

# IEMS 304 Lecture 1: Introduction to Statistical Learning

---

Yiping Lu

yiping.lu@northwestern.edu

*Industrial Engineering & Management Sciences*  
*Northwestern University*



NORTHWESTERN  
UNIVERSITY

# Logistics

---

**Textbook:** James G, Witten D, Hastie T, et al. *An introduction to statistical learning.*

**Time and Location:** Monday, Wednesday and Friday, 9.00 A.M.- 9.50 A.M.  
Tech L251

**Office Hour:** Friday: 1 P.M. Tech M237

**TA Office Hour:**

## Pre-requisite and Pre-test

This is a **mathematically intense** course. But that's why it's exciting and rewarding!

**Pre-requisite:** A previous course in statistics at the level of IEMS 303 plus a course in matrix analysis. Comfort with programming (we will be programming in R) is also necessary.

**Pre-test:** Passing the pretest is worth 3% of your final course grade. You must achieve a passing score of 70% or higher by

Monday, October 2 at 11:59 p.m. This deadline will be firmly enforced.

# Honor Code

## Do's

- form study groups (with arbitrary number of people); discuss and work on homework problems in groups
  - write down the solutions independently
  - write down the names of people with whom you've discussed the homework
  - use ChatGPT as a TA
- 

## Don'ts

- It is an honor code violation to copy, refer to, or look at written or code solutions from a previous year, including but not limited to: official solutions from a previous year, solutions posted online, solutions you or someone else may have written up in a previous year, and solutions for related problems.
- Directly copy the answer from ChatGPT

# Lab Session

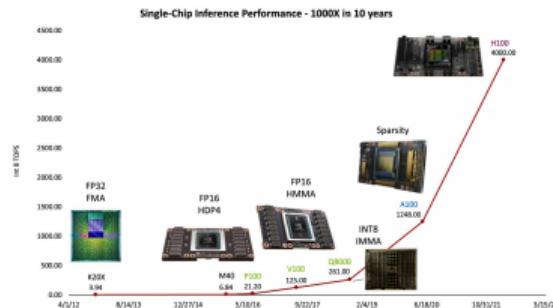
---

# Massive Data

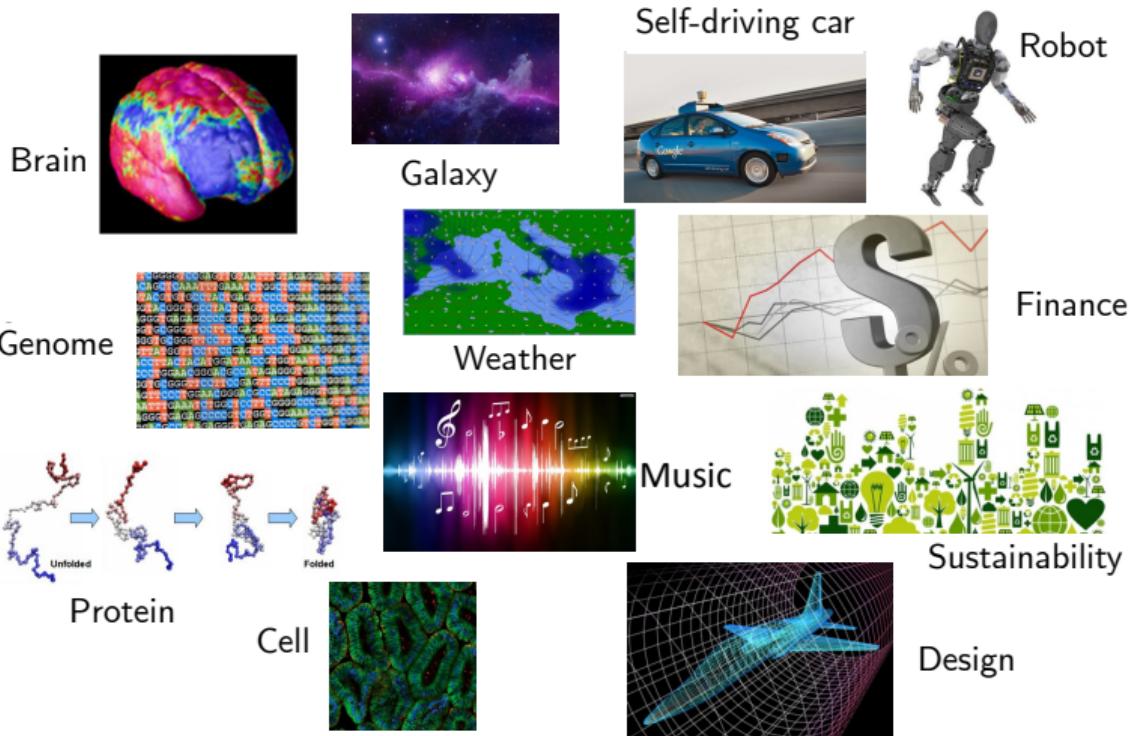
Massive complex data : Images, Acoustic signals, Text, ...

- Wikipedia pages: 13 millions (2014), 57 million (2022)
- Facebook users: 800 million (2014), 2.96 billion (2022)
- Flickr photos: 6 billion (2014), 10 billion (2022)
- Twitter tweets/day: 340 million (2014), 500 million (2022)
- Youtube video/min: 24 hours (2014), 500 hours (2022)
- Google pages:  $\geq 1$  trillion (2014),  $\geq 130$  trillions (2016)

Massive Computing : Huang's Law



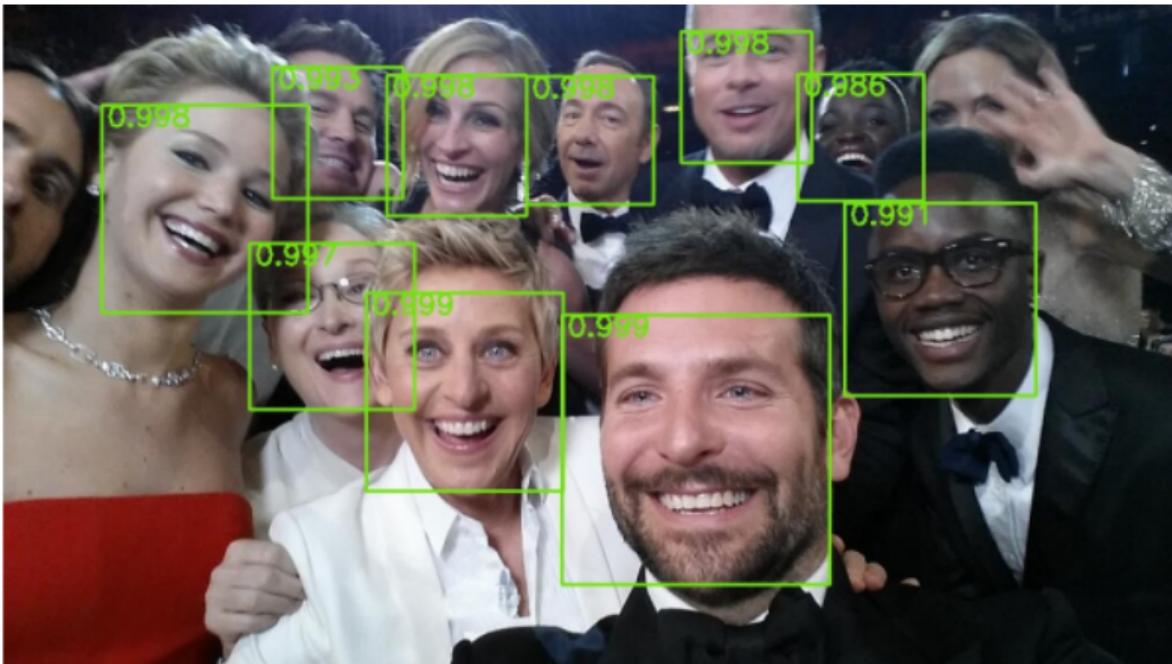
# Broad Applications in Science and Engineering



# Image Classification

			
<b>mite</b> mite black widow cockroach tick starfish	<b>container ship</b> container ship lifeboat amphibian fireboat drilling platform	<b>motor scooter</b> go-kart moped bumper car golfcart	<b>leopard</b> leopard jaguar cheetah snow leopard Egyptian cat
			
<b>grille</b> convertible grille pickup beach wagon fire engine	<b>mushroom</b> agaric mushroom jelly fungus gill fungus dead-man's-fingers	<b>cherry</b> dalmatian grape elderberry ffordshire bullterrier currant	<b>Madagascar cat</b> squirrel monkey spider monkey titi indri howler monkey

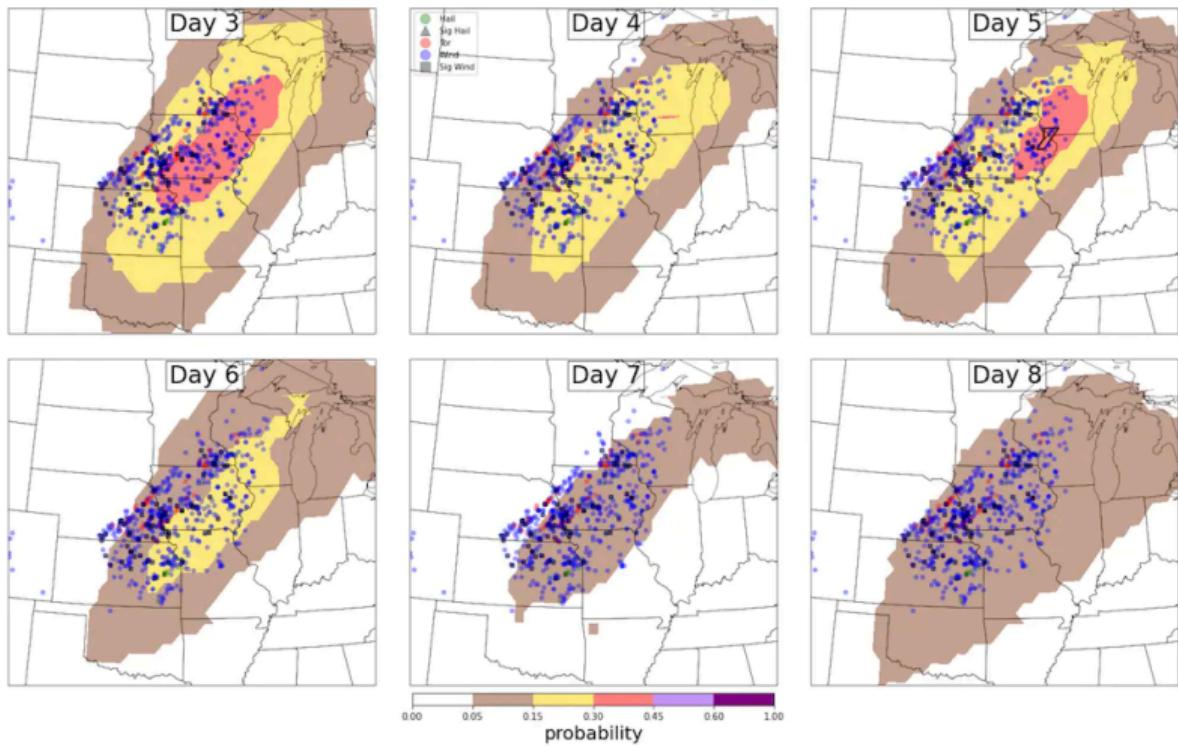
# Face Detection



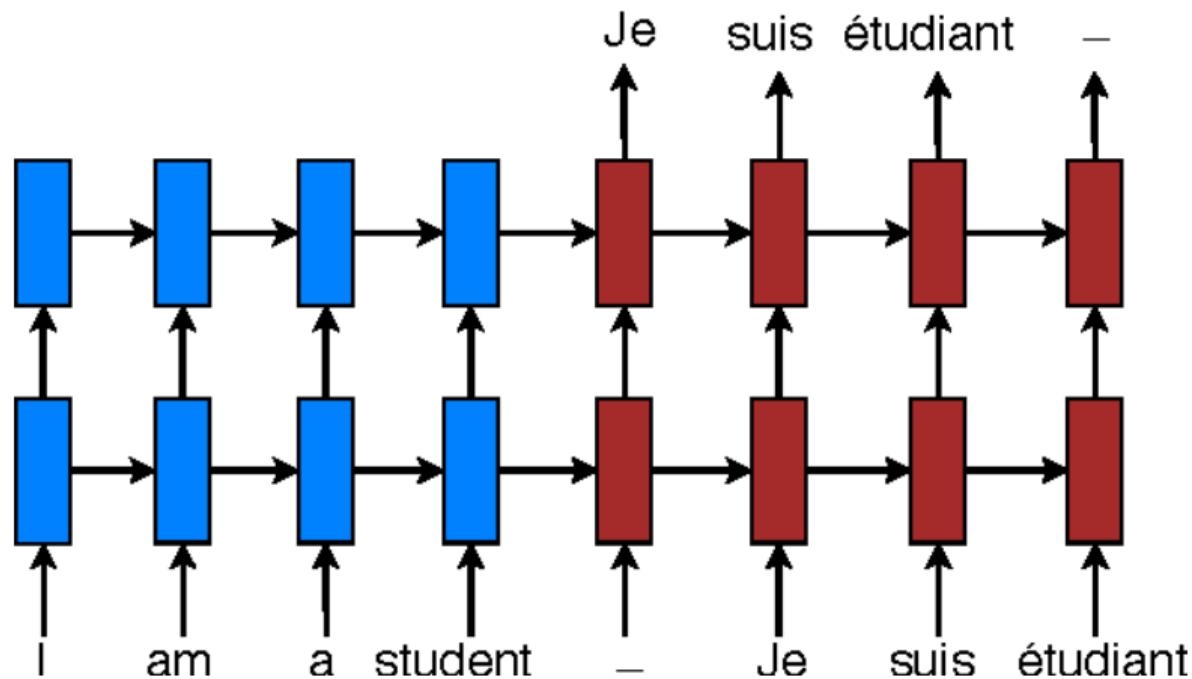
# Spam Detection

	Subject	Correspondents	Date
✉	URGENT RFQ	↳ ← AL WALEED EQUIPMENTS	03/13/2017 06:55
✉	New Order Attached **KINDLY SEND INVOICE	↳ ← starsescorts@gmail.com	03/15/2017 01:27
✉	We're sad to let you know that our delivery was unsuccessful....	↳ ← Amr Hassan	03/15/2017 19:30
✉	47929 username2	↳ ← FedEx Expedited Express	03/16/2017 02:53
✉	Delivery Status Notification	↳ ← pkeith@gejlaw.com	03/16/2017 05:29
✉	Formal Inquiry	↳ ← webmaster@stroy-exp...	03/16/2017 05:47
✉	We have delivery problems with your parcel #7104543	↳ ← vowsbyjudy@shaw.ca	03/16/2017 14:38
✉	INQUIRY	↳ ← "Anaïs VANACKER" <Va...	03/16/2017 21:16
✉	54343 username	↳ ← webmaster@whfarm2....	03/17/2017 00:57
✉	Item Delivery Notification	↳ ← Saigon Offshore	03/17/2017 03:47
✉	UPS courier can not deliver parcel #004287245 to you	↳ ← dava@ac-lyon.fr	03/17/2017 14:25
✉	Parcel Delivery Notification	↳ ← juanro5554@hotmail.c...	03/17/2017 14:48
✉	Visa Card Award	↳ ← alifeof8@server.alifeofj...	00:34
✉	Problems with item delivery, n.4930349	↳ ← webmaster@stroy-exp...	06:23
✉	Package Delivery Notification	↳ ← abidjanbateau@vps286...	06:52
✉	Delivery Status Notification	↳ ← info@visa.com	07:21
		↳ ← Apache	09:54
		↳ ← Apache	10:06
		↳ ← contrav8@box980.blue...	17:05

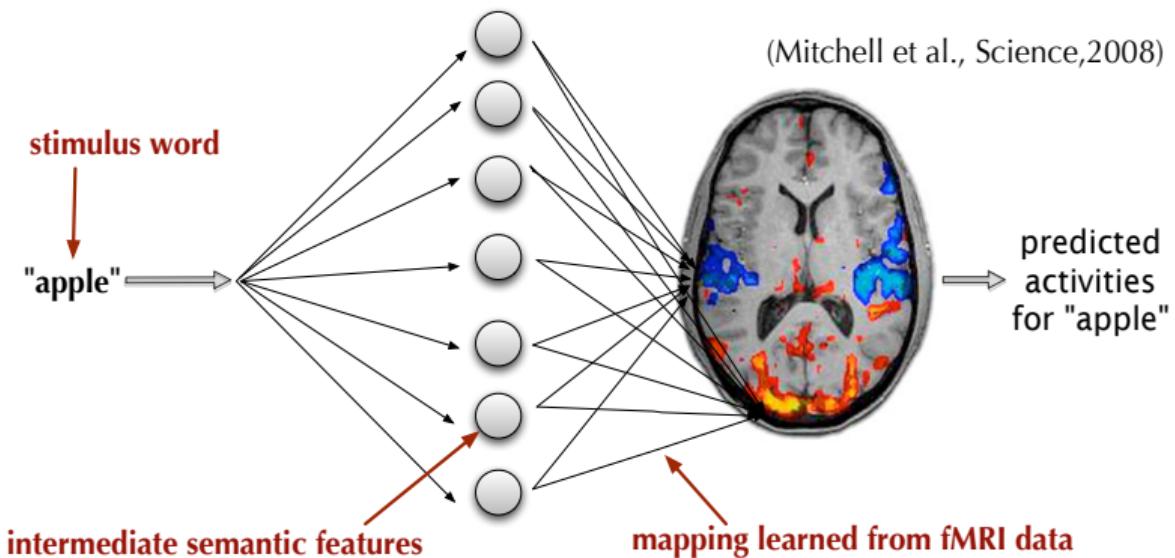
# Weather Forecasting



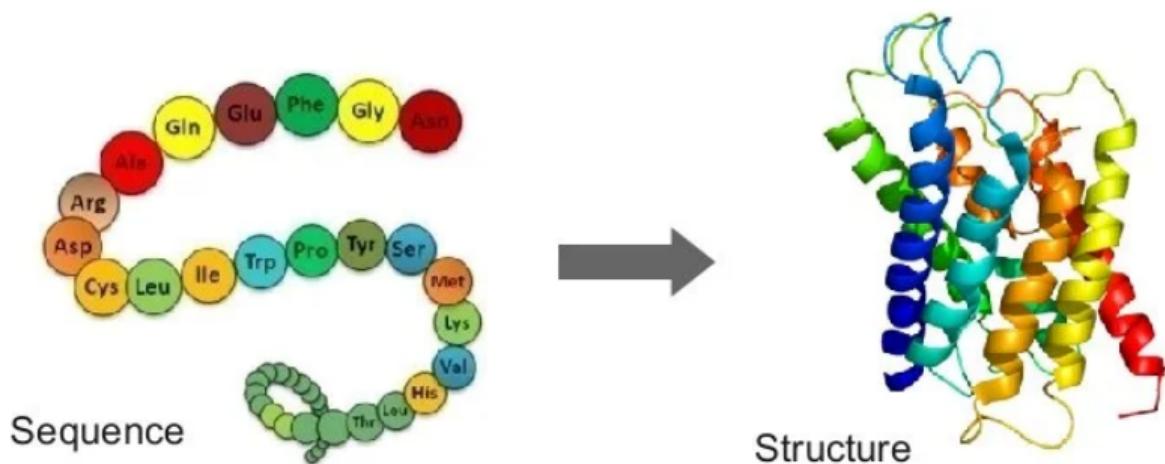
# Machine Translation



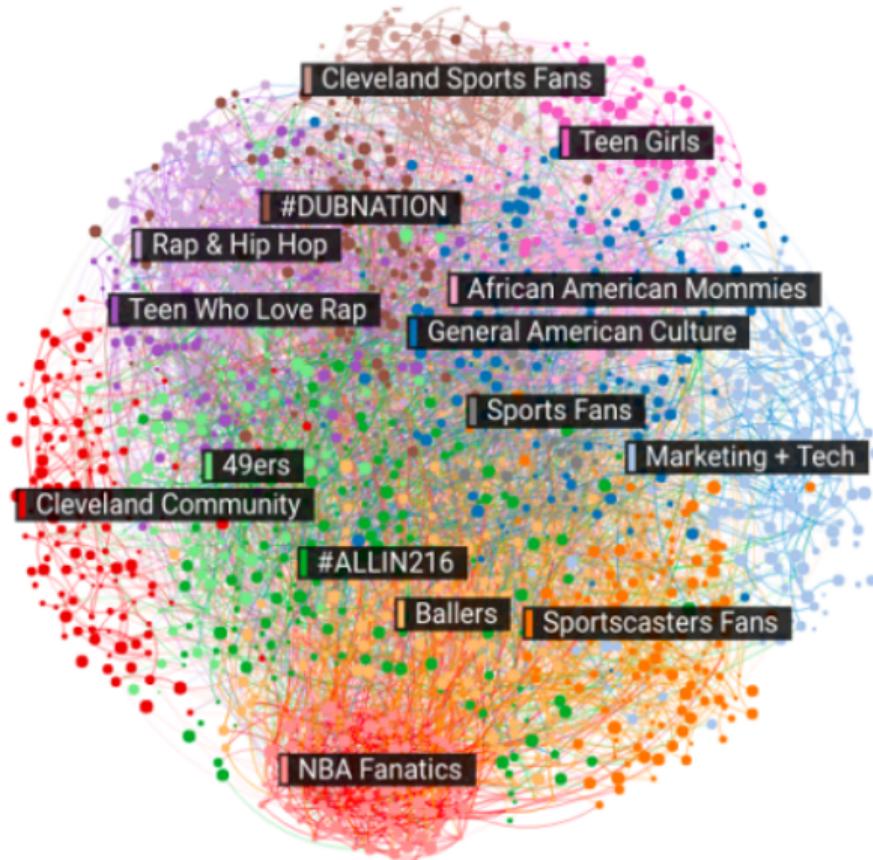
# Neural Semantic Discovery



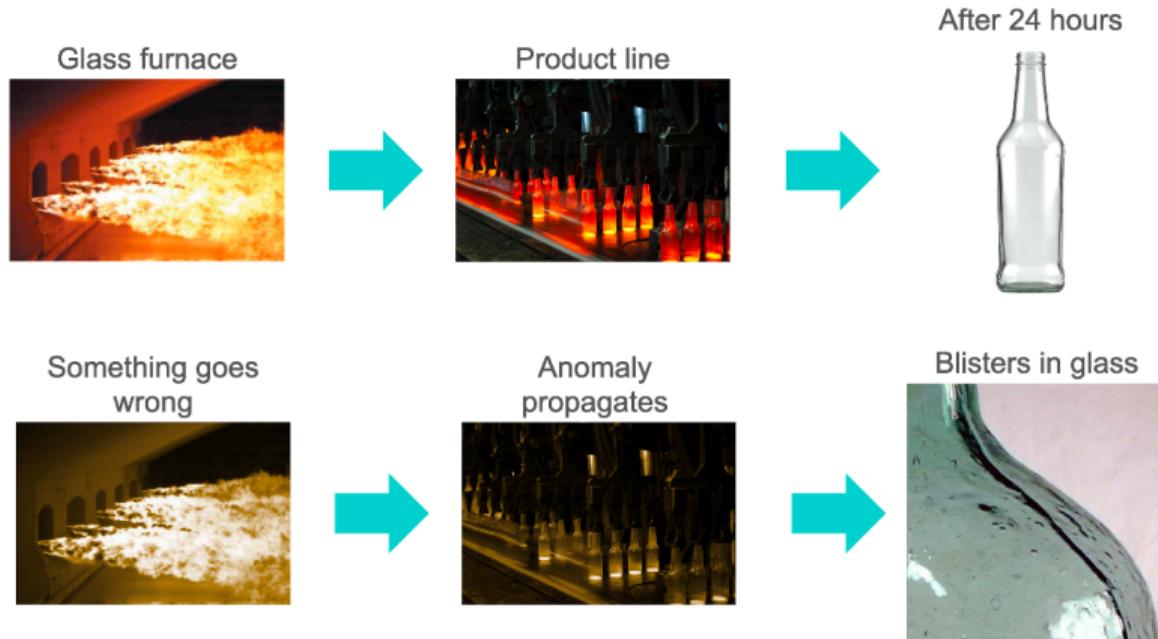
# Protein Structure Prediction



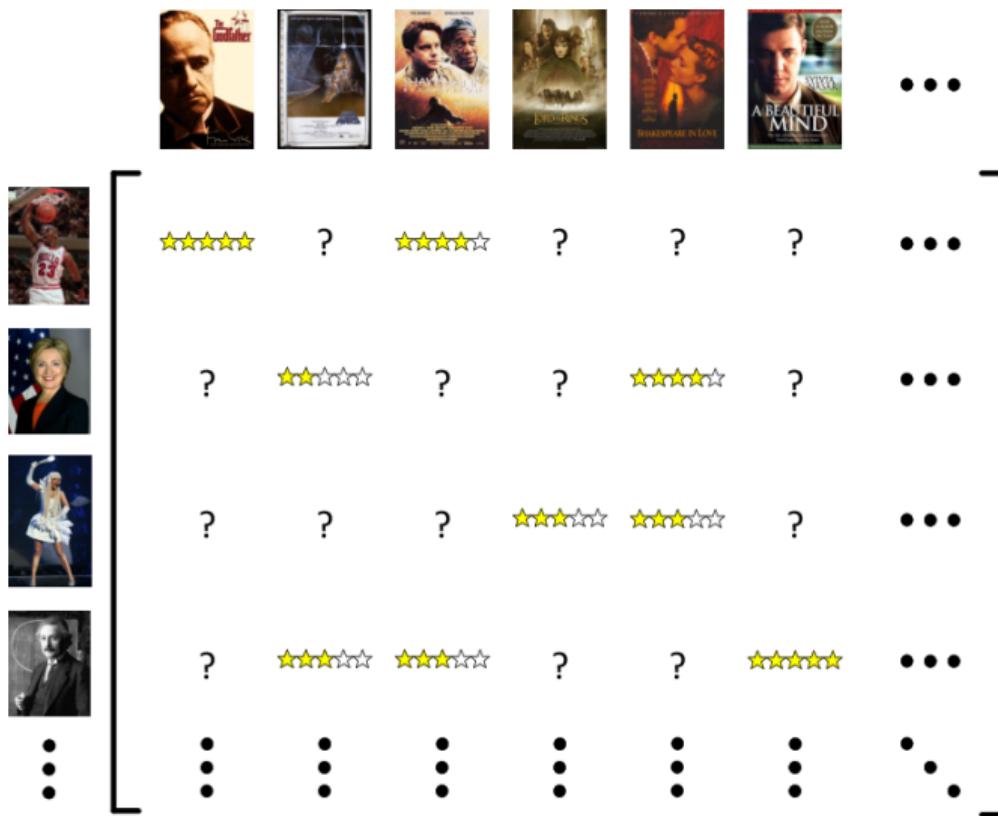
# Community Detection

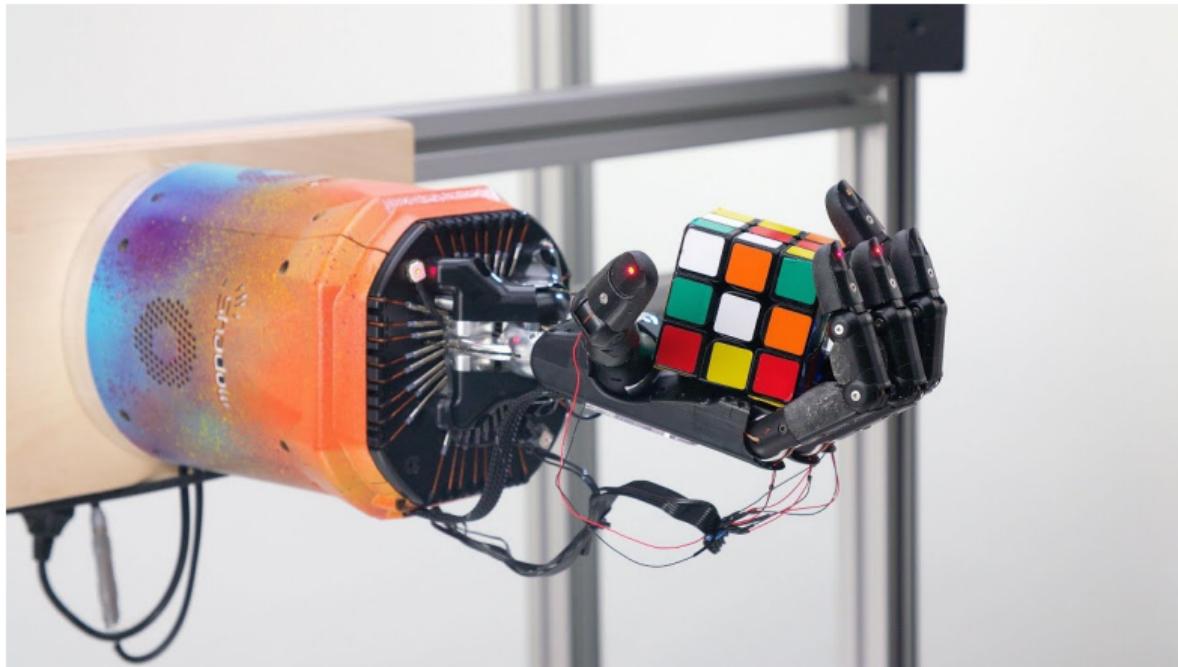


# Anomaly Detection



# Movie Recommendation

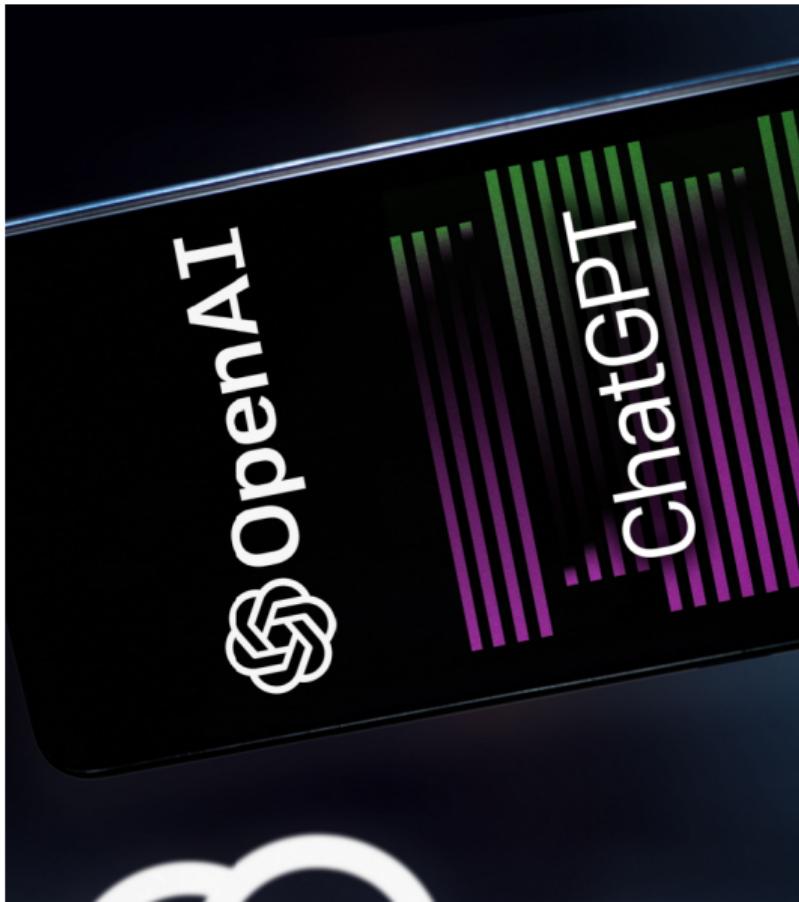




# Autonomous Driving



# Chatbot



# MIDJOURNEY

AI

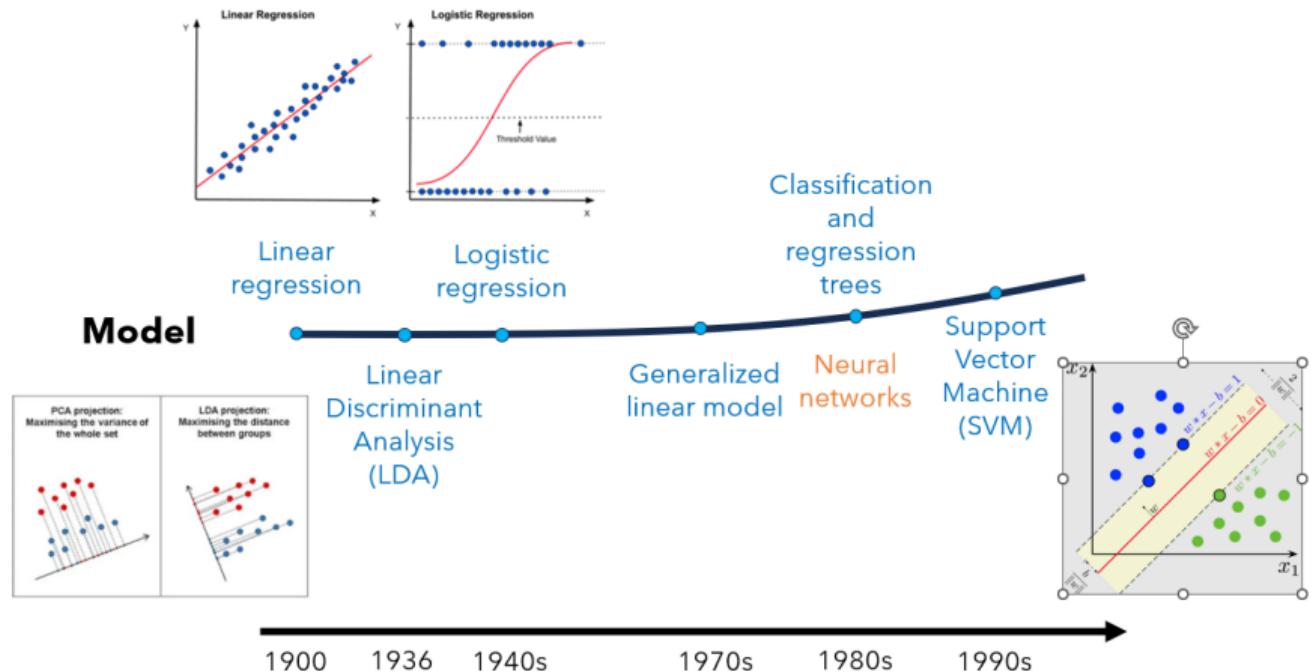


# Introduction: Machine Learning

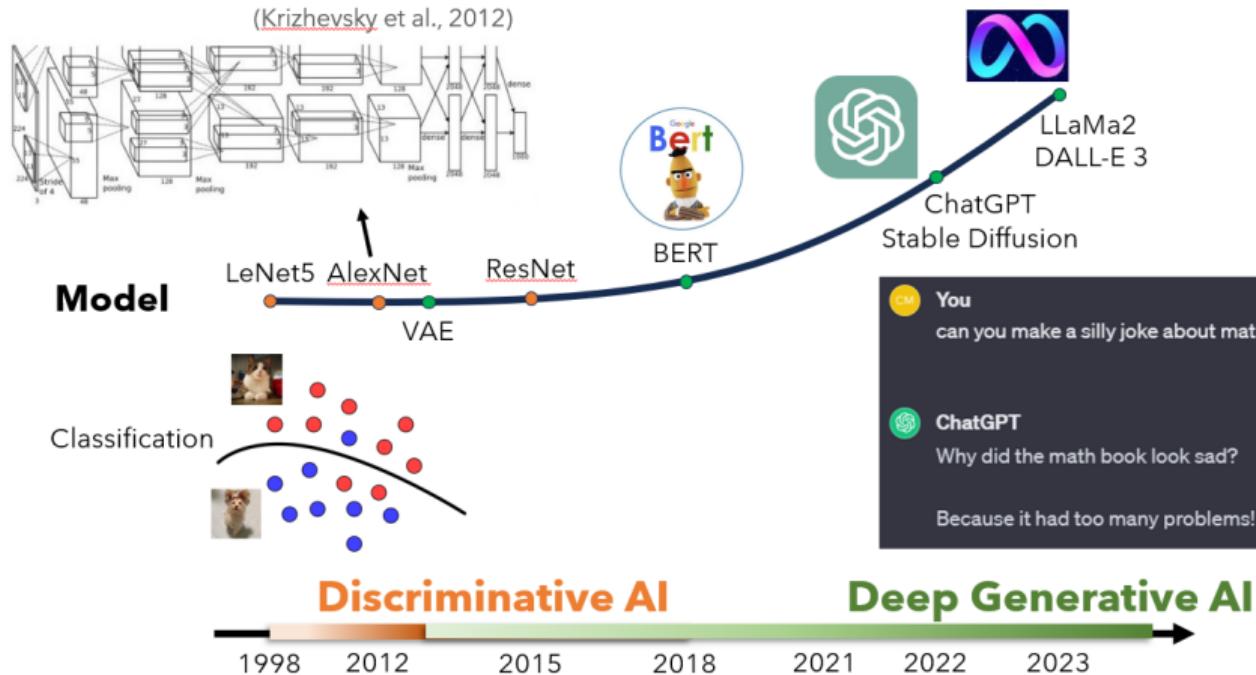
**Tom Mitchell (1998):** a computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.

- Experience (data): games played by the program (with itself)
  - Performance measure: winning rate
- 
- (⌚) We want to provide clear, interpretable models. These models allow you to understand the direct influence of each predictor on the outcome, which is essential in fields where insight into relationships (rather than just prediction) is needed.
  - (⌚) No confidence interval estimation
  - (⌚) In cases where data is scarce, simpler parametric models used in statistical learning can perform better. (**Why?**)

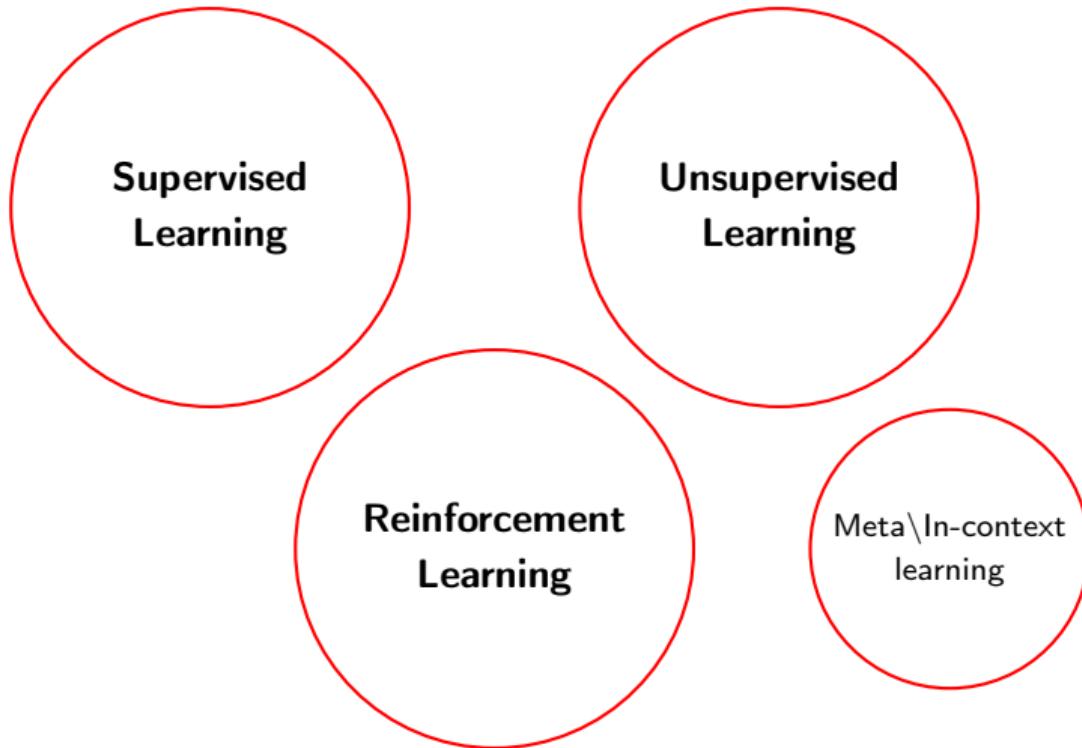
# Early History



# Contemporary Developments



# Taxonomy of Machine Learning



## Supervised Learning (Regression)

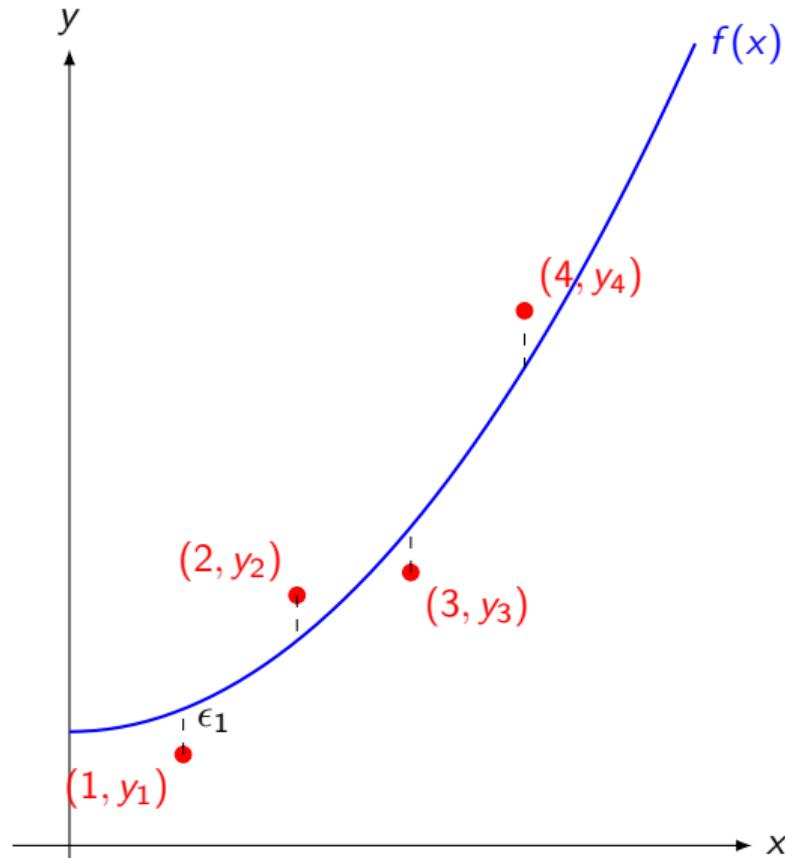
**Supervised Learning:** a set of observed data points  $\{(x_i, y_i)\}_{i=1}^n$ , where  $x_i$  represents the predictor (or vector of predictors) and  $y_i$  represents the response variable. Regression is the process of modeling the relationship between  $x$  and  $y$  by assuming:

$$y_i = f(x_i) + \epsilon_i,$$

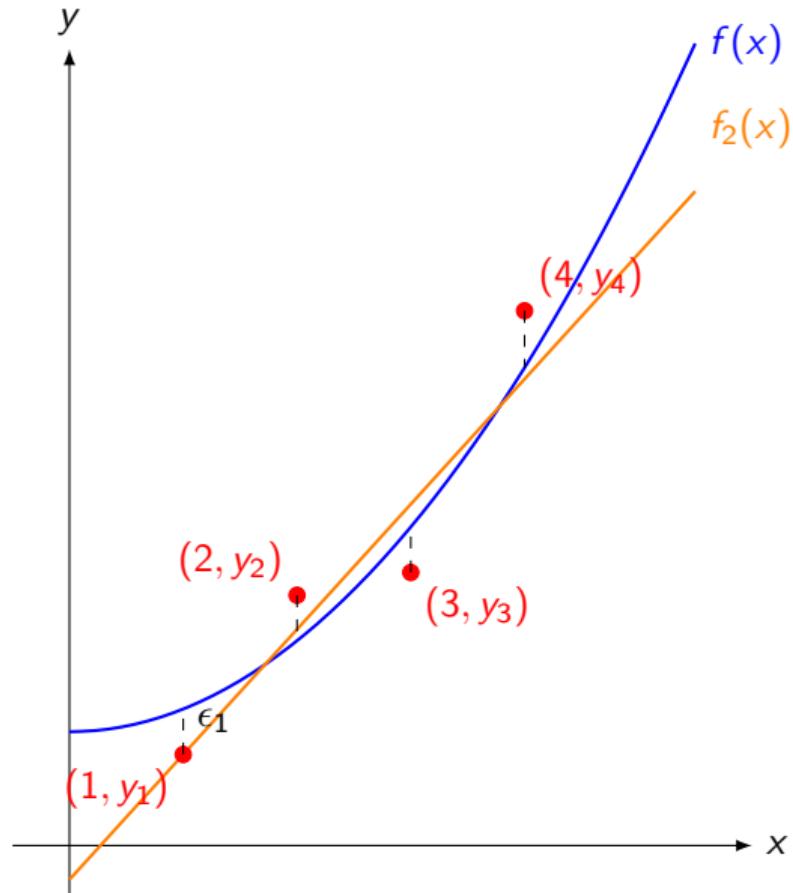
where:

- $f(x_i)$  is an unknown function that describes the systematic component of the relationship
- $\epsilon_i$  is a random error term.

# Regression



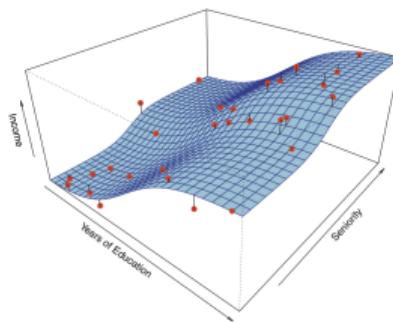
# Regression



# High Dimensional Features

□  $x \in \mathbb{R}^d$

$$x = \begin{bmatrix} x_1 & \text{--- living size} \\ x_2 & \text{--- lot size} \\ x_3 & \text{--- \# floors} \\ \vdots & \text{--- condition} \\ x_d & \text{--- zip code} \end{bmatrix} \longrightarrow y \text{ --- price}$$



## Different Prediction

---

- Point Prediction : retrun  $\hat{f}(x)$  since it returns a number.
- Interval Prediction , e.g.,  $Y$  will be within an interval  $[l, u]$  with probability  $1 - \alpha$
- distributional prediction , e.g.  $Y$  will follow an  $N(m, v)$  distribution.

## Bias and Variance Trade-off

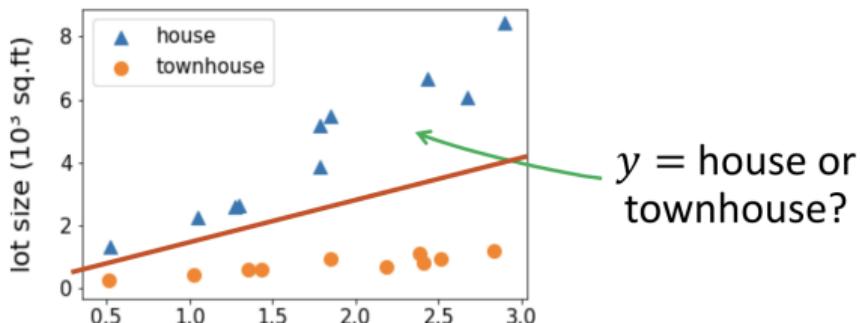
$$\mathbb{E}[(y - \hat{f}(x))^2] = \underbrace{\left(f(x) - \mathbb{E}[\hat{f}(x)]\right)^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}\left[\left(\hat{f}(x) - \mathbb{E}[\hat{f}(x)]\right)^2\right]}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{Irreducible}}$$

- 
- ⌚ An unbiased estimator could still make systematic mistakes – for example, if it overestimates 99% of the time, and underestimates 1% of the time \*by a lot\*, in expectation it could be unbiased.
  - ⌚ An unbiased estimator is \*\*not\*\* necessarily better than a biased estimator, because the total error depends on both the bias and variance of the estimator.

# Classification

- Regression : if  $y \in \mathbb{R}$  is a continuous variable
- classification : the label is a discrete variable

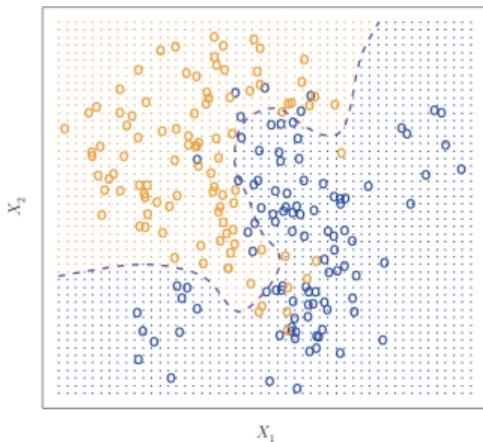
(size, lot size) → house or townhouse ?



## Classification as Regression: Bayes Classifier

training error rate:  $\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$

Here the function  $I(y_i \neq \hat{y}_i)$  is an indicator variable that equals 1, if  $y_i \neq \hat{y}_i$  and 0 otherwise. If  $y_i \neq \hat{y}_i$ , then the  $i$ -th observation was classified incorrectly; otherwise it was not misclassified.



Consider random label:  $\mathbb{P}(Y = j | X = x_0)$ .  
The Bayes classifier returns

$$1 - \max_j \mathbb{P}(Y = j | X = x_0)$$

produces the lowest possible test error rate,  
called the *Bayes error rate* is given by

$$\underbrace{1 - \mathbb{E}\left[\max_j \mathbb{P}(Y = j | X)\right]}_{\text{Irreducible}}$$

## Prediction Accuracy and Model Interpretability

Why would we ever choose to use a more restrictive method instead of a very flexible approach?

# $x$ and $y$ in Computer Vision

## Task. Image Classification

$x = ?, y = ?$

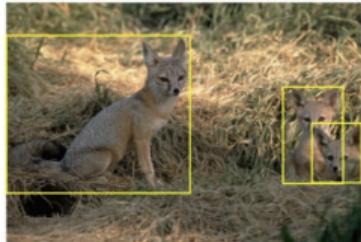
ILSVRC



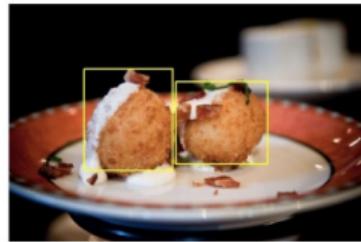
# $x$ and $y$ in Computer Vision

## Task. Object localization and detection

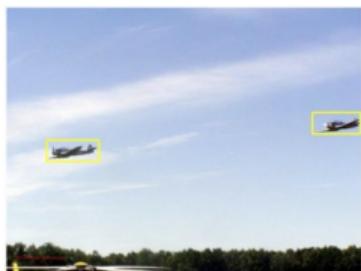
$x = ?, y = ?$



kit fox



croquette



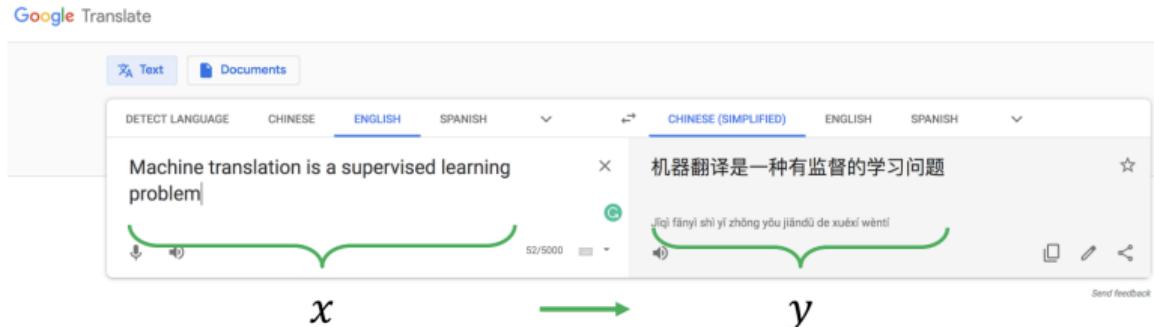
airplane



frog

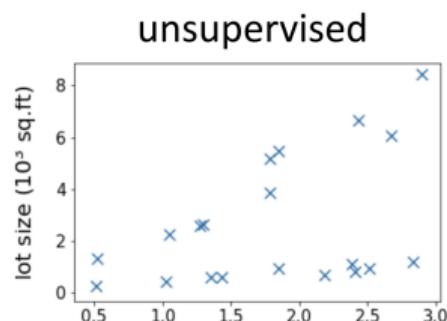
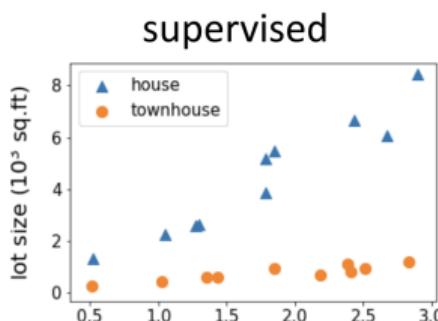
# $x$ and $y$ in Natural Language

Task. Machine Translation d  $x = ?$ ,  $y = ?$



# Unsupervised Learning (Clustering)

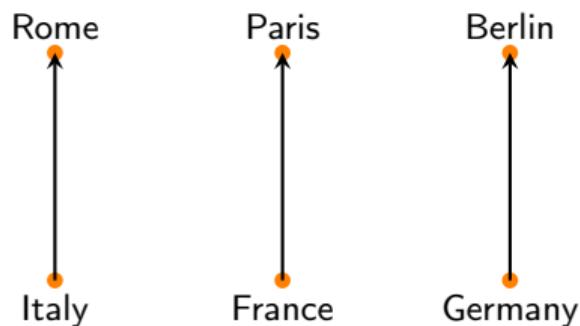
- ❑ Dataset contains **no** labels:  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$
- ❑ Goal (**vaguely-posed**): to find interesting structures in the data



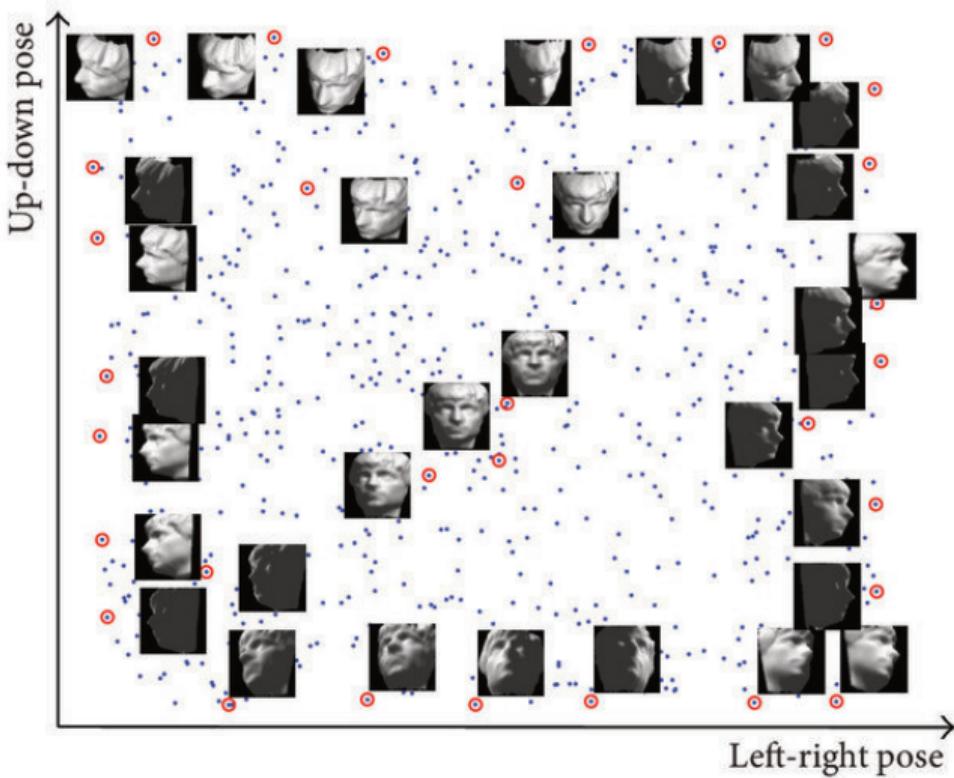
# Unsupervised Learning (Feature Extraction)

- Word : Encode as vectors
- Relationship : represent as direction

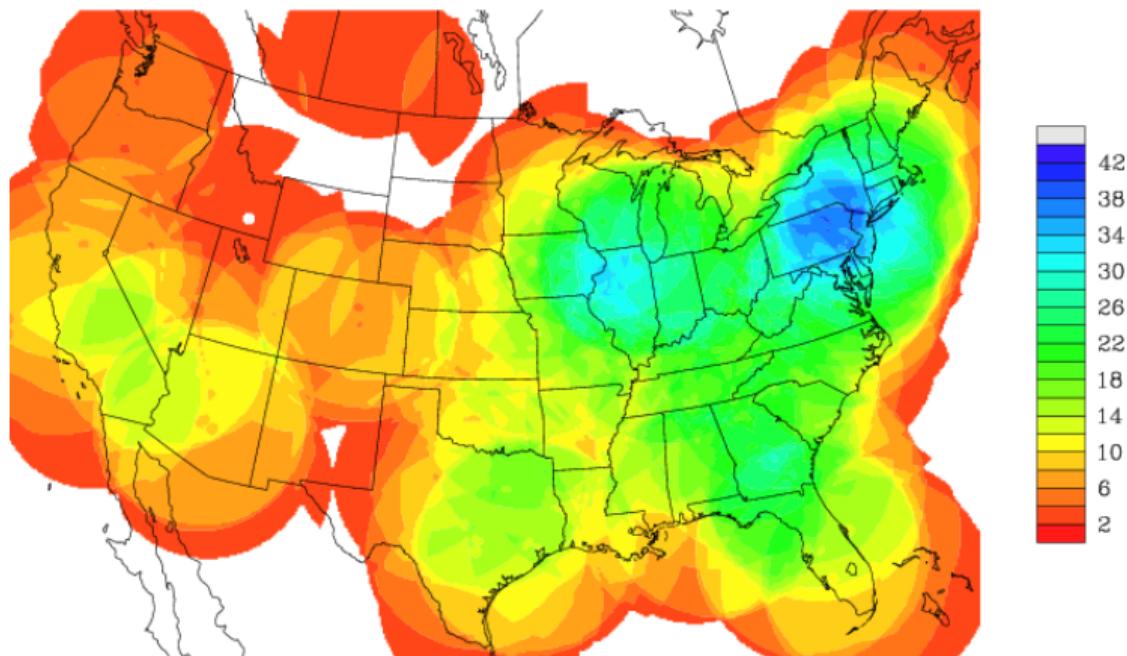
▷ word → encode → vector  
▷ relation → encode → direction



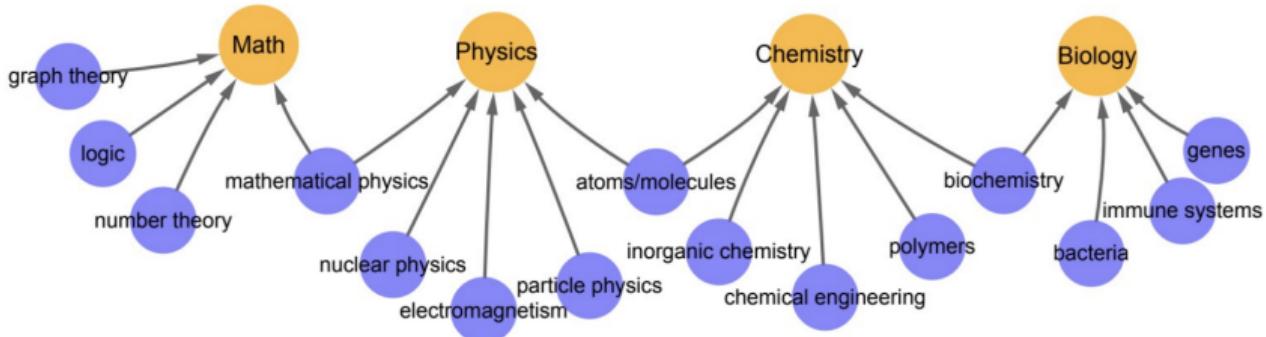
# Unsupervised Learning (Feature Extraction)



# Unsupervised Learning (Density Estimation)



# Unsupervised Learning



	logic deductive propositional semantics	graph subgraph bipartite vertex	boson massless particle higgs	polyester polypropylene resins epoxy	acids amino biosynthesis peptide
tag	<i>logic</i>	<i>graph theory</i>	<i>particle physics</i>	<i>polymer</i>	<i>biochemistry</i>

# Reinforcement Learning

Learning to make sequential decisions

