A 3D rendered white robot with blue eyes is sitting at a desk, reading a book. The robot has a rounded head and a friendly appearance. The background is plain white.

Machine Learning

Info on the course

Paola Velardi

What is machine learning (in a nutshell)

- A set of methodologies to find regularities in data
- These findings are use to predict future outcomes and/or to classify unseen data

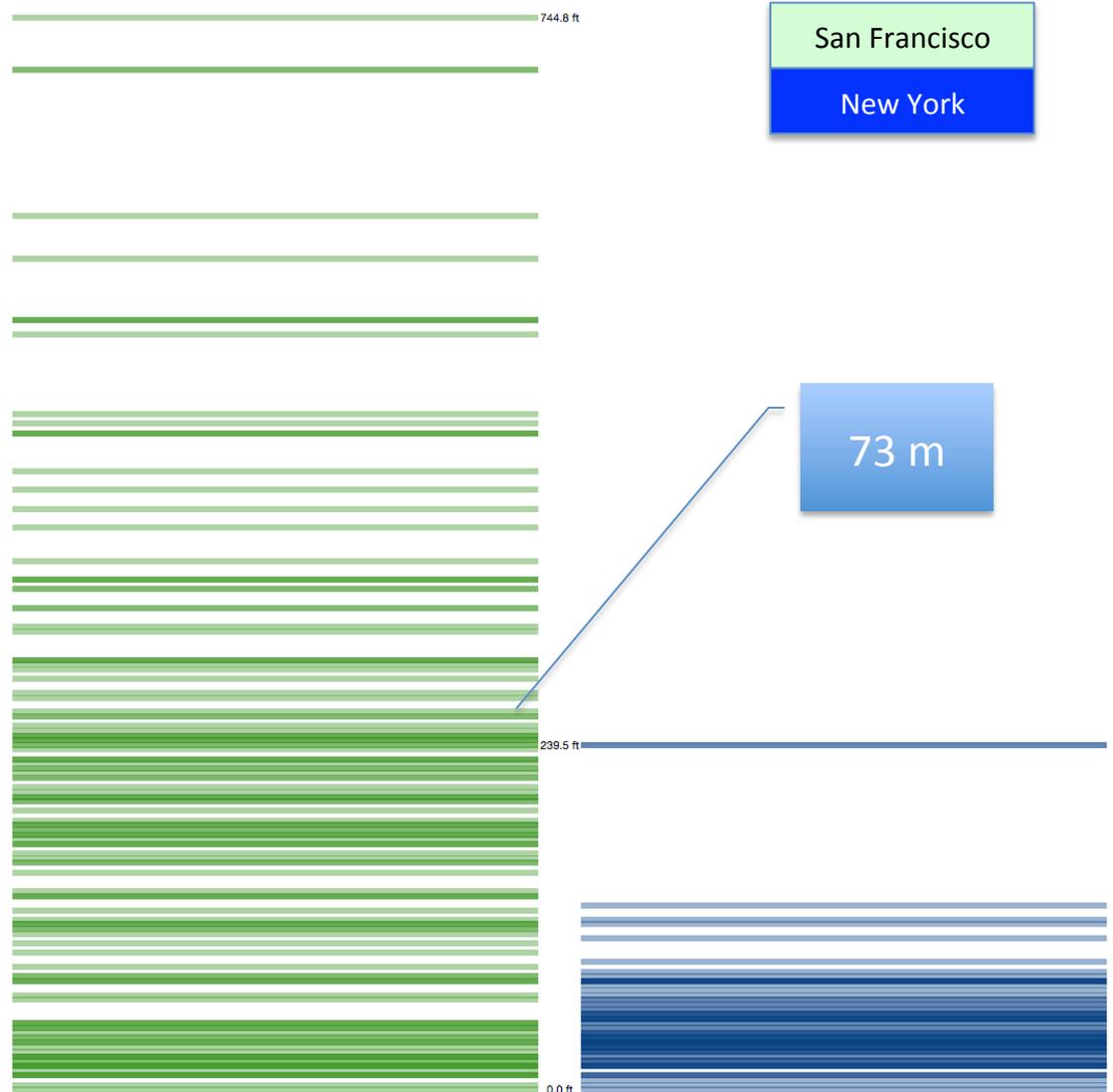
A (VERY) simple example

- Using a data set about homes, we wish to create a machine learning model to distinguish homes in New York from homes in San Francisco.
- Let's say you had to determine whether a home is in **San Francisco** or in **New York**. In machine learning terms, categorizing data points is a **classification** task.

To begin, we consider building elevation

- Since San Francisco is relatively hilly, the elevation of a home may be a good way to distinguish the two cities.
- Based on the home-elevation data to the right, you could argue that a home above 73 meters should be **classified** as one in San Francisco.
- So we can infer the following rule:

IF elevation > 73m THEN city = SF

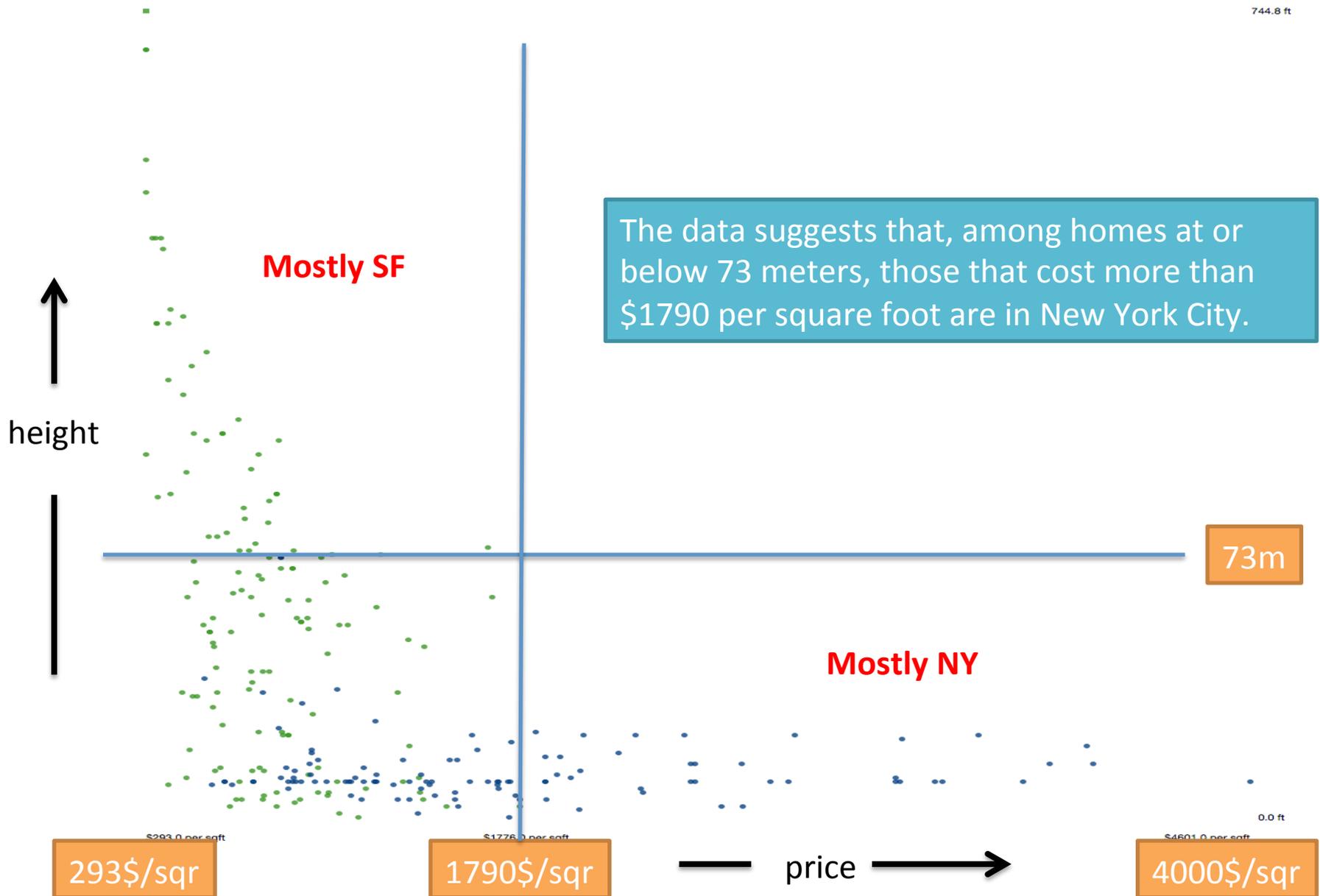


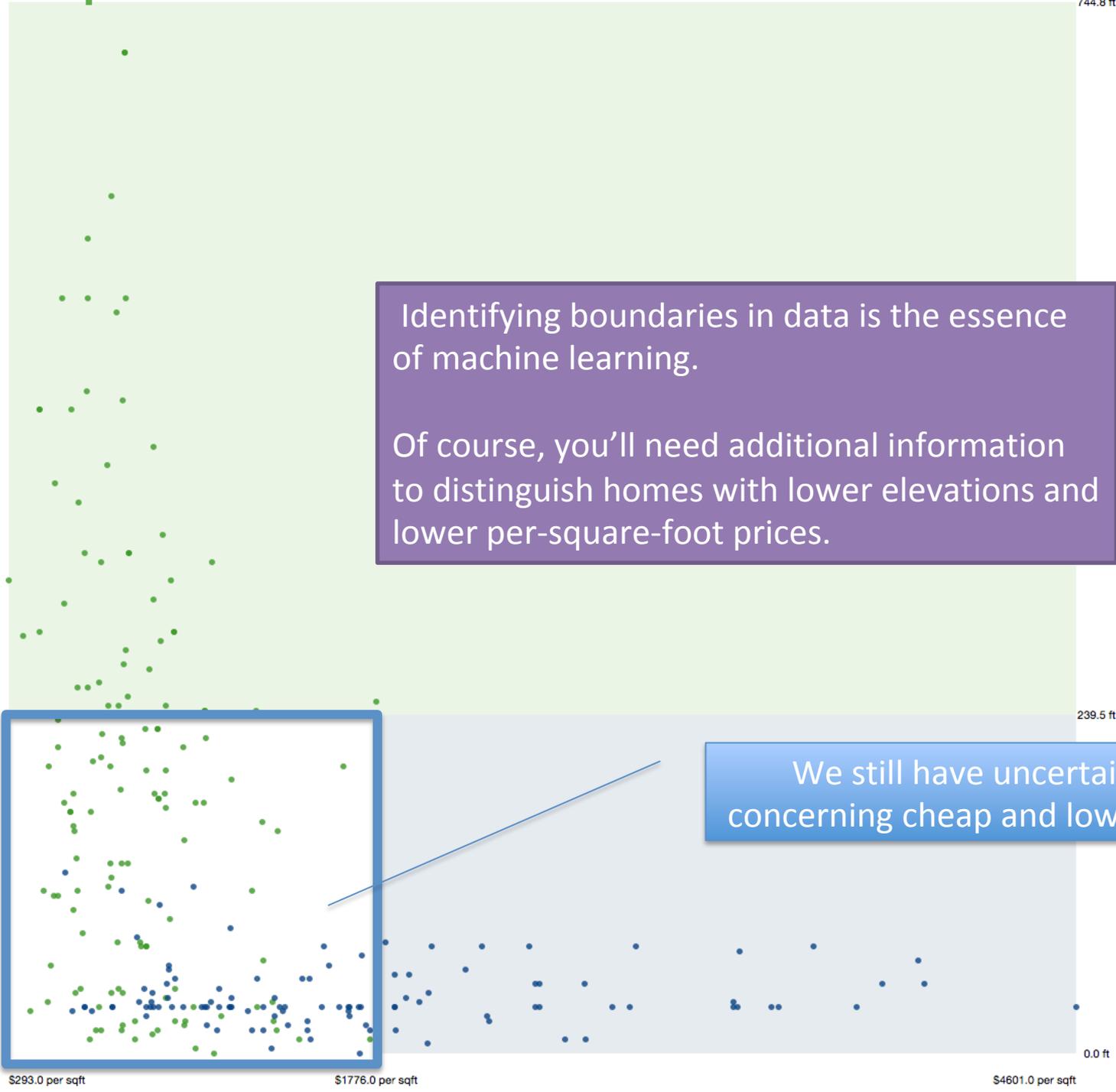
Is this enough?

- What if elevation < 73 ? We cannot tell..
- Adding another **dimension** allows for more nuance. For example, New York apartments can be extremely expensive per square foot.
- Dimensions (such as height, price) are denoted as **FEATURES** in machine learning (= what may characterize the objects we wish to classify). They are also called attributes, descriptors, dimensions..

Elevation *and* price per square foot scatterplot

744.8 ft





Identifying boundaries in data is the essence of machine learning.

Of course, you'll need additional information to distinguish homes with lower elevations and lower per-square-foot prices.

We still have uncertainty concerning cheap and low homes

\$293.0 per sqft

\$1776.0 per sqft

\$4601.0 per sqft

239.5 ft

0.0 ft

elevation

elevation

Year built

bathrooms

bedrooms

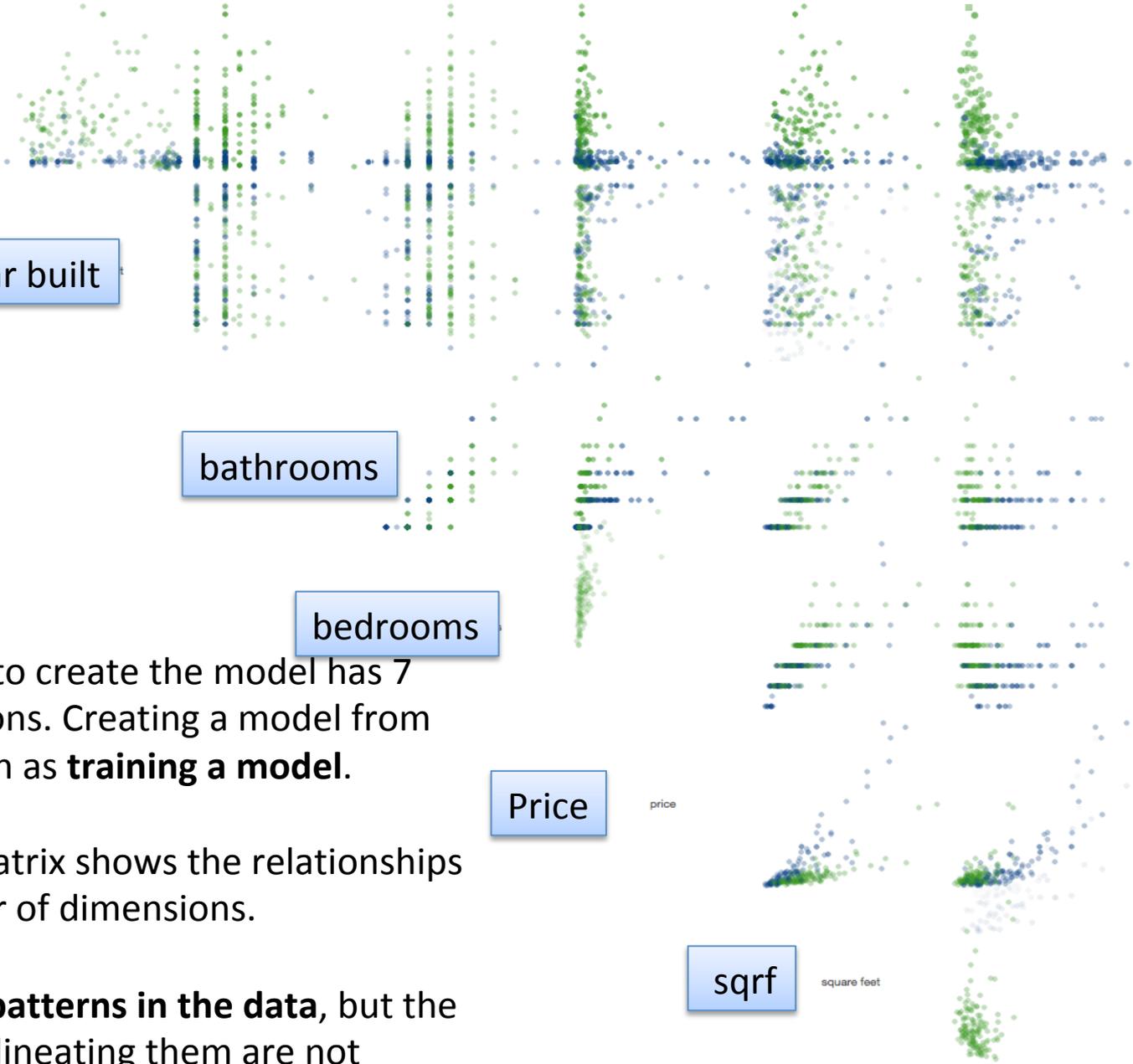
Price

price

sqrf

square feet

Price per sqrf



The dataset used to create the model has 7 different dimensions. Creating a model from data is also known as **training a model**.

The scatterplot matrix shows the relationships between each pair of dimensions.

There are clearly **patterns in the data**, but the boundaries for delineating them are not obvious.

So, machine learning is...

- Finding patterns in data
- Machine learning methods use learning algorithms to identify regularities and boundaries
- “Boundaries” may come in several forms, they can be probabilities, separating lines, or rules (if height >73 then SF)

Issues in Machine Learning

So, what is Machine Learning?

What is Machine Learning?

- Machine Learning is a scientific discipline that addresses the following question: *'How can we program systems to automatically learn and to improve with experience?'*
- **What** is learning?
- **How** do we learn?
- **What** can we learn?
- **How** can we “improve”, and over what??
- **What** is “experience”??

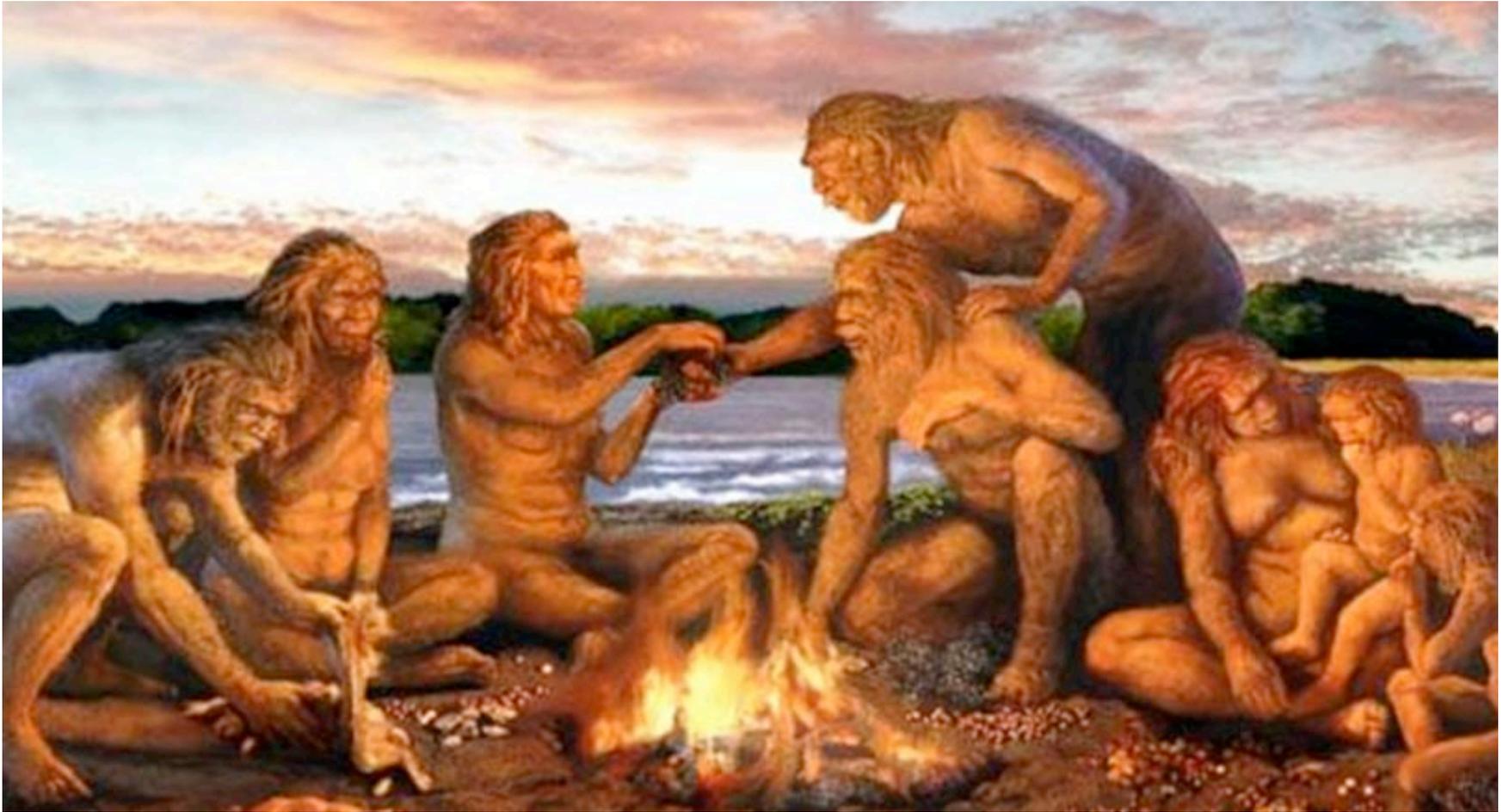
What is learning??



Fire burns!



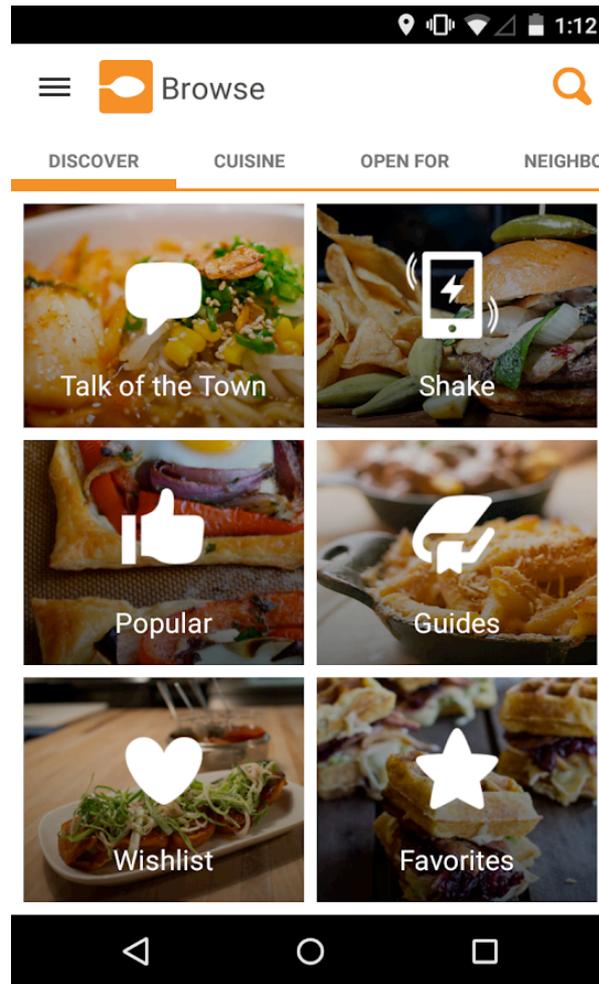
But we learned using it



you can study (learn) Machine learning



An then build an app to learn your favorite food and advise you on restaurants



So, what is learning?

- **Make sense** of a subject, event or feeling by interpreting it into our own words or actions.
- **Use** our newly acquired ability or knowledge - in conjunction with skills and understanding we already possess - **to do something useful** with the new knowledge or skill.
- <http://www.skillsyouneed.com/general/learning.html#ixzz3QtQQTKvg>

What is learning?

**UNDERSTAND + GAIN KNOWLEDGE +
USE NEW KNOWLEDGE TO DO
SOMETHING**

But, how do we learn??



How do humans learn?

- Someone tell us (**teacher**, or **whatching others**)
- Try and test (**learning by doing**) - As in the fire example -



There is only one thing more painful than learning from experience, and that is not learning from experience.

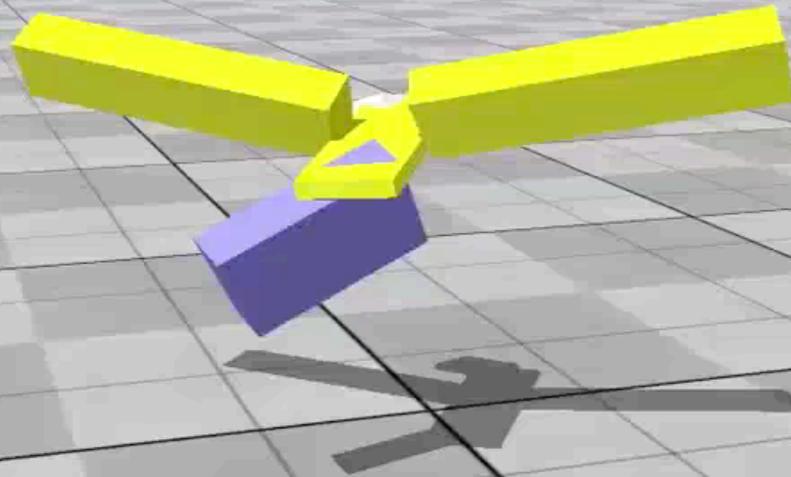
Laurence J. Peter

Is there something humans cannot learn??

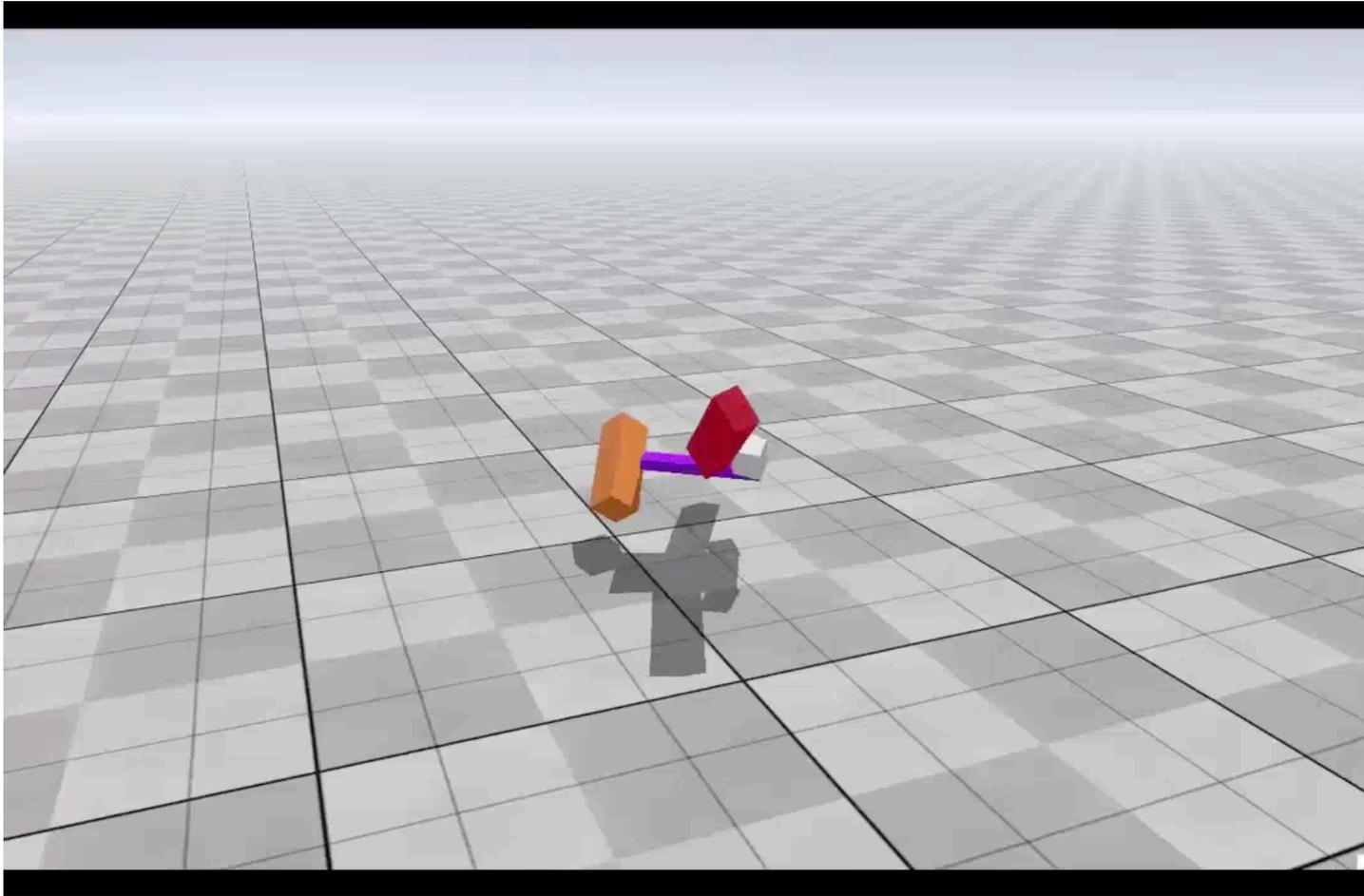


Machines that learn to fly

This is the story of a 2946 generation virtually evolved creature.
Evolved in the program "3D virtual creature evolution".



.. After some time..

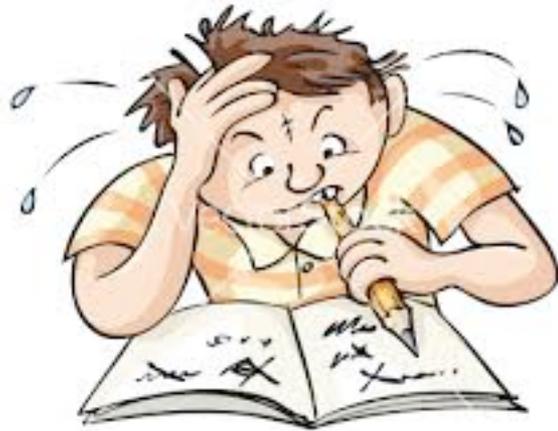


After some more time..

Here is the final creature. This is the best creature of generation 2946.

Besides things that we cannot learn,
there are others that are either..

- Difficult to learn
- Difficult to teach



When is it difficult for humans to learn? (1)

- If there are **too many data**, humans cannot easily make sense of them (e.g. finding regularities in human genoma, learning to recognize one among millions of objects, market analysis and forecasts)



When is it difficult for humans to learn? (2)

- If data **change too frequently**, humans might be unable to continuously adapt their knowledge (e.g. personalized recommendations, market analysis forecast)



The image shows a blurred screenshot of a stock market data table. The table has multiple columns, including stock prices, changes in price, and percentages. Green arrows indicate an increase in price, while red arrows indicate a decrease. The text is mostly illegible due to the blur, but some values are visible.

Symbol	Price	Change	% Change	Volume	Market Cap	PE Ratio	Dividend Yield
(N/A)	15258.07	↑ 1.89	↑ 0.02%	3798	11	100	0 (N/A)
(N/A)	158.62	↑ 2.44	↑ 0.01%	2403	95	100	0 (N/A)
(N/A)	14756.29	↑ 0.11	↑ 1.54%	2308	-20	100	0 (N/A)
(N/A)	7456.28	↑ 0.10	↑ 0.00%	3745	100	100	0 (N/A)
(N/A)	13856.03	↓ -0.15	↓ -0.00%	3075	107	100	0 (N/A)
(N/A)	7057.36	↑ 1.18	↑ 0.02%	3348	0	100	0 (N/A)
(N/A)	26516.93	↑ 2.93	↑ 0.01%	1074	90	100	0 (N/A)
(N/A)	308705.97	↑ 0.97	↑ 0.00%	473	90	100	0 (N/A)
(N/A)	16994.26	↓ -5.74	↓ -0.03%	1621	175	100	0 (N/A)
(N/A)	13666.09	↑ 0.09	↑ 0.00%	47	0	100	0 (N/A)
(N/A)	33145.62	↓ -4.38	↓ -0.01%	3679	90	100	0 (N/A)
(N/A)	16666.33	↑ 0.33	↑ 0.00%	4428	-17	100	0 (N/A)
(N/A)	513.81	↑ 1.81	↑ 0.35%	5856	-21	100	0 (N/A)
(N/A)	666.03	↓ -2.97	↓ -0.45%	-5	-20	100	0 (N/A)
(N/A)	581.81	↑ 1.81	↑ 0.31%	179	100	100	0 (N/A)
(N/A)	820.70	↑ 8.70	↑ 1.06%	42	-6	100	0 (N/A)
(N/A)	132.40	↑ 4.40	↑ 3.32%	229	95	100	0 (N/A)
(N/A)	998.26	↓ -0.74	↓ -0.07%	32	-48	100	0 (N/A)
(N/A)	59.26	↑ 5.26	↑ 0.05%	348	-2	100	0 (N/A)

Stock market values
And quotes

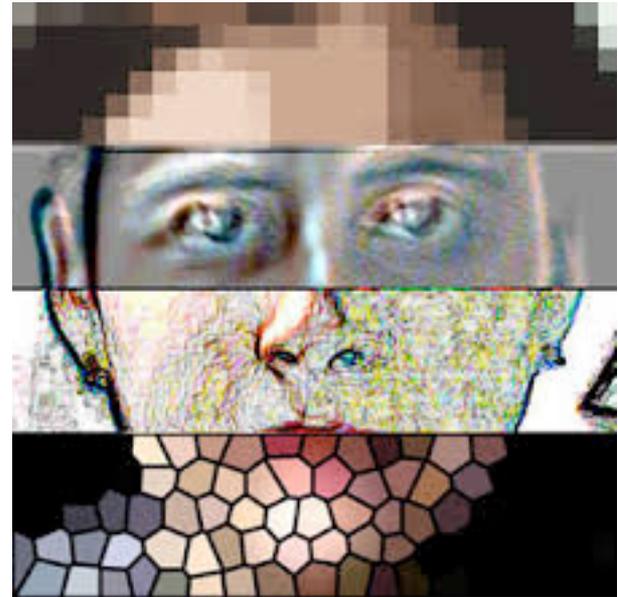
When is it difficult for humans to learn? (3)

- If the environment is dangerous, “learning by doing” cannot be applied (e.g. rescue systems)



When is it difficult for humans to teach? (4)

- If there is not enough information or previous expertise to “understand and gain knowledge” (we actually **do not understand** the image and speech recognition process by humans – it is not “teachable”)



So when is it advisable to use Machine Learning? (1)

- ML is used when:
 - No expertise
 - Human expertise does not exist (navigating on Mars), or there is a danger
 - Humans are unable to explain their expertise (speech/image recognition)
 - Too many data, data change frequently:
 - solution changes in time (market data for market forecast)
 - Solution needs to be adapted to particular cases (personalized systems for recommendation, diagnosis, etc.)

So when is it advisable to use Machine Learning? (2)

- Develop systems that are too difficult/expensive to construct manually because they require specific detailed skills or knowledge tuned to a specific task (*knowledge engineering bottleneck*).
 - Expert systems
- Develop systems that can automatically adapt and **customize** themselves to individual users.
 - Personalized news or mail filter
 - Personalized tutoring
 - Recommenders
- Discover new knowledge from large databases (*data mining*).
 - Market basket analysis (e.g. diapers and beer)
 - Medical text mining (e.g. migraines to calcium channel blockers to magnesium)
 - Twitter mining

An interdisciplinary topic: many related disciplines!

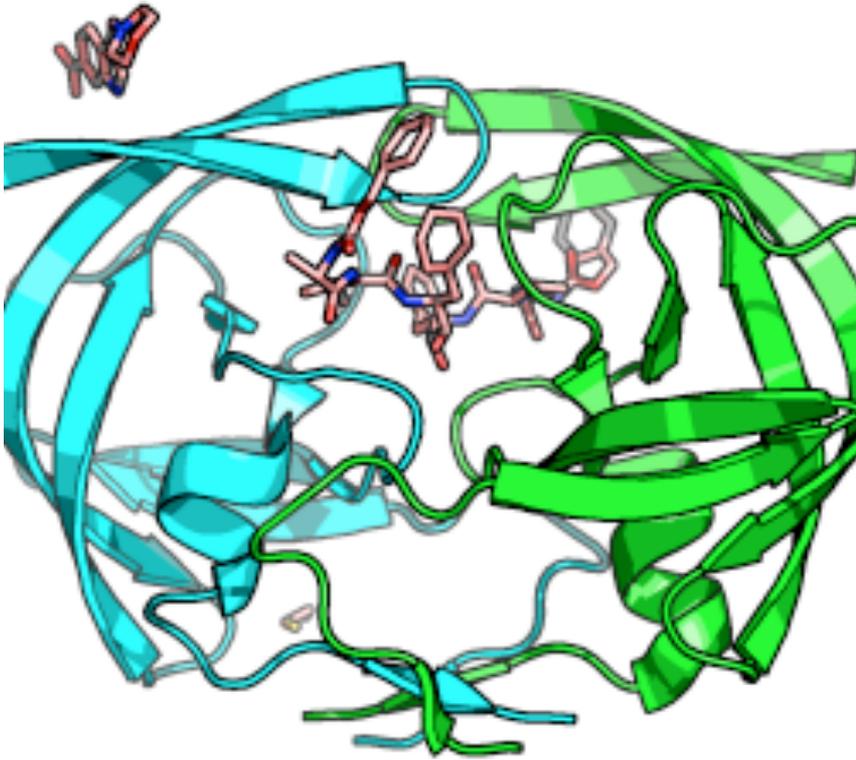
- Artificial Intelligence
- Data Mining
- Probability and Statistics
- Information theory
- Numerical optimization
- Computational complexity theory
- Control theory (adaptive)
- Psychology (developmental, cognitive)
- Neurobiology
- Linguistics
- Philosophy

ML is perhaps the most interdisciplinary of CS areas!!

Some “real hot” ML applications

- It is really hard to find a problem where machine learning is not already applied -- machine learning is practically everywhere, in business applications and science!
- Here is a list of “hot” applications:

Computational Biology & E-health



- Predicting diseases and complications from patient's health records
- Predicting disease genes through the analysis of biological networks (e.g. interactions between proteins)
- Predicting epidemics through the analysis of human interaction data (e.g., population density, data on population moves, climatic data, etc.)

Web Search and Recommendation Engines

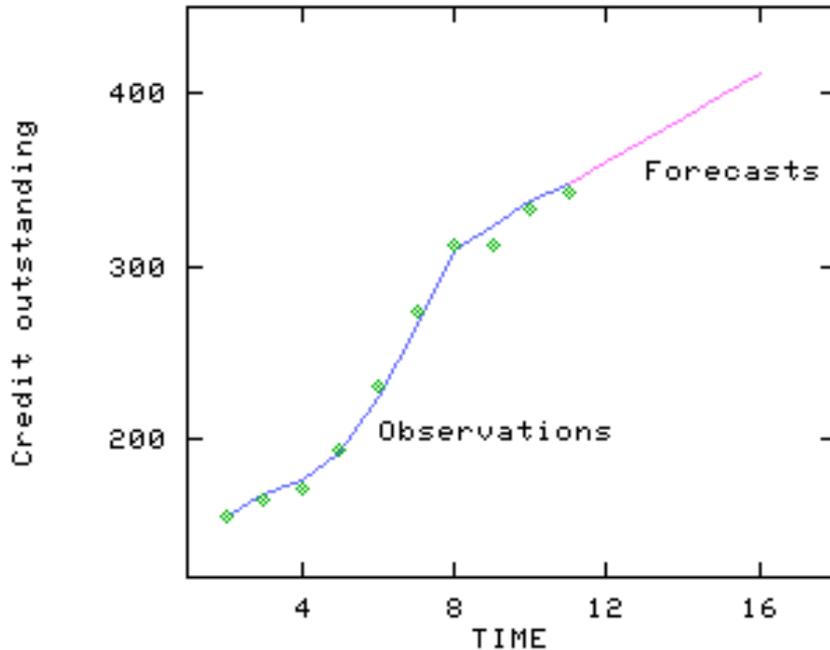
The screenshot shows a web design tool interface for configuring product recommendations. The main preview area displays a woman in a grey jacket with the text "We Think You'll Love" and "NEW MARKETGOODS". Below this is a grid of three product cards:

- Earthtone Crop Top**: Lightweight, handspun cotton with blocked panels. \$44.50. Buy Now
- Dhusara Shawl**: Oversized summer shawl in handwoven chambray. \$44.50. Buy Now
- Khadi Pullover**: Handspun and handwoven cotton khadi. \$54.50. Buy Now

The right sidebar is titled "Product Recommendations" and contains settings for content, style, and settings. It shows "Number of recommendations" set to 3, "Range to display" set to "from 2 - 4", and "Optional Details" including Name, Price, and Button (set to "Buy Now").

- find relevant searches, predict which results are most relevant to us, return a ranked output (Google)
- recommend similar products (e.g., Netflix, Amazon, etc.)

Finance



- predict if an applicant is credit-worthy
- detect credit card fraud
- find promising trends on the stock market

Text and Speech Recognition

- handwritten digit and letter recognition at the post office
- voice assistants (Siri)
- language translation services

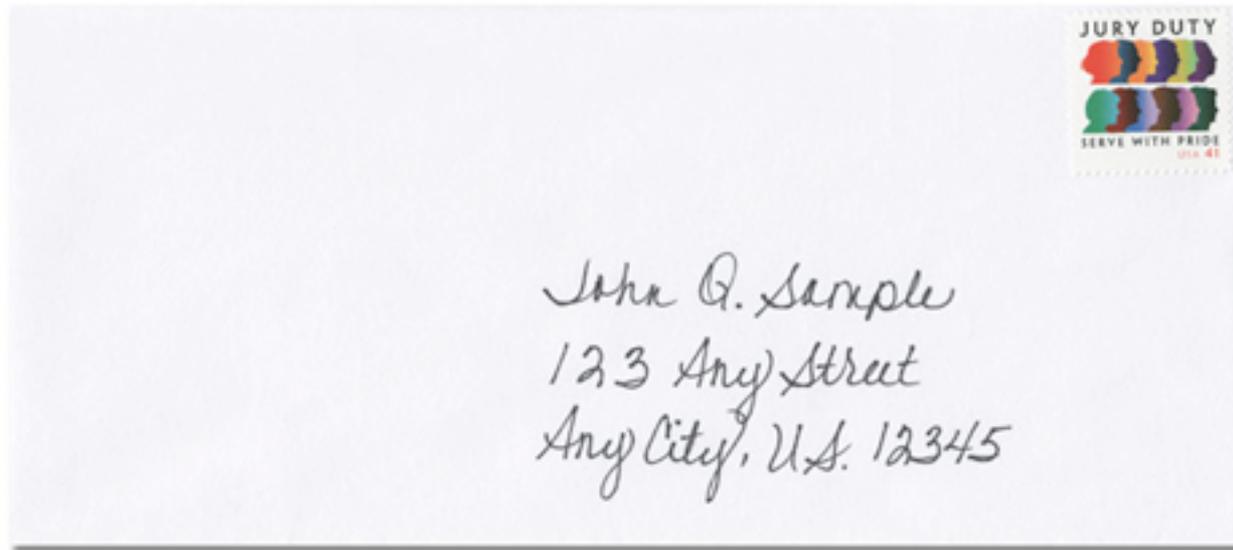


Image Understanding and Robotics



- Identification of relevant information (objects) in large amounts of Astronomy data
- Robotics for industry, energy saving, and smart cities
- Self-driving cars

Social Networks and Advertisement

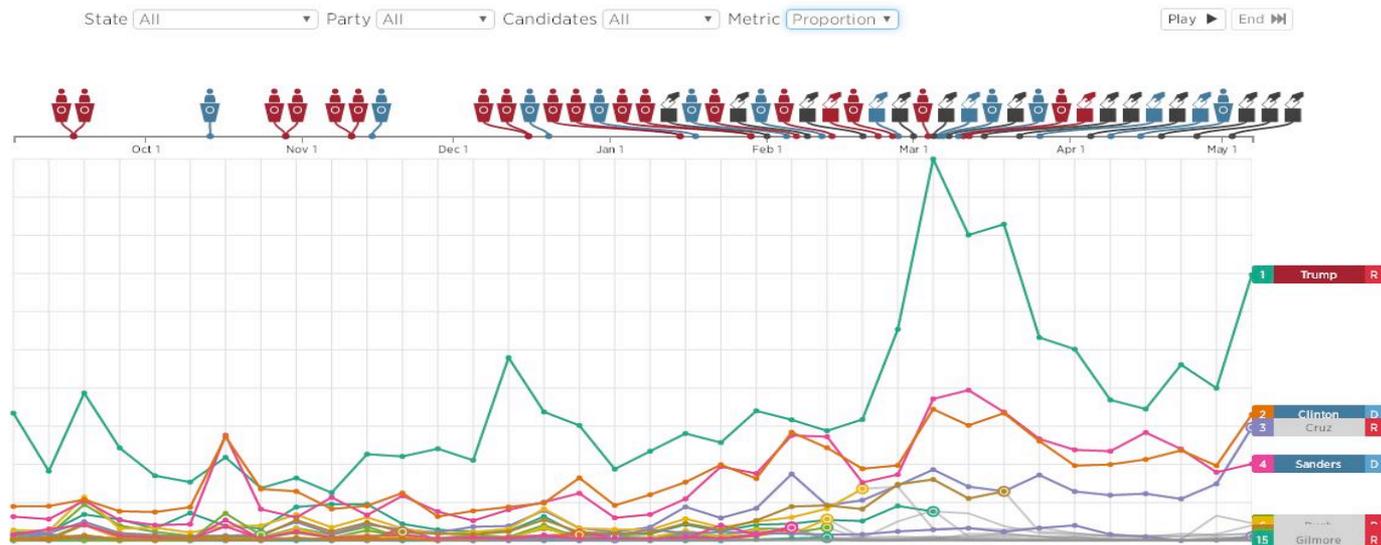
- Social data mining
 - data mining of personal information
 - Predict/analyze opinions, political choices, purchase behaviors

#interactive

Tweet Embed

#Election2016: US Presidential Candidate Twitter Buzz

As the fortunes of the 2016 US presidential candidates rise and fall throughout the campaign, so does the amount of conversation about them on Twitter. Below is an interactive graphic that allows you to take a look back at the amount of buzz each presidential candidate received on Twitter since September. By default, the graphic ranks all candidates using national data, but you can filter by party, state and status of candidacy or order it proportionally.



COURSE ORGANIZATION AND SYLLABUS

Course material

- Slides (partly) from:
<http://www.cs.utexas.edu/users/mooney/cs391L/> and many other sources
- Textbook: Tom Mitchell, Machine Learning, McGraw Hill, 1997 (new 2017 chapters on <http://www.cs.cmu.edu/~tom/NewChapters.html>)
- Introduction to machine learning ETHEM ALPAYDIN (on line book)
- Deep learning (MIT press):
<https://www.deeplearningbook.org/>
- Course twiki
<http://twiki.di.uniroma1.it/twiki/view/ApprAuto>

**THERE IS PLENTY OF MATERIAL ON THE
WEB, AND PLENTY OF DATASETS AND
LIBRARIES**

Course Syllabus

1. Concept learning and Decision Tree Learning
2. Practical Machine learning: feature selection and feature engineering
3. Evaluation methods: experimental and theoretical methods
4. Artificial Neural Networks and Deep learning (convolutional and denoising)
5. Support Vector Machines
6. Probabilistic Learning: Naive Bayes
7. Ensemble Methods
8. Unsupervised learning:
 1. Rule learning (Apriori, FP-Growth)
 2. Clustering methods
9. Reinforcement learning: Q-learning and Genetic Algorithms

Course labs

- Algorithms experimented on Weka toolkit <http://www.cs.waikato.ac.nz/ml/weka/> (only for those really weak with programming) AND on <http://scikit-learn.org/stable/> (a python platform). Other libraries can be used (e.g. TensorFlow for deep methods)
- Objective of labs is learning **practical ML**: selecting features, processing real life datasets, choosing algorithms, hyper-parameter tuning, evaluation experiments.

Caveat: Coverage of ML topics is limited!

- This is a first-level “basic” ML course
- On year 2 there is an advanced course
- ML algorithms for specific applications (NLP, security..) are also taught in other courses

Exam and Project

- Written exam on course material (60%)
- + Course Project: free choice of a topic (40%)
- Projects can be carried on by teams of 2 students
- How should the project be:
 - Not a trivial problem. Choose a **real life** problem, or invent one
 - Sufficiently **large data** set (many repositories available). Best if merging different datasets
 - Dataset must need some **feature engineering** (some non trivial pre-processing of data)
 - **More than one algorithm** tested, hyper-parameter tuning
 - **Evaluation** and analysis of results must show that you understand why you get a given result
 - I don't care if you get very good performances (in complex problems results might not be so good) but rather that you **understand what is going on**

How is the course organized

- Theoretical lessons + labs
- After every (or so) lesson, self-assessment are sent
- **PLEASE DO SUBSCRIBE TO GOOGLE GROUP** (use your Sapienza email and don't forgive to check it often or use redirect –don't miss my mails!)
- Self assessment useful to test your understanding of the subject. Very useful to pass the written test
- Google group **also useful to discuss self-assessment solutions among students!** (peer evaluation)