# Introduction to Machine Learning

## CS 534: Machine Learning

# Why Machine Learning?

"We are drowning in information and starving for knowledge." — John Naisbitt

- Big data era

- Use algorithms to discover new relationships, scale tasks, and perform decision making under uncertainty

# Machine learning workflow

**Unsupervised**

**Feature extraction**

**Machine learning algorithm**

**Grouping of objects**

based on some common characteristics

**Training set**

**Supervised**

**Predictive model**
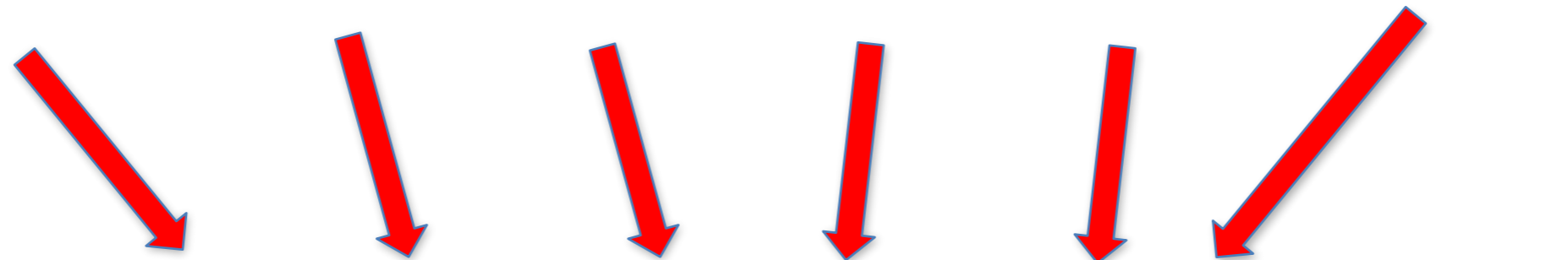
**New Data**

**Annotated data**

# ML, DM, or AI?

- Machine Learning (ML) - The study, design and development of algorithms that endow learning capabilities

- Data Mining (DM) – Using ML and statistical techniques to learn something new from data

- Artificial Intelligence (AI) – Broad study of developing intelligent agents (think Turing Test)

# Diverse Community

# Journals & Conferences

- Journals

  - Transactions on Pattern Analysis and Machine Intelligence (IEEE)

  - Journal of Machine Learning Research (ACM)

  - Machine Learning (Springer)

  - Foundations and Trends in Machine Learning (ACM)

- Conferences

  - ICML – International Conference on Machine Learning

  - NIPS – Neural Information Processing Systems

  - UAI – Uncertainty in Artificial Intelligence

  - CVPR – Computer Vision Pattern Recognition

**Complete list at Microsoft Research Rankings**

# All About Benchmarking

- Abundant data available to compare algorithms

- Required for publication

- Makes ML more of a science

- Still difficult to make fair comparisons

  - What parameters to use

  - Difficult to sweep parameter space

# Tools for Machine Learning

- Python

  Pros: free, fast, many algorithms available

  Cons: can be slow

- R

  Pros: free, standard in bioinformatics & statistics, great vector graphics

  Cons: extremely slow, poorly documented, bad language conventions

- Matlab

  Pros: fast, large user community & codebase, well documented

  Cons: not free

Code examples will be provided in Python

# Course Logistics

# Course Website

http://joyceho.github.io/cs534-s17/index.html

- Lectures

- Assignments

- Example code (when applicable)

# About Me (Joyce Ho)

- Undergraduate / MEng from MIT

- PhD from University of Texas at Austin

- Research interests:

  - Data Mining / Machine Learning

  - Healthcare Informatics

- More information: http://joyceho.github.io

# Contact Information

- Email: joyce.c.ho@emory.edu

- Office Hours @ MSC W414

  - M 1:00 pm - 3:30 pm

  - W 9:30 am - 12:00 pm

# Communication

- Piazza: http://piazza.com/emory/spring2017/cs534

  - Announcements

  - Questions + Discussions

  - Assignment Clarifications + Slide Corrections

- Office Hours

- By Appointment!

# Course Textbook

- Elements of Statistical Learning

  - PDF available online

- Machine Learning: a Probabilistic Perspective, by Kevin Murphy (Optional)

- Pattern Recognition and Machine Learning, by Christopher Bishop (Optional)

# Evaluation

- 4-5 assignments (40%)

  - Both theory & programming

- Midterm (15%)

- Project (40%)

- Participation (5%)

# Collaboration Policy

- Try the assignments on your own first

- Discuss with others if necessary

- Write-up solutions on your own

- List the people you collaborated with

# Project

- Work in groups of 1-2

- Emphasis on public data sets (e.g., Kaggle competitions, MovieLens, KDD Cup, etc.)

- Project proposal due by spring break for feedback

- Goal is to either develop a new algorithm or try multiple algorithms to achieve good performance

# Prerequisites

- REQUIRED:

  - Probability theory

  - Linear algebra

- STRONGLY RECOMMENDED:

  - Statistics

  - Good programming skills

# Jupyter [IPython] Notebooks

- Designed to make lectures more 'interactive' — in-class activity to learn together

- Provide a means for more exercises to better understand the material

- Learn some Python along the way

# Preview of Topics

# Linear Regression

- How can we build linear models to predict continuous-valued outcomes?

- How can we analyze these models to understand the importance of features?

# Linear Regression: Stock Market



Apple Inc.
NASDAQ: AAPL - Dec 20, 7:59 PM EST

**116.93** USD ↑0.29 (0.25%)
After-hours: 116.96 ↑0.03%

| 1 day | 5 day | 1 month | 3 month | 1 year | 5 year | max |

| Open | 116.74 | | Mkt cap | 627.11B |
| High | 117.50 | | P/E ratio | 14.13 |
| Low | 116.68 | | Div yield | 1.95% |

# Linear Regression: Weather

| Hour | Weather | | Temp. | Precip. | Wind |
|------|---------|--|-------|---------|------|
| 10pm | | Mostly Clear | 41°F | 0 in | NW - 5 mph |
| 12am | | Mostly Clear | 39°F | 0 in | NW - 3 mph |
| 02am | | Mostly Clear | 39°F | 0 in | NW - 3 mph |
| 04am | | Mostly Clear | 37°F | 0 in | NW - 3 mph |
| 06am | | Mostly Clear | 36°F | 0 in | NW - 3 mph |
| 08am | | Mostly Sunny | 43°F | 0 in | WNW - 3 mph |
| 10am | | Mostly Sunny | 50°F | 0 in | W - 2 mph |
| 12pm | | Mostly Sunny | 55°F | 0 in | SW - 2 mph |
| 02pm | | Mostly Sunny | 57°F | 0 in | S - 3 mph |
| 04pm | | Mostly Sunny | 57°F | 0 in | S - 3 mph |
| 06pm | | Mostly Clear | 54°F | 0 in | SSE - 3 mph |
| 08pm | | Mostly Clear | 50°F | 0 in | S - 3 mph |
| 10pm | | Partly Cloudy | 46°F | 0 in | S - 4 mph |

# Linear Classifiers

- What if the response variable is categorical or discrete?

- What are strategies for selecting the best features and dealing with noise?

# Linear Classifiers: Spam Filtering



spam
vs
not spam

# Linear Classifiers: Weather Prediction

| Hour | Weather | | Temp. | Precip. | Wind |
|------|---------|---|-------|---------|------|
| 10pm | | Mostly Clear | 41°F | 0 in | NW - 5 mph |
| 12am | | Mostly Clear | 39°F | 0 in | NW - 3 mph |
| 02am | | Mostly Clear | 39°F | 0 in | NW - 3 mph |
| 04am | | Mostly Clear | 37°F | 0 in | NW - 3 mph |
| 06am | | Mostly Clear | 36°F | 0 in | NW - 3 mph |
| 08am | | Mostly Sunny | 43°F | 0 in | WNW - 3 mph |
| 10am | | Mostly Sunny | 50°F | 0 in | W - 2 mph |
| 12pm | | Mostly Sunny | 55°F | 0 in | SW - 2 mph |
| 02pm | | Mostly Sunny | 57°F | 0 in | S - 3 mph |
| 04pm | | Mostly Sunny | 57°F | 0 in | S - 3 mph |
| 06pm | | Mostly Clear | 54°F | 0 in | SSE - 3 mph |
| 08pm | | Mostly Clear | 50°F | 0 in | S - 3 mph |
| 10pm | | Partly Cloudy | 46°F | 0 in | S - 4 mph |

# Learning Theory

- How can we gauge the accuracy of a hypothesis on unseen data?

- How do we quantify our ability to **generalize** as a function of the amount of training data and the hypothesis space?

- How do we find the best hypothesis?

# Occam's Razor Principle

- William of Occam: Monk living in the 14th century

- Principle of parsimony:
  "One should not increase, beyond what is necessary, the number of entities required to explain anything"

- When many solutions are available for a given problem, we should select simplest one
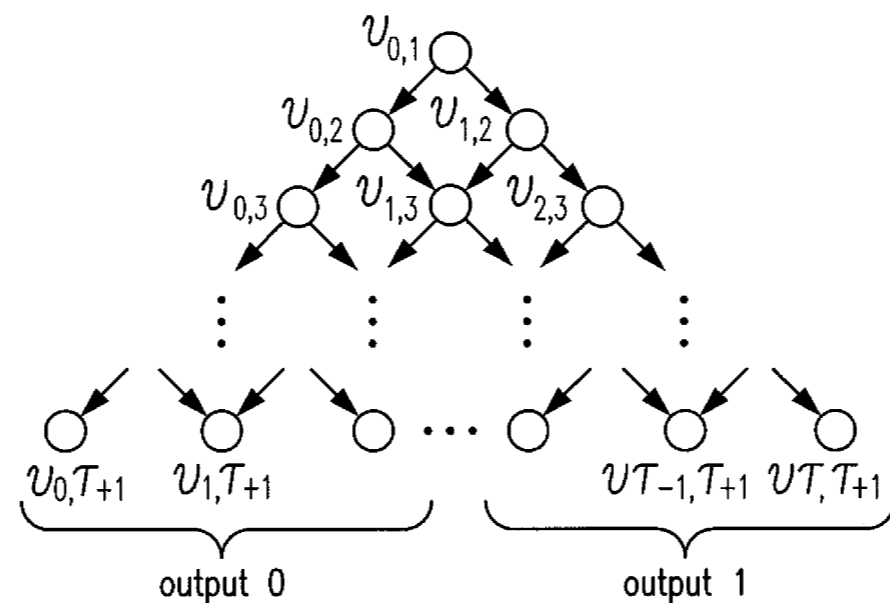
But what does simplest really mean?

**Samy Bengio**

# Validation

- How do we objectively measure performance of ML algorithms?

- How do we select the best algorithms and parameter values?

- Probably the most important topic — most abused and neglected
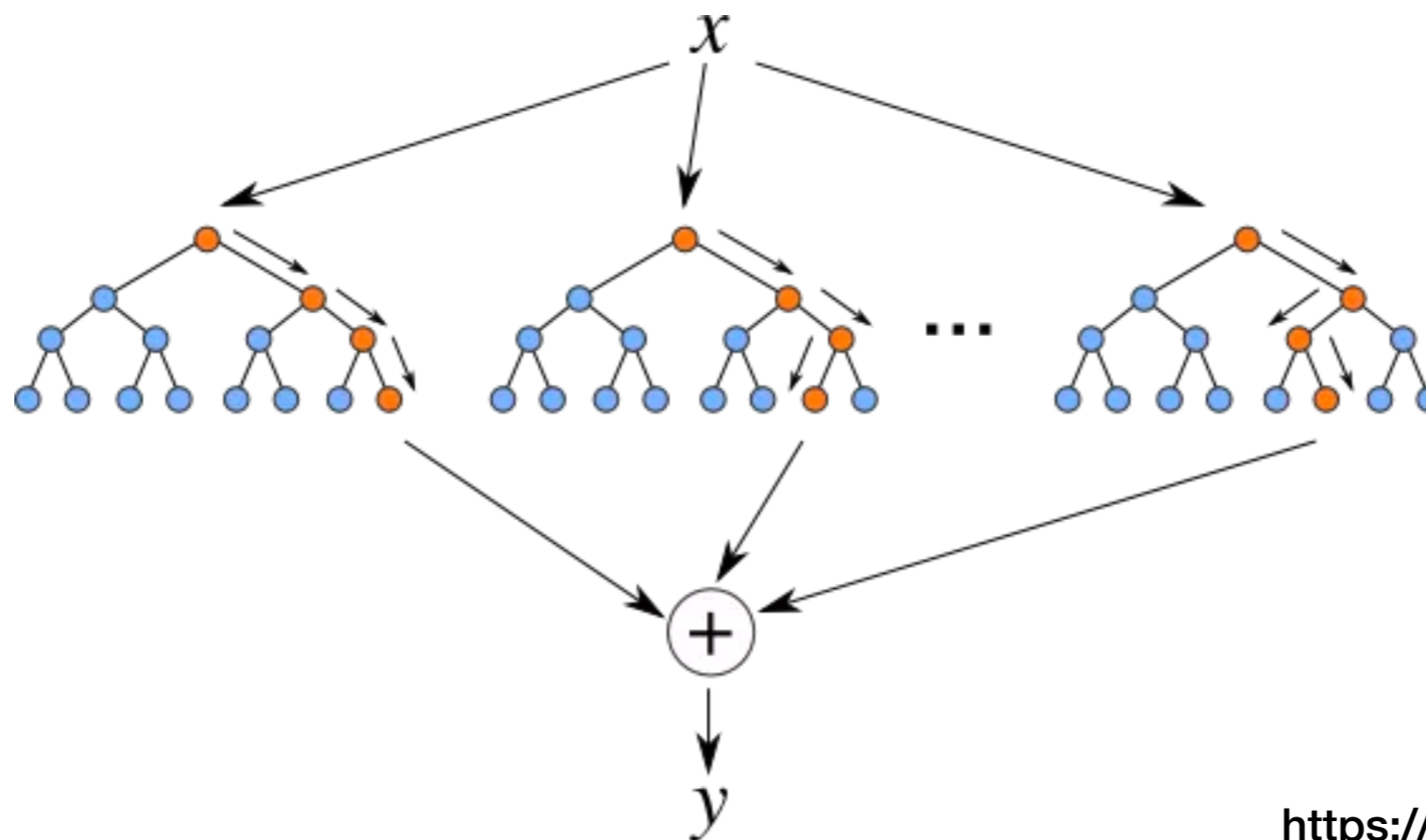
# Boosting, Trees, & Additive Models

- Can we achieve good performance by combining many primitive models?

- Can we build a strong learning from many weak learners?
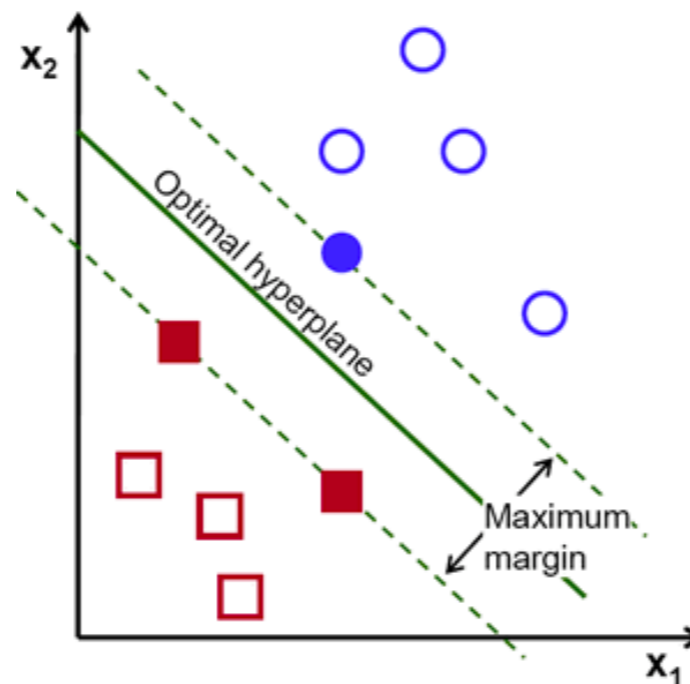
# Ensembles & Random Forests

- One of the most popular methods around

- Combining trees in a special way to get powerful classifier



https://kgpdag.wordpress.com/
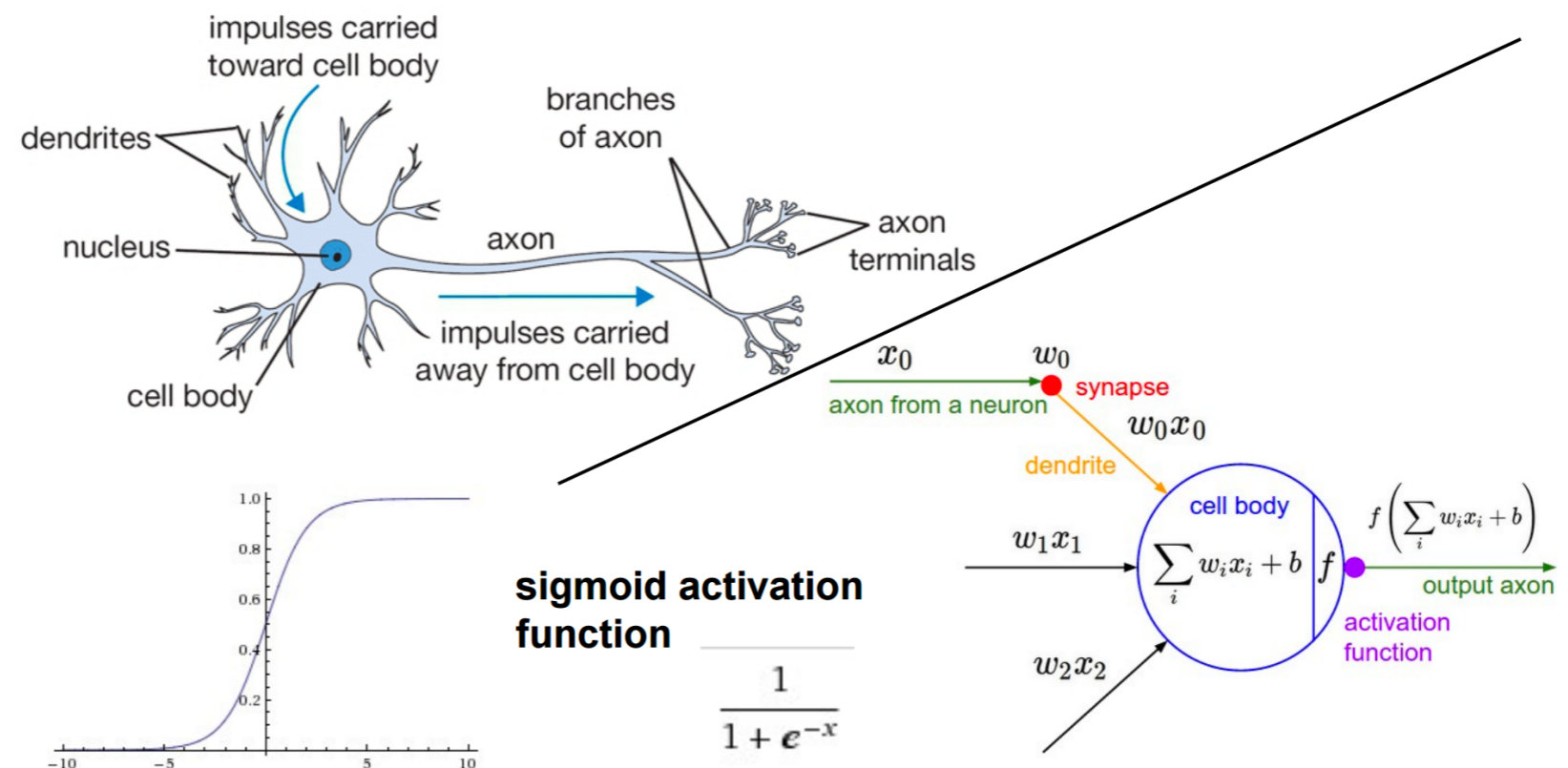
# Support Vector Machines

- Widely used classifier

- How can we implement non-linear classifiers by automatically transforming the data?

# Neural Networks

- Can we use biology as an inspiration for classification?

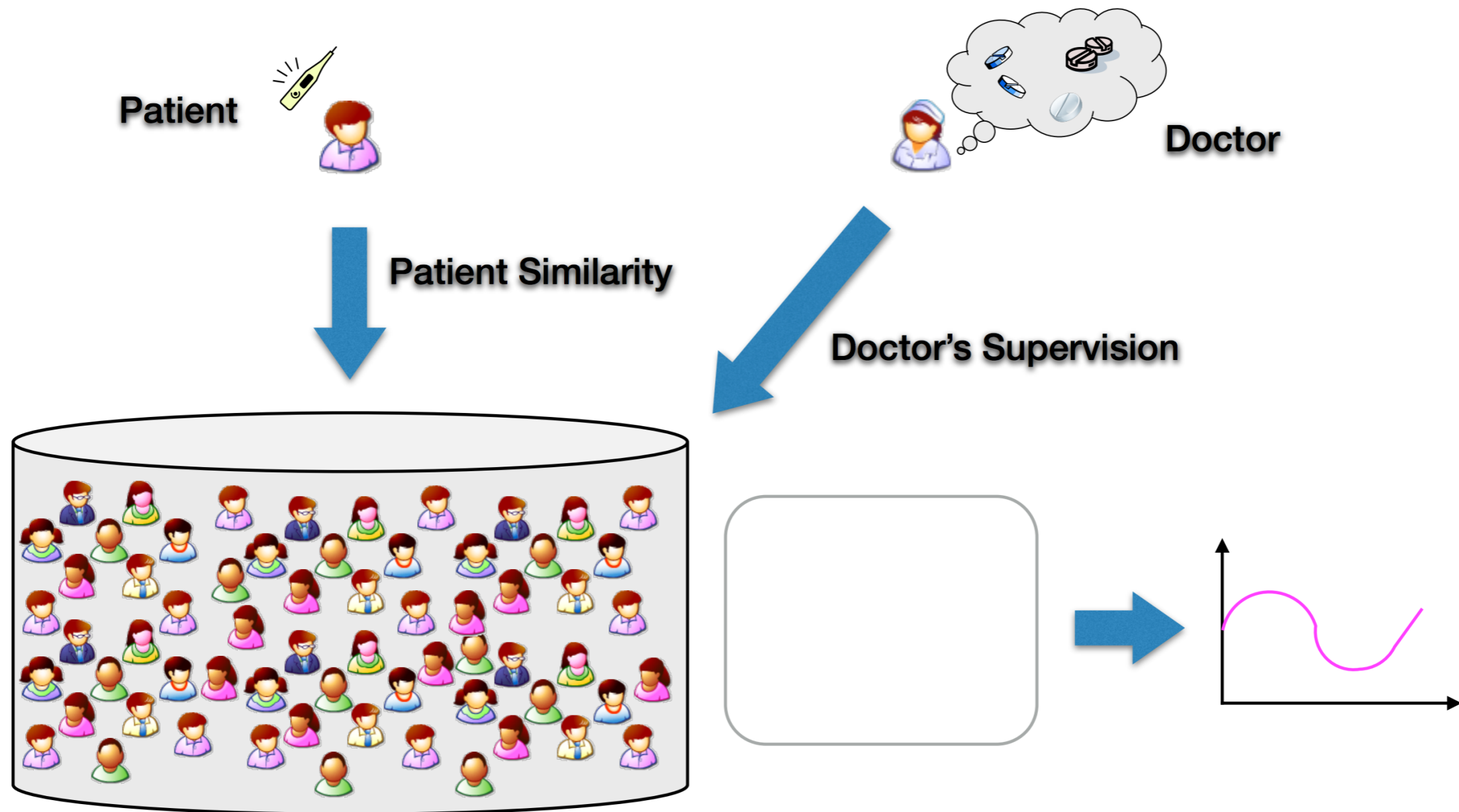- Connect virtual neurons in "natural" architectures and train to make decisions



sigmoid activation function

$$\frac{1}{1 + e^{-x}}$$

http://vision.stanford.edu/teaching/cs231n/slides/lecture5.pdf

# Prototype Methods

- What if we don't have an underlying model of how the data is shaped?

- Can we use the relationships between data points to develop good classifiers?
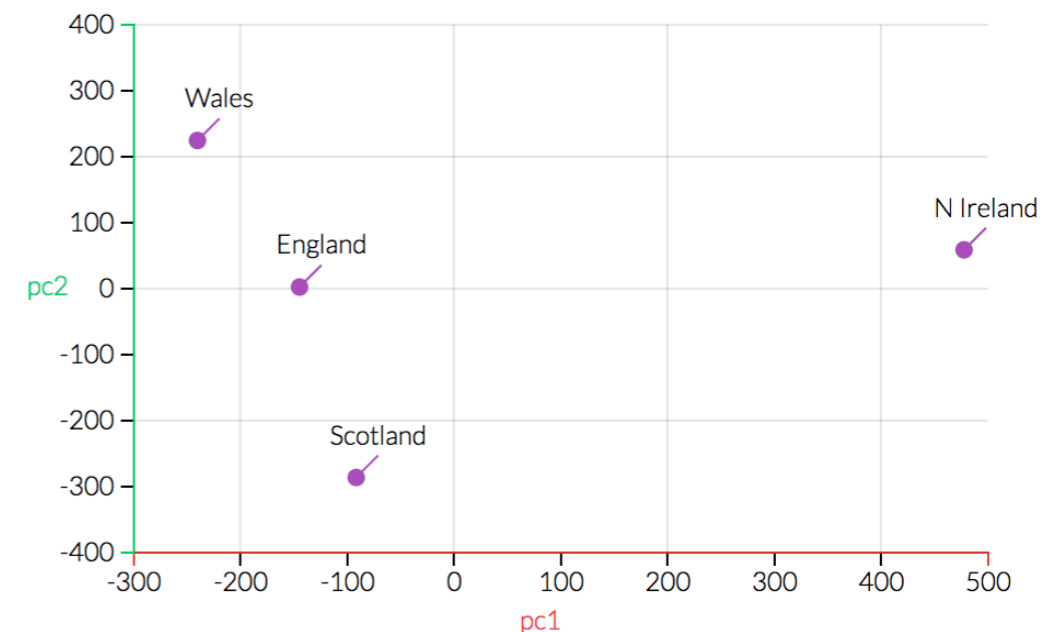
# Prototype Methods: kNN



**Patient**

**Doctor**

**Patient Similarity**

**Doctor's Supervision**

# Unsupervised Learning

- What happens when we don't know the outcome or have classes?

- How to explore data to look for structure and patterns?

# Unsupervised Learning: Visualization

| | England | N Ireland | Scotland | Wales |
|---|---|---|---|---|
| Alcoholic drinks | 375 | 135 | 458 | 475 |
| Beverages | 57 | 47 | 53 | 73 |
| Carcase meat | 245 | 267 | 242 | 227 |
| Cereals | 1472 | 1494 | 1462 | 1582 |
| Cheese | 105 | 66 | 103 | 103 |
| Confectionery | 54 | 41 | 62 | 64 |
| Fats and oils | 193 | 209 | 184 | 235 |
| Fish | 147 | 93 | 122 | 160 |
| Fresh fruit | 1102 | 674 | 957 | 1137 |
| Fresh potatoes | 720 | 1033 | 566 | 874 |
| Fresh Veg | 253 | 143 | 171 | 265 |
| Other meat | 685 | 586 | 750 | 803 |
| Other Veg | 488 | 355 | 418 | 570 |
| Processed potatoes | 198 | 187 | 220 | 203 |
| Processed Veg | 360 | 334 | 337 | 365 |
| Soft drinks | 1374 | 1506 | 1572 | 1256 |
| Sugars | 156 | 139 | 147 | 175 |



http://setosa.io/ev/principal-component-analysis/

# Unsupervised Learning: Clustering

# Unsupervised Learning: Topic Models

**Personal Finance:** (money, 0.15), (retire, 0.10), (risk, 0.03) ...

**Politics:** (President Obama, 0.10), (congress, 0.08), (government, 0.07), ...

## Parceling Out a Nest Egg, Without Emptying It
By **PAUL SULLIVAN**

What clients often forget are fixed costs — homes, cars, insurance — that must come down but take time to reduce, she said. Beyond that is her clients' skittish approach to risk; putting all of their money in cash may make them feel safe, she said, but it probably will not support the lifestyle they want for decades.

A generational disconnect is at work here: most people plan to retire at 65, the retirement age established for Social Security in 1935, when the average life expectancy was 61. Today the average is over 80 for men and women with a college degree.

So the $5.12 million gift exemption — created in a compromise between President Obama and Congress in 2010 — presents the well-off with a decision laden with short- and long-term consequences. How much should they give heirs now — and thus avoid giving the government in estate taxes later — while maintaining their lifestyle over a probably longer but still unpredictable remaining life span?

# Graphical Models

- Marriage between graph theory and probability theory
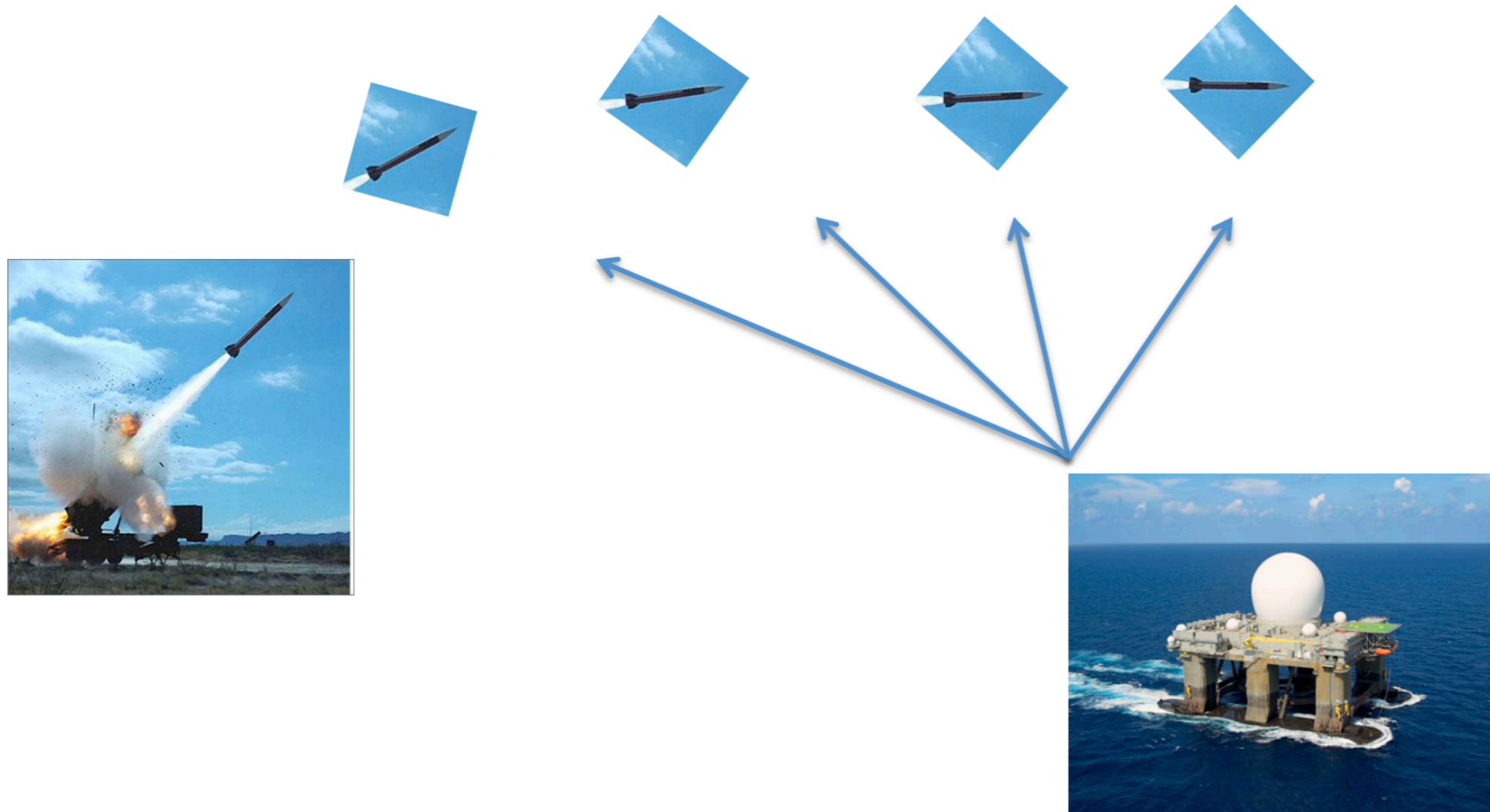
- Great for modeling sequential data (e.g., time series, speech processing)



https://en.wikipedia.org/wiki/Hidden_Markov_model

CS 534 [Spring 2017] - Ho

# Graphical Models: Tracking

Observe noisy measurements of missile location



Where is the missile now? Where will it be in 1 minute?

# Deep Learning

- Form of representation learning

- Aimed at learning feature hierarchies

- Features from higher levels of the hierarchy are formed by lower level features

- Each hidden layer allows for more complex features of input



http://www.deeplearningbook.org/contents/intro.html

# Recommendation Systems



How to build a system that provides or suggests items to the end users?

# iPython Setup

# Jupyter [iPython] Notebook

- Interactive computational environment which save output in a nice notebook format

  - Combine code execution, rich text, math, plots, and may other things

  - Supports markdown, LaTeX, HTML, etc

  - Popular in Data Science and can be easily shared with others

- More information: http://jupyter.org/

# Jupyter Notebook Setup

- Suggestion: Anaconda, an open source data science platform powered by Python: https://www.continuum.io/downloads

  - Contains Windows, OS X, Linux installations

  - Use either Python 2.7 or Python 3.5
    (Note: class will use 2.7 syntax primarily)

# Jupyter Notebook Setup (2)

- Once installed, to start, just open a terminal and run jupyter notebook



- A browser should open with jupyter running and you can import the .ipynb notebook from today's activity