

# CPSC 340 and 532M: Machine Learning and Data Mining

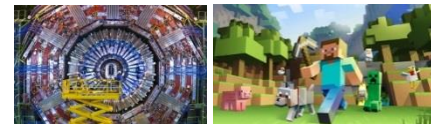
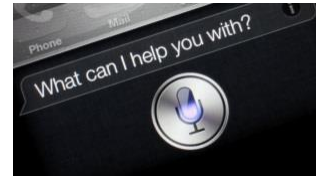
Mark Schmidt

University of British Columbia, Fall 2019

[www.cs.ubc.ca/~schmidtm/Courses/340-F19](http://www.cs.ubc.ca/~schmidtm/Courses/340-F19)

# Big Data Phenomenon

- We are **collecting and storing data** at an unprecedented rate.
- Examples:
  - YouTube, Facebook, MOOCs, news sites.
  - Credit cards transactions and Amazon purchases.
  - Transportation data (Google Maps, Waze, Uber)
  - Gene expression data and protein interaction assays.
  - Maps and satellite data.
  - Large hadron collider and surveying the sky.
  - Phone call records and speech recognition results.
  - Video game worlds and user actions.

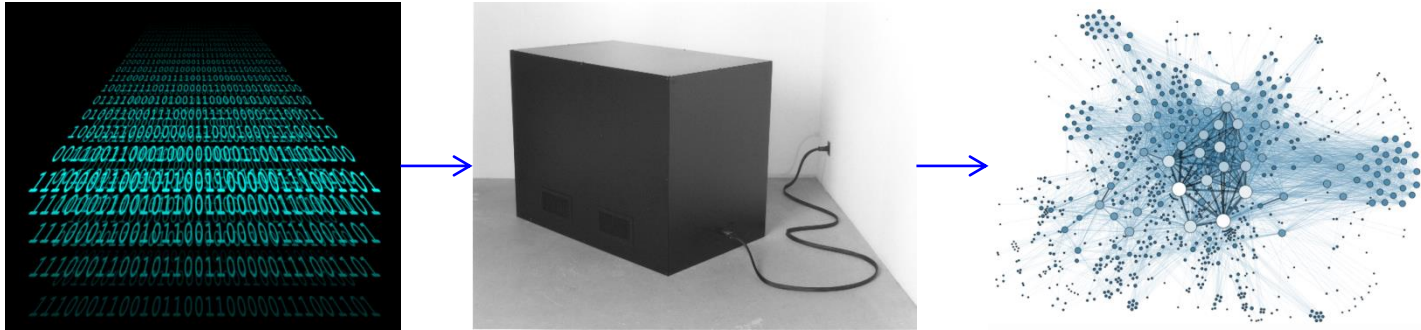


# Big Data Phenomenon

- What do you do with all this data?
  - Too much data to search through it manually.
- But there is valuable information in the data.
  - How can we use it for fun, profit, and/or the greater good?
- Data mining and machine learning are key tools we use to make sense of large datasets.

# Data Mining

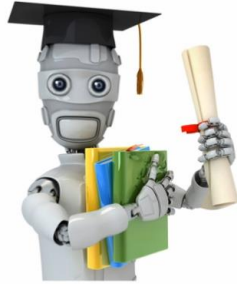
- Automatically **extract useful knowledge** from large datasets.



- Usually, to help with human decision making.

# Machine Learning

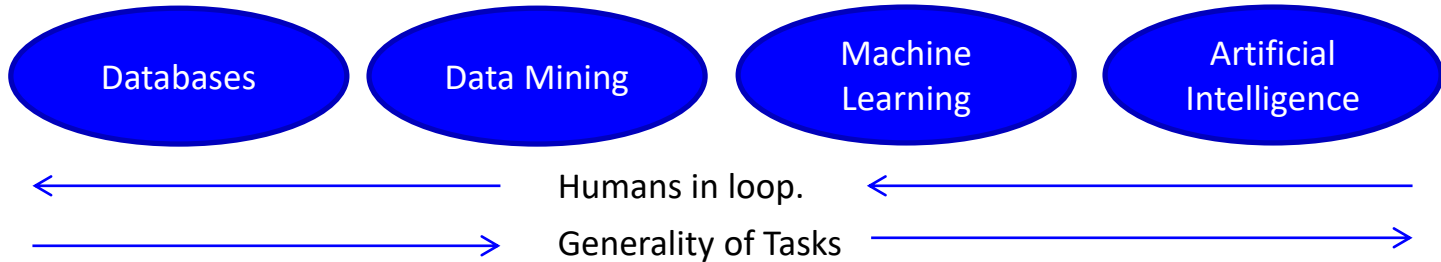
- Using computer to automatically **detect patterns in data and use these to make predictions** or decisions.



- Most useful when:
  - We want to automate something a human can do.
  - We want to do things a human can't do (look at 1 TB of data).

# Data Mining vs. Machine Learning

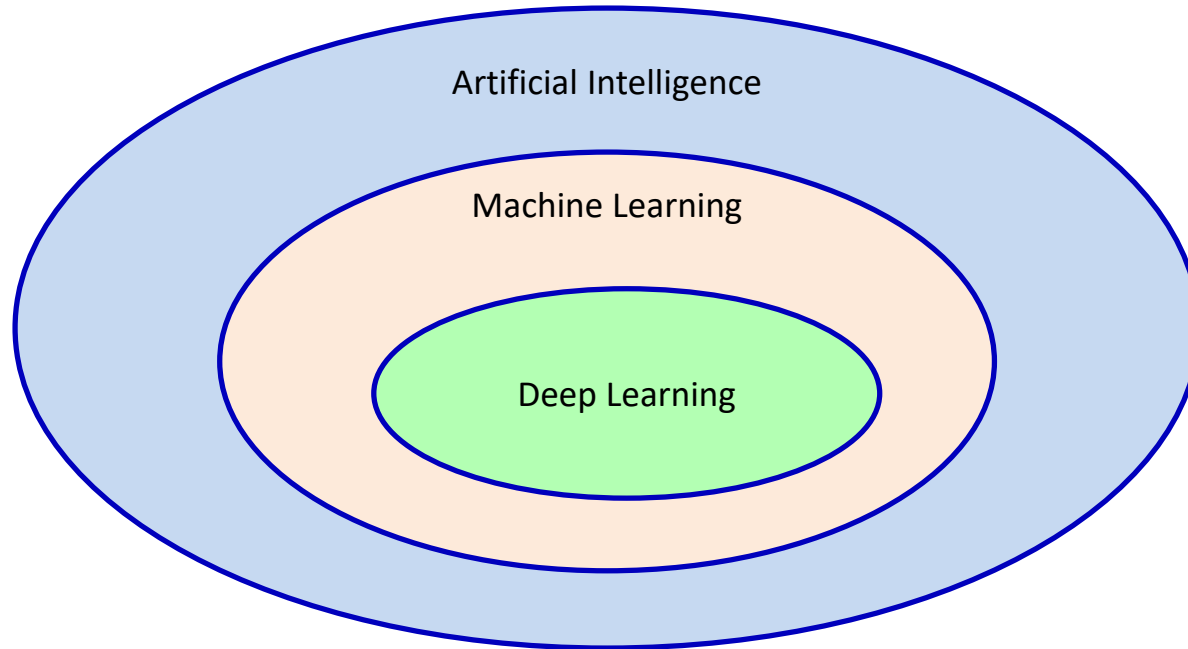
- Data mining and machine learning are very similar:
  - Data mining often viewed as closer to databases.
  - Machine learning often viewed as closer AI.



- Both are similar to statistics, but more emphasis on:
  - Large datasets and computation.
  - Predictions (instead of descriptions).
  - Flexible models (that work on many problems).

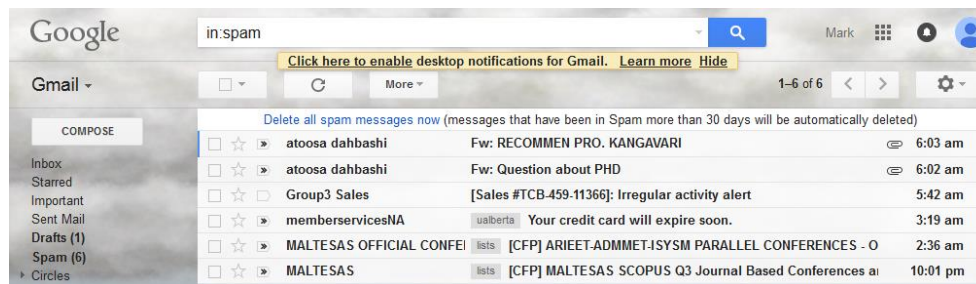
# Deep Learning vs. Machine Learning vs. AI

- Traditional we've viewed ML as a subset of AI.
  - And “deep learning” as a subset of ML.



# Applications

- Spam filtering:
- Credit card fraud detection:
- Product recommendation:



Transaction Date	Posted Date	Transaction Details	Debit	Credit
Aug. 27, 2015	Aug. 28, 2015	BEAN AROUND THE WORLD VANCOUVER, BC	\$10.95	

Customers Who Bought This Item Also Bought

Page 1 of 20

Pattern Recognition and Machine Learning (Information Science and...)   
 Christopher Bishop   
 ★★★★★☆ 115   
 Hardcover   
 \$60.76 ✓Prime

Learning From Data   
 Yaser S. Abu-Mostafa   
 ★★★★★☆ 88   
 Hardcover

The Elements of Statistical Learning: Data Mining, Inference, and Prediction...   
 Trevor Hastie   
 ★★★★★☆ 50   
 Hardcover   
 \$62.82 ✓Prime

Probabilistic Graphical Models: Principles and Techniques (Adaptive...)   
 Daphne Koller   
 ★★★★★☆ 28   
 Hardcover   
 \$91.66 ✓Prime

Foundations of Machine Learning (Adaptive Computation and...)   
 Mehryar Mohri   
 ★★★★★☆ 8   
 Hardcover   
 \$65.68 ✓Prime



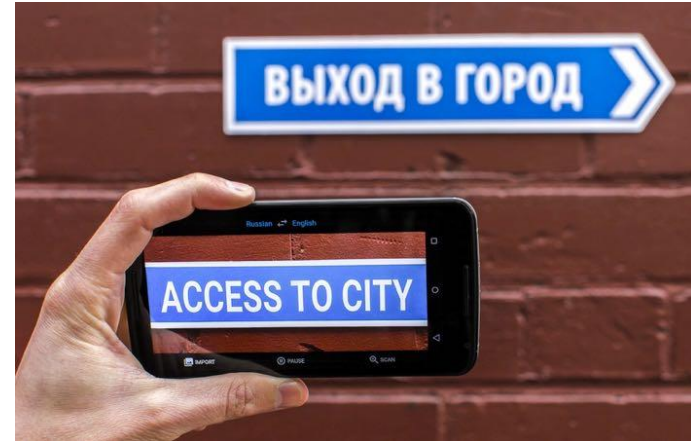
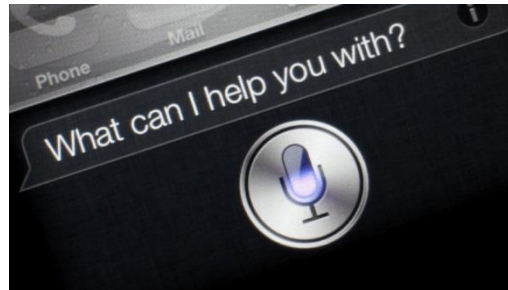
# Applications

- Motion capture:



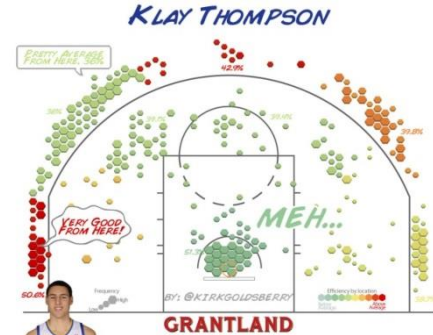
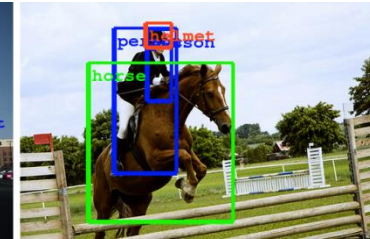
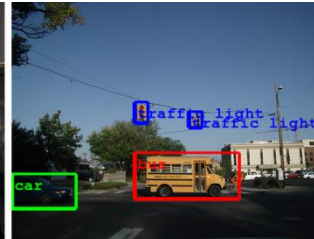
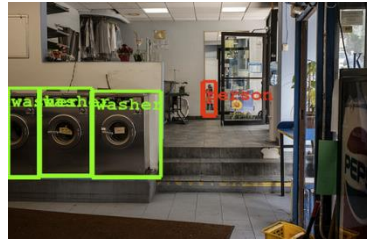
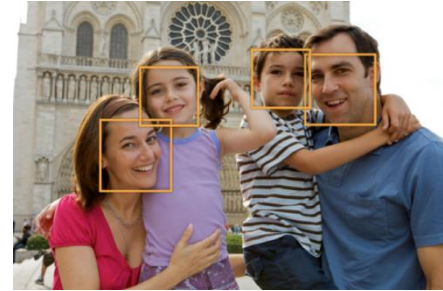
- Optical character recognition and machine translation:

- Speech recognition:



# Applications

- Face detection:
- Object detection:
- Sports analytics:

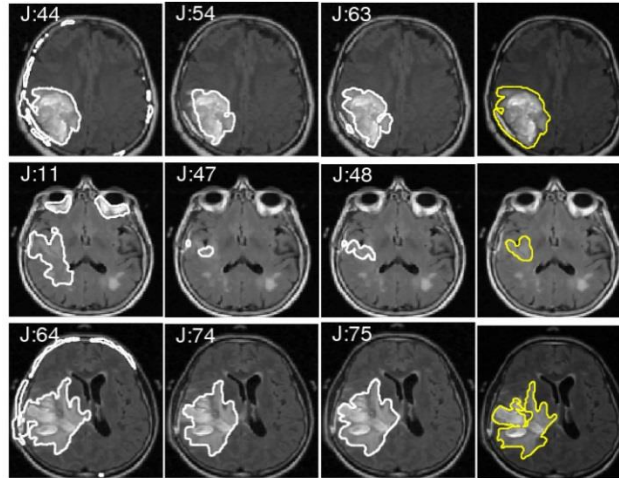


# Applications

- Personal Assistants:



- Medical imaging:

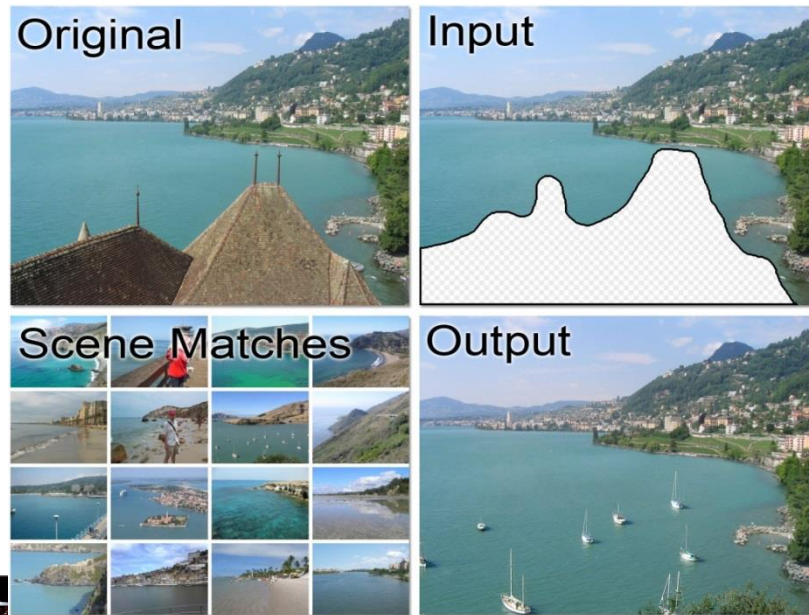


- Self-driving cars:



# Applications

- Scene completion:



- Image annotation:



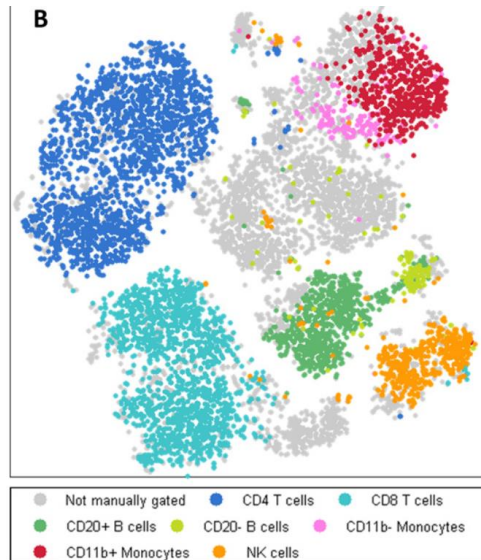
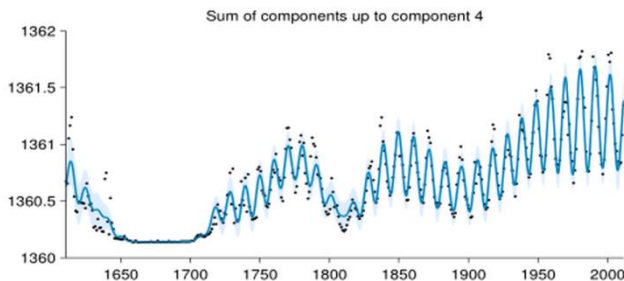
# Applications

- Discovering new cancer subtypes:

- Automated Statistician:

**2.4 Component 4 : An approximately periodic function with a period of 10.8 years. This function applies until 1643 and from 1716 onwards**

This component is approximately periodic with a period of 10.8 years. Across periods the shape of this function varies smoothly with a typical lengthscale of 36.9 years. The shape of this function within each period is very smooth and resembles a sinusoid. This component applies until 1643 and from 1716 onwards.





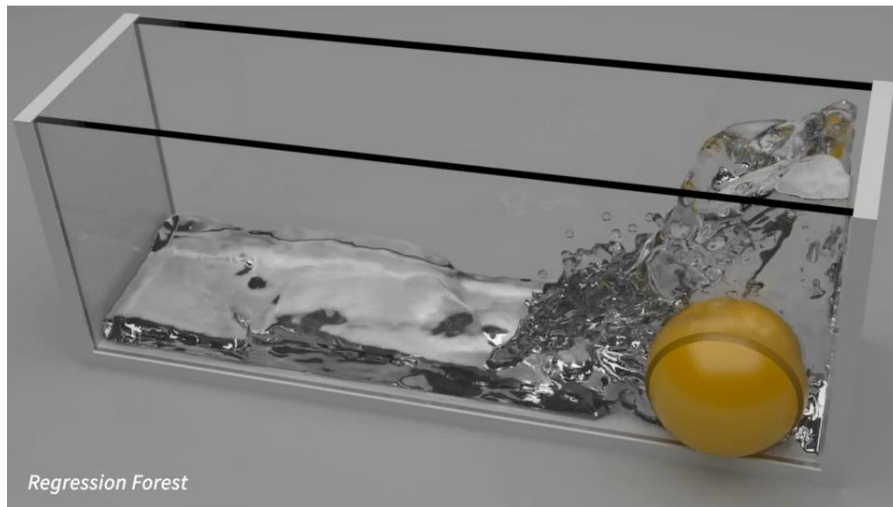
# Applications

- Mimicking artistic styles:



# Applications

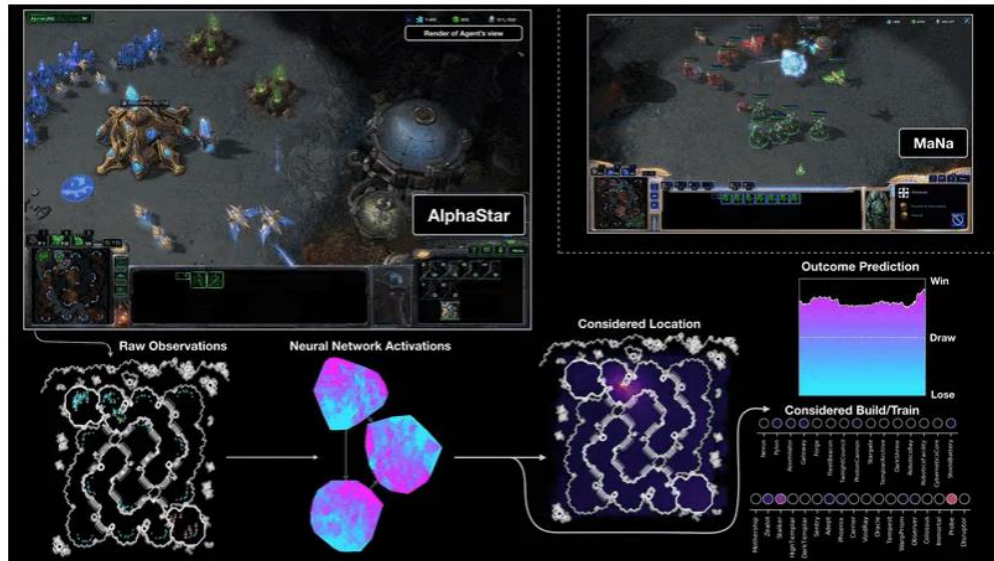
- Fast physics-based animation:



- Mimicking art style in [video](#).
- Recent work on generating text/music/voice/poetry/dance.

# Applications

- Beating humans in Go and Starcraft:





- Summary:
  - There is a lot you can do with a bit of statistics and a lot data/computation.
- We are in exciting times.
  - Major recent progress in fields like speech recognition and computer vision.
  - Things are changing a lot on the timescale of 3-5 years.
  - NeurIPS conference sold out in ~11 minutes last year.
  - A bubble in ML investments (most “AI” companies are just doing ML).
- But it is important to know the **limitations** of what you are doing.
  - “The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data.” – John Tukey
  - A huge number of people applying ML are just “**overfitting**”.
    - Or don’t understand the assumptions needed for them to work.
    - Their **methods do not work** when they are released “into the wild”.

# Bonus Slides

- I will include a lot of “bonus slides”.
  - May mention advanced variations of methods from lecture.
  - May overview big topics that we don’t have time for.
  - May go over technical details that would derail class.
- You are **not expected to learn** the material on these slides.
  - But they’re useful if you want to take 540 or work in this area.
- I’ll use this colour of background on bonus slides.

# Course Outline

- Next class discusses “exploratory data analysis”.
- After that, the remaining lectures focus on five topics:
  - 1) Supervised Learning.
  - 2) Unsupervised learning.
  - 3) Linear prediction.
  - 4) Latent-factor models.
  - 5) Deep learning.
- [“What is Machine Learning?”](#) (overview of many class topics)

Photo I took in the UK on the way home from the “Optimization and Big Data” workshop:

