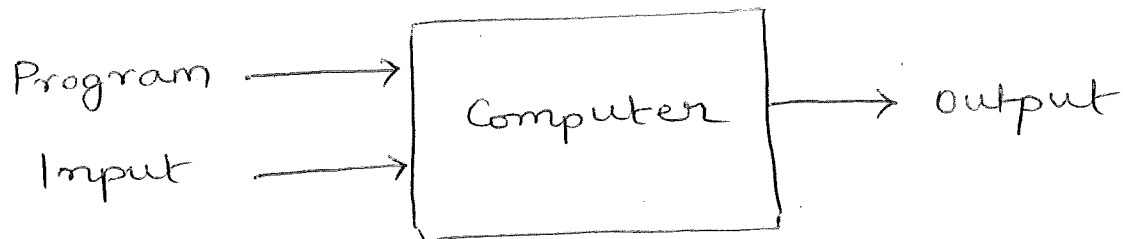1. What is machine learning?

ML is the branch of engineering that develops technology for automated inference (prediction)

- It combines
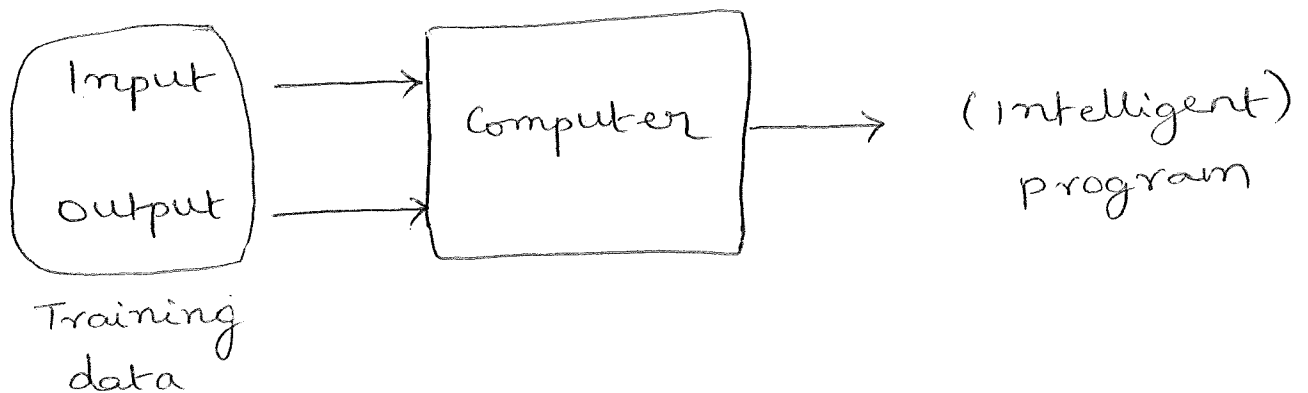
Probability + Statistics + Optimization + Algos

2. ML = Automating automation

## Traditional programming

Program $\longrightarrow$

Input $\longrightarrow$ Computer $\longrightarrow$ output

## Machine learning

Input

output

Computer $\longrightarrow$ (intelligent) program

Training data

2. Learning paradigms

   Supervised learning $\longrightarrow$ main focus of our class

   Semi-Supervised learning

   Active learning

   Reinforcement learning

3. Supervised Learning

   $x \rightarrow$ input

   $y \rightarrow$ output

   <u>Classification</u> :  $y$  is a discrete label

       Binary : 2 Labels (positive/negative)

       E.g: Spam vs. Non-Spam & Male vs. Female

       Multi-class : More than 2 Labels (Say K)

       E.g: Face recognition, Document classification

   <u>Regression</u> :  $y$  is a continuous Label

       E.g: Stock market price as a functions of financial specs.

   <u>Ranking</u> : $y$ is an ordering of a set of objects

       E.g: Search engines rank documents based on keywords

# 4. Learning a Classifier

Training Data
$(x_1, y_1)$
$(x_2, y_2)$
$\vdots$
$(x_n, y_n)$

Learning Algorithm

$\theta$

$F(x, \theta)$

x ----- feature vector

Testing or Inference

Y ----- class label

Training

We will study algorithms:
- Perceptron
- K-nearest neighbor
- Support Vector Machines
- Decision Trees

## 5. Semi-Supervised Learning

Small amount of labeled data and large amount of unlabeled data

- find a classifier that separates labeled points & unlabelede points "well"

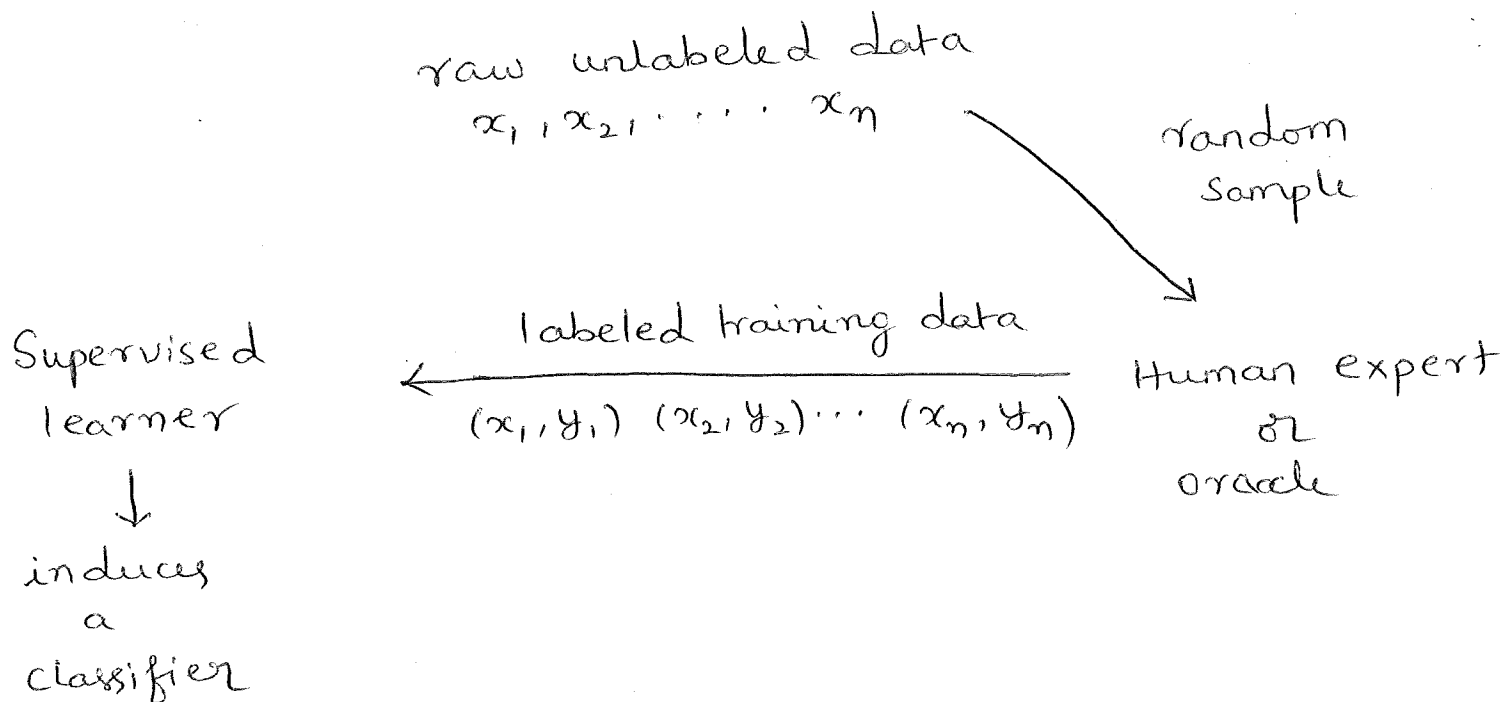- Co-traing: leverage diversity in learners to learn from each other

Eg: Diversity: multiple redundant views

webpage classification:
1) words
2) links that point to the page.

## 6. Active learning

(passive) Supervised learning

raw unlabeled data
$x_1, x_2, \ldots x_n$ → random Sample

Supervised learner ← labeled training data / $(x_1, y_1) \ (x_2, y_2) \ldots (x_n, y_n)$ → Human expert or oracle

↓
induces a classifier

# Active learning

raw unlabeled data

$x_1, x_2, x_3, \ldots$

inspect unlabeled data

request labels for selected data

$(x_1, ?)$

$\longrightarrow$

Human Expert
or
Oracle

Active learner

$\longleftarrow$ $(x_1, y_1)$

$x_2 ?$ $\longrightarrow$

$\downarrow$

induces a classifier

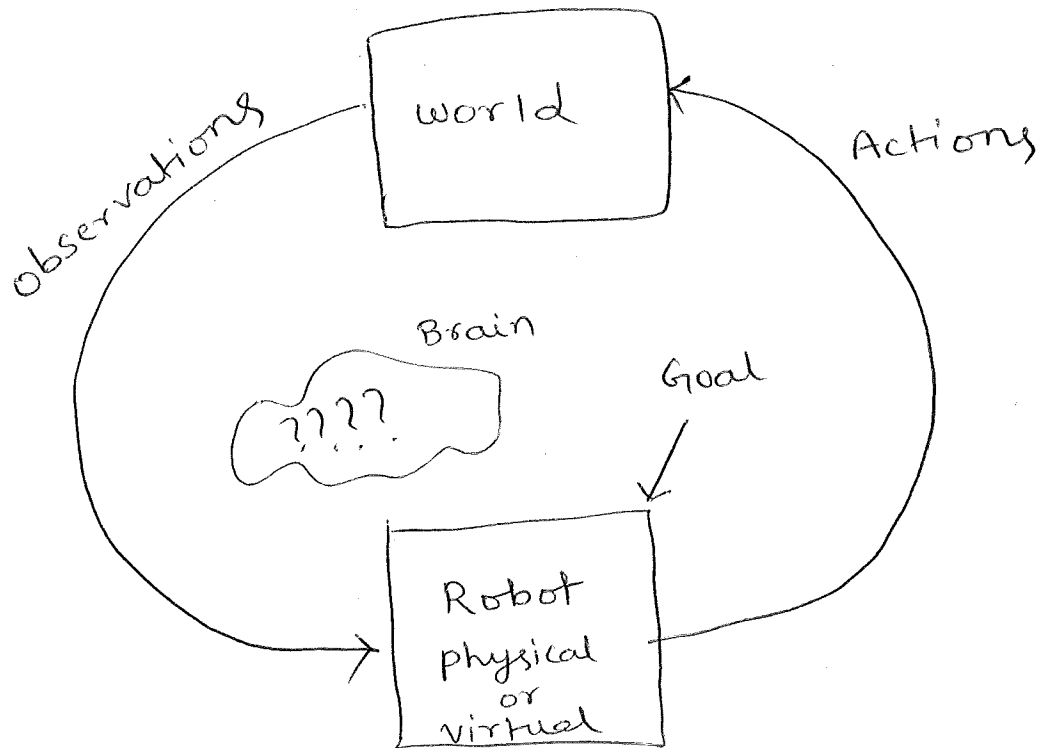$\longleftarrow$ $(x_2, y_2)$

Why ??

- Labeling is expensive
- Want to learn a highly-accurate classifier with few labeled examples
- Intelligently select the examples to get labels

Advange: Exponential efficiency

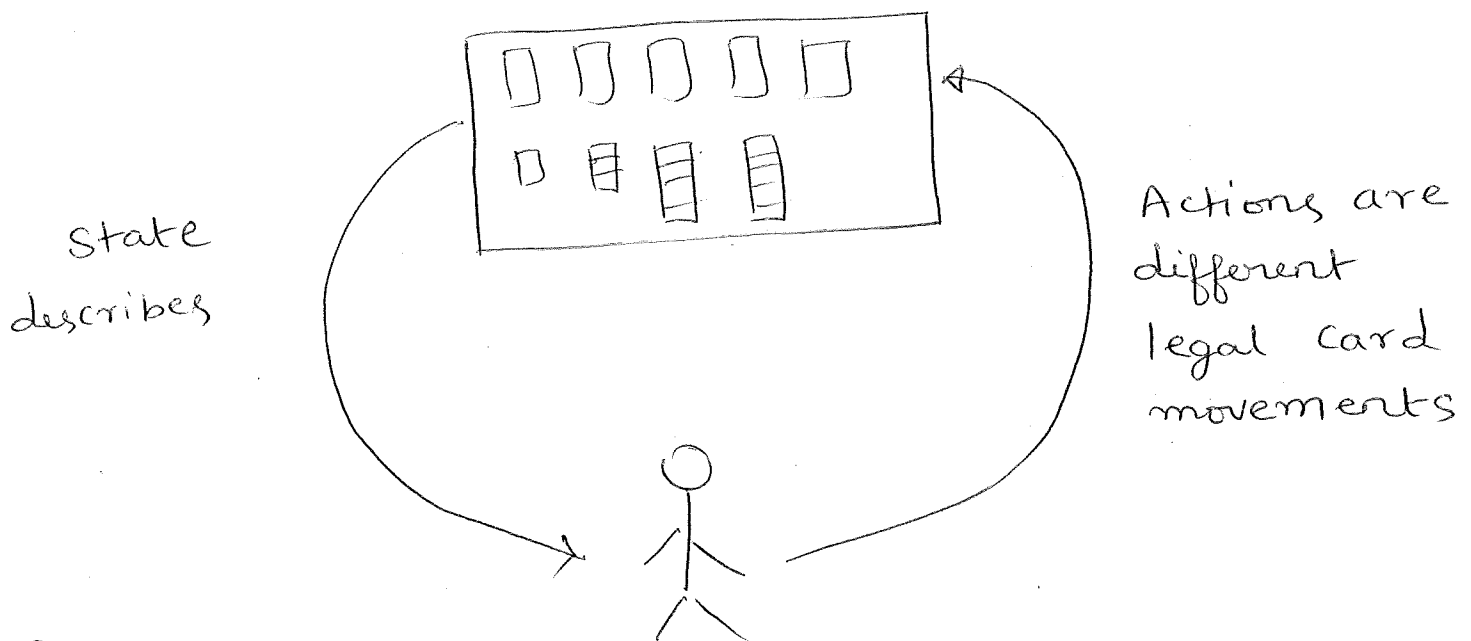$O(n)$ examples for Supervised passive learning

$O(\log n)$ " " Active learning

$\hookrightarrow$ learn accurate classifier

# 7. Reinforcement learning



Brain ????

world

observations

Actions

Goal

Robot
physical
or
virtual

**Goal:** maximize expected reward over lifetime

**Example:** card game



State
describes

Actions are
different
legal card
movements

**Goal:** win the game or play max # of cards

Alphag GO :    Deep learning + Monte-Carlo Tree Search

AlphaZero :    No Search needed


Note :  Supervised learning is often used in the
         inner-loop of different learning paradigms.


8.   Input Examples : Representation

Input examples ( emails, text documents, images)
are often represented as real-valued vectors

$$x \in \mathbb{R}^d$$

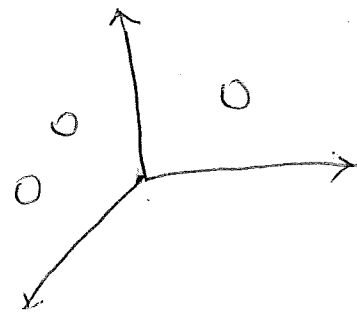   - each co-ordinate of "$x$" is called a
     feature.


   Examples :
              Bag of words representation of text
              Histograms of colors in image
              Sound frequency histogram


Bag-of-words model :          1. To be, or not to be,
   - Sentences to points        2. To be a woman,
                                3. To not be a man

| To | be | or | not | woman | a | man |
|----|----|----|----|----|----|----|
| 2 | 2 | 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 | 1 | 1 |

9. Every machine learning algorithm has three components

    1. Representation

    2. Evaluation

    3. Optimization

Rep Examples:

Linear hyperplanes
Decision trees
Neural networks.

$$x \in R^d$$
$$w \in R^d$$

$\boxed{Sign \ (w \cdot x)}$ classifier

Evaluation examples:

Accuracy
Precision and Recall
likelihood
Entropy

Optimization Examples:

    Combinatorial : greedy Search
                  dynamic Programming

Convex optimization : gradient descent

Constrained     " : linear programming
                        Quadratic programming

# Supervised Learning:

Given: a set of training examples $\{(x, y^*)\}$ drawn from an unknown target distribution $D$.

Find: a function $F$ that maps input $(x)$ to output $(y)$ such that predictions have ~~accur~~ high accuracy on <u>unseen</u> inputs from $D$.

<u>Learning goal</u>: Generalization

(not memorization)

## <u>Two types of learning algorithms</u>:

1. Online learning
   - processes one training example at a time incrementally. $(f_t \rightarrow f_{t+1})$
   - Game between teacher & student

2. Batch learning
   - processes all training examples at once and produces a globally optimized function $F$

# Formal Setting: Classification

1. Instances input $x$

   e.g: Emails

2. Output Labels $y \in \{+1, -1\}$

   e.g: Spam vs. Non-Spam

3. Prediction rule $f(x) = \hat{y}$

   e.g: linear prediction rule

4. Loss $\ell(y^*, \hat{y}) \in [0, \infty)$

   e.g: Zero-one error

# Linear classifier:

$$\hat{y} = \text{Sign}(f(x))$$

$$= \text{Sign}(w \cdot x)$$

weights $\in R^d$      features $\in R^d$

Confidence $= |f(x)| = |w \cdot x|$

# Perceptron Algorithm:

Simple and most popular ML algorithm

Online learning algorithm

## Online learning Framework:

Initialize classifier $f_1(x)$

Algorithm works in rounds

on round "$t$", the online algorithm:

- Receives input $x$
- Outputs a prediction $\hat{y} = f_t(x)$
- Receives a feedback $y^*$
    Label
- Computes Error $\quad l(\hat{y}, y^*)$

- if loss $> 0$, update rule

$$\boxed{f_t \rightarrow f_{t+1}}$$

Learning : Suffer small cumulative loss/Error
goal

$$\sum_{t=1}^{Tmax} l(\hat{y}, y^*)$$

# Why online learning?

- Fast
- Memory efficient : process one example at a time
- Simple to implement : less than 30 mins 😊
- Online to Batch Conversions
- Adaptive

# Design Principle :

If the learner suffers non-zero loss at any round, then we want to balance two goals :

1. <u>Corrective</u> : Update function so that we don't make this error again

2. <u>Conservative</u> : Don't change the function too much

\* Different online learning algorithms make different trade-off's between these two competing goals.

# Perceptron algorithm :

The rule to update function

$$f_t \rightarrow f_{t+1}$$

# Linear Classifiers :

Find $w_{t+1}$ from $w_t$ based on the training example $(x, y^*)$

# Algorithm :

If no mistake : $y^*(w_t \cdot x) > 0$

then <u>DO nothing</u>

$$\Rightarrow w_{t+1} = w_t$$

If mistake : $y^*(w_t \cdot x) \leq 0$

Update weights

$$w_{t+1} \leftarrow w_t + y^* x$$

# Running Example :

Training data :

|  | $x_1$ | $y_1$ |
|---|---|---|
|  | $(4, 0)$ | $+1$ |
|  | $x_2$ | $y_2$ |
|  | $(1, 1)$ | $-1$ |
|  | $x_3$ | $y_3$ |
|  | $(0, 1)$ | $-1$ |
|  | $x_4$ | $y_4$ |
|  | $(-2, -2)$ | $+1$ |

$w_1 = 0$

For $t = 1$,

$$y^* \cdot (w_1 \cdot x)$$

$$= 0 \qquad \text{// mistake}$$

$\Rightarrow$

$$w_2 = w_1 + 1 \cdot (4, 0)$$

$$= (4, 0)$$

For $t = 2$

$$y^* (w_2 \cdot x) < 0$$

$\Rightarrow \quad w_3 = w_2 + (-1) \cdot (1, 1)$

$$= (4, 0) - (1, 1)$$

$$= (3, -1)$$

For $t = 3$

$$y^*.(w_3 \cdot x) > 0 \qquad \text{// correct}$$

$$\Rightarrow \qquad w_{34} = w_3 = (3, -1)$$

For $t = 4$

$$y^*(w_4 \cdot x) < 0$$

$$\Rightarrow \qquad w_5 = w_4 + 1 \cdot (-2, -2)$$

$$= (3, -1) + (-2, -2)$$

$$= (1, -3)$$

## When does Perceptron Converge:

linear Separability : If there exists a weight vector that can separate positive and negative points.



linearly
Separable

Not linearly
Separable.

# Measure of Separability :

**Margin :** For a weight vector $w$ & training set $S = \{(x, y^*)\}$, margin of $w$ with respect to $S$ is defined as follows :

$$\gamma(w) = \min_{(x, y^*) \in S} y \cdot (w \cdot x)$$

+ + + x

+ + +
 + +
 -
  - -
    -

High-margin data

low-margin data

- - -
 - - -

**Convergence :** If training set is linearly separable with margin $\gamma$, then perceptron makes $\leq \dfrac{1}{\gamma^2}$ mistakes

1. lower margin implies more mistakes

2. May need more than one pass over the training data to get a classifier with no mistakes

# what if data is not linearly Separable?

Perceptron Still works
- there will be few mistakes close to the decision boundary
- will never converge to a single "w" as we make more passes.

## Voted Perceptron:

Initialization: $m=1$; $w_1 = 0$; $c_m = 1$

Training examples: for $t = 1, 2, 3, \cdots$

△ If mistake, update weights

$$w_{m+1} = w_m + y^*_b x_b$$

$$m = m+1$$

$$c_m = 1$$

△ Else

$$c_m = c_m + 1 \quad // \text{Counting how long } w_m \text{ survived}$$

Output: $(w_1, c_1), (w_2, c_2), (w_3, c_3) \cdots$

# Voted Perceptron classifier

weighted majority vote of all the classifiers.

## Drawbacks:

1. we have to store many classifiers (space)

2. we need to make many predictions (time)

## Averaged Perceptron:

$$W = \frac{\sum_{i=1}^{K} c_i * w_i}{K}$$

Averaging $\Rightarrow$ robustness & regularization (leading to better generalization)

## Some Practical tricks:

1. Shuffling: shuffle the training examples in each iteration

2. Variable learning rate: decrease as learning progresses

# Learning Curve:

- Training iterations vs. no of mistakes

- You want to see that mistakes decrease as we increase the no of iterations.

- Very useful in debugging & seeing the learning behaviour.

# Hyper-parameter Optimization:

- Split the training data: Sub-train
  +
  validation

- Tune hyper-parameters (e.g., no of iterations) on the validation data

- The learner should not look at the test data.

# Multi-class classification:

Suppose we have $(K > 2)$ classes.

$K$ weight vectors : $w_1, w_2, \ldots, w_K \in R^d$

input instance $x \in R^d$

Score (label $r$) = $w_r \cdot x$

| Class $r$ | $w_r \cdot x$ |
|-----------|---------------|
| 1         | $-1.08$       |
| 2         | $1.66$        |
| 3         | $0.37$        |
| 4         | $-2.09$       |

Prediction: output label (class) with highest score.

Learning:

$$w_{y^*} = w_{y^*} + x$$

$$w_{\hat{y}} = w_{\hat{y}} - x$$

# Regression Learning:

$y$ is continuous value.

Prediction Rule: $F(x) = w \cdot x$

# Widrow - Hoff Algorithm:

- Initialize $w_1 = 0$

for $t = 1$ to $T$ do

- get $x_t \in R^d$

- predict $\hat{y}_t = w_t \cdot x_t$

- observe $y^*_t$

- Incur loss of $(\hat{y}_t - y^*_t)^2$

- $w_{t+1} = w_t - \eta (w_t \cdot x_t - y^*_t) * x_t$

end