



# CS60050: Machine Learning

## Autumn 2023

Sudeshna Sarkar

**Introduction to ML**

3-4 Aug 2023

# Course Website

<https://sites.google.com/view/cs60050-fall-2023/home>

## Course Timings

- Wed 11-12 PM, Thu 12-1 PM, Fri 8-9 AM
- **Section 1 (Dr. Somak Aditya, NC442)**
  - Even Roll Numbers
- **Section 2 (Dr. Sudeshna Sarkar, NC 243)**
  - Odd Roll Numbers

Contact

MS Teams

Email: [sudeshna@cse.iitkgp.ac.in](mailto:sudeshna@cse.iitkgp.ac.in) Subject “CS60050 ML23”

# Course Evaluation Plan

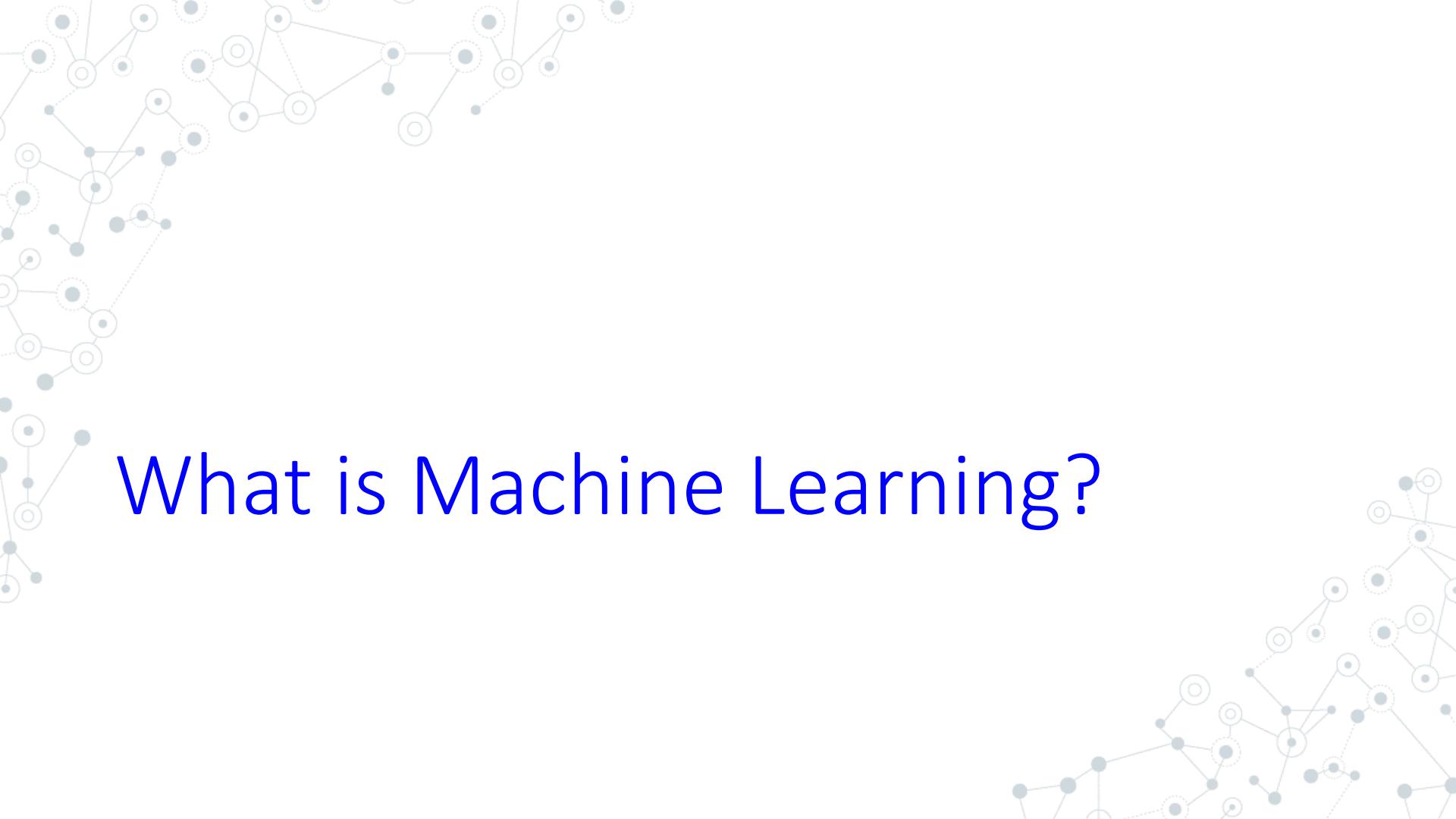
- Mid Term - 25%
- Final Exam - 40%
- Assignments (3-4) - 20%
- Class-Tests (Two) - 15%

# Attendance Policy

Students must attend all classes

Attendance Policy: Attendance record will be maintained and uploaded on the website.

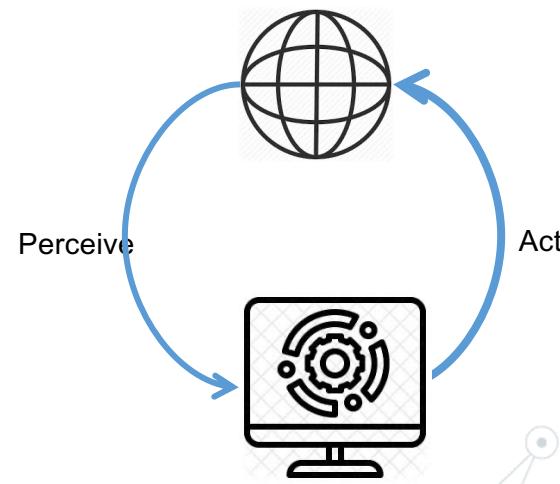
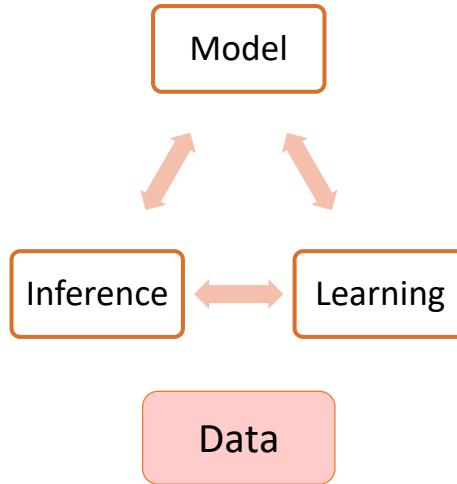
All students must maintain 75% attendance to continue the course.



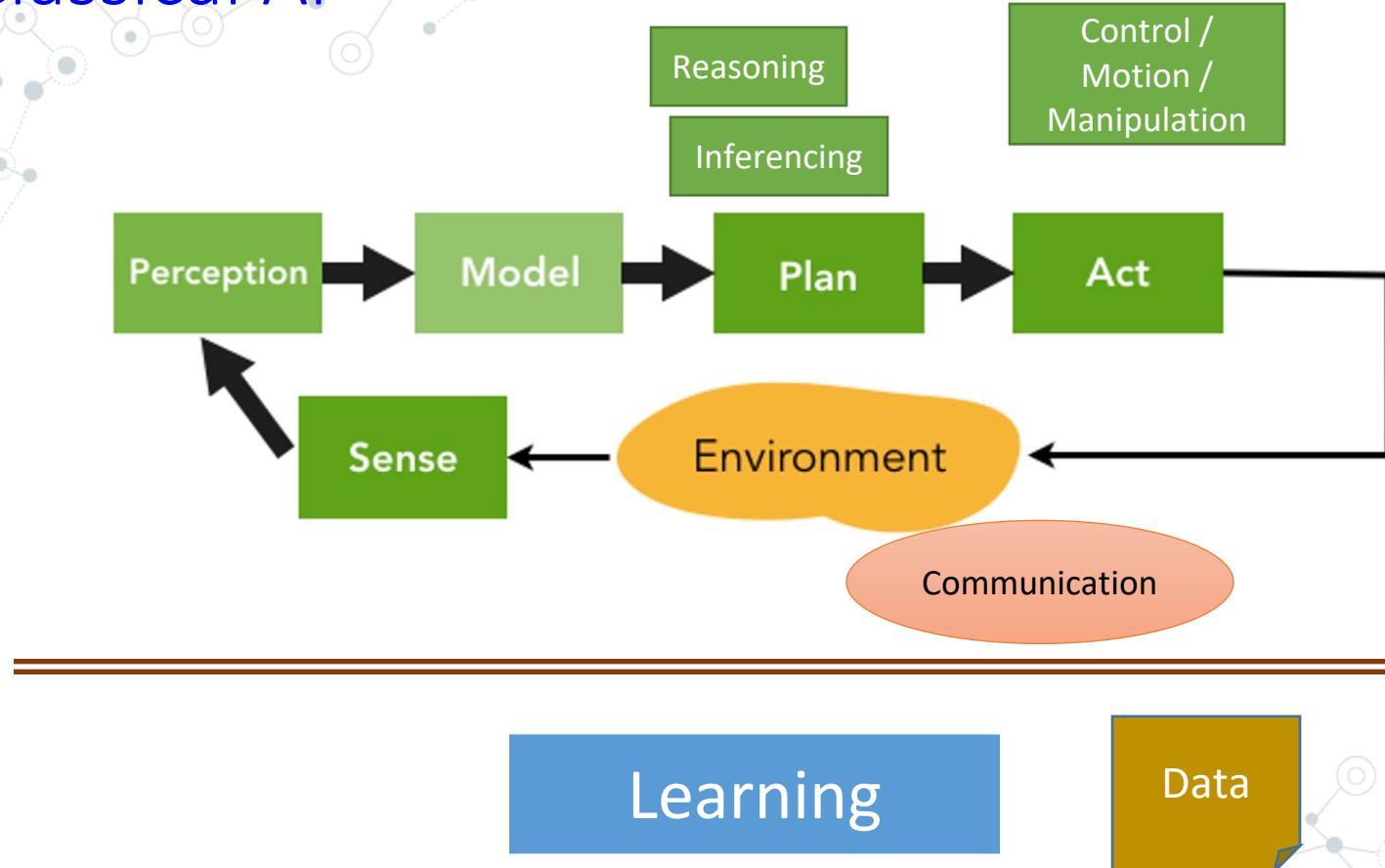
# What is Machine Learning?

# Artificial Intelligence

- The quest for building Intelligent Agents that can think and act rationally (or humanlike)



# Classical AI



# Machine Learning Toolbox

- Machine Learning
- Statistics
- Probability
- Computer Science
- Optimization

# What is ML?

## Speech Recognition

### 1. Learning to recognize spoken words

“...the SPHINX system (e.g., Lee 1989) learns speaker-specific strategies for recognizing the primitive sounds (phonemes) and words from the observed speech signal...neural network methods...hidden Markov models...”

(Mitchell, 1997)



Figure from <https://btpenguin.com/alexa-vs-siri-vs-google-assistant/>

## Computer Vision

### 4. Learning to recognize images

#### THEN

“...The recognizer is a convolution network that can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors....”



(LeCun et al., 1995)

#### NOW



Figure from <https://proceedings.neurips.cc/paper/2012/file/c399861d3b9d6b76c8436e924a6c45b-Paper.pdf>

## Robotics

### 2. Learning to drive an autonomous vehicle

“...the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars...”

(Mitchell, 1997)



9

## Games / Reasoning

### 3. Learning to beat the masters at board games

“...the world's top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself...”

(Mitchell, 1997)



15

## Learning Theory

### • 5. In what cases and how well can we learn?

#### Sample Complexity Results

Definitions: The sample complexity of a learning algorithm is the number of examples required to learn a hypothesis with respect to the hypothesis space, with high probability (i.e. close to 1).

Your Case: we care about...

Feasible ( $H$ )

Infeasible ( $H'$ )

For  $\epsilon = 5\%$

Q: That few (aka required #d) (aka Generalization Error)  
 $P(\hat{E} \leq \epsilon) \geq 1 - \delta$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $P(\hat{E} \leq \epsilon) = P(\hat{E} \leq E + \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )

① Then  $E = \frac{1}{n} \sum_{i=1}^n E_i$  (aka empirical error)  
 $R(n) = P_{\hat{E}}(E \leq \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = P_{\hat{E}}(E \leq E + \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = P_{\hat{E}}(E \leq E + \epsilon) = P_{\hat{E}}(E \leq E + \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = P_{\hat{E}}(E \leq E + \epsilon) = \frac{1}{n} \sum_{i=1}^n P_{\hat{E}_i}(E_i \leq E + \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = \frac{1}{n} \sum_{i=1}^n P_{\hat{E}_i}(E_i \leq E + \epsilon) = \frac{1}{n} \sum_{i=1}^n P(E_i \leq E + \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = \frac{1}{n} \sum_{i=1}^n P(E_i \leq E + \epsilon) = \frac{1}{n} \sum_{i=1}^n (1 - P(E_i > E + \epsilon))$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = \frac{1}{n} \sum_{i=1}^n (1 - P(E_i > E + \epsilon)) = 1 - \frac{1}{n} \sum_{i=1}^n P(E_i > E + \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = 1 - \frac{1}{n} \sum_{i=1}^n P(E_i > E + \epsilon) = 1 - \frac{1}{n} \sum_{i=1}^n (1 - P(E_i \leq E + \epsilon))$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = 1 - \frac{1}{n} \sum_{i=1}^n (1 - P(E_i \leq E + \epsilon)) = \frac{1}{n} \sum_{i=1}^n P(E_i \leq E + \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = \frac{1}{n} \sum_{i=1}^n P(E_i \leq E + \epsilon) = \frac{1}{n} \sum_{i=1}^n (1 - \epsilon)$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = \frac{1}{n} \sum_{i=1}^n (1 - \epsilon) = 1 - \frac{1}{n} \sum_{i=1}^n \epsilon$  (i.e.  $\hat{E}$  is close to  $E$ )  
 $R(n) = 1 - \frac{1}{n} \sum_{i=1}^n \epsilon = 1 - \epsilon$  (i.e.  $\hat{E}$  is close to  $E$ )

1. How many examples do we need to learn?
2. How do we quantify our ability to generalize to unseen data?
3. Which algorithms are better suited to specific learning settings?

Learning

Data

Model

- Requires a leap of faith: generalization

# ML Big Picture

## Learning Paradigms:

*What data is available and when? What form of prediction?*

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

## Theoretical Foundations:

*What principles guide learning?*

- probabilistic
- information theoretic
- evolutionary search
- ML as optimization

## Problem Formulation:

*What is the structure of our output prediction?*

boolean	Binary Classification
categorical	Multiclass Classification
ordinal	Ordinal Classification
real	Regression
ordering	Ranking
multiple discrete	Structured Prediction
multiple continuous	(e.g. dynamical systems)
both discrete & cont.	(e.g. mixed graphical models)

## Application Areas

*Key challenges?*

NLP, Speech, Computer Vision, Robotics, Medicine, Search

## Facets of Building ML Systems:

*How to build systems that are robust, efficient, adaptive, effective?*

1. Data prep
2. Model selection
3. Training (optimization / search)
4. Hyperparameter tuning on validation data
5. (Blind) Assessment on test data

## Big Ideas in ML:

*Which are the ideas driving development of the field?*

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

# ML and Ethics

*What ethical responsibilities do we have as machine learning experts?*

**Question:** What are the possible societal impacts of machine learning for each case below?

**Answer:**

- 1) Search results for news are optimized for ad revenue.



- 2) An autonomous vehicle is permitted to drive unassisted on the road.



- 3) A doctor is prompted by an intelligent system with a plausible diagnosis for her patient.



# Define Machine Learning

# Well-Posed Learning Problem

***Three components***  $\langle T, P, E \rangle$ :

1. Task,  $T$
2. Experience,  $E$
3. Performance measure,  $P$

***Definition of learning:***

A computer program **learns** if its performance at task  $T$ , as measured by  $P$ , improves with experience  $E$ .

# Example Learning Problem(s)

*Learning to respond to voice commands (Siri)*

- 1. Task,  $T$ :
- 2. Performance measure,  $P$ :
- 3. Experience,  $E$ :

# Example Learning Problems

Learning to respond to voice commands (Siri)

1. Task, T: predicting action from speech
2. Performance measure, P: percent of correct actions taken in user pilot study
3. Experience, E: examples of (speech, action) pairs

# Capturing the Knowledge of Experts



## Solution #1: Expert Systems

- Over 20 years ago, we had rule-based systems:
  1. Put a bunch of linguists in a room
  2. Have them think about the structure of their native language and write down the rules they devise

Give me directions to Starbucks

If: "give me directions to X"  
Then: directions(here, nearest(X))

How do I get to Starbucks?

If: "how do i get to X"  
Then: directions(here, nearest(X))

Where is the nearest Starbucks?

If: "where is the nearest X"  
Then: directions(here, nearest(X))

# Capturing the Knowledge of Experts



## Solution #2: Annotate Data and Learn

- Experts:
  - **Very good at** answering questions about specific cases
  - **Not very good at** telling **HOW** they do it
- 1990s: So why not just have them tell you what they do on **SPECIFIC CASES** and then let **MACHINE LEARNING** tell you how to come to the same decisions that they did

# Capturing the Knowledge of Experts



## Solution #2: Annotate Data and Learn

1. Collect raw sentences  $\{x^{(1)}, \dots, x^{(n)}\}$
2. Experts annotate their meaning  $\{y^{(1)}, \dots, y^{(n)}\}$

$x^{(1)}$ : How do I get to Starbucks?

$y^{(1)}$ : directions (here,  
nearest (Starbucks))

$x^{(2)}$ : Show me the closest Starbucks

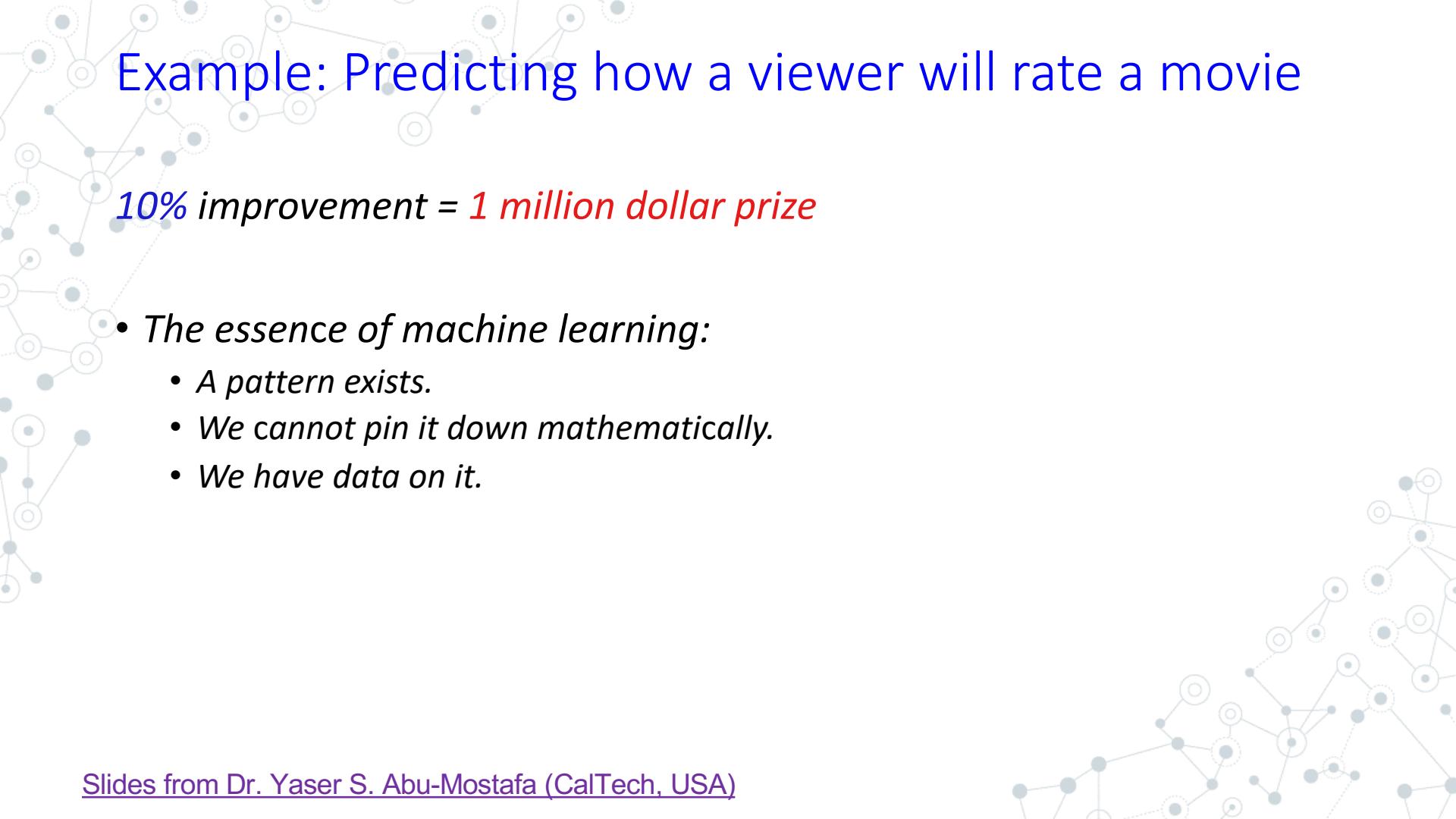
$y^{(2)}$ : map (nearest (Starbucks))

$x^{(3)}$ : Send a text to John that I'll be late

$y^{(3)}$ : txtmsg (John, I'll be late)

$x^{(4)}$ : Set an alarm for seven in the morning

$y^{(4)}$ : setalarm (7:00AM)

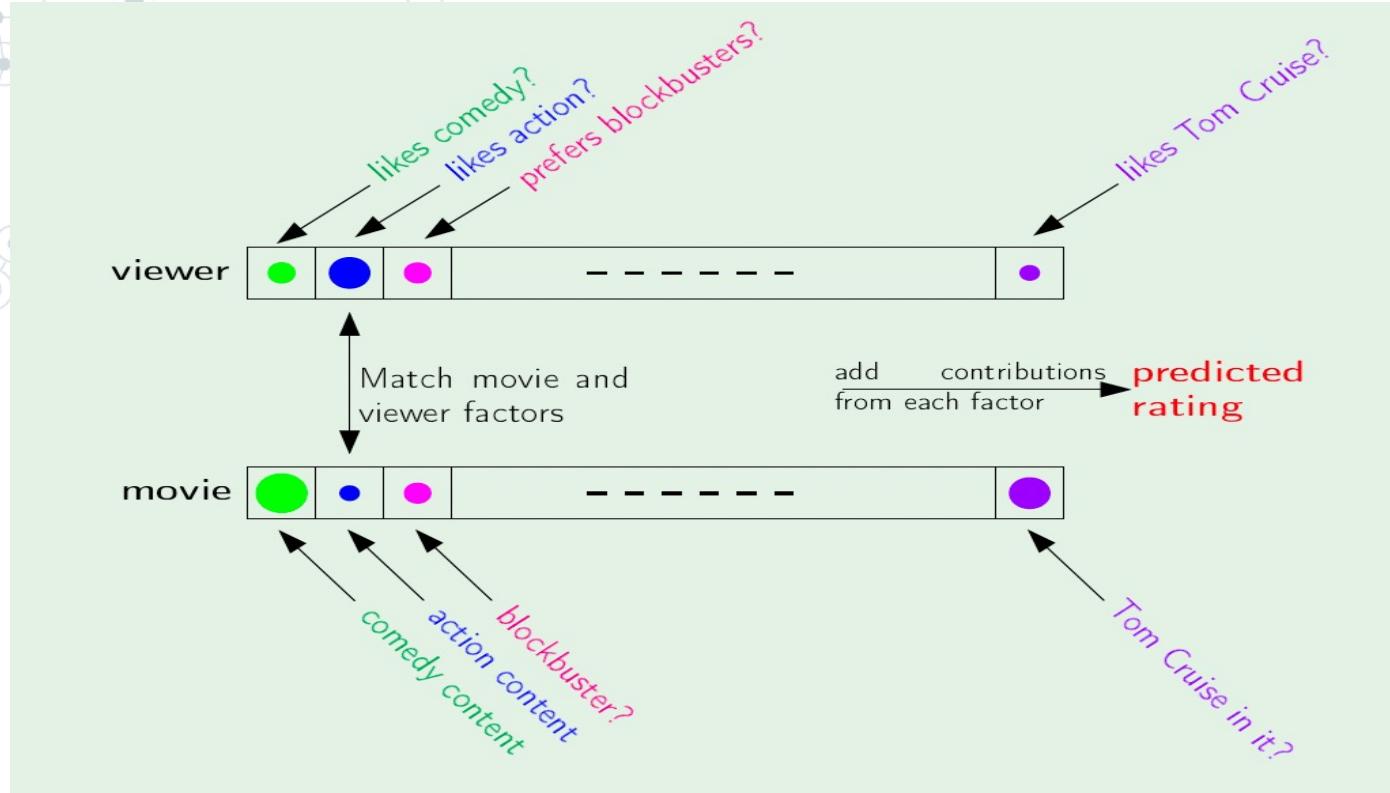


# Example: Predicting how a viewer will rate a movie

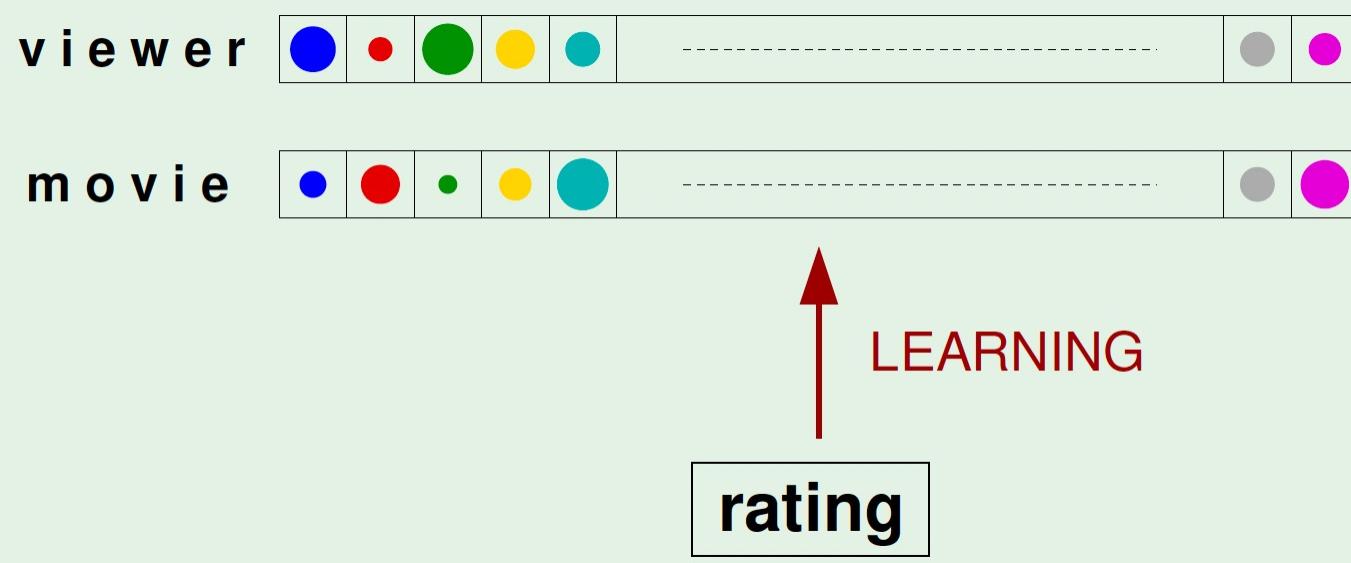
*10% improvement = 1 million dollar prize*

- *The essence of machine learning:*
  - A pattern exists.
  - We cannot pin it down mathematically.
  - We have data on it.

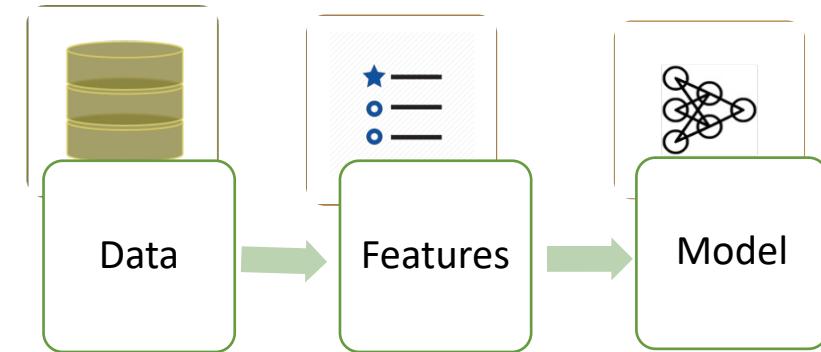
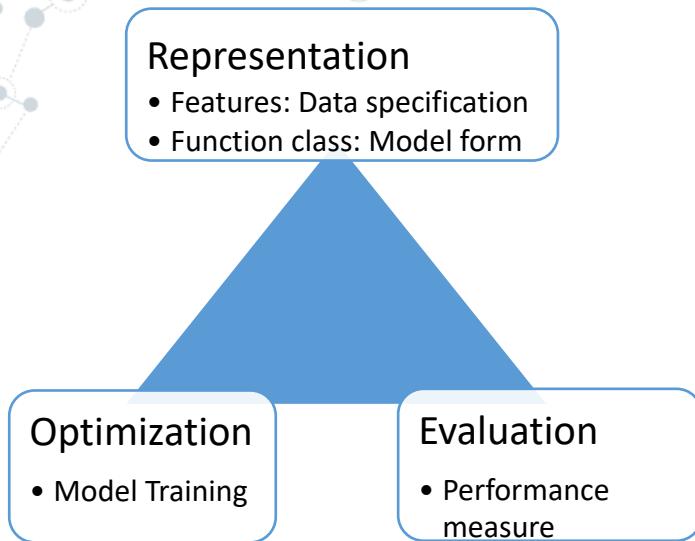
# Movie Rating – A Solution



# The Learning approach



# Components of a ML application



# A. Representation of Data

How is the data specified?

## A. Features

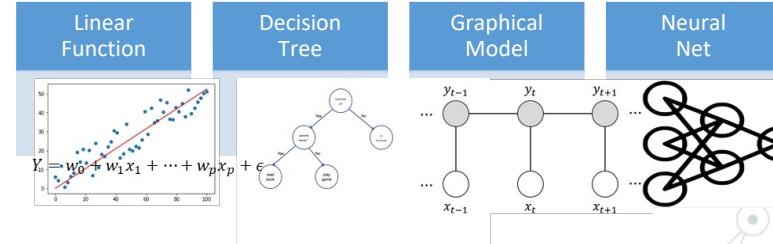
- Feature vector of  $n$  features

$$\bar{x} = (x_1, x_2, \dots, x_n)$$

The richer the representation, the more useful it is for subsequent problem solving.

The richer the representation, the more difficult it is to learn.

Hypothesis Set  
 $\mathcal{H} = \{h\}$



## 2. Evaluation

1. Accuracy =  $\frac{\# \text{ correctly classified}}{\# \text{ all test examples}}$

2. Logarithmic Loss:

$$L_i = -\log(P(Y = y_i | X = x_i))$$

$$L = \sum_{c=1}^M y_{oc} \log(p_{oc})$$

3. Mean Squared error

$$MSE = \frac{1}{m} \sum (y_{pred} - y_{true})^2$$

### 3. Optimization

- Define loss function
- Optimize loss function

- Stochastic Gradient Descent  
(Convex functions)
- Combinatorial optimization
  - E.g.: Greedy search
- Constrained optimization
  - E.g.: Linear programming

- Task: Credit approval
- Applicant information:

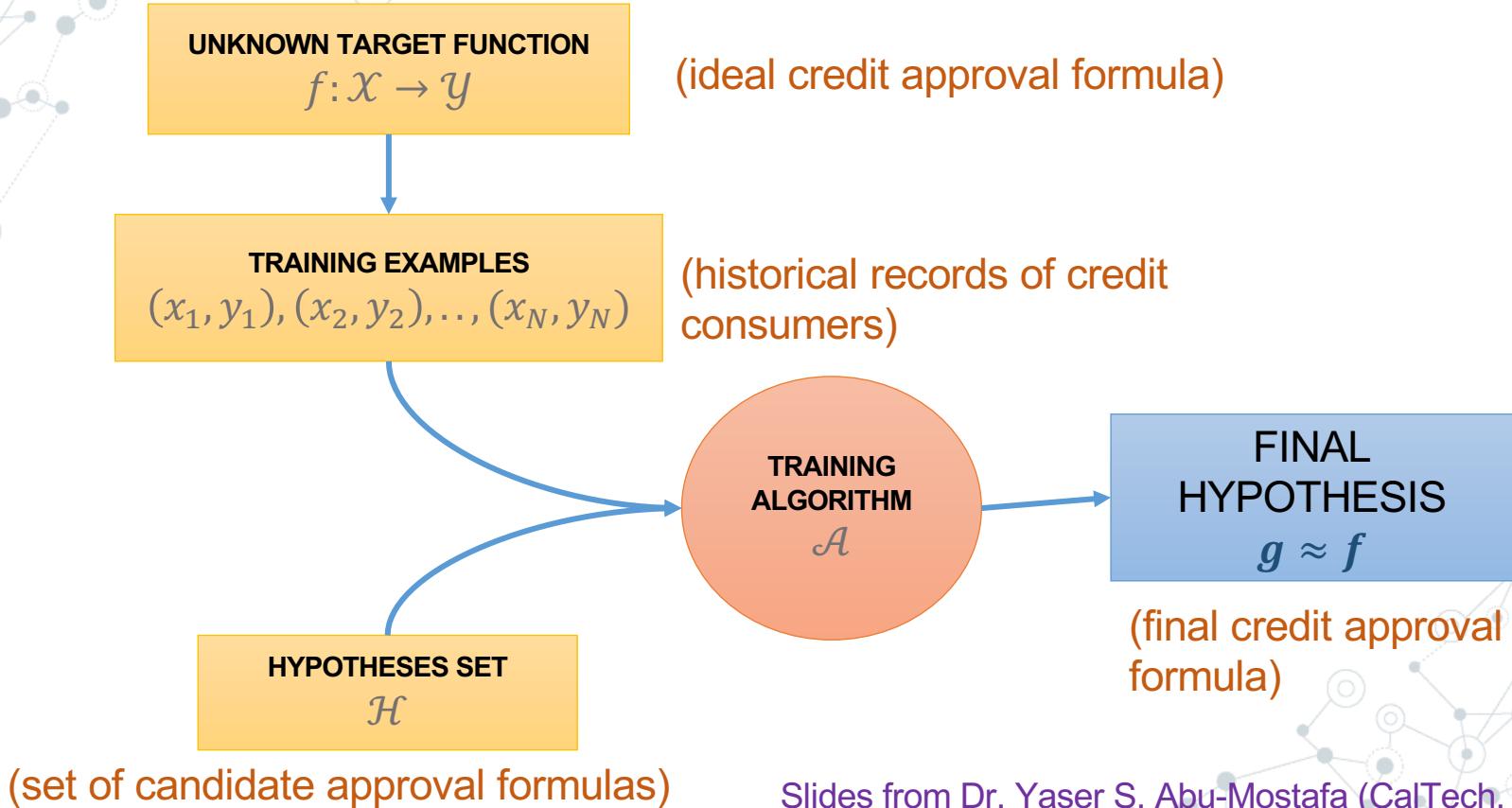
Age	23 years
Gender	male
Annual salary	\$30000
Years in residence	1 year
Years in job	1 year
Current debt	\$15000
...	...

- Approve Credit?

## Formalization

- Input  $x$  (customer application)
- Output  $y$  (good/bad customer?)
- Target function  $f: \mathcal{X} \rightarrow \mathcal{Y}$  (ideal credit approval formula)
- Data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$  (historical records)
- Hypothesis  $g: \mathcal{X} \rightarrow \mathcal{Y}$  (formula to be used)

# Components of Learning

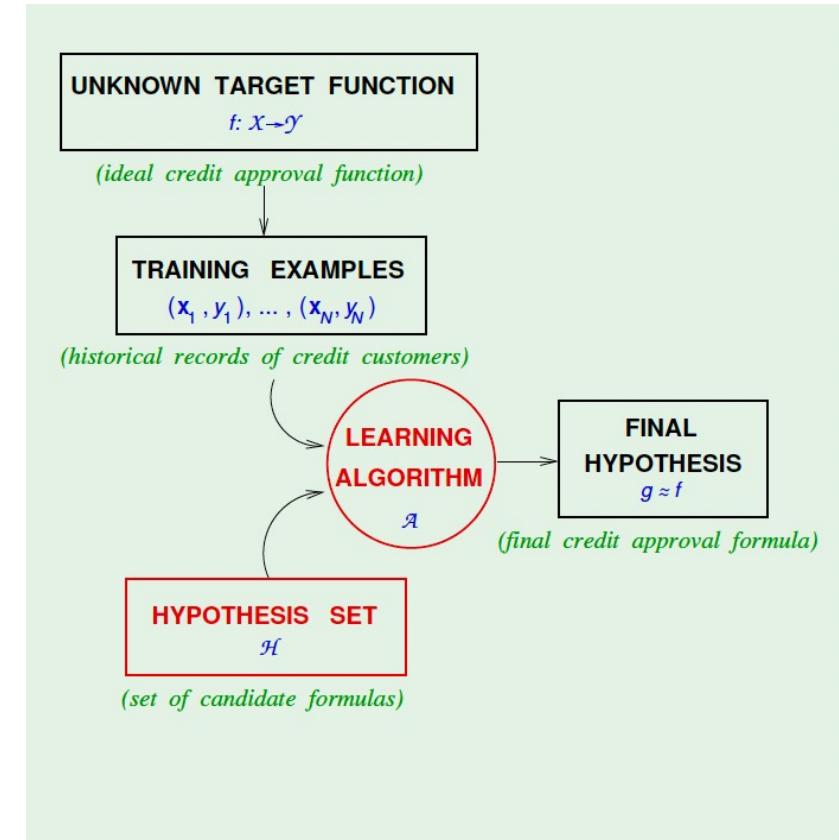


# Solution Components

2 solution components of the learning problem

- The Hypothesis Set  
 $\mathcal{H} = \{h\}, g \in \mathcal{H}$
- The Learning Algorithm

Together, they are referred to as the *learning model*.



# A simple hypothesis set – the ‘perceptron’

For input  $x = (x_1, \dots, x_d)$  ‘attributes of a customer’

Approve credit if  $\sum_{i=1}^d w_i x_i > \text{threshold}$ ,

Deny credit if  $\sum_{i=1}^d w_i x_i < \text{threshold}$ ,

Linear formula  $h \in \mathcal{H}$  can be written as

$$h(x) = \text{sign}(\sum_{i=1}^d w_i x_i - \text{threshold})$$

# Separability

The perceptron implements

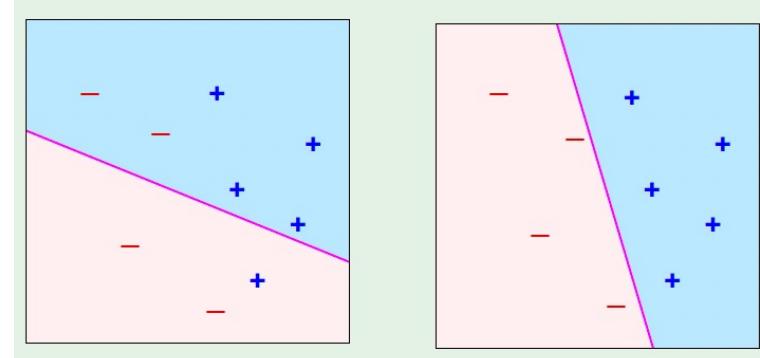
$$h(x) = \text{sign}(\sum_{i=1}^d w_i x_i + w_0)$$

Introducing an artificial coordinate  $x_0 = 1$

$$h(x) = \text{sign}(\sum_{i=0}^d w_i x_i)$$

In vector form, the perceptron implements

$$h(x) = \text{sign}(w^T x)$$



'linearly separable' data

# A simple Learning Algo - PLA

The perceptron implements

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$$

Given the training set:

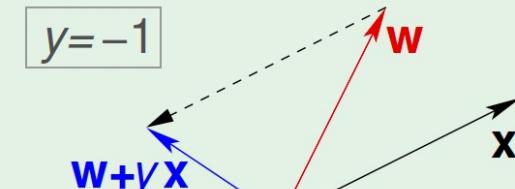
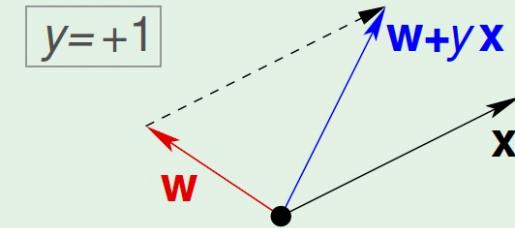
$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$$

pick a **misclassified** point:

$$\text{sign}(\mathbf{w}^\top \mathbf{x}_n) \neq y_n$$

and update the weight vector:

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$



# Iterations of PLA

- One iteration of the PLA:

$$\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$$

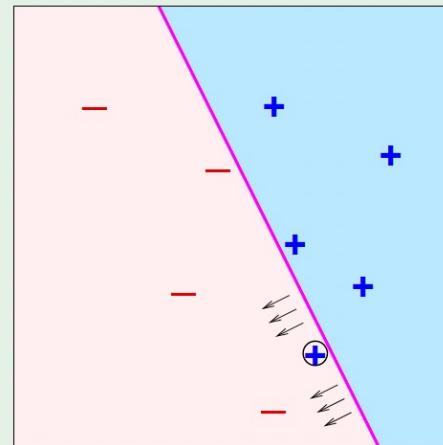
where  $(\mathbf{x}, y)$  is a misclassified training point.

- At iteration  $t = 1, 2, 3, \dots$ , pick a misclassified point from

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$$

and run a PLA iteration on it.

- That's it!



# Broad types of machine learning

- Supervised Learning
  - Training Data with labels:  $X, y$  (pre-classified)
  - Given an observation  $x$ , what is the best label for  $y$ ?
- Unsupervised learning
  - Training Data without labels:  $X$
  - Given a set of  $x$ 's, find hidden structure
- Semi-supervised Learning
  - Training Data + some Labels
- Reinforcement Learning
  - Given: observations and periodic rewards as the agent takes sequential action in an environment
  - Determine optimum policy

# Supervised Learning

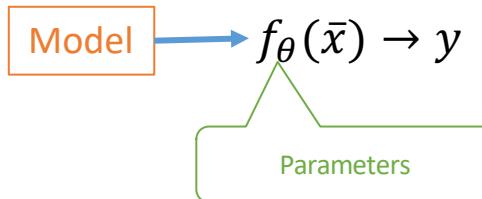
- Given data containing the inputs and outputs:

Training Data:

$$\{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_m, y_m)\}$$

- Learn a function  $f(x)$  to predict  $y$  given  $x$

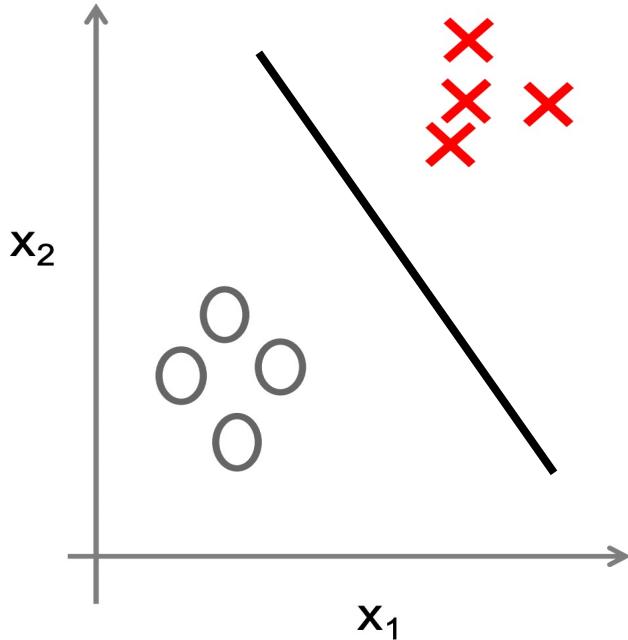
$\bar{x}$	$y$
$\bar{x}_1$	$y_1$
$\bar{x}_2$	$y_2$
...	..
$\bar{x}_m$	$y_m$



Training: Learn the model from the Training Data

Given Test instance  $\bar{x}'$ , predict  $y' = f_{\theta}(\bar{x}')$

# Supervised Learning

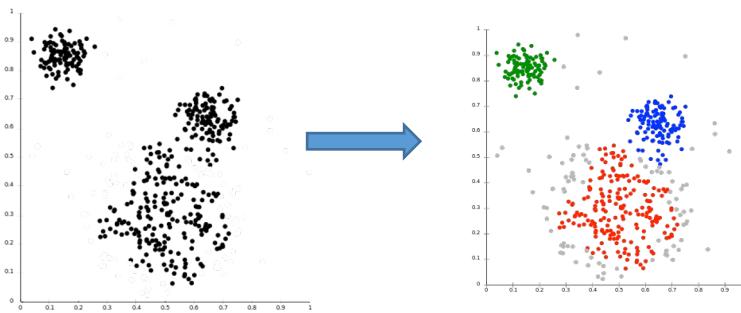


## Classification (one example)

- Input: (input-features, correct output)
  - $\langle \text{size}, \# \text{rooms} \rangle, \langle \text{cheap/costly} \rangle$
- Output of learning algorithm
  - Function maps features to output
  - $F(\langle \text{size}, \# \text{rooms} \rangle) = \text{cheap/costly}$

# Unsupervised Learning (Clustering)

- Given  $\{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m\}$  without labels
- Find hidden structure in the data
  - Clustering
  - Dimensionality Reduction
- Clustering: Grouping similar objects

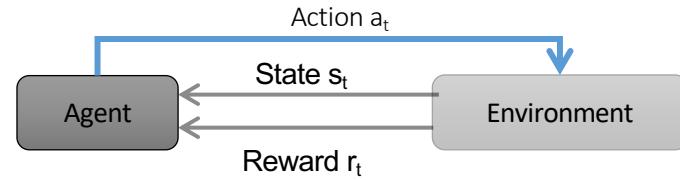


# Reinforcement Learning

- Given a sequence of states and actions with (delayed) rewards, output a policy.

- Receive feedback in the form of **rewards**
- Agent's utility is defined by the reward function
- Must (learn to) act so as to **maximize expected rewards**

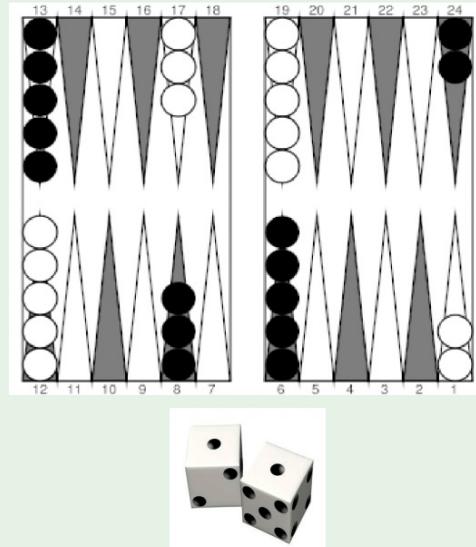
**Goal:** Constantly learn to make 'optimal' predictions based on real-time feedback from past predictions



- Examples:
  - Game playing (Go)
  - Robot grasping
  - Controlling aircraft and robotic motion

# Reinforcement Learning

Instead of (**input, correct output**),  
we get (**input, *some* output, grade for this output**)



The world champion was  
a neural network!

# Why Machine Learning?

- Human expertise not sacrosanct
  - Pricing, Promotions
- Human expertise cannot be coded
  - face/handwriting/speech recognition
  - driving a car, flying a plane
- Rapidly changing situation
  - Credit scoring, financial modeling
  - Fraud detection
- Need for personalization
  - Product recommendation
- Too much Data