

• Customize an email spam detection system.

• Identify the numbers in a handwritten zip code.

• Classify a tissue sample into one of several cancer classes, based on a gene expression profile.

• Establish the relationship between salary and demographic variables in population survey data.

• Classify the pixels in a LANDSAT image, by usage.

## Phoneme Examples

Log-periodogram vs Frequency

Legend: aa, ao

## Phoneme Classification: Raw and Restricted Logistic Regression

Logistic Regression Coefficients vs Frequency

---

# Statistical Learning Problems

- Identify the risk factors for prostate cancer.

- Classify a recorded phoneme based on a log-periodogram.

- **Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.**

- Customize an email spam detection system.

- Identify the numbers in a handwritten zip code.

- Classify a tissue sample into one of several cancer classes, based on a gene expression profile.

- Establish the relationship between salary and demographic variables in population survey data.

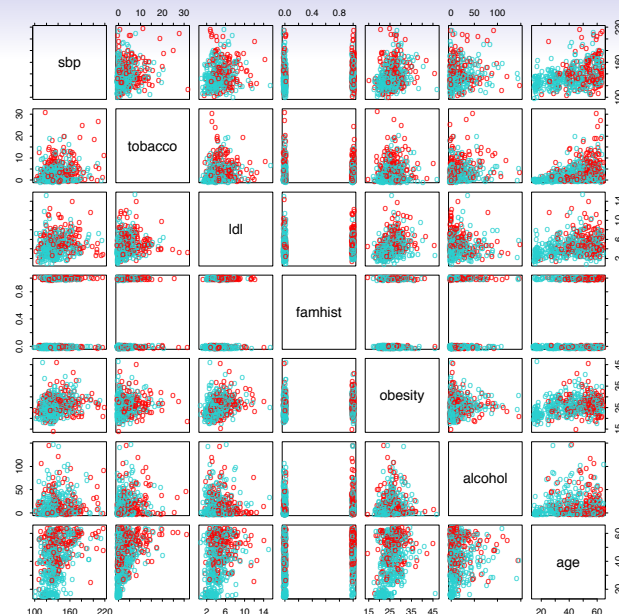- Classify the pixels in a LANDSAT image, by usage.

---

Scatterplot matrix of variables: sbp, tobacco, ldl, famhist, obesity, alcohol, age

---

# Statistical Learning Problems

- Identify the risk factors for prostate cancer.

- Classify a recorded phoneme based on a log-periodogram.

- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.

- **Customize an email spam detection system.**

- Identify the numbers in a handwritten zip code.

- Classify a tissue sample into one of several cancer classes, based on a gene expression profile.

- Establish the relationship between salary and demographic variables in population survey data.

- Classify the pixels in a LANDSAT image, by usage.

## Spam Detection

- data from 4601 emails sent to an individual (named George, at HP labs, before 2000). Each is labeled as *spam* or *email*.
- goal: build a customized spam filter.
- input features: relative frequencies of 57 of the most commonly occurring words and punctuation marks in these email messages.

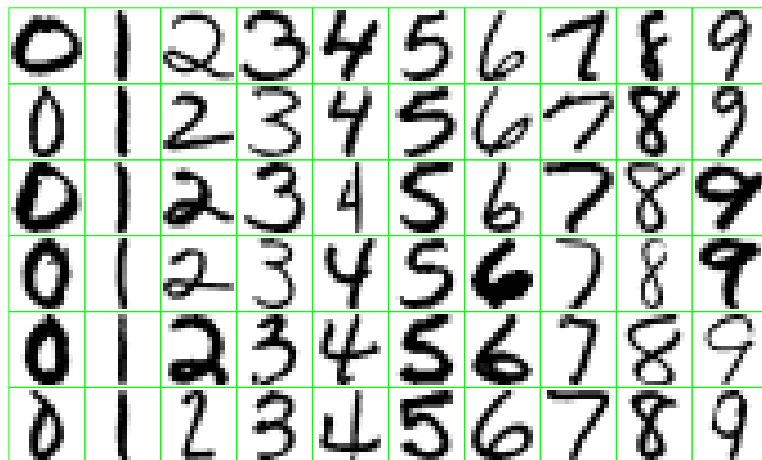|       | george | you  | hp   | free | !    | edu  | remove |
|-------|--------|------|------|------|------|------|--------|
| spam  | 0.00   | 2.26 | 0.02 | 0.52 | 0.51 | 0.01 | 0.28   |
| email | 1.27   | 1.27 | 0.90 | 0.07 | 0.11 | 0.29 | 0.01   |

*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.

## Statistical Learning Problems

- Identify the risk factors for prostate cancer.
- Classify a recorded phoneme based on a log-periodogram.
- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.
- Customize an email spam detection system.
- Identify the numbers in a handwritten zip code.
- Classify a tissue sample into one of several cancer classes, based on a gene expression profile.
- Establish the relationship between salary and demographic variables in population survey data.
- Classify the pixels in a LANDSAT image, by usage.

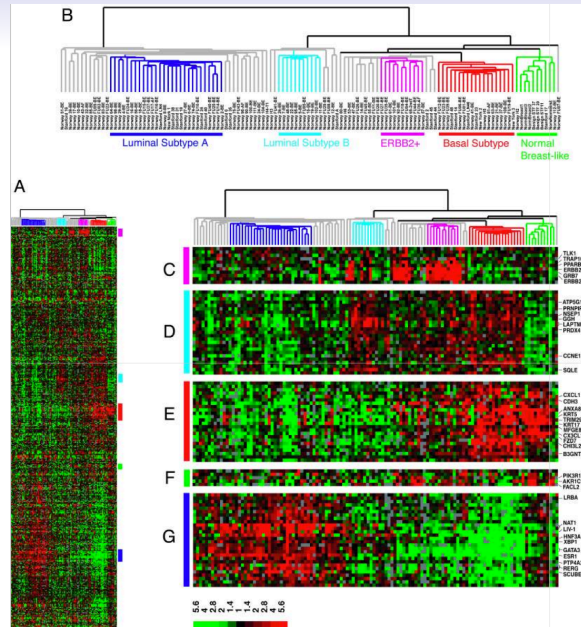## Statistical Learning Problems

- Identify the risk factors for prostate cancer.
- Classify a recorded phoneme based on a log-periodogram.
- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.
- Customize an email spam detection system.
- Identify the numbers in a handwritten zip code.
- Classify a tissue sample into one of several cancer classes, based on a gene expression profile.
- Establish the relationship between salary and demographic variables in population survey data.
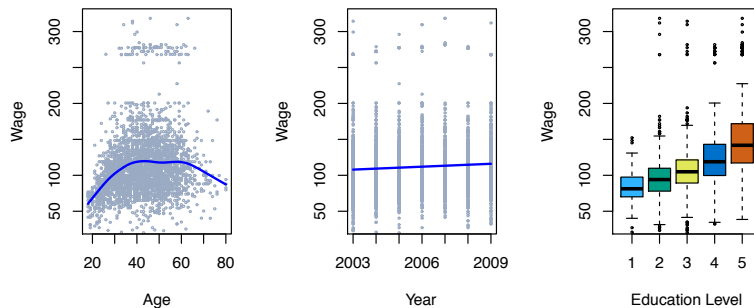- Classify the pixels in a LANDSAT image, by usage.

## Statistical Learning Problems

- Identify the risk factors for prostate cancer.

- Classify a recorded phoneme based on a log-periodogram.

- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.

- Customize an email spam detection system.

- Identify the numbers in a handwritten zip code.

- Classify a tissue sample into one of several cancer classes, based on a gene expression profile.

- Establish the relationship between salary and demographic variables in population survey data.
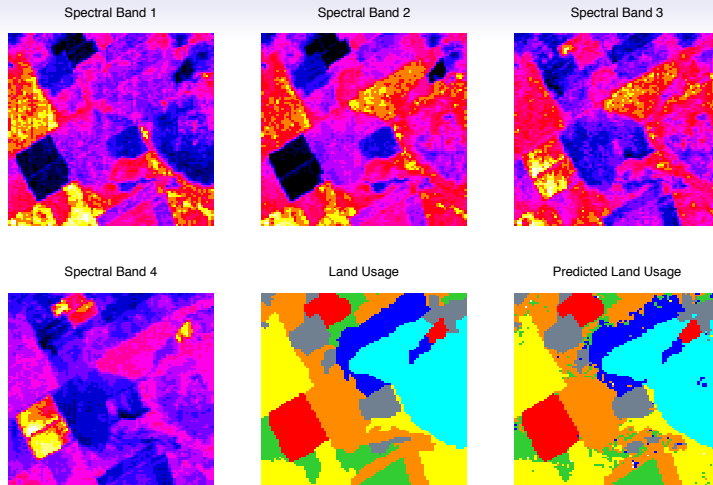
- Classify the pixels in a LANDSAT image, by usage.

Income survey data for males from the central Atlantic region of the USA in 2009.

## Statistical Learning Problems

- Identify the risk factors for prostate cancer.

- Classify a recorded phoneme based on a log-periodogram.

- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.

- Customize an email spam detection system.

- Identify the numbers in a handwritten zip code.

- Classify a tissue sample into one of several cancer classes, based on a gene expression profile.

- Establish the relationship between salary and demographic variables in population survey data.

- Classify the pixels in a LANDSAT image, by usage.

Spectral Band 1    Spectral Band 2    Spectral Band 3

Spectral Band 4    Land Usage    Predicted Land Usage

*Usage* ∈ {*red soil, cotton, vegetation stubble, mixture, gray soil, damp gray soil*}

# The Supervised Learning Problem

*Starting point:*

- Outcome measurement $Y$ (also called dependent variable, response, target).
- Vector of $p$ predictor measurements $X$ (also called inputs, regressors, covariates, features, independent variables).
- In the *regression problem*, $Y$ is quantitative (e.g price, blood pressure).
- In the *classification problem*, $Y$ takes values in a finite, unordered set (survived/died, digit 0-9, cancer class of tissue sample).
- We have training data $(x_1, y_1), \ldots, (x_N, y_N)$. These are observations (examples, instances) of these measurements.

# Objectives

On the basis of the training data we would like to:

- Accurately predict unseen test cases.
- Understand which inputs affect the outcome, and how.
- Assess the quality of our predictions and inferences.

# Philosophy

- It is important to understand the ideas behind the various techniques, in order to know how and when to use them.
- One has to understand the simpler methods first, in order to grasp the more sophisticated ones.
- It is important to accurately assess the performance of a method, to know how well or how badly it is working [simpler methods often perform as well as fancier ones!]
- This is an exciting research area, having important applications in science, industry and finance.
- Statistical learning is a fundamental ingredient in the training of a modern *data scientist.*

## Unsupervised learning

- No outcome variable, just a set of predictors (features) measured on a set of samples.
- objective is more fuzzy — find groups of samples that behave similarly, find features that behave similarly, find linear combinations of features with the most variation.
- difficult to know how well your are doing.
- different from supervised learning, but can be useful as a pre-processing step for supervised learning.

## The Netflix prize

- competition started in October 2006. Training data is ratings for $18,000$ movies by $400,000$ Netflix customers, each rating between 1 and 5.
- training data is very sparse— about 98% missing.
- objective is to predict the rating for a set of 1 million customer-movie pairs that are missing in the training data.
- Netflix's original algorithm achieved a root MSE of 0.953. The first team to achieve a 10% improvement wins one million dollars.
- is this a supervised or unsupervised problem?

| Rank | Team Name | Best Test Score | % Improvement | Best Submit Time |
|---|---|---|---|---|
| \multicolumn{5}{|c|}{Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos} |
| 1 | BellKor's Pragmatic Chaos | 0.8567 | 10.06 | 2009-07-26 18:18:28 |
| 2 | The Ensemble | 0.8567 | 10.06 | 2009-07-26 18:38:22 |
| 3 | Grand Prize Team | 0.8582 | 9.90 | 2009-07-10 21:24:40 |
| 4 | Opera Solutions and Vandelay United | 0.8588 | 9.84 | 2009-07-10 01:12:31 |
| 5 | Vandelay Industries ! | 0.8591 | 9.81 | 2009-07-10 00:32:20 |
| 6 | PragmaticTheory | 0.8594 | 9.77 | 2009-06-24 12:06:56 |
| 7 | BellKor in BigChaos | 0.8601 | 9.70 | 2009-05-13 08:14:09 |
| 8 | Dace_ | 0.8612 | 9.59 | 2009-07-24 17:18:43 |
| 9 | Feeds2 | 0.8622 | 9.48 | 2009-07-12 13:11:51 |
| 10 | BigChaos | 0.8623 | 9.47 | 2009-04-07 12:33:59 |
| 11 | Opera Solutions | 0.8623 | 9.47 | 2009-07-24 00:34:07 |
| 12 | BellKor | 0.8624 | 9.46 | 2009-07-26 17:19:11 |

BellKor's Pragmatic Chaos wins, beating The Ensemble by a narrow margin.

## Statistical Learning versus Machine Learning

- Machine learning arose as a subfield of Artificial Intelligence.
- Statistical learning arose as a subfield of Statistics.
- *There is much overlap* — both fields focus on supervised and unsupervised problems:
  - Machine learning has a greater emphasis on *large scale* applications and *prediction accuracy*.
  - Statistical learning emphasizes *models* and their interpretability, and *precision* and *uncertainty*.
- But the distinction has become more and more blurred, and there is a great deal of "cross-fertilization".
- Machine learning has the upper hand in *Marketing!*