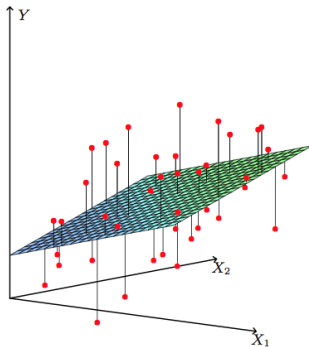# Data Mining and Machine Learning

Kuangnan Fang

Department of Statistics, Xiamen University
*Email: xmufkn@xmu.edu.cn*

# 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.

*"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.

"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."

**Multimedia**



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.

| Cat | Dog |

Deep learning

Example:
MLPs

Example:
Shallow
autoencoders

Representation learning

Example:
Logistic
regression

Example:
Knowledge
bases

Machine learning

AI

# 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!*

# 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.
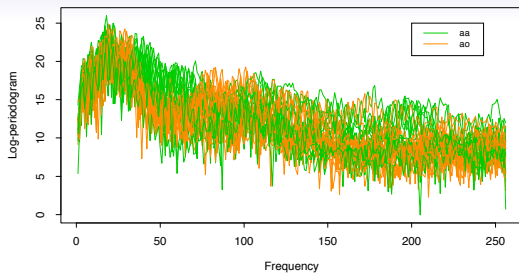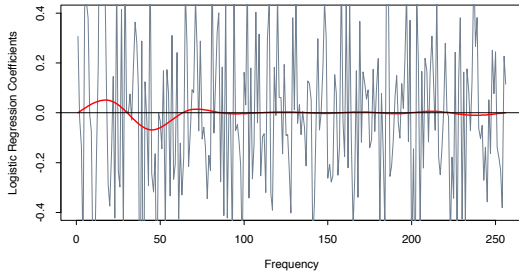
# Statistical Learning Problems

- 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.

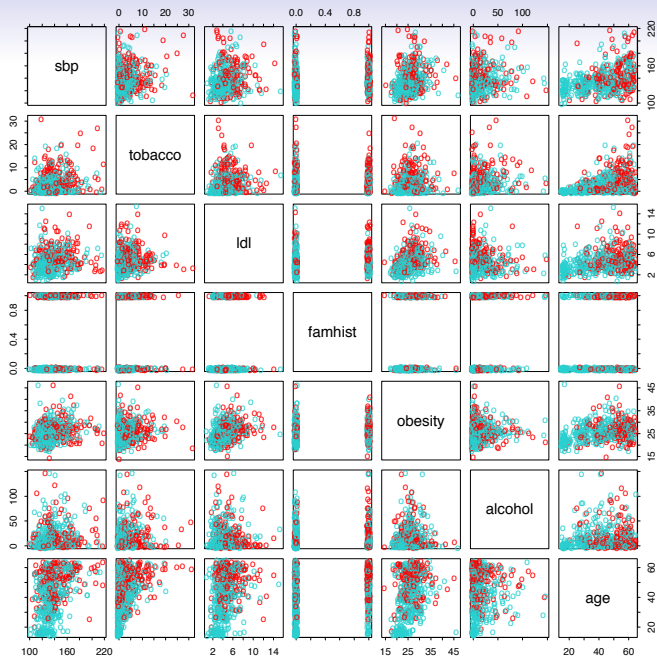# Statistical Learning Problems

- 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.

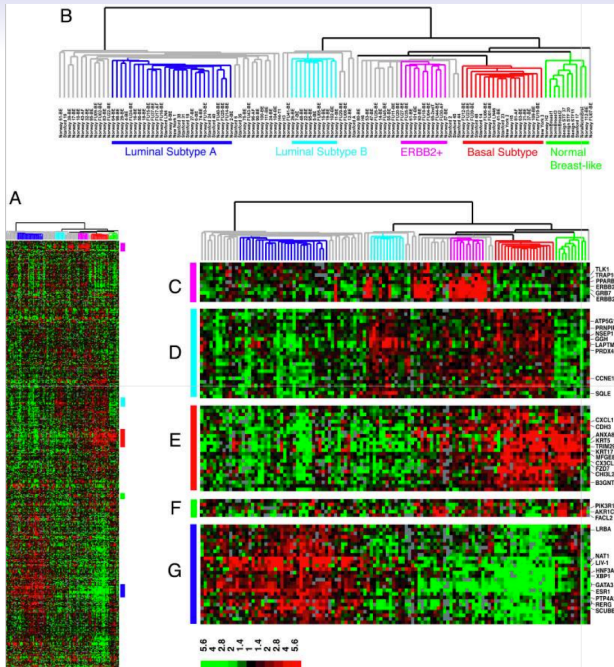- Classify the pixels in a LANDSAT image, by usage.

# 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.

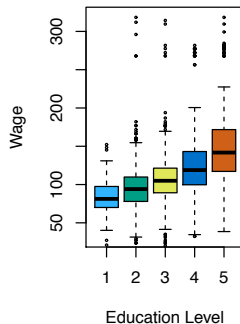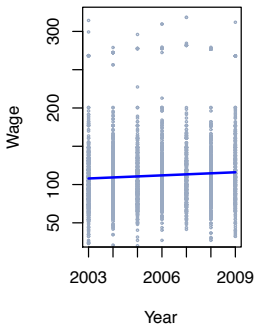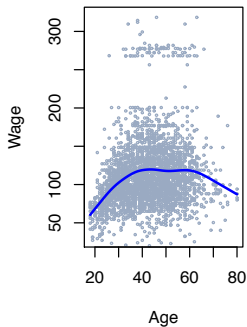- Classify the pixels in a LANDSAT image, by usage.

# 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.

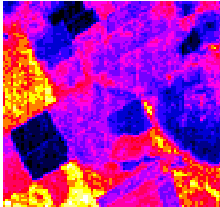- Classify the pixels in a LANDSAT image, by usage.

# 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`.
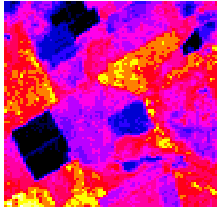
# Statistical Learning Problems

- 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.

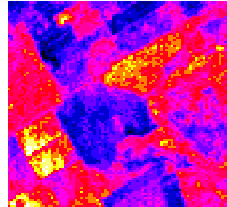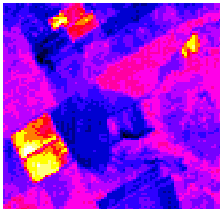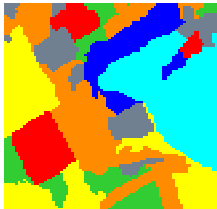- Classify the pixels in a LANDSAT image, by usage.

# 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

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

# Graphical network

17,214 gene expressions and 22,247 CNVs are available from TCGA(X.Fan,K.Fang,S.Ma, Q. Zhang. Assisted Graphical Model for Gene Expression Data Analysis. Statistics in Medicine. 2019).



Top left panel: AGM; Top right panel: GM; Bottom left panel: difference between AGM and GM; Bottom middle: difference between moderate connections ; Bottom right: difference between strong connections.

The FICO score was first introduced in 1989 by Fair, Isaac, and Company, whic is used by the vast majority of banks and credit grantors, and is based on consumer credit files of the three national credit bureaus: Experian, Equifax, and TransUnion.



**30%** Amounts Owed
**10%** New Credit
**FICO'SCORE**
**15%** Length of Credit History
**35%** Payment History
**10%** Credit Mix

### ACCOUNT SUMMARY

| | | |
|---|---|---|
| Previous Balance | | $42.68 |
| Payments and Credits | - | $712.83 |
| Purchases | + | $670.15 |
| Balance Transfers | + | $0.00 |
| Cash Advances | + | $0.00 |
| Fees Charged | + | $0.00 |
| Interest Charged | + | $0.00 |
| New Balance | | $0.00 |

See Interest Charge Calculation section following the Transactions section for detailed APR information

| | |
|---|---|
| Credit Line | $1,900 |
| Credit Line Available | $1,900 |
| Cash Advance Credit Line | $400 |
| Cash Advance Credit Line Available | $400 |

You may be able to avoid interest on Purchases. See reverse for details.

**FICO 724**
Your FICO® Credit Score on 8/27/15.

Track recent scores on your FICO® page in this statement.

Please pay online at www.Discover.com or make checks payable to Discover. Phone or internet payment# Pay before midnight ET on your payment due date for same day credit.

# Credit score in China

# Interesting problem for Credit scoring

- Variable selection
- reject inference
- reject option
- imbalance data
- label noise
- missing data
- fraud detection

# Data integratoin

Two type of Data integration (Data fushion): Sample integration and Variable integration

Table 1: 3 datasets on breast cancer from GEO

| Breast | Gene | Sample | Case | Control |
|---|---|---|---|---|
| GSE9574 | 20995 | 29 | 14 | 15 |
| GSE21947 | 20995 | 30 | 15 | 15 |
| GSE5364 | 20995 | 196 | 183 | 13 |

Variable integration: Gene and CNA
Integration of supervised and unsupervised learning

# CSA Breath Data



1. Tidal volume breathing for 5 minutes;

# CSA Breath Data



1. Tidal volume breathing for 5 minutes;
2. Exhaled breath drawn over the sensor array;

# CSA Breath Data



1. Tidal volume breathing for 5 minutes;
2. Exhaled breath drawn over the sensor array;
3. Images were converted to numerical values in the red, green, blue spectra, and 4 ultraviolet spectra.

# CSA Breath Data



1. Tidal volume breathing for 5 minutes;
2. Exhaled breath drawn over the sensor array;
3. Images were converted to numerical values in the red, green, blue spectra, and 4 ultraviolet spectra.
4. Totally 128 (the number of colorants) $\times$ 7( changes in the red, green, blue, and 4 ultra-color spectrum of each colorant) $= 896$ groups.

# CSA Breath Data

# Object detection



**CF图片项目背景**

人工检查

缺点：

- 人工成本高
- 检测精度低
- 人员培训困难
- 工作单调枯燥

# Object detection



## CF图片项目背景

人工智能

- 优化产业链
- 智能制造
- 提升产品竞争力

# Object detection

🔍 标注示例——FI312



输出结果包含
- 分类标签FI312
- 分类置信度
- 目标所在的位置框

# Network experience Index

## II. Construction of the composite NEI/NPI index

**Objective:** Develop a new strategy for constructing the composite NEI/NPI index which can "inherit" the advantages of the existing approaches while overcoming their limitations.

**VaR**

Value at risk is a measure of the risk of loss for investments (with a given probability).

**Feature screening**

- From Step I analysis: connectivity and Laplacian
- Experts' opinion

**Feature selection and estimation**

PCA (and other feature selection methods): identify relevant features and estimate their loadings.

**Bayesian update**

**Delphi method**

Collect prior information from experts through questionnaire.

**Extreme value theory**

It assesses the probability of an event that is more extreme than previously observed.
POT(Peaks Over Threshold) and BMM(Block Maxima Method) are used to find the extreme values.

PCA        Sparse PCA

Design of the loadings

Upper and lower bounds of indexes

**Deliverable:** a composite NEI/NPI index

# 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?

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

# Kaggle

# Software