

Machine Learning - Intro

Pradeep Ravikumar

Machine Learning 10-701
Jan 18, 2017



MACHINE LEARNING DEPARTMENT

Carnegie Mellon.
School of Computer Science

SOME LOGISTICS

Machine Learning Class webpage

- We will be using Piazza as a class webpage
- <https://piazza.com/class#spring2017/10701>
- Those who have been admitted will be added onto the Piazza webpage
- Those on waitlist can request to be added as well

Auditing

- To satisfy the auditing requirement, you must:
 - Get a passing grade in the course
 - Please send the instructors an email
 - Saying that you will be auditing the class
 - An audit plan of what you plan to do.

Prerequisites

- Probabilities
 - Distributions, densities, marginalization...
- Basic statistics
 - Moments, typical distributions, regression...
- Algorithms
 - Dynamic programming, basic data structures, complexity...
- Programming
 - Mostly your choice of language, but Matlab will be very useful
- We provide some background, but the class will be fast paced
- Ability to deal with “abstract mathematical concepts”

Recitations

- Strongly recommended
 - Brush up pre-requisites
 - Review material (difficult topics, clear misunderstandings, extra new topics)
 - Ask questions
- Basics of Probability, Dan Schwartz
- Thursday, Jan 19, 6:00 – 7:00 Tomorrow!
- GHC 6115

Textbooks

- Recommended Textbook:
 - Pattern Recognition and Machine Learning; Chris Bishop
- Secondary Textbooks:
 - The Elements of Statistical Learning: Data Mining, Inference, and Prediction; Trevor Hastie, Robert Tibshirani, Jerome Friedman
 - Machine Learning; Tom Mitchell
 - Machine Learning: A probabilistic perspective; Kevin Murphy.

Grading

- 5 Homeworks (50%)
 - Start early, Start early
- Final project (25%)
 - Kaggle style data analysis project, with leaderboard, and possibly discussion board
 - Projects done individually
 - Stay tuned for further details
- Midterm (25%)
 - Apr 5, in class

Homeworks

- Homeworks are hard, start early 😊
- Due in the beginning of class
- No late days
- Will have three components
 - Problem solving
 - Programming assignment
 - Using off-the-shelf packages
 - Explore the consequences of ML algorithms rather than implementing ML algorithms
 - Multiple choice

Homeworks

- Collaboration
 - You may **discuss the questions**
 - Each student writes their own answers
 - Each student must write their own code for the programming part
 - **Please don't search for answers on the web, Google, previous years' homeworks, etc.**
 - please ask us if you are not sure if you can use a particular reference

First Point of Contact for HWs

- To facilitate interaction, TA(s) will be assigned to each homework question
- These will be your “first point of contact” for this question
 - But, you can always ask any of us

Communication Channel

- For e-mailing instructors, always use:
 - 10701-instructors@cs.cmu.edu
- We highly recommend using Piazza to ask questions
- For announcements, we will be using Piazza

Instructors

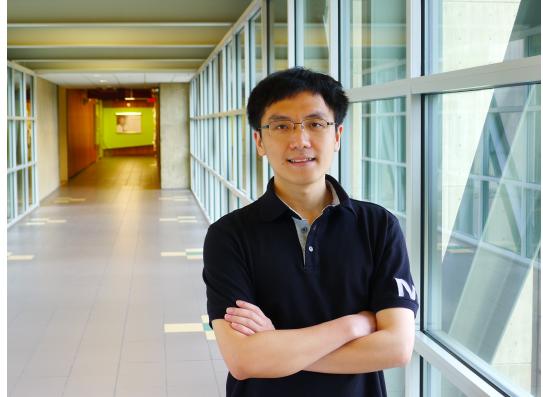


Aarti Singh



Pradeep Ravikumar

Your saviors - TAs



(Adams) Wei Yu



Calvin Murdock



Yichong Xu



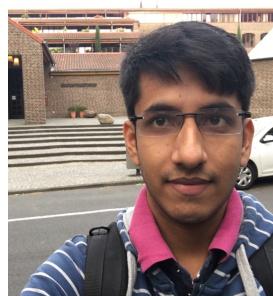
Hao Li



Dan Schwartz



Weixiang Ding



Danish Danish



Yiting Hao



Prakhar Naval

Great resources for learning, Interact with them!

(Adams) Wei Yu

- PHD in Machine Learning Department
- Research interests:
 - Large Scale Optimization
 - Statistical Machine Learning
 - Deep Learning Algorithms and Applications.
- More about me:
 - www.cs.cmu.edu/~weiyu/

Calvin Murdock

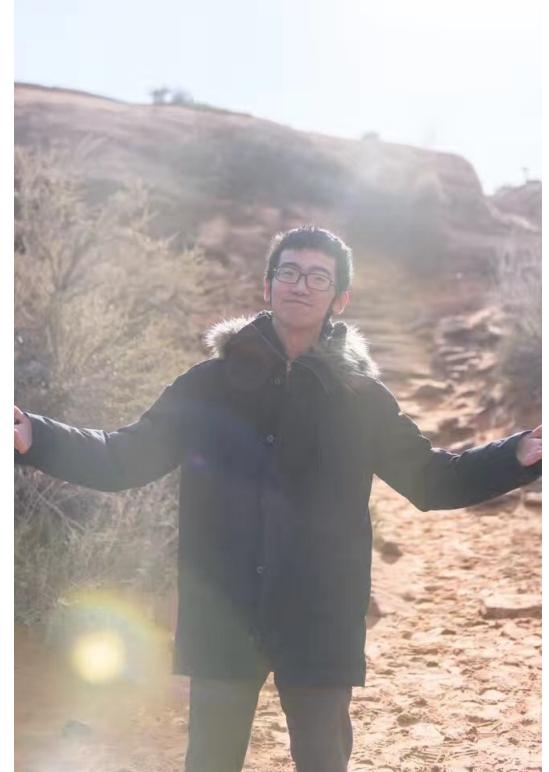
(cmurdock@andrew.cmu.edu)

- 4th year PhD student in MLD
- Broadly, interested in **computer vision**
- Specifically, **unsupervised visual learning**:
 - Representation learning (component analysis, manifold learning, deep neural networks)
 - Priors for encouraging interpretability
 - Geometric methods (iterative projection algorithms, learning on manifolds)



Yichong Xu

- Research Interests:
 - Active & Interactive Learning
 - Statistical Learning Theory
 - Deep Learning, Game theory and other stuffs
- Office: GHC 8215
- yichongx@cs.cmu.edu



YITING HAO



- Second Year Master Student from LTI
- Natural Language Processing, especially in semantics analysis

Danish's research interests



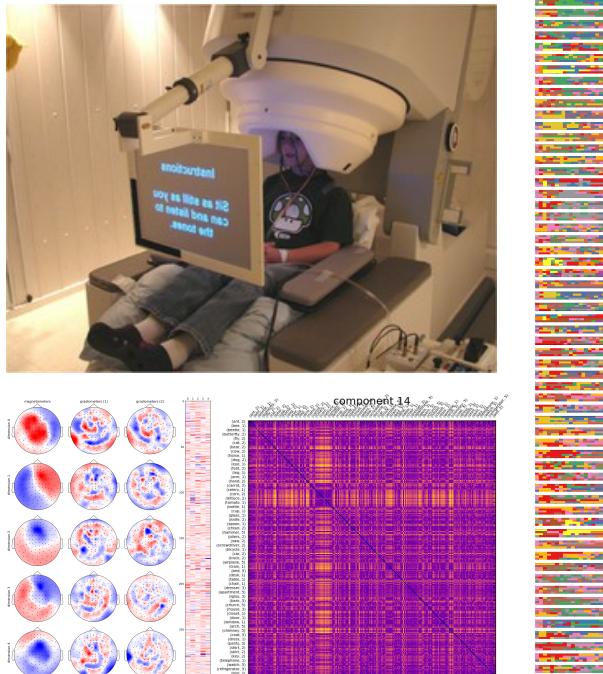
- **Representation Learning** — how to obtain neural representations that are not just expressive but also *interpretable*?
- **Machine Reading** — how to build curated readers, for biomedical domain, that *capture various events, activities and interactions* amongst different proteins?
- **Prediction with pairwise comparisons** — how can one benefit from *additional inexpensive pairwise comparisons*?

Weixiang Ding (Vincent)

- 2nd-year Master student (MSBIC) from LTI
- Experience in:
 - Natural Language Processing
 - Computer Vision (Pattern Recognition)
- Interest:
 - Neural Network



Dan Schwartz



Research Interests

- Brain activity analysis
 - Using machine learning techniques to understand the brain
- Semantics in the human brain
 - Understanding the basis of meaning representations in the brain
- Interpretable embeddings of time-series data



Experience and Interests:

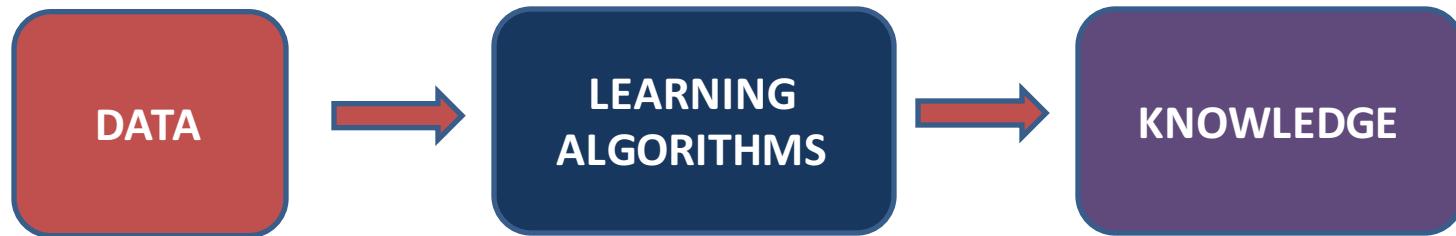
- Applications of Machine Learning techniques in technology industry
 - e.g. How e-commerce companies use ML techniques to auto-investigate countless online commodities?
- Machine Learning algorithms performance (efficiency) tuning in C++/Python
 - e.g. How to write more efficient and manageable programs/algorithms in C++/Python?



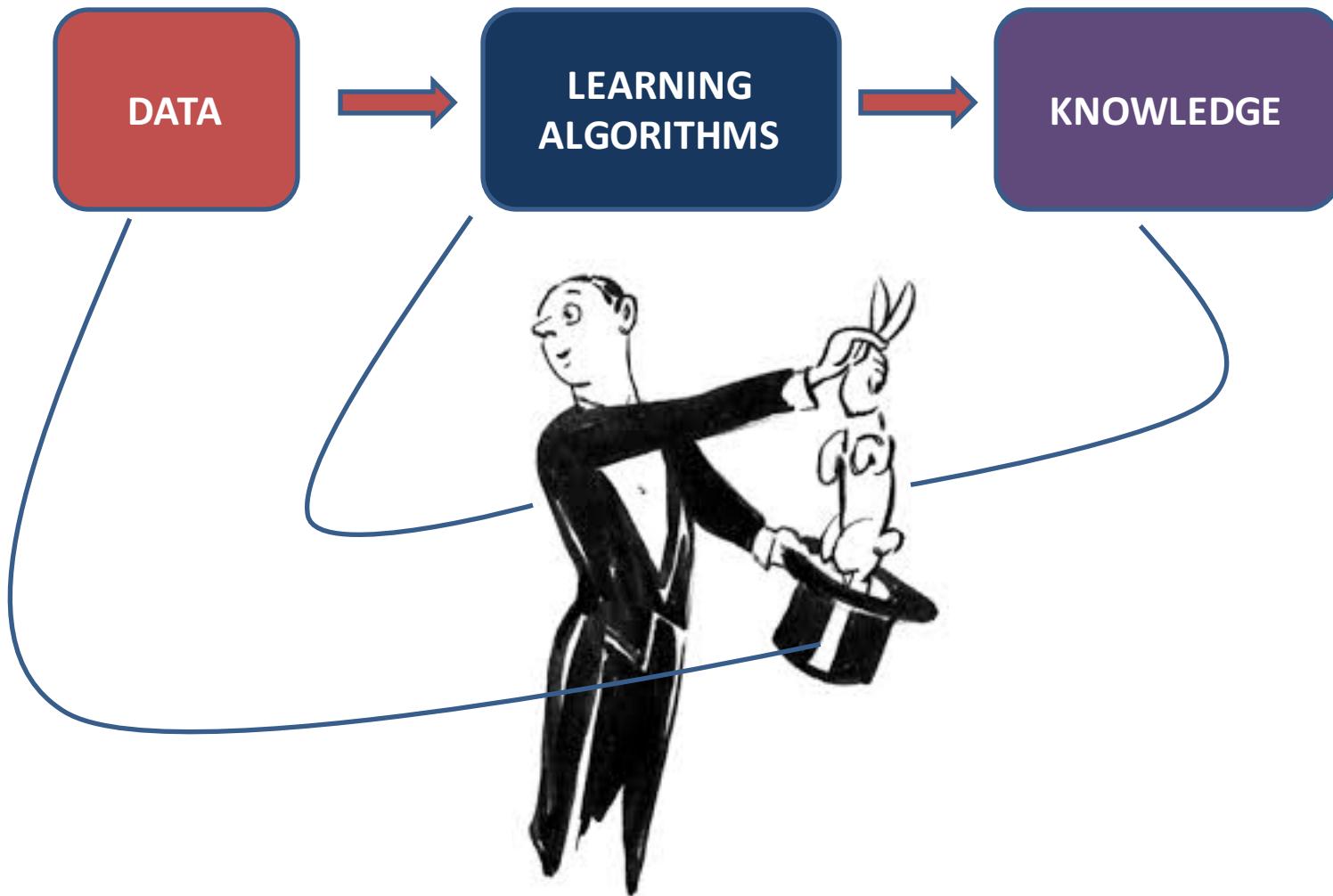
Hao Li

What is Machine Learning?

What is Machine Learning?



What is Machine Learning?



Machine Learning



- Algorithms that improve their knowledge towards some task with data
- How is it different from Statistics?
 - Same, but with better PR?
 - Statistics + Computation?
- What is its relationship with AI, Data Science, Data Mining?

Machine Learning

- It is useful to differentiate these different fields by their goals
- The goal of machine learning is the underlying mechanisms and algorithms that allow improving our knowledge with more data
 - Data construed broadly, e.g. “experiences”
 - Knowledge construed broadly e.g. possible actions

The Justice League

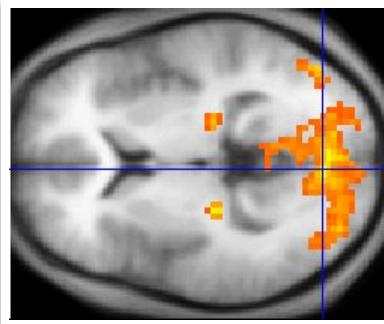
- Statistics: the goal is the understanding of the data at hand
- Artificial Intelligence: the goal is to build an intelligent agent
- Data Mining: the goal is to extract patterns from large-scale data
- Data Science: the science encompassing collection, analysis, and interpretation of data

From Data to Understanding ...

Machine Learning in Action

Machine Learning in Action

- Decoding thoughts from brain scans



Rob a bank ...

Home » Health & Wellness

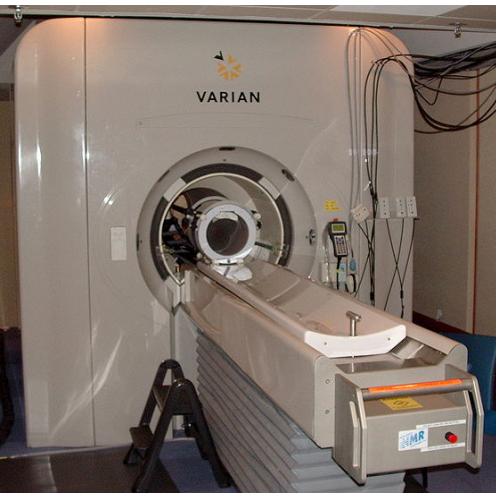
Brain Scans: Are You a Criminal?



Published February 07, 2007 by:
Andrea Okrentowich
[View Profile](#) | [Follow](#) | [Add to Favorites](#)

More: [Brain Scans](#) | [Brain Scan](#) | [Disposition](#) | [Defendant](#) | [Criminal Behavior](#)

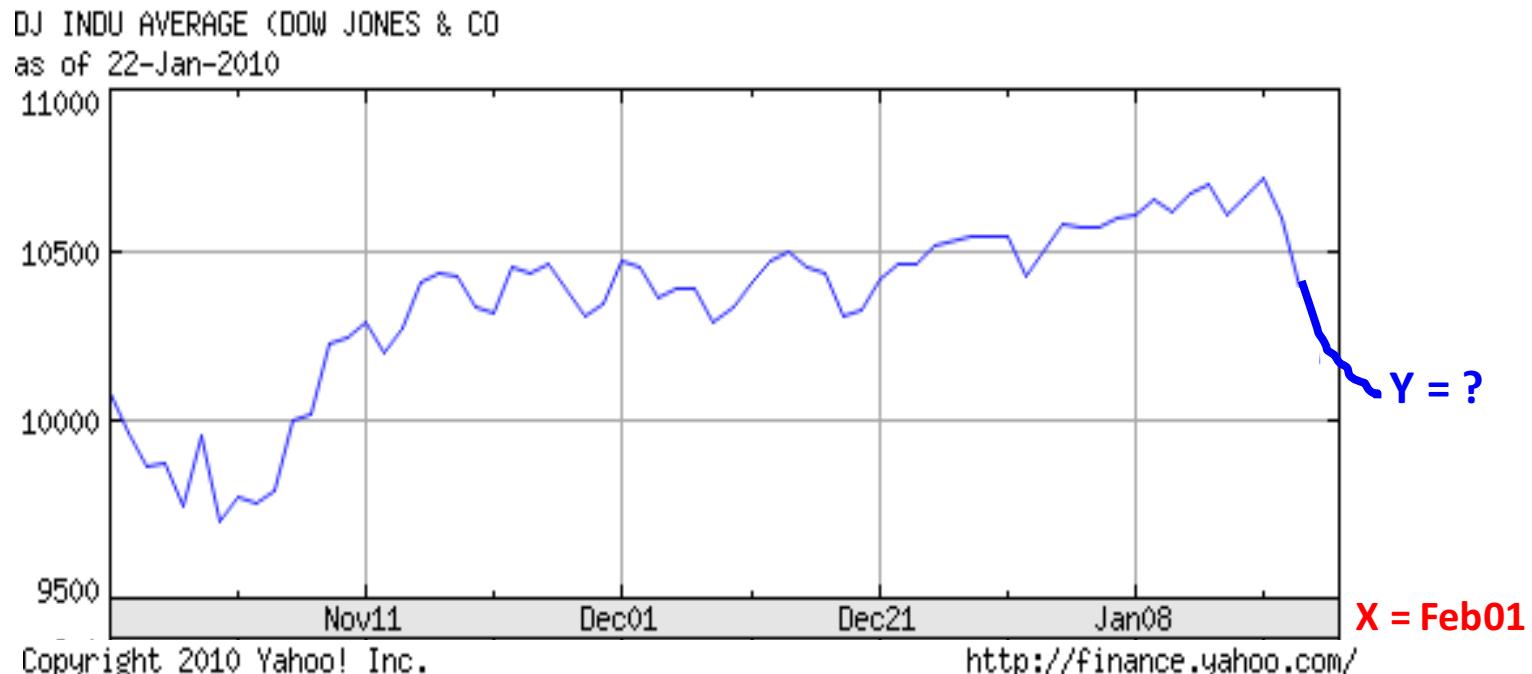
MRI Scans as Courtroom Evidence



The average Joe's MRI scan can show a brain abnormality, do we proceed to check him into the nearest mental institution or prison? That would make about as much sense as trying to prove a defendant innocent of a violent

Machine Learning in Action

- Stock Market Prediction



Machine Learning in Action

- Document classification



Sports
Science
News

Machine Learning in Action

- Spam filtering

Welcome to New Media Installation: Art that Learns

Hi everyone,

Welcome to New Media Installation:Art that Learns

The class will start tomorrow.

Make sure you attend the first class, even if you are on the Wait List.

The classes are held in Doherty Hall C316, and will be Tue, Thu 01:30-4:20 PM.

By now, you should be subscribed to our course mailing list: 10615-announce@cs.cmu.edu.

**Natural _LoseWeight SuperFood Endorsed by Oprah Winfrey, Free Trial 1 bottle,
pay only \$5.95 for shipping mfw rlk**



Spam/
Not spam

==== Natural WeightLOSS Solution ===

Vital Acai is a natural WeightLOSS product that Enables people to lose weight and cleanse their bodies faster than most other products on the market.

Here are some of the benefits of Vital Acai that You might not be aware of. These benefits have helped people who have been using Vital Acai daily to Achieve goals and reach new heights in their dieting that they never thought they could.

- * Rapid WeightLOSS
- * Increased metabolism - BurnFat & calories easily!
- * Better Mood and Attitude

Machine Learning in Action

- Cars navigating on their own



Boss, the self-driving SUV
1st place in the DARPA Urban
Challenge.

Photo courtesy of Tartan Racing.



Machine Learning in Action

- The **best** helicopter pilot is now a computer!
 - it runs a program that learns how to fly and make acrobatic maneuvers by itself!
 - no taped instructions, joysticks, or things like that ...



Machine Learning in Action

- Robot assistant? [http://stair.stanford.edu/]



Machine Learning in Action

- Many, many more...

Speech recognition, Natural language processing

Computer vision

Web forensics

Medical outcomes analysis

Computational biology

Sensor networks

Social networks

...

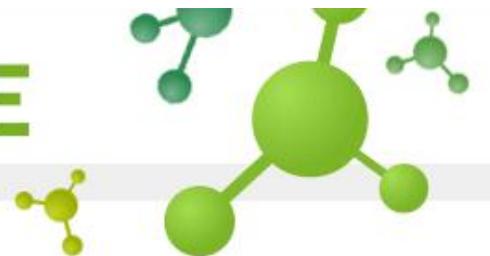
ML is trending!

- Wide applicability
- Very large-scale complex systems
 - Internet (billions of nodes), sensor network (new multi-modal sensing devices), genetics (human genome)
- Huge multi-dimensional data sets
 - 30,000 genes x 10,000 drugs x 100 species x ...
- Improved machine learning algorithms
- Improved data capture (Terabytes, Petabytes of data), networking, faster computers
- New York Times is regularly talking about machine learning

ML has a long way to go ...

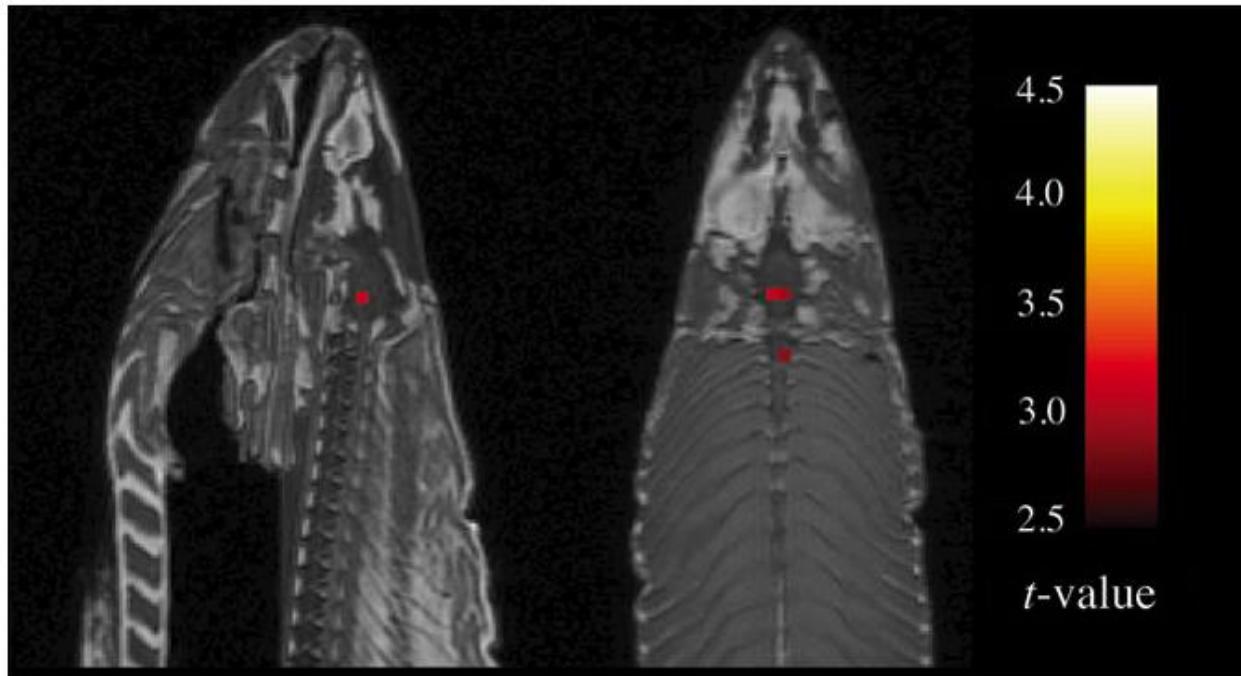
WIRED SCIENCE

NEWS FOR YOUR NEURONS



Scanning Dead Salmon in fMRI Machine Highlights Risk of Red Herrings

By Alexis Madrigal September 18, 2009 | 5:37 pm | Categories: Brains and Behavior



Unique Challenges with increased use of Machine Learning

Google researchers develop a test for machine learning bias



BY MIKE WHEATLEY
(HTTP://SILICONANGLE.COM/BLOG/AUTHOR/MIKEWHEATLEY/)

UPDATED 19:05 EST . 23 DECEMBER 2016

2 MIN READ



A team of researchers at Google Inc. has developed a method for testing whether or not machine learning algorithms inject bias, such as gender or racial bias, into their decision-making processes.

For some time, concerns have been raised about the possibility that machine learning algorithms are injecting bias into applications such as advertising, credit, education, employment and justice. Recent examples include a crime prediction algorithm that targeted black neighborhoods (<http://uk.businessinsider.com/predictive-policing-discriminatory-police-crime-2016-10?r=US&IR=T>) and an online advertising platform that was found to show highly paid executive jobs to men more often than women (<http://www.independent.co.uk/life-style/gadgets-and-tech/news/googles-algorithm-shows-prestigious-job-ads-to-men-but-not-to-women-10372166.html>) .

“Decisions based on machine learning can be both incredibly useful and have a profound impact on our lives,” said Moritz Hardt, a senior research scientist at Google, who co-authored the paper, “Equality of Opportunity in Supervised Learning.” “Despite the demand, a vetted methodology for avoiding discrimination against protected attributes in machine learning is lacking.”

What this course is about

- Covers a wide range of Machine Learning techniques
 - from basic to state-of-the-art
- You will learn about the methods you heard about:
 - Naïve Bayes, logistic regression, nearest-neighbor, decision trees, boosting, neural nets, overfitting, regularization, dimensionality reduction, PCA, error bounds, VC dimension, SVMs, kernels, margin bounds, K-means, EM, mixture models, semi-supervised learning, HMMs, graphical models, active learning, reinforcement learning...
- Covers algorithms, theory and applications
- **It's going to be fun and hard work ☺**

Machine Learning

- The goal of machine learning is identifying the underlying mechanisms and algorithms that allow improving our knowledge with more data
- Algorithms that improve their knowledge towards some task with data

Three axes of ML

- Data
- Tasks i.e. kind of knowledge
- Algorithms

First Axis: Data

- Fully observed
- Partially observed
 - Some variables systematically not observed
 - e.g. “topic” of a document
 - Some variables missing some of the time
 - “missing data”
- Actively collect/sense data

Second Axis: Algorithms

- Model-based Methods
 - Probabilistic Model of the data
 - Parametric Models
 - Nonparametric Models
- Model-free Methods

Model-based ML



Model-based ML



- Learning: From data to model
 - A model thus is a summary of the data
 - But also a story of how the data was generated
 - Could thus be used to describe how future data can be generated
 - **E.g. given (symptoms, diseases) data, a model explains how symptoms and diseases are related**
- Inference: From model to knowledge
 - Given the model, how can we answer questions relevant to us
 - **E.g. given (symptom, disease) model, given some symptoms, what is the disease?**

Model-based ML

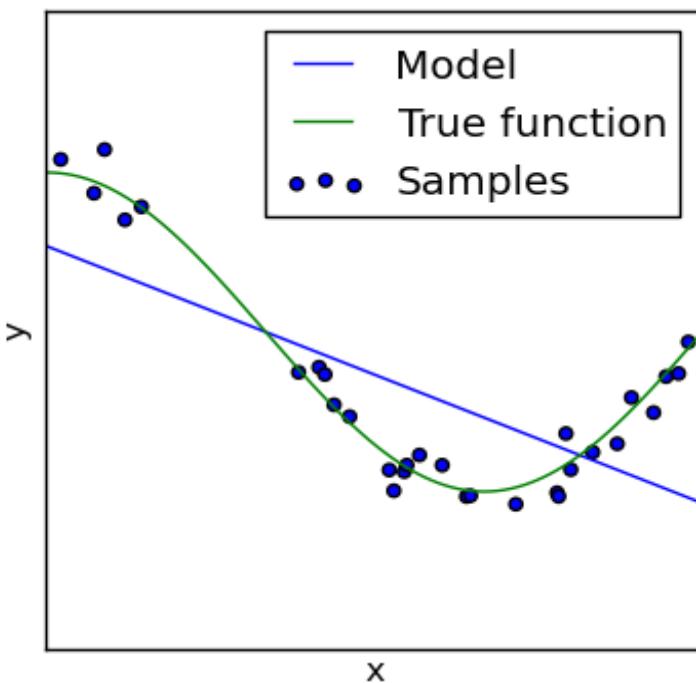


- Learning: From data to model
 - A model thus is a summary of the data
 - But also a story of how the data was generated
 - Could thus be used to describe how things work
 - **E.g. given (symptoms, diseases) model, what is the disease?**
symptoms and diseases are related
- Inference: From model to knowledge
 - Given the model, how can we answer questions?
 - **E.g. given (symptom, disease) model, what is the disease?**



Parametric Models

- “Fixed-size” models that do not “grow” with the data
- More data just means you learn/fit the model better

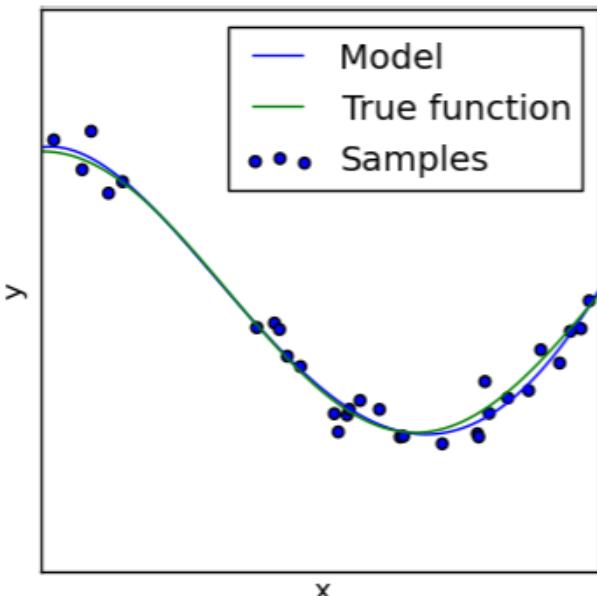


Fitting a simple line (2 params)
to a bunch of one-dim. samples

Model: data = point on line + noise

Nonparametric Models

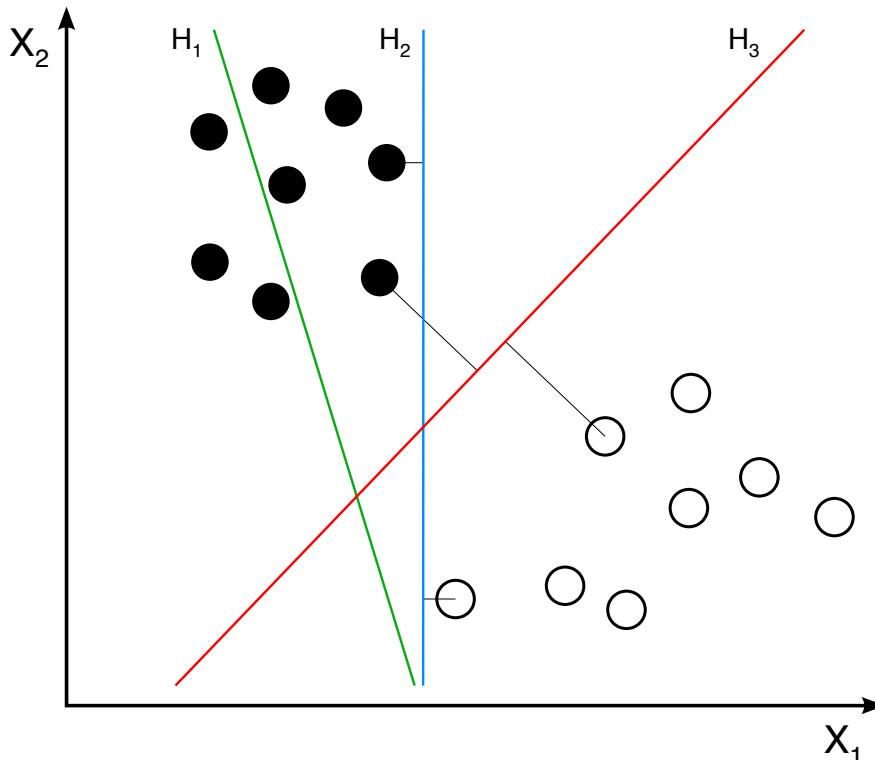
- Models that grow with the data
- More data means a more complex model



Fitting a “smooth function”
to a bunch of one-dim. samples

Model: data = point on smooth curve + noise

Model-free Methods



- Find best line that separates black from white points
- No modeling assumption e.g. that data generated from some point on line + noise

Third Axis: Knowledge/Tasks

- Prediction:
 - Estimate output given input

Prediction Problems

Feature Space \mathcal{X}



Words in a document

Label Space \mathcal{Y}

“Sports”
“News”
“Science”

...



Share Price
“\$ 24.50”



Task: Given $X \in \mathcal{X}$, predict $Y \in \mathcal{Y}$.

Prediction - Classification

Feature Space \mathcal{X}

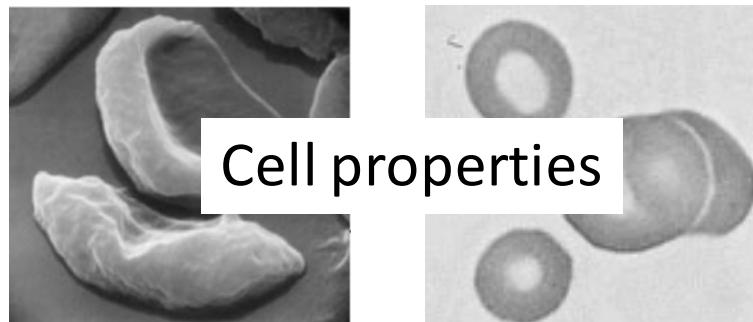


Words in a document

Label Space \mathcal{Y}

“Sports”
“News”
“Science”

...



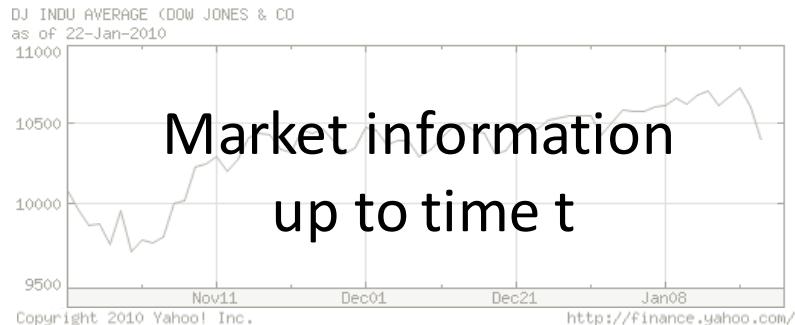
“Anemic cell”
“Healthy cell”



Discrete Labels

Prediction - Regression

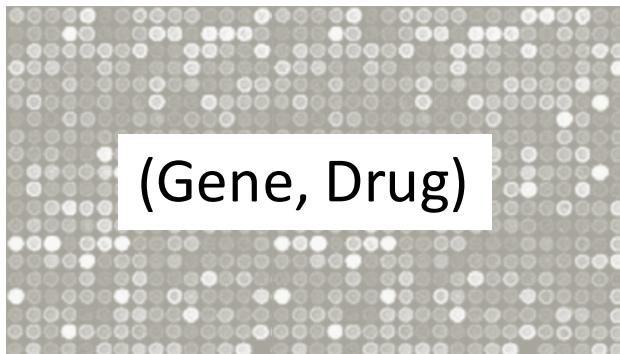
Feature Space \mathcal{X}



Label Space \mathcal{Y}



Share Price
“\$ 24.50”



Expression level
“0.01”

Continuous Labels

Prediction problems

Features?

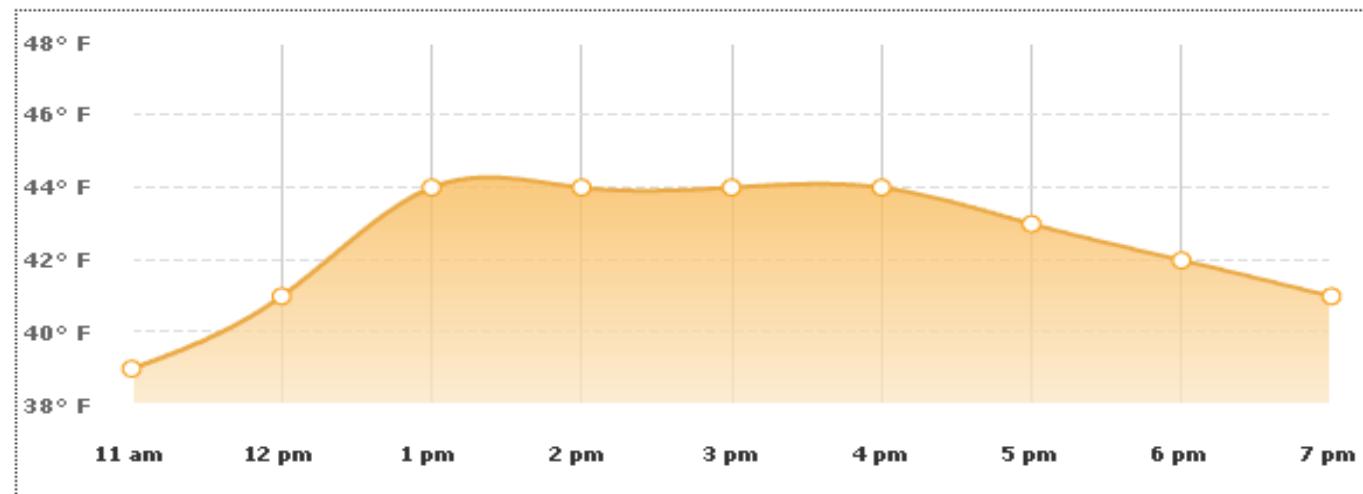
Labels?

Classification/Regression?

11 am	12 pm	1 pm	2 pm	3 pm	4 pm	5 pm	6 pm
39° F	41° F	44° F	44° F	44° F	44° F	43° F	42° F

Precip:
10%

Precip:
0%



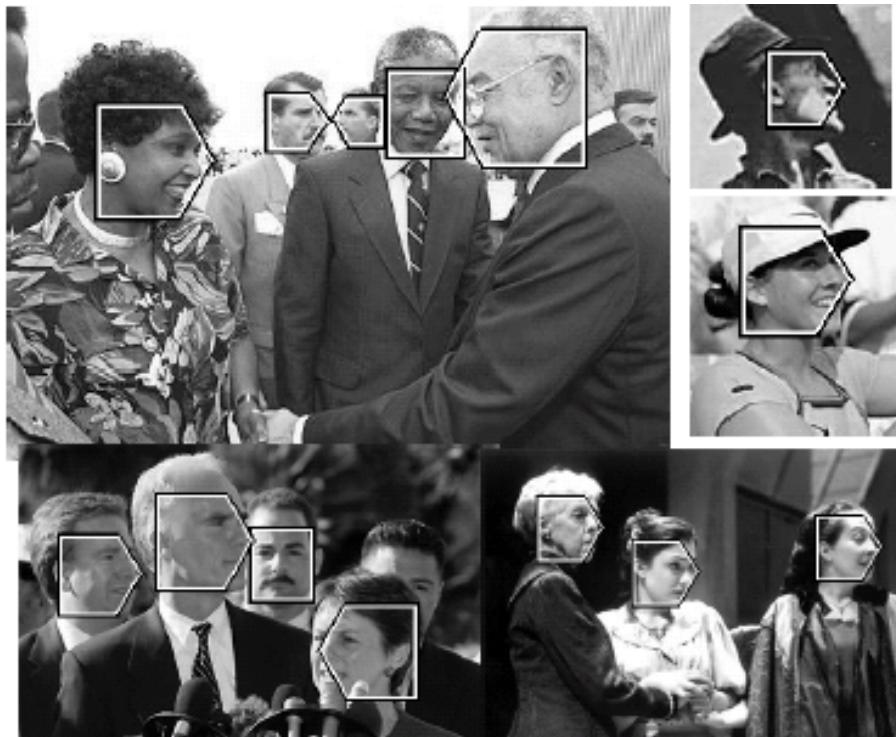
Temperature/Weather prediction

Prediction problems

Features?

Labels?

Classification/Regression?



Face Detection

Prediction problems

Features?

Labels?

Classification/Regression?



Environmental Mapping

Prediction problems

Features? Labels? Classification/Regression?



Robotic Control

Third Axis: Tasks

- Other than prediction problems, another class of tasks are **description** problems
- Examples:
 - Density estimation
 - Clustering
 - Dimensionality reduction
- Also called **unsupervised learning**
 - When first axis (data) consists only of inputs
 - No “supervision” in data as to the descriptive outputs

Unsupervised Learning

Aka “learning without a teacher”

Feature Space \mathcal{X}



Words in a document

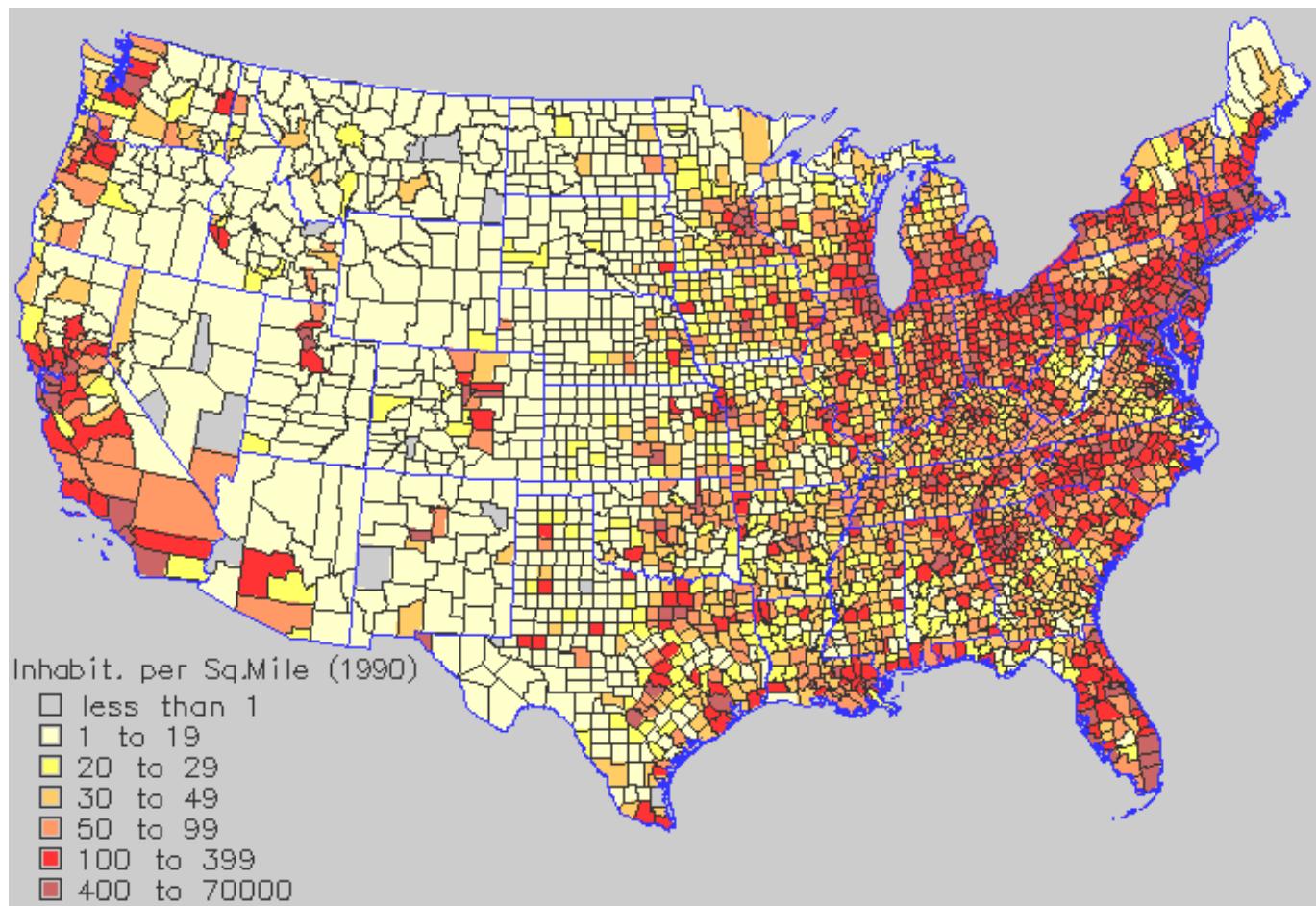


Word distribution
(Probability of a word)

Task: Given $X \in \mathcal{X}$, learn $f(X)$.

Unsupervised Learning – Density Estimation

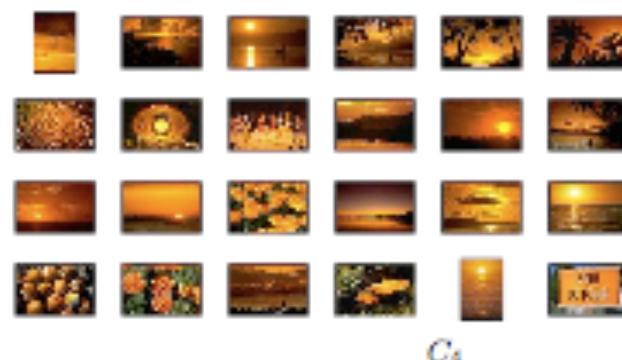
Population density



Unsupervised Learning – Clustering

Group similar things e.g. images

[Goldberger et al.]



Unsupervised Learning – clustering web search results

The screenshot shows the Clusty search interface. At the top, there's a navigation bar with links for 'web', 'news', 'images', 'wikipedia', 'blogs', 'jobs', and 'more'. Below that is a search bar containing the word 'race'. To the right of the search bar are 'Search' and 'advanced preferences' buttons. The main content area starts with a message: 'Cluster Human contains 8 documents.' On the left, there's a sidebar with tabs for 'clusters', 'sources', and 'sites'. The 'clusters' tab is selected, showing a list of categories: 'All Results (238)', 'Car (28)', 'Race cars (1)', 'Photos, Races Scheduled (5)', 'Game (4)', 'Track (3)', 'Nascar (2)', 'Equipment And Safety (2)', 'Other Topics (7)', 'Photos (22)', 'Game (14)', 'Definition (13)', 'Team (18)', 'Human (8)', 'Classification Of Human (2)', 'Statement, Evolved (2)', 'Other Topics (4)', 'Weekend (8)', 'Ethnicity And Race (7)', and 'Race for the Cure (8)'. The categories 'Race cars (1)', 'Photos, Races Scheduled (5)', 'Game (4)', 'Track (3)', 'Nascar (2)', 'Equipment And Safety (2)', 'Other Topics (7)', 'Game (14)', 'Definition (13)', 'Team (18)', and 'Ethnicity And Race (7)' are circled in red. The main content area lists the 8 documents found in the 'Human' cluster:

- Race (classification of human beings) - Wikipedia, the free ...**
The term **race** or racial group usually refers to the concept of dividing **humans** into populations or groups on the basis of visible traits (especially skin color, cranial or facial features and hair texture), and self-identified by culture and over time, and are often controversial for scientific as well as social and political reasons. History · McGraw-Hill · en.wikipedia.org/wiki/Race_(classification_of_human_beings) - [cache] - Live, Ask
- Race - Wikipedia, the free encyclopedia**
General. **Racing** competitions The **Race** (yachting **race**), or La course du millénaire, a no-rules round-the-world sail of **human** beings) **Race** and ethnicity in the United States Census, official definitions of "race" used by the US Census Bureau, genetics. Historical definitions of **race**; **Race** (bearing), the inner and outer rings of a rolling-element bearing. **RACE** · Literature · Video games
en.wikipedia.org/wiki/Race - [cache] - Live, Ask
- Publications | Human Rights Watch**
The use of torture, unlawful rendition, secret prisons, unfair trials, ... Risks to Migrants, Refugees, and Asylum Seekers
...
www.hrw.org/backgrounder/usa/race - [cache] - Ask
- Amazon.com: Race: The Reality Of Human Differences: Vincent Sarich ...**
Amazon.com: **Race: The Reality Of Human Differences**: Vincent Sarich, Frank Miele: Books ... From Publishers Weekly
www.amazon.com/Race-Reality-Differences-Vincent-Sarich/dp/0813340861 - [cache] - Live
- AAPA Statement on Biological Aspects of Race**
AAPA Statement on Biological Aspects of **Race** ... Published in the American Journal of Physical Anthropology, vol. 111, pp. 1-12, 2000. evolution and variation, ...
www.physanth.org/positions/race.html - [cache] - Ask
- race: Definition from Answers.com**
race n. A local geographic or global **human** population distinguished as a more or less distinct group by genetically similar characteristics.
www.answers.com/topic/race-1 - [cache] - Live

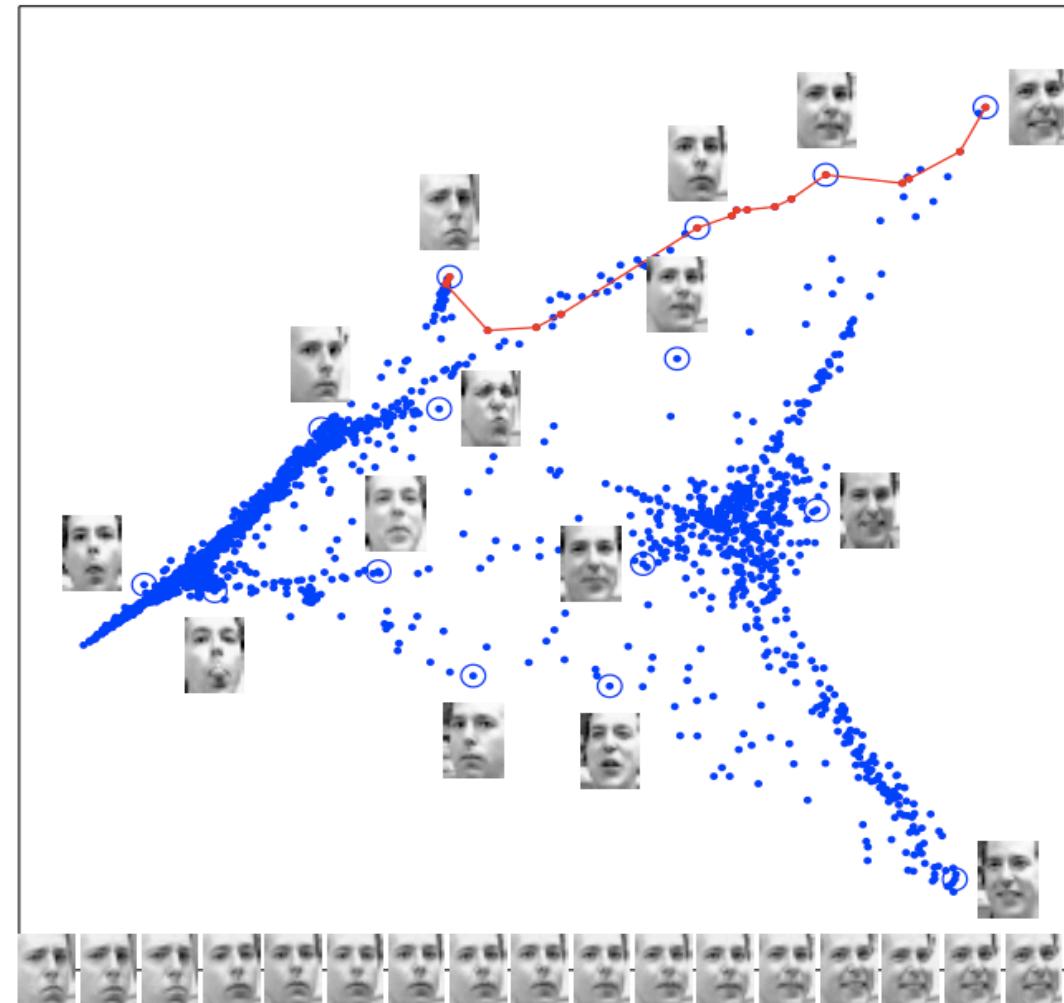
Unsupervised Learning - Embedding

Dimensionality Reduction

[Saul & Roweis '03]

Images have thousands or millions of pixels.

Can we give each image a coordinate,
such that similar images are
near each other?



Unsupervised Learning - Embedding

Dimensionality Reduction - words

[Joseph Turian]



Unsupervised Learning - Embedding

billmark mary
bob jack stephen elizabeth
tony edward
miss jimmike brian richard alexander
steve chris rich william charles
joe tom harry robinson francis maria
mr. andrew frank paul davida james louis
mr. sam arthur george jean thomas
don ray martin
simon howard

dr. ben al
r. a. lee
m. e. h. j. scott lewis bush

c. s. w. wilson jackson fox
b. d. p. taylor johnson fox
von smith williams
van jones davis ford grant
bell

van

los angeles

la los
et dad el san
santa

des hong
core

cape

east

june august
february november
january september
april october
december march

amkong

usa philippines

virginia
columbia missouri
indiana maryland
colorado tennessee
washington wisconsin
oregon kansas
california carolina
houston philadelphia
philadelphia pennsylvania
new jersey georgia
detroit toronto
chicago ontario
massachusetts york
sydney new zealand
boston melbourne
montreal cambridge
london manchester
victoria quebec
moscow quebec
mexico scotland
walengland
canada ireland britain
australia sweden
singapore spain
america norway france
europe austria
asia germany
africa russia
india japan rome
korea china
pak egypt
israel
vietnam

Machine Learning Tasks

Broad categories -

- **Prediction tasks**
 - Classification, Regression
- **Description tasks**
 - Typically synonymous with unsupervised learning
 - Density estimation, Clustering, Dimensionality reduction

Machine Learning Subfields

- Topics in machine learning can be categorized by these three axes
 - Data
 - Tasks i.e. kind of knowledge
 - Algorithms

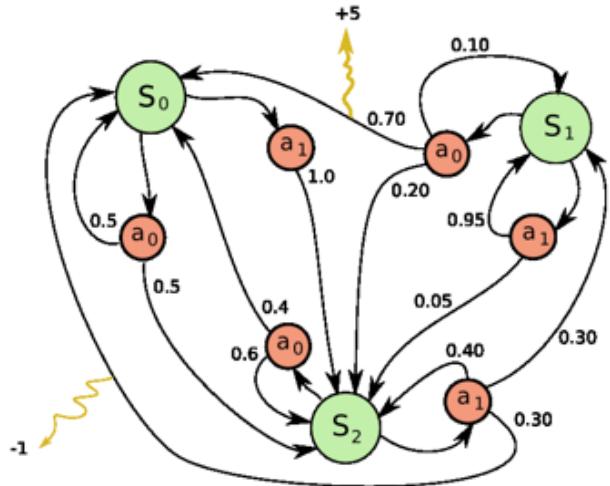
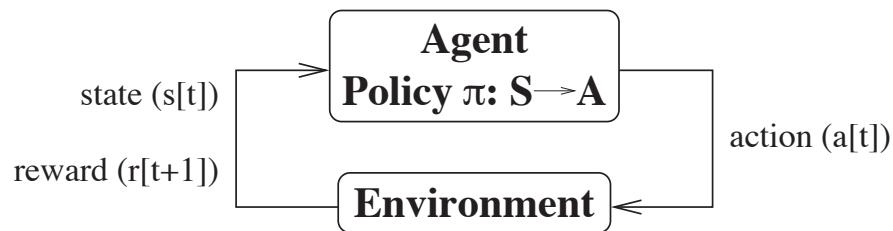
Machine Learning Subfields

- **Particular choices of the three axes**
- **Supervised learning**
 - First axis (data) consists of both inputs and outputs
 - Third axis (tasks) consists of prediction
- **Semi-supervised learning**
 - First axis (data) consists of inputs and only some of them with outputs
 - Third axis (tasks) consists of prediction
- Many more (active learning,)

Some other Machine Learning Subfields

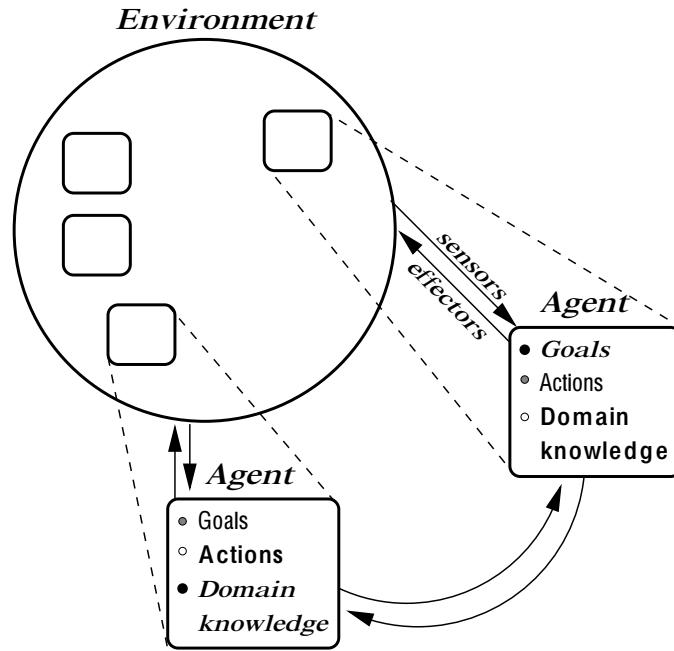
- Reinforcement learning
- Multi-agent systems
 - In addition to observations, we also have “actions” we can take

Reinforcement Learning



- Data consists of rewards that come through taking **actions** that has a feedback i.e. affects future observations
- Task: maximize reward

Multi-agent Systems



- Multiple agents
- Same setup as in reinforcement learning
- But now, the data also consists of other agents' actions

Enjoy!

- ML is becoming ubiquitous in science, engineering and beyond
- This class should give you the basic foundation for applying ML and developing new methods
- The fun begins...