# Machine Learning for OR & FE

## Introduction and Course Overview

### Martin Haugh

Department of Industrial Engineering and Operations Research
Columbia University
Email: martin.b.haugh@gmail.com

## What is Data Mining / Machine Learning / "Big Data"?

Data mining / machine learning / "big data" are all related to the process of discovering new patterns from large data sets.

There is substantial overlap between all three "fields"

– they all use ideas from statistics and algorithms.

According to Hand, Mannila and Smyth:

*"Data mining is the analysis of (often large) observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful to the data owner."*

But isn't all this just statistics?!

## Differences Between Machine Learning and Statistics?

Traditional statistics:

- First hypothesize, then collect data, then analyze.
- Often model-oriented with an emphasis on parametric models.
- Focus on understanding and hypothesis testing.

Machine learning (ML):

- Few if any a priori hypotheses.
- Typically data has already been collected.
- Analysis is typically data-driven not hypothesis-driven.
- Often algorithm-oriented rather than model-oriented.
- Focus on prediction.
- ML is more associated with artificial intelligence than statistics.

But statistical ideas and algorithms are extremely useful in ML

- e.g. inference, bootstrapping, over-fitting solutions, regression and classification algorithms.

"Big data" also encompasses data-management technology

- e.g. cloud computing, NoSQL databases, map-reduce and Hadoop etc.

### Glossary

| Machine learning | Statistics |
|---|---|
| network, graphs | model |
| weights | parameters |
| learning | fitting |
| generalization | test set performance |
| supervised learning | regression/classification |
| unsupervised learning | density estimation, clustering |
| large grant = $1,000,000 | large grant= $50,000 |
| nice place to have a meeting: Snowbird, Utah, French Alps | nice place to have a meeting: Las Vegas in August |

But "uncool" jargon isn't unique to statistics; witness "neural networks".

# Why the Massive Popular Interest in "Big Data"?

There are many reasons including:

- The explosion of data from many domains such as biology, the internet, astrophysics, telecommunications, engineering and sensor networks.
- The availability of ever cheaper storage
    - developments in database technology such as NoSQL databses
- Faster and cheaper computation to analyze the data.
- Competitive pressure in business: data has value!
- Software developments:
    - commercial products: SAS, SPSS, Google Analytics, IBM, Oracle, Amazon, S-Plus
    - open source products: R, Python libraries such as scikit learn and pandas, Weka.

Don't necessarily need to be a ML expert to do ML

– but many pitfalls so some competence is required!

## Is "Big Data" Over-Hyped?

But "big data" is perhaps over-hyped in the popular media
– and startup scene seems focussed on less important matters:



Budd Aldrin, *MIT Technology Review*, October 2012.

# But There Are Many Data Mining Successes ...

- Google: the entire company!
- Market Basket (WalMart).
- Recommender Systems: Amazon.com, Netflix.
- Fraud Detection in Telecommunications: AT&T.
- Targeted Marketing: Target.
- Financial Markets: Hedge funds using text mining, Bloomberg.
- DNA Microarray analysis.
- Smart grid: dynamic demand response.

# Taxonomy of Machine Learning Problems

**Supervised Learning**

- Observe training samples, and learn a function mapping inputs to outputs.
- Examples include regression, classification, ranking, etc.

**Unsupervised Learning**

- Observe data and construct a low complexity description of the data.
- Clustering, dimensional reduction techniques such as principal components analysis, non-negative matrix factorization etc.

**Reinforcement Learning (RL)**

- RL attempts to learn a map from system states to actions with the goal of maximizing total cumulative reward.
- A given action in a given state influences what the next state will be
    - so need to trade off exploitation (of current knowledge) with exploration (to learn new information).
- Closely related to dynamic programming except that we do not know the parameters of the model.

Will focus on supervised and unsupervised problems in this course.

## Data Types

Data comes in many forms ...

- Flat files or vector data.
- Text data; e.g. an article or collection of articles.
- Transactional data: web logs, phone calls.
- Relational data from relational databases.
- Time series data; e.g. financial data, auction data from EBay.
- Image data and video data.
- Spatial temporal data e.g. geophysical data.
- Network data: physical networks, social networks etc.

Regardless of original data-type, typically need to convert it into quantitative data

– hopefully without losing important information.

## Exploratory Data Analysis (EDA)

**Goal**: To obtain a general sense and understanding of the data
- EDA is a vital and necessary first step for any machine learning / statistics task.

EDA tasks:
- What do the marginal distributions look like? Normal, skewed, etc.
- What is the quality of the data? Is some of it missing?
    - If so, why is it missing? Is it missing at random (MAR), missing completely at random (MCAR) or missing not at random (MNAR)?
    - A very important issue!
- Are there outliers? What should we do with them?
- What are the dependencies / correlation between different variables?
- What subsets of the variables are of particular interest?
- What sort of (functional) relationships are we looking for?

General guidelines for EDA:
- Make it data-driven and typically model-free ...
- Think interactive and visual ...
- Many dimensions to play with: x, y, z, space, color, time ...

## Some EDA Strategies

Summary statistics:

- mean, median, variance, skewness, kurtosis

Single variable visualization:

- Histograms
- Density estimates
- Box plots

Two variable visualization:

- Scatter plot
- Binning
- Transparent plots
- Contour plots
- Bar charts: categorical variables

**Bottom line:** it is always well worth looking at your data!

- EDA will not be a focus of this course but it should always be the start of any analysis.

## Data Visualization

Data visualization is an increasingly important subfield of "big data"

  – if "a picture is worth a thousand words" what is a video worth?

See Hans Rosling's famous talk on the joy of statistics (and global health dynamics) at http://www.youtube.com/watch?v=jbkSRLYSojo.



Data visualization also becoming increasingly important in journalism

  – see http://www.theguardian.com/data if you need convincing.

## Domain Knowledge

An expertise in machine learning / statistics generally not enough to apply these methods successfully.

Domain-specific knowledge generally also required:

- – How can we apply techniques to data from biology, the internet, astrophysics, telecommunications, engineering and sensor networks, etc. without understanding these domains?
- – How do we know what the right questions are?
- – How can we make sense of the results?
- – Or figure out if the results even make sense?

The most successful ML applications almost always based on good domain knowledge

- – also helps in the EDA process.

Should avoid the temptation of blindly applying your favorite ML tool

- – a hammer looking for a problem is not a good recipe for success!

# Course Highlights I: Classification and Reduced-Rank LDA
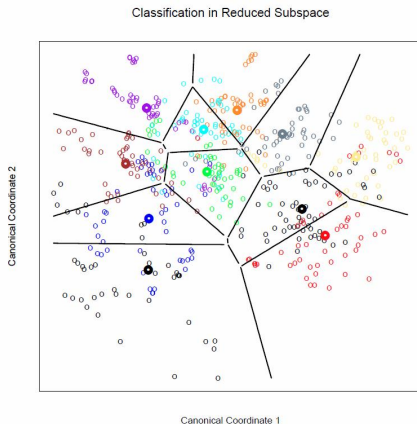


Classification in Reduced Subspace

**Figure 4.11 from HTF**: Decision boundaries for the vowel training data, in the two-dimensional subspace spanned by the first two canonical variates. Note that in any higher-dimensional subspace, the decision boundaries are higher-dimensional affine planes, and could not be represented as lines.

## Course Highlights II: Clustering and Vector Quantization



**FIGURE 14.9.** *Sir Ronald A. Fisher (*$1890 - 1962$*) was one of the founders of modern day statistics, to whom we owe maximum-likelihood, sufficiency, and many other fundamental concepts. The image on the left is a $1024 \times 1024$ grayscale image at 8 bits per pixel. The center image is the result of $2 \times 2$ block VQ, using 200 code vectors, with a compression rate of 1.9 bits/pixel. The right image uses only four code vectors, with a compression rate of 0.50 bits/pixel*

*Figure 14.9 taken from HTF*

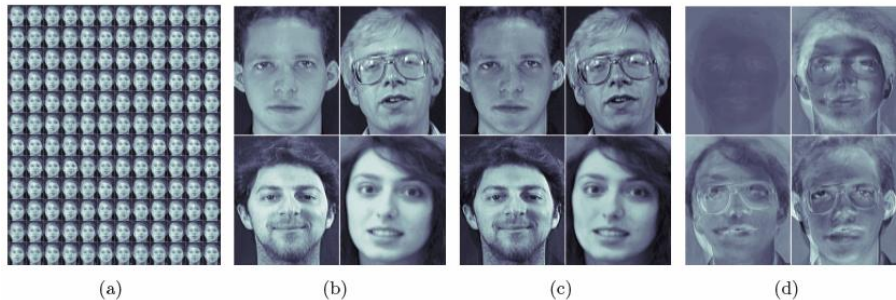## Course Highlights III: Dimension Reduction Techniques



**Figure 15.15 from Barber**: (a) Training data, consisting of a positive (convex) combination of the base images. (b): The chosen base images from which the training data is derived. (c): Basis learned using conditional PLSA on the training data. This is virtually indistinguishable from the true basis. (d): Eigenbasis (sometimes called 'eigenfaces').

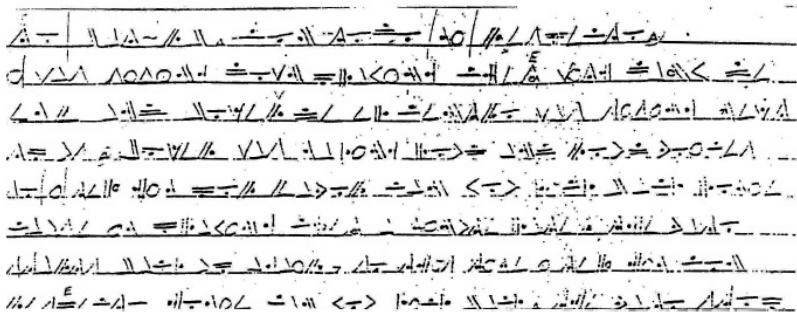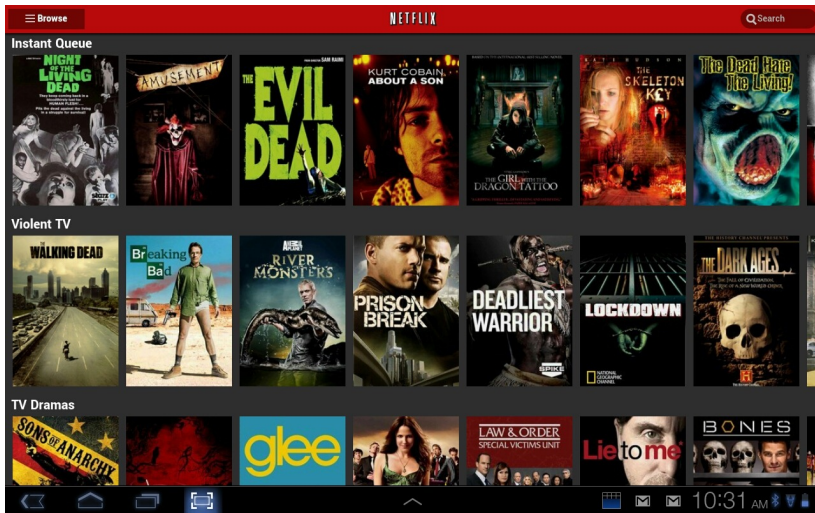## Course Highlights IV: MCMC and Cryptography



Figure taken from "The Markov Chain Monte Carlo Revolution", by Persi Diaconis in the *Bulletin of the American Mathematical Society* (2008).

- An unusual application of MCMC: decoding the coded message.

**Course Highlights V: Collaborative Filtering and Recommender Systems**



Source: *Netflix*
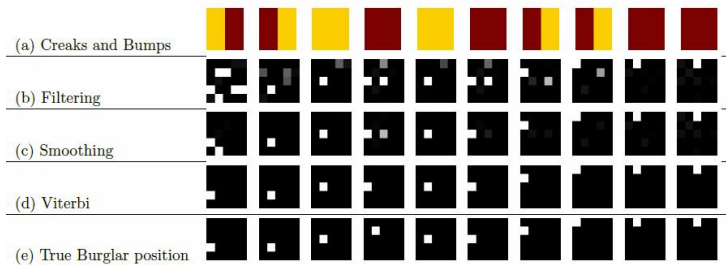
# Course Highlights VI: HMMs and Burglar Tracking



**Figure 23.7 from Barber**: Localising the burglar through time for 10 time steps. (a): Each panel represents the visible information $v_t = \left(v_t^{\text{creak}}, v_t^{\text{bump}}\right)$, where $v^{\text{creak}} = 1$ means that there was a 'creak in the floorboard' ($v^{\text{creak}} = 2$ otherwise) and $v^{\text{bump}} = 1$ meaning 'bumped into something' (and is in state 2 otherwise). There are 10 panels, one for each time $t = 1, \ldots, 10$. The left half of the panel represents $v_t^{\text{creak}}$ and the right half $v_t^{\text{bump}}$. The yellow shade represents the occurrence of a creak or bump, the red shade the absence. (b): The filtered distribution $p(h_t \mid v_{1:t})$ representing where we think the burglar is. (c): The smoothed distribution $p(h_t \mid v_{1:10})$ that represents the distribution of the burglar's position given that we know both the past and future observations. (d): The most likely (Viterbi) burglar path $\text{argmax}_{h_{1:10}} \, p(h_{1:10} \mid v_{1:10})$. (e): The actual path of the burglar.

Figure 4: Difference between conditional probabilities given evidence and prior probabilities. Evidence is observed in prospect A, and follows the explanations in section 3.4 . Figure shows the effect of Evidence 1 (left) and 2 (right).

Figure taken from "Strategies for Petroleum Exploration Based on Bayesian Networks: a Case Study", by Martinelli et al. (2012).
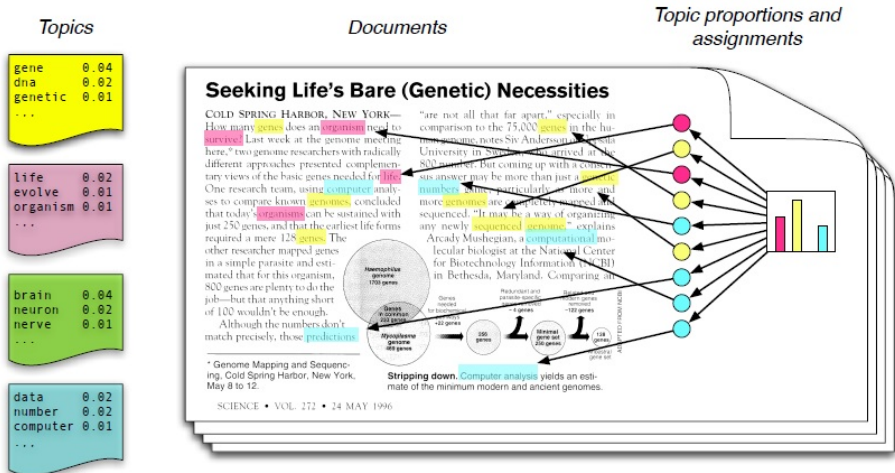
Figure taken from "Introduction to Probabilistic Topic Models", by D.M. Blei. (2011).

**Lots of Machine Learning in Finance but . . .**



Source: *Bloomberg*

# The *Google* Smart Car (Not a Course Highlight!)



**Under the bonnet**
How a self-driving car works

Signals from **GPS (global positioning system)** satellites are combined with readings from tachometers, altimeters and gyroscopes to provide more accurate positioning than is possible with GPS alone

**Lidar (light detection and ranging)** sensors bounce pulses of light off the surroundings. These are analysed to identify lane markings and the edges of roads

**Video cameras** detect traffic lights, read road signs, keep track of the position of other vehicles and look out for pedestrians and obstacles on the road

**Radar sensor**

**Ultrasonic sensors** may be used to measure the position of objects very close to the vehicle, such as curbs and other vehicles when parking

The information from all of the sensors is analysed by a **central computer** that manipulates the steering, accelerator and brakes. Its software must understand the rules of the road, both formal and informal

**Radar sensors** monitor the position of other vehicles nearby. Such sensors are already used in adaptive cruise-control systems

Source: *The Economist*

Source: *The Economist*