# Stat315a: Statistical Learning

# Statistics in the news

## How IBM built Watson, its "Jeopardy"-playing supercomputer by Dawn Kawamoto DailyFinance 02/08/2011
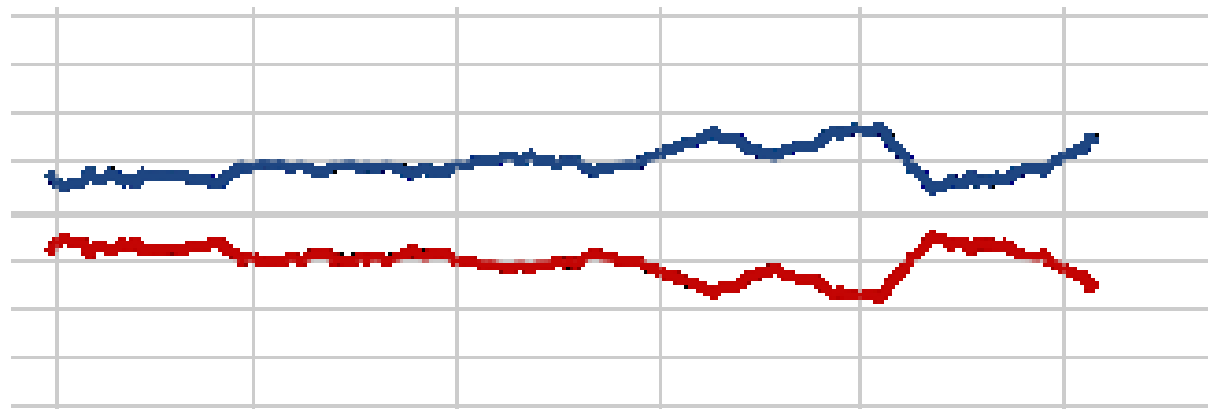


**Learning from its mistakes** According to David Ferrucci (PI of Watson DeepQA technology for IBM Research), Watson's software is wired for more that handling natural language processing.

"Its machine learning allows the computer to become smarter as it tries to answer questions — and to learn as it gets them right or wrong."

Click to LOOK INSIDE!

the signal and the noise
and the noise and the
the noise and the
noise and the no

why so many and
predictions fail—
but some don't th
and the noise an
the noise and the

nate silver noise
noise and the no

# 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."

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.
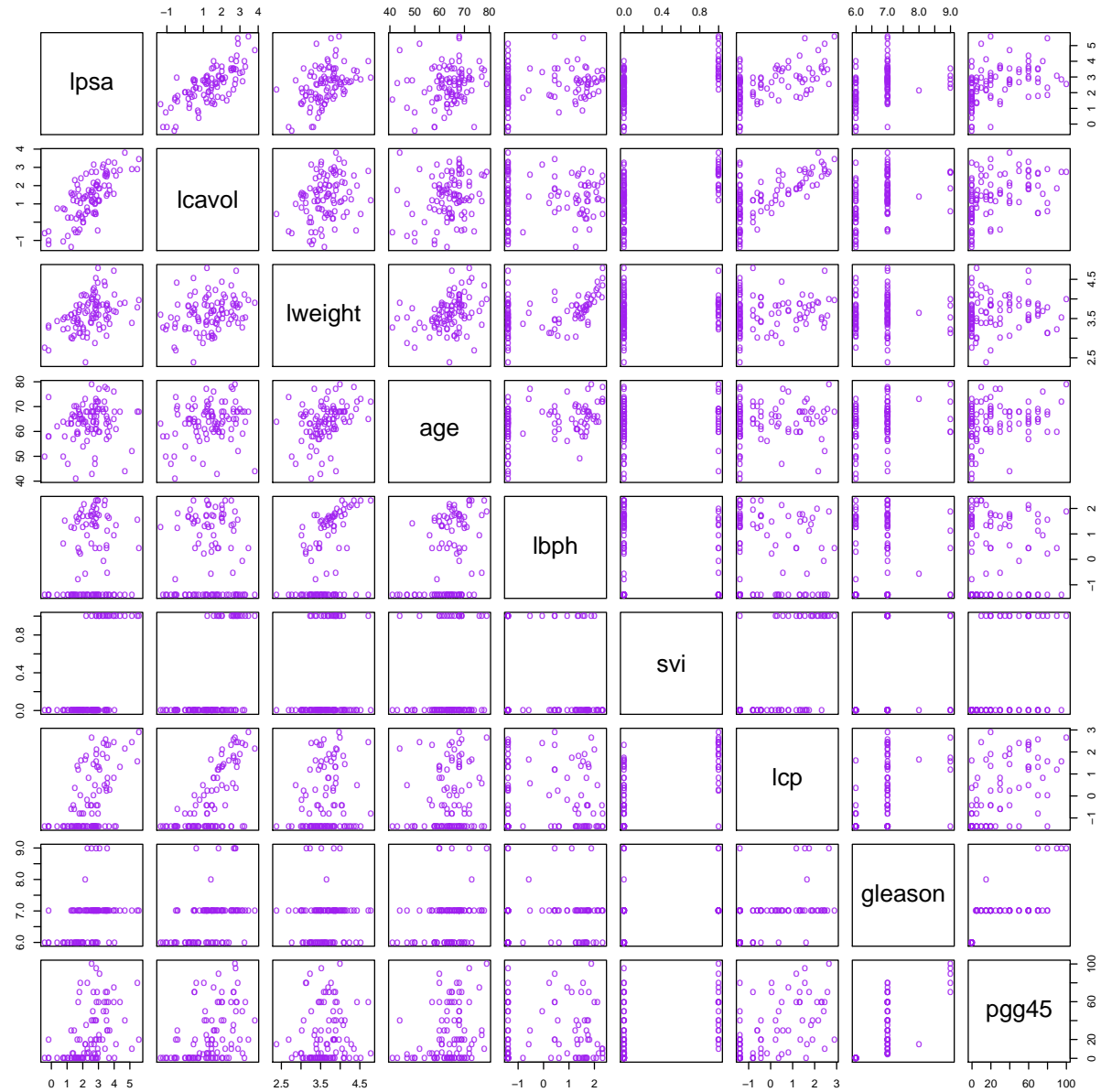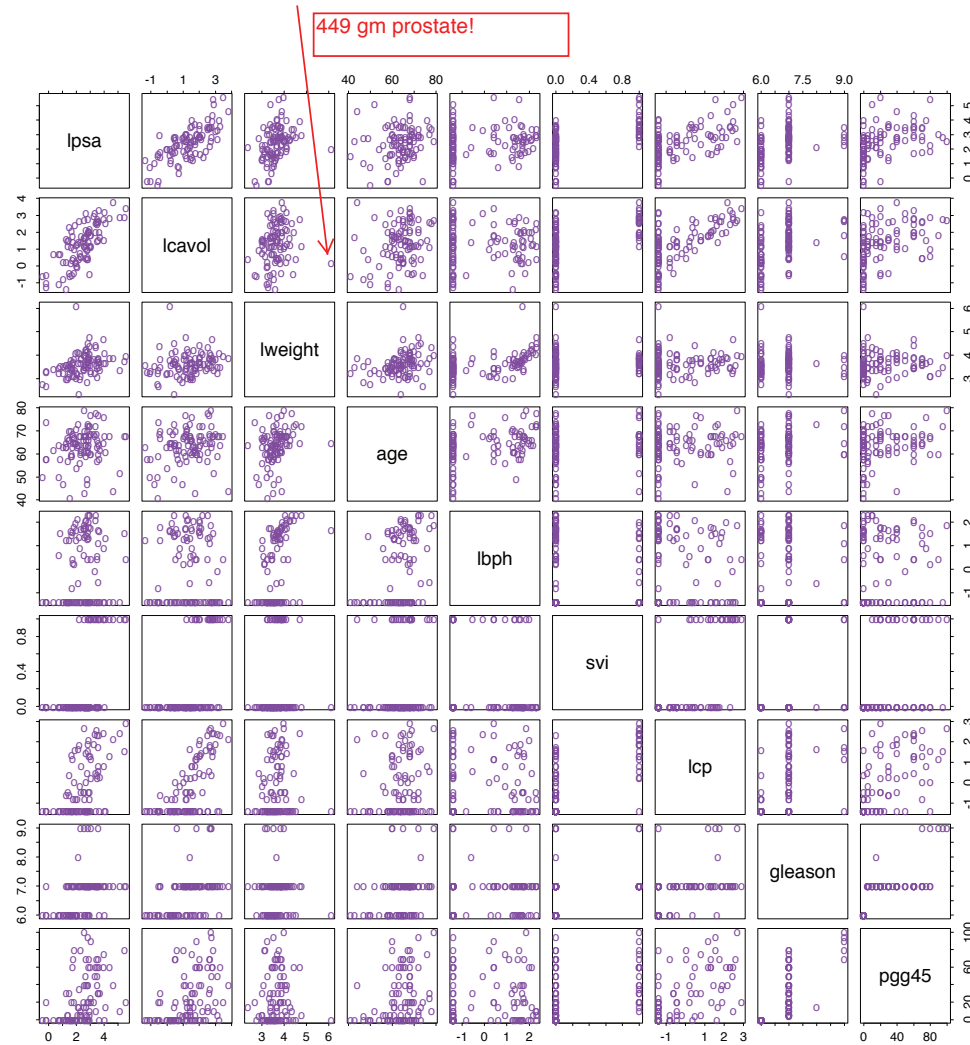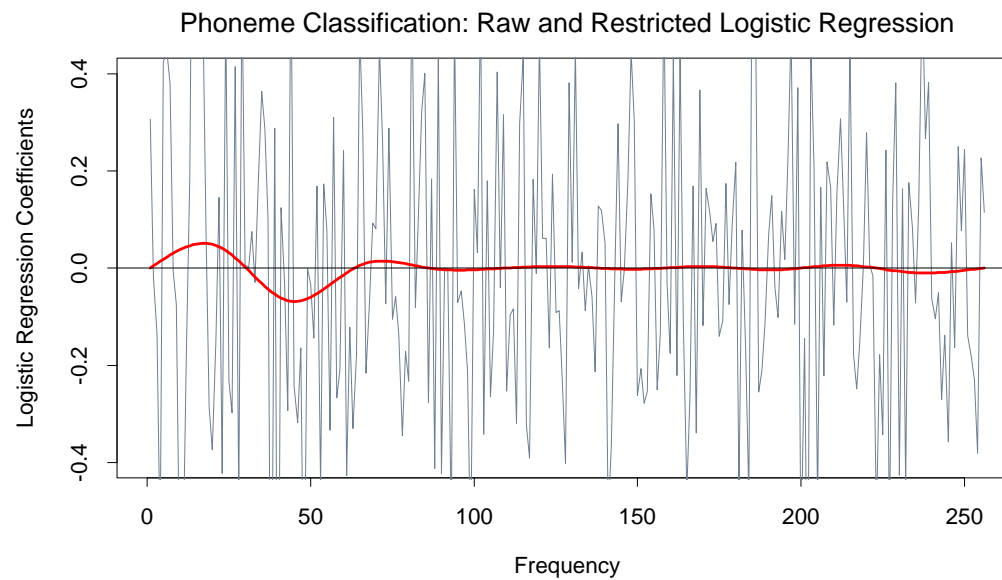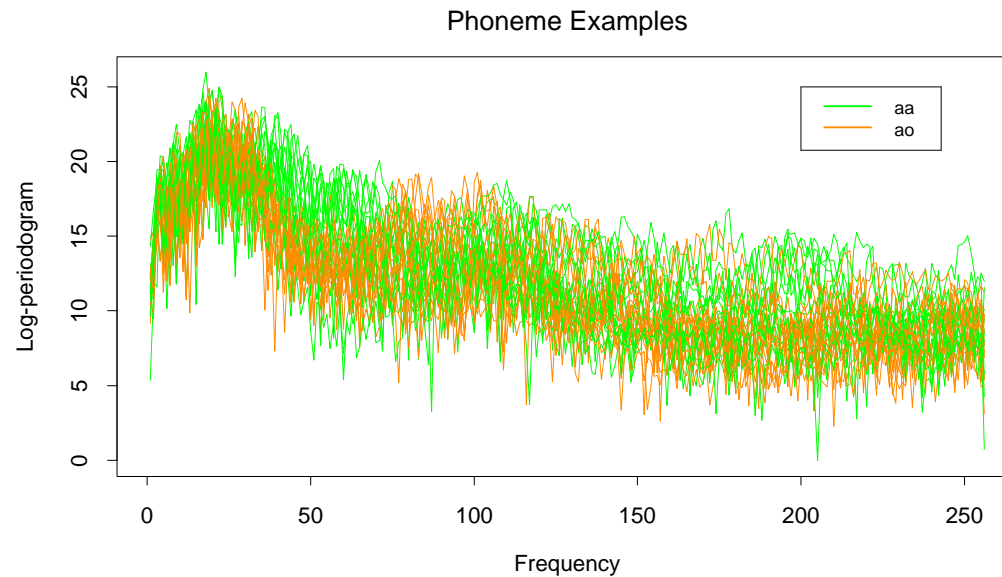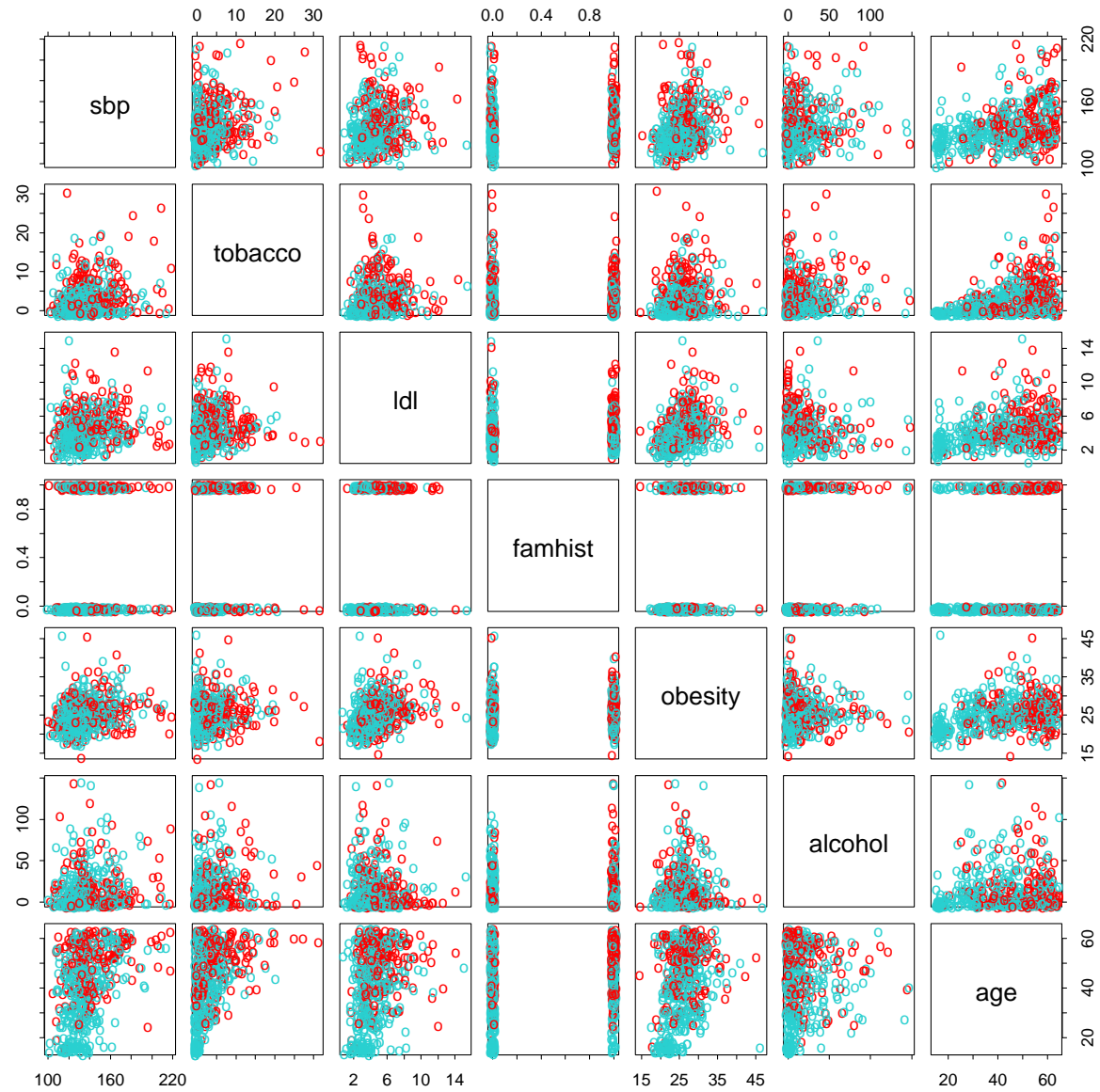
5

OLD



NEW

# Statistical Learning Problems

- Identify the risk factors for prostate cancer (Fig 1.1).

- Classify a recorded phoneme (Fig 5.5) based on a log-periodogram.

- Predict whether someone will have a heart attack (Fig 4-12) on the basis of demographic, diet and clinical measurements

- Customize an email spam (Tab 1.1) detection system.

- Identify the numbers in a handwritten zip code (Fig 1.2), from a digitized image

- Classify a tissue sample into one of several cancer classes, based on a gene expression (Fig 1.3) profile

- Classify the pixels in a LANDSAT (Fig 13.6) image, by usage:

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

449 gm prostate!

Error pointed out by retired urologist Dr Stephen W. Link

Phoneme Examples



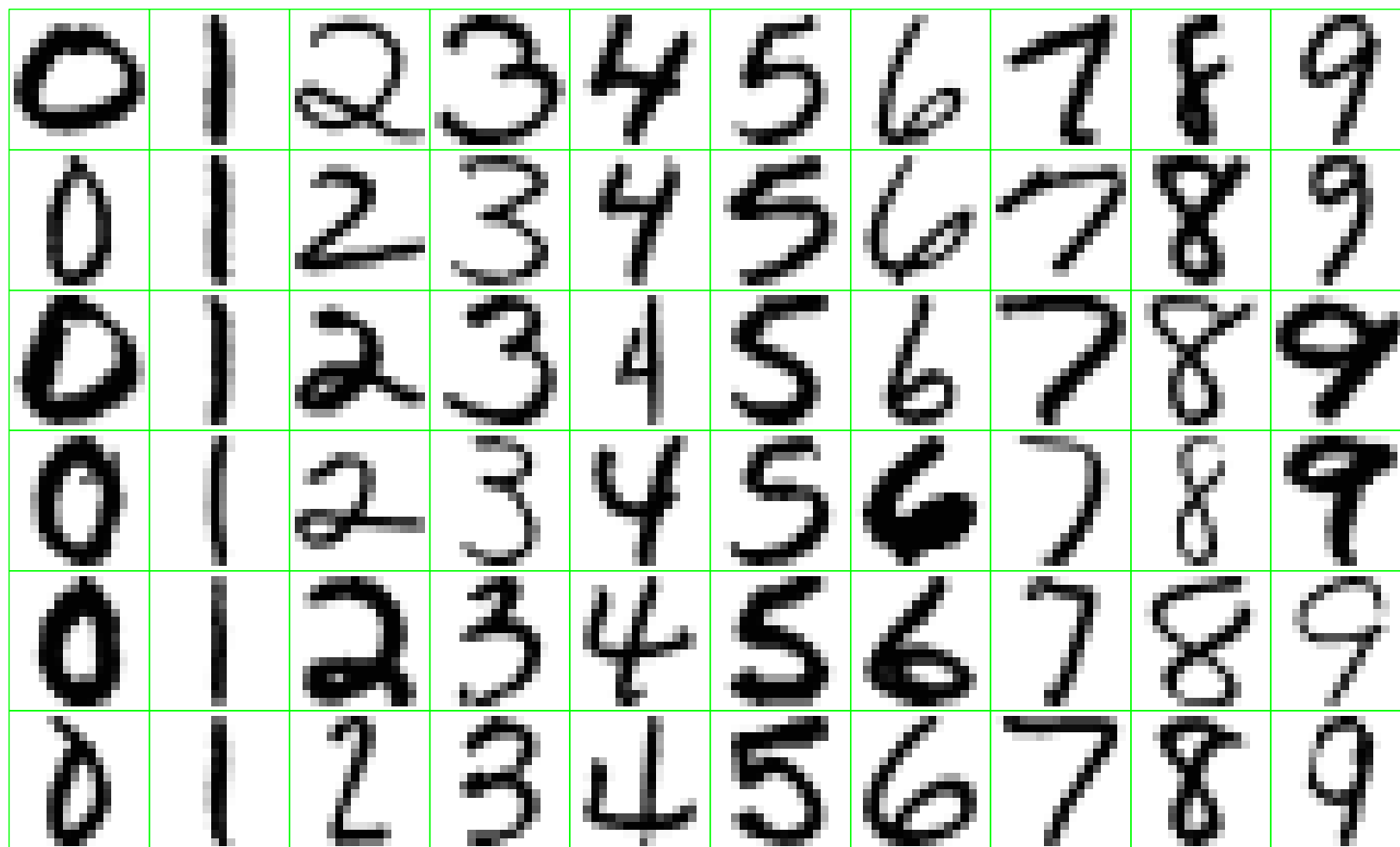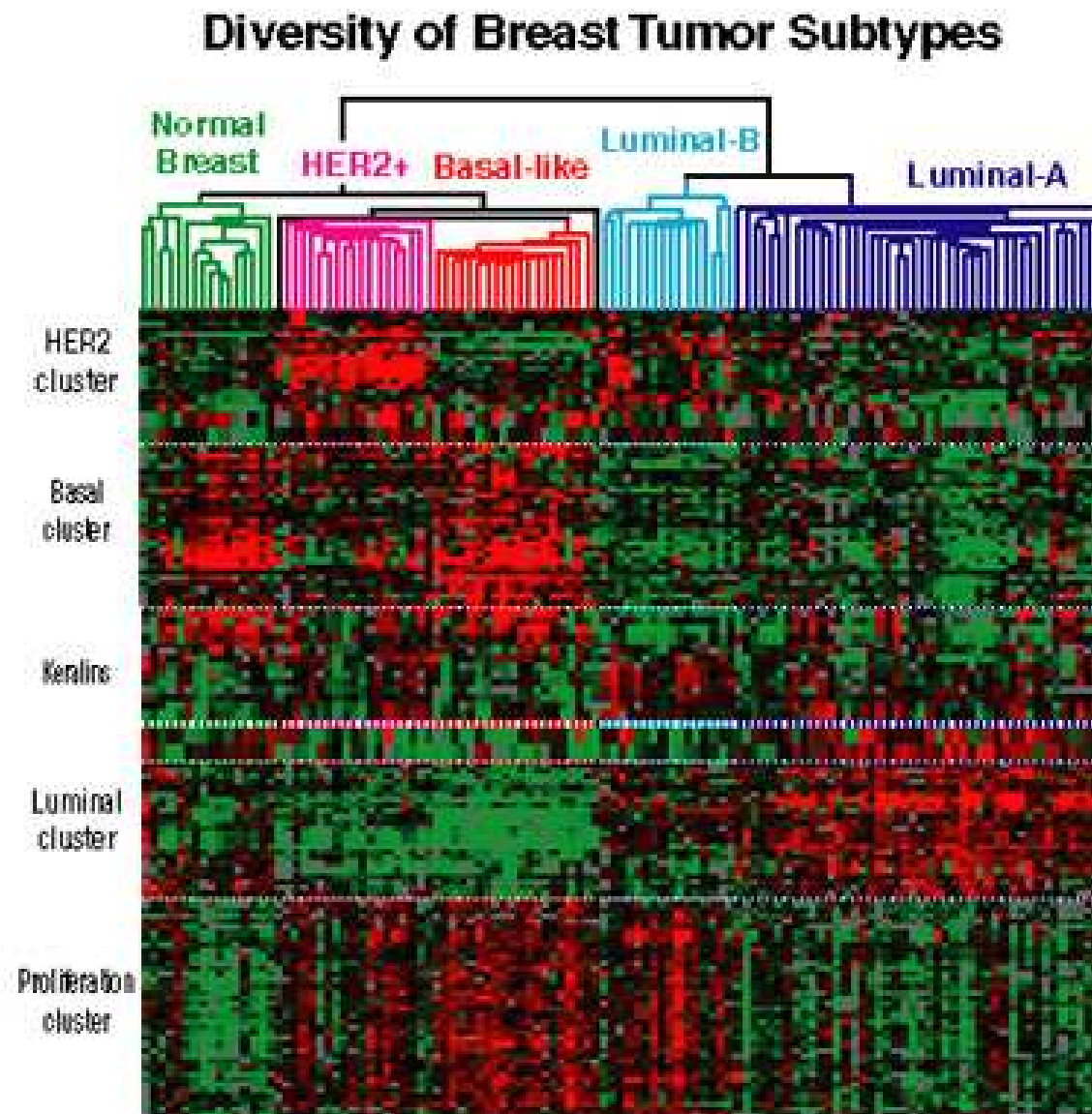Phoneme Classification: Raw and Restricted Logistic Regression

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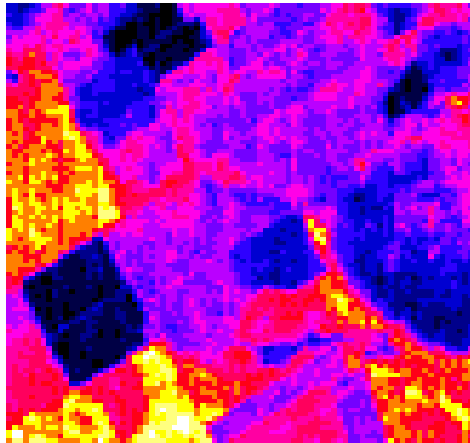
| | george | you | hp | free | ! | edu | remove |
|---|---|---|---|---|---|---|---|
| spam | 0.00 | 2.26 | 0.02 | 0.52 | 0.51 | 0.01 | 0.28 |
| email | 1.27 | 1.27 | 0.90 | 0.07 | 0.11 | 0.29 | 0.01 |

*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.
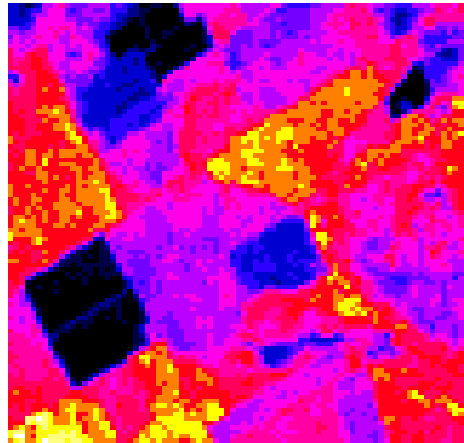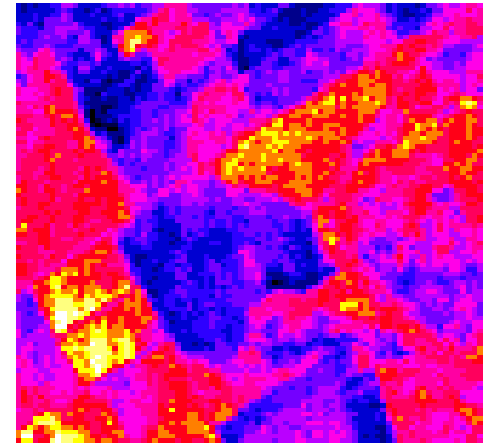
Diversity of Breast Tumor Subtypes
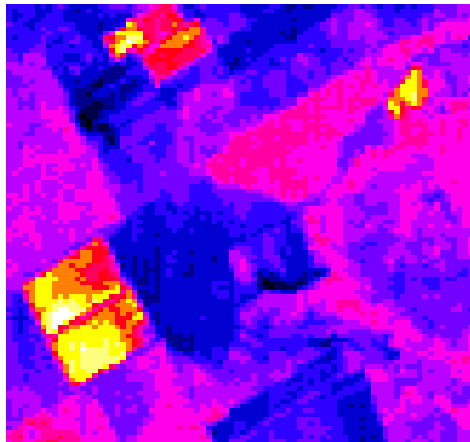
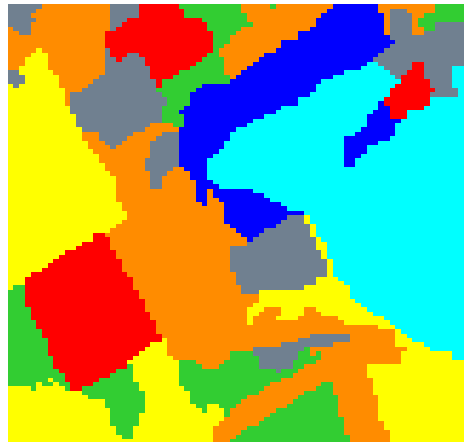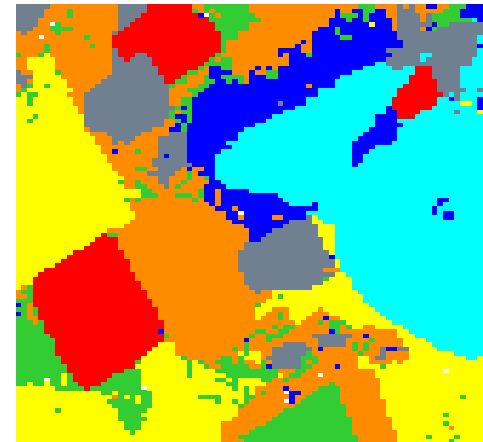Spectral Band 1    Spectral Band 2    Spectral Band 3

Spectral Band 4    Land Usage    Predicted Land Usage

# 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
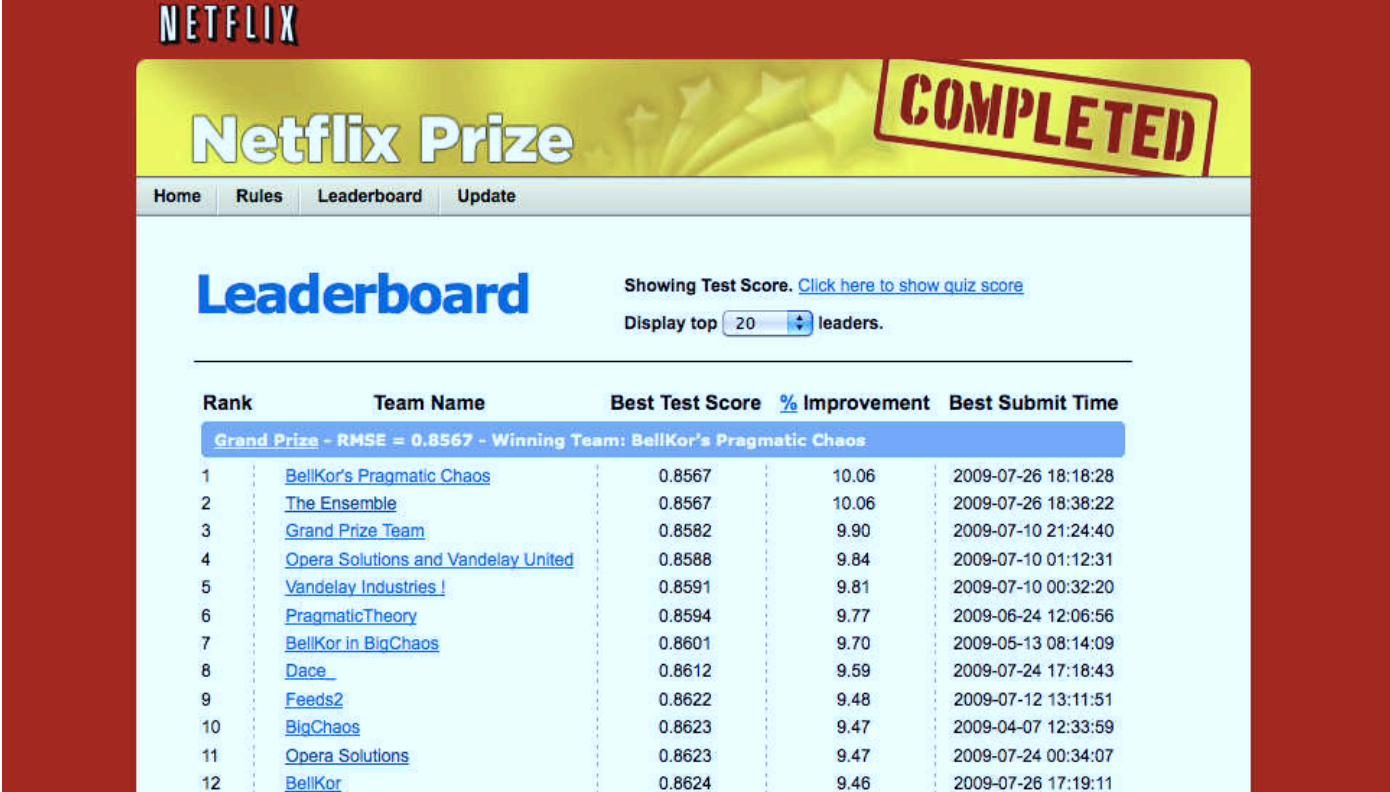
# Philosophy

- 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.

# 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 1 million dollars.

- is this a supervised or unsupervised problem?