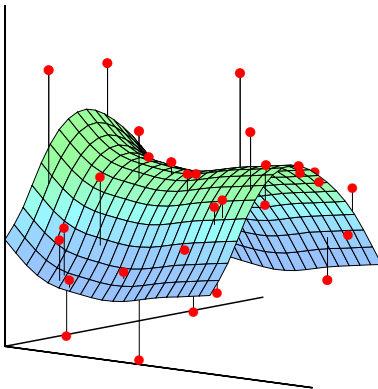


# Statistical Learning



*Trevor Hastie and Robert Tibshirani*

## Statistics in the news

How IBM built Watson, its *Jeopardy!*-playing supercomputer by [Dawn Kawamoto](#) DailyFinance

02/08/2011



**Learning from its mistakes** According to David Ferrucci (PI of Watson DeepQA technology for IBM Research), Watson's software is wired for more than handling natural language processing.

*"It's [machine learning](#) allows the computer to become smarter as it tries to answer questions — and to learn as it gets them right or wrong."*

## For Today's Graduate, Just One Word: Statistics

By [STEVE LOHR](#)

Published: August 5, 2009

MOUNTAIN VIEW, Calif. — At Harvard, Carrie Grimes majored in anthropology and archaeology and ventured to places like Honduras, where she studied Mayan settlement patterns by mapping where artifacts were found. But she was drawn to what she calls “all the computer and math stuff” that was part of the job.

[Enlarge This Image](#)



Thor Swift for The New York Times  
Carrie Grimes, senior staff engineer at Google, uses statistical analysis of data to help improve the company's search engine.

### Multimedia



“People think of field archaeology as Indiana Jones, but much of what you really do is data analysis,” she said.

Now Ms. Grimes does a different kind of digging. She works at [Google](#), where she uses statistical analysis of mounds of data to come up with ways to improve its search engine.

Ms. Grimes is an Internet-age statistician, one of many who are changing the image of the profession as a place for dronish number nerds. They are finding themselves increasingly in demand — and even cool.

“I keep saying that the sexy job in the next 10 years will be statisticians,” said Hal Varian, chief economist at Google. “And I’m not kidding.”

[SIGN IN TO  
RECOMMEND](#)

[SIGN IN TO  
E-MAIL](#)

[PRINT](#)

[REPRINTS](#)

[SHARE](#)



QUOTE OF THE DAY,  
NEW YORK TIMES,  
AUGUST 5, 2009

“I keep saying that the sexy job in the next 10 years will be statisticians. And I’m not kidding.”  
— HAL VARIAN, chief economist at Google.

# FiveThirtyEight

Nate Silver's Political Calculus

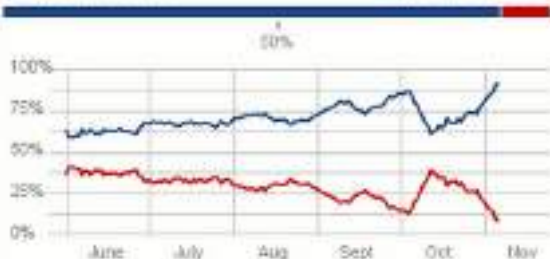
90.9%

+13.5 since Oct. 30

Chance of  
Winning

9.1%

-13.5 since Oct. 30

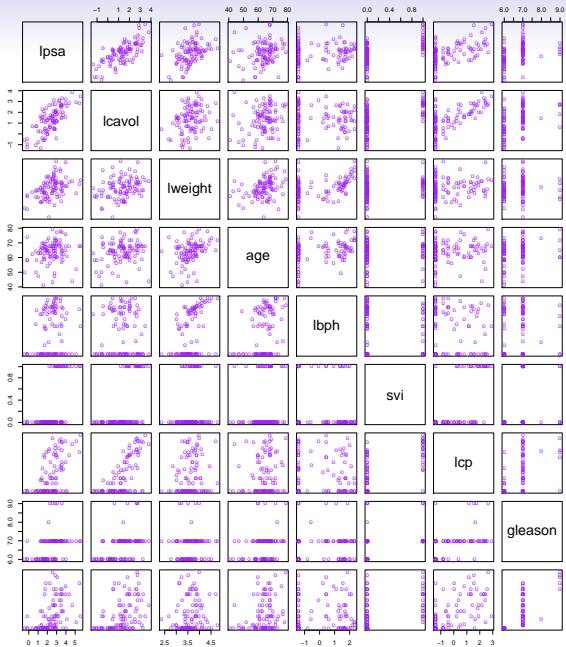


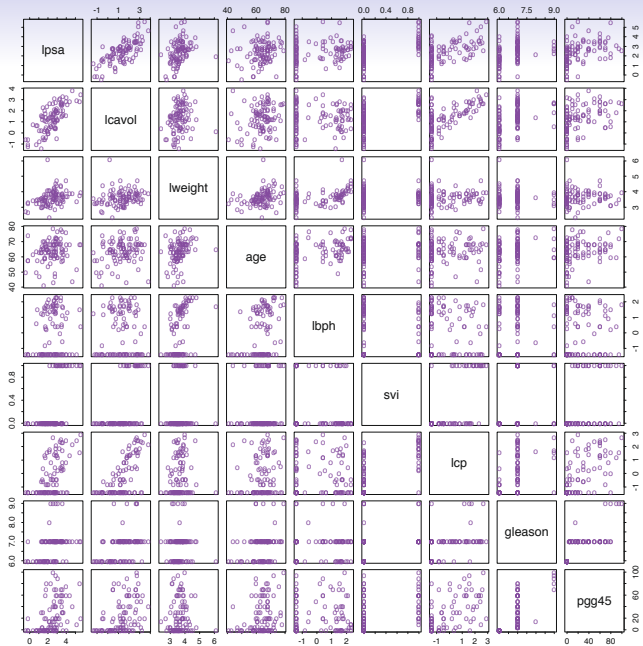
ask a LIDOK model

the signal and the noise  
the noise and the signal  
why so many predictions fail -  
but some don't  
and the noise and the  
nate silver

# 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.



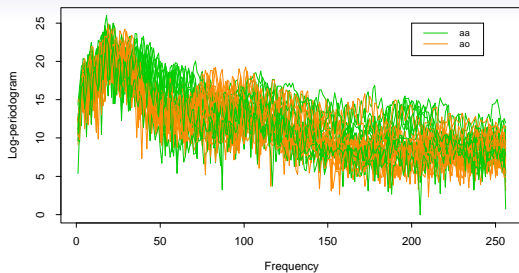


# Statistical Learning Problems

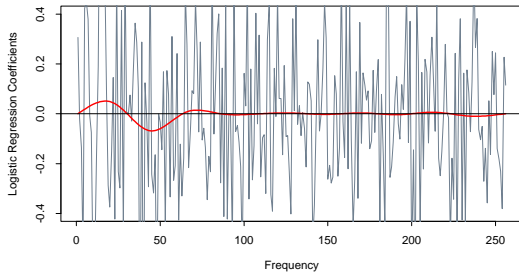
- Identify the risk factors for prostate cancer.
- Classify a recorded phoneme based on a log-periodogram.
- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.
- Customize an email spam detection system.
- Identify the numbers in a handwritten 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

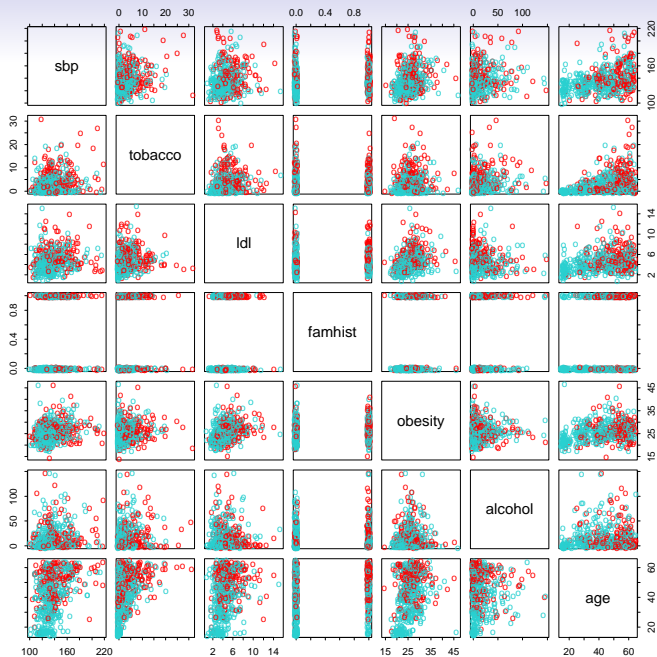


Phoneme Classification: Raw and Restricted Logistic Regression



# 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.

## 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.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

# Statistical Learning Problems

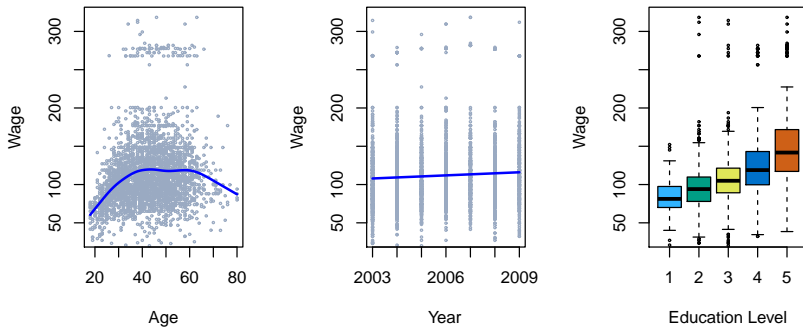
- Identify the risk factors for prostate cancer.
- Classify a recorded phoneme based on a log-periodogram.
- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements.
- Customize an email spam detection system.
- Identify the numbers in a handwritten 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.

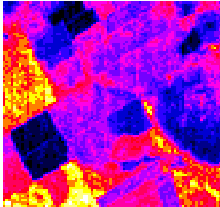


Income survey data for males from the central Atlantic region of the USA in 2009.

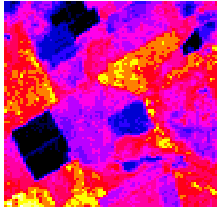
# Statistical Learning Problems

- 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.

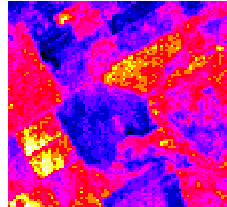
Spectral Band 1



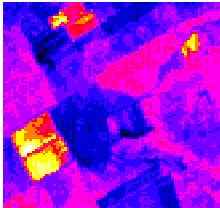
Spectral Band 2



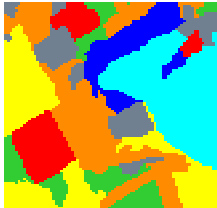
Spectral Band 3



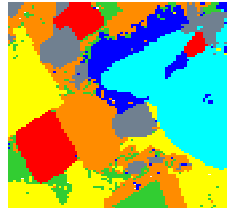
Spectral Band 4



Land Usage



Predicted Land Usage



$Usage \in \{\text{red soil, cotton, vegetation stubble, mixture, gray soil, damp gray soil}\}$

# 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  $(x_1, y_1), \dots, (x_N, y_N)$ . 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?

## Netflix Prize

COMPLETED

[Home](#)
[Rules](#)
[Leaderboard](#)
[Updates](#)

## Leaderboard

Showing Test Scores. [Click here to show our 2009](#)Display top  leaders.

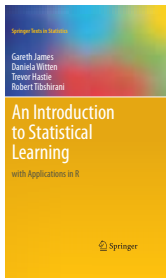
Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize: \$1,000,000 - \$2,000 - Winning Team: BellKor's Pragmatic Chaos				
1	<a href="#">BellKor's Pragmatic Chaos</a>	0.8567	10.00	2009-07-26 18:18:28
2	<a href="#">The Ensemble</a>	0.8567	10.00	2009-07-26 18:38:22
3	<a href="#">Grand Prize Team</a>	0.8582	0.06	2009-07-10 21:24:40
4	<a href="#">Coresolutions and Verbelive United</a>	0.8586	0.04	2009-07-10 01:12:31
5	<a href="#">Vandelay Industries</a>	0.8591	0.01	2009-07-10 08:52:20
6	<a href="#">PhishKitTives</a>	0.8594	0.17	2009-06-24 12:08:58
7	<a href="#">BellKor's RealChaos</a>	0.8601	0.16	2009-05-13 08:14:09
8	<a href="#">Cores...</a>	0.8612	0.58	2009-07-26 17:18:43
9	<a href="#">Ensemble</a>	0.8623	0.48	2009-07-12 13:11:51
10	<a href="#">RealChaos</a>	0.8623	0.47	2009-04-07 12:33:59
11	<a href="#">Cores Solutions</a>	0.8623	0.47	2009-07-26 08:34:57
12	<a href="#">BellKor</a>	0.8624	0.48	2009-07-26 17:18:11

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

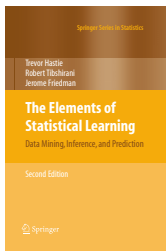
# Statistical Learning versus Machine Learning

- 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!*

## Course Texts



The course will cover most of the material in this Springer book (ISLR) published in 2013, which the instructors coauthored with Gareth James and Daniela Witten. Each chapter ends with an R lab, in which examples are developed. By January 1st, 2014, an electronic version of this book will be available for free from the instructors' websites.



This Springer book (ESL) is more mathematically advanced than ISLR; the second edition was published in 2009, and coauthored by the instructors and Jerome Friedman. It covers a broader range of topics. The book is available from Springer and Amazon, a free electronic version is available from the instructors' websites.