

AI

23-01-24

Artificial Intelligence : Computers mimicking human behaviour

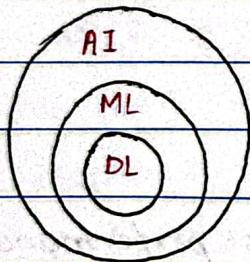
Data ↗ continuous (regression)  
└ discrete (classification)

Tesla cars, ChatGPT all are examples of weak AI

AI ↗ weak AI (narrow AI) → task specific (all AI tech we use)

AI ↗ strong AI machine is as smart as humans

super AI machine surpasses human behavior



supervised learning : Labeled training data

unsupervised learning : given unlabeled data

- Agent : software or system designed to act autonomously & perform specific tasks.
- sensors → perceives environment      actuators → act through using them
- Rational agent → does the right thing
- percept sequence → the prior knowledge, history
- PEAS → Sensors  
 ↓  
 performance measure      environment      ↓  
 actuators
- Fully observable : complete knowledge of environment
- Partially observable
- Unobservable → dishwasher
- Deterministic → next state of environment is determined by current state & action executed by agent otherwise stochastic



- stochastic → certain degree of randomness or probability

- strategic → refers to decision making process

- single-agent vs multi-agent

solving crossword  
puzzle

chess playing

competitive Partially competitive

co-operative

- episodic vs sequential

↓      ↓  
single action current decision affects future decisions

doesn't affect

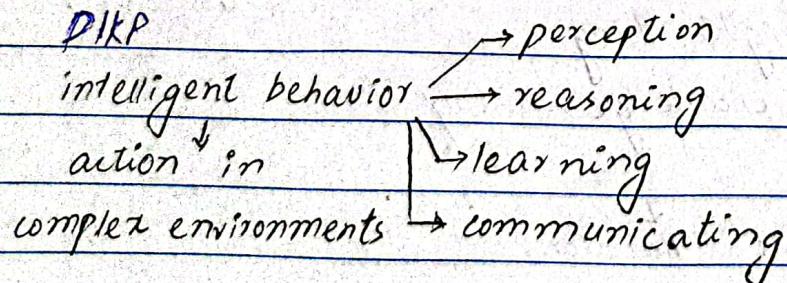
any other action

- static vs dynamic → env can change

semi-dynamic → env itself doesn't change but agent's performance score does

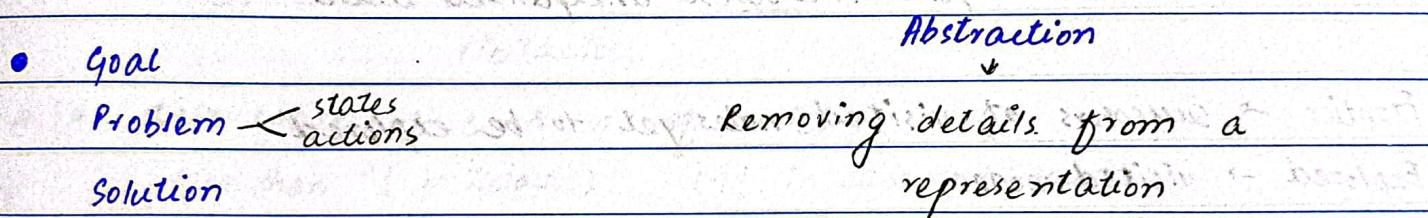
- Discrete vs continuous

- Known vs unknown → focused on agent's state of knowledge



Agent function → maps the percept into an action

- Simple Reflex Agent → simple "if then" rules
- Model-based → percept history (previously-observed history) →
- Goal-based → goal-oriented → focuses on overall goal. Takes into account the actions at a particular instance of time irrespective of the overall goal.
- Utility-based → efficiency, desirability of a particular state or outcome
- Learning agent → improve their performance over time



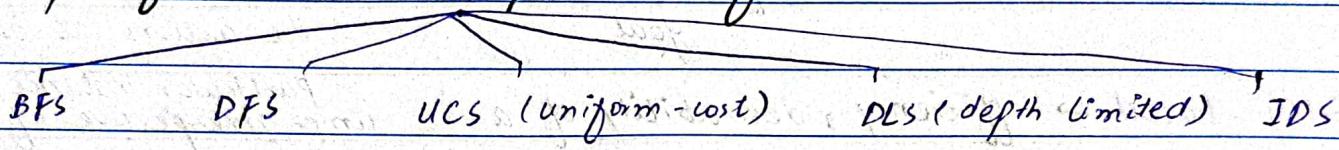
- Toy problem: illustrating some problem solving method using a less complex problem than real world scenarios.
- $n^2$  → # of states in state space graph.
- 8-Queen Problem : Place 8 Queens on chessboard such that no Queen attack any other
  - ↪ complete state
  - ↪ Incremental
    - all queens placed on the board
    - adding one queen at a time
- Search strategy defined by picking order of node expansion

predicate calculus

31-01-2024

Uniformed Search Strategies : No prior knowledge about the no of nodes etc. Only "goal node" is known.

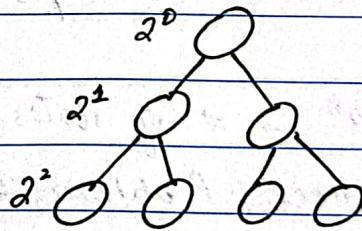
→ only info available in problem definition



- Breadth-first search → Queue → FIFO Expand Shallowest unexpanded node

Branching factors → # of successors of a node

the branching factor  $\leftarrow b^d \rightarrow$  depth



- Shallowest node → the node closest to root node.
- Time complexity & space complexity both  $\rightarrow O(b^d)$
- unequally weighted & unweighted graphs → BFS optimal

• Depth-first search →

stack → LIFO

Expand deepest unexpanded node in the current fringe.

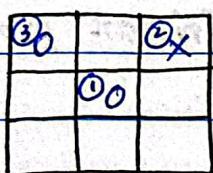
Not complete, not optimal

→ Problem solving & State space:-

- initial state
- current state
- goal state

State space → All possible states

Duplicate states problem → Because sometimes diff paths can lead to the same state



different order & diff branches of tree.

Tic Tac Toe

- So, state space's size ≠ no of nodes

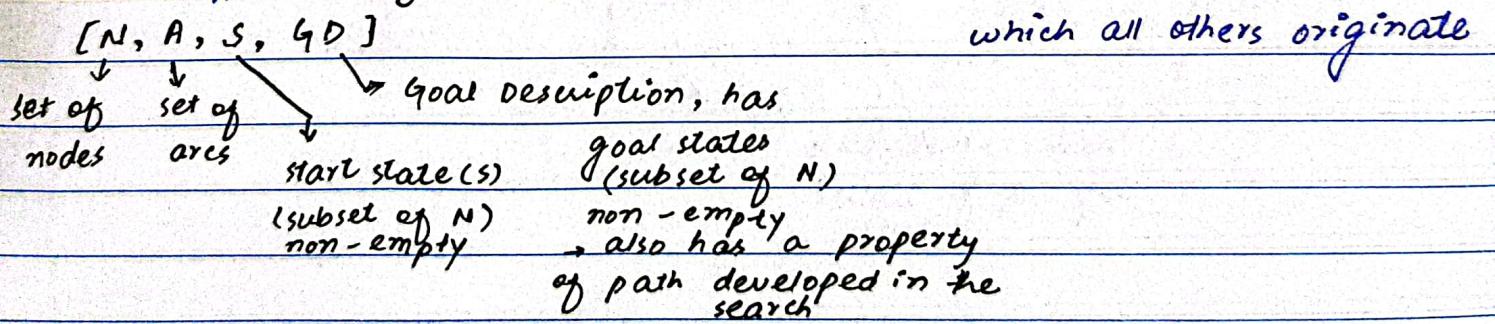
(less than  $9!$  in Tic Tac Toe) ( $9!$  in case of Tic Tac Toe)

- we cannot do exhaustive search in case of this many nodes in tree
- travelling salesman problem complexity →  $(n-1)!$

Graph : nodes & arcs

- labeled graph , directed , undirected , rooted graph

- state space → four-tuple



- Strategies for state space search

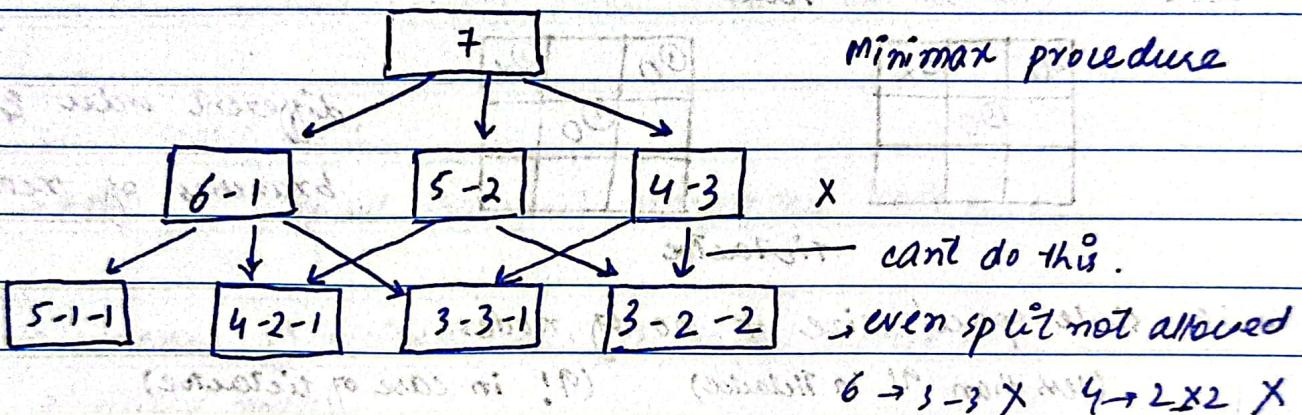
Data-Driven  
Forward chaining

Goal-Driven  
Backward-chaining

$S \rightarrow A \rightarrow B \rightarrow D \quad 7$ 

$\hookrightarrow \{ [S \rightarrow A \rightarrow C \rightarrow G, 4], [S \rightarrow A \rightarrow C \rightarrow D \rightarrow G, 6], [S \rightarrow A \rightarrow B \rightarrow D, 7], [S \rightarrow G, 12] \}$

UGS  $\rightarrow$  path cost 0  $\rightarrow$  stuck in recursive loop.

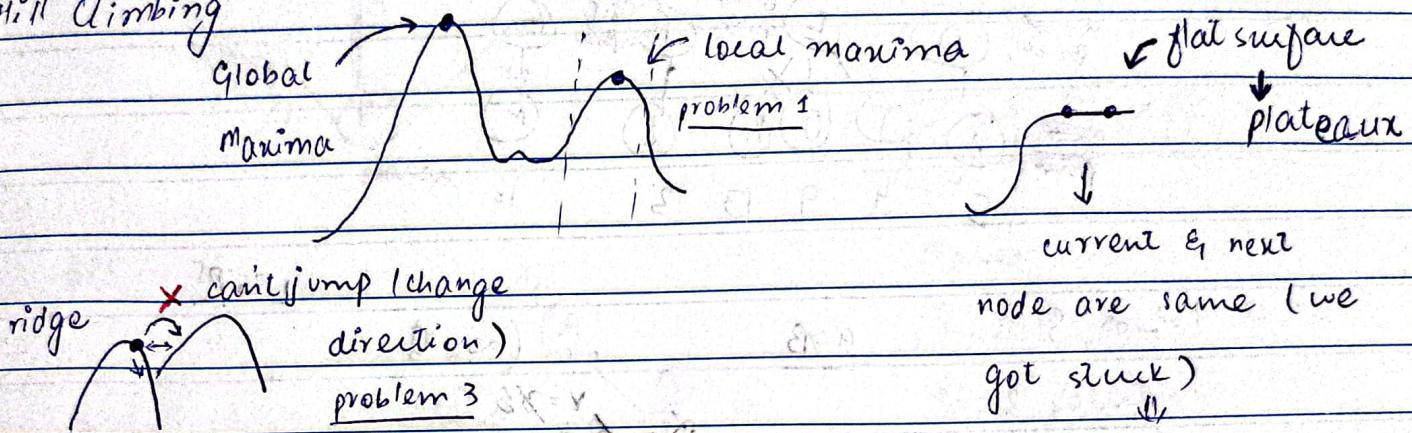


04-03-24

## Local Search Algorithms :

- ↳ path to goal irrelevant, goal state itself is the solution
- ↳ can't backtrack
- ↳ memory efficient

### Hill Climbing



No previous info retained, like climbing everest in thick fog and amnesia.

node are same (we got stuck)

take a bad move & move to next node.

⇒ can be a plato or a shoulder

### Variants

• Simple hill climbing  
examines one neighbour

• Steepest Ascent  
examines all neighbour

• stochastic hill climbing  
chooses at random from among uphill moves

• Random-restart  
restart if you get stuck

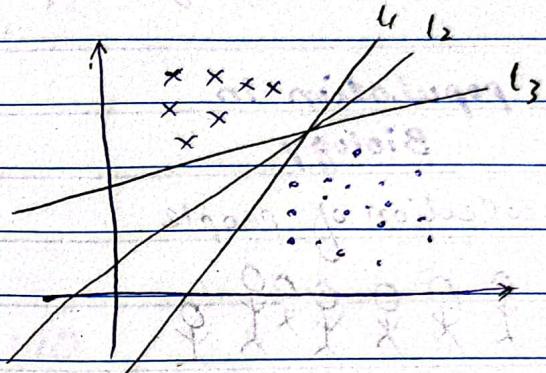
## Evolutionary Algorithms

### → Genetic Algorithm

→ biologically inspired model of intelligence

→ survival of the fittest

→ genetic operators → crossover,  
mutation



### Attributes :-

- genetic representation of candidate solutions (chromosomes)

each one ↓  
is unique

phenotype → genotype  
conversion (encoding)



$$y = mx + c$$

↑      ↑

index = col val = row

	X		
	.		X
X	.	*	
*	X	.	

- Population size

- fitness func / evaluation func / score

- Genetic operators (crossover, mutation)

- selection algorithm

- generation gap (how many children produced and replaced by parent)

- amount of elitism used

- # of duplicates allowed

4-Queen problem

### Chromosome :

Gene → basic unit, represents one characteristic of individual

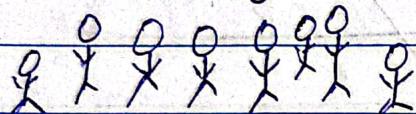
Allele → value of each gene (values it can take i.e. eye color gene can be green, black, brown etc).

String of genes, represents a point in search space

**population:** collection of chromosomes

**population in  
Biology**

collection of people



**population in  
Algorithms**

collection of states

x	1	0	
0	x		x
		0	

**fittest chromosome  
in Biology**

more healthy less prone  
to disease

**fittest chromosome  
in algorithms**

closest to final  
solution

**Crossover:-**

P1  
11.0000

P2  
10111.0

A	1 0 0 1 0 1	0 0 1 0 1 0	fitness
B	1 1 0 1 1 1	1 1 0 1 1 1	1
C	1 1 0 0 0 0	1 1 0 0 0 0	2
D	1 0 1 1 1 0	1 0 1 1 1 0	4

for ③

③

so, next generation doesn't guaranteedly gives optimal or better score than parents.

don't care terms = # or X

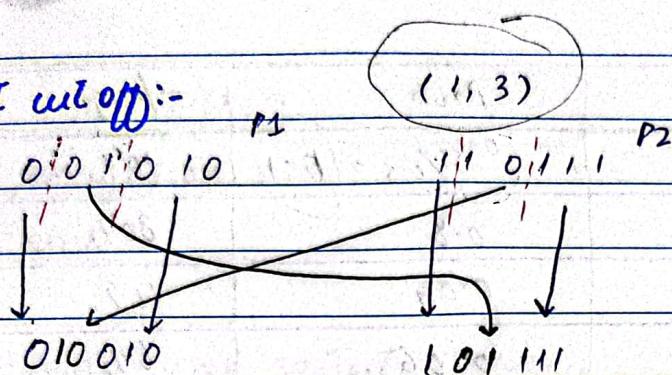
001010 sol

11#010 fitness score = 4

001010 sol

110010 f.s = 3

Two point cut off :-



multipoint crossover :-

01111      001110

Randomly select bits from both parents.

- crossover mask :- 

0	1	0	0	1	0
---	---	---	---	---	---

This means when there is 0, pick bit from P1 and when 1, pick bit from P2.

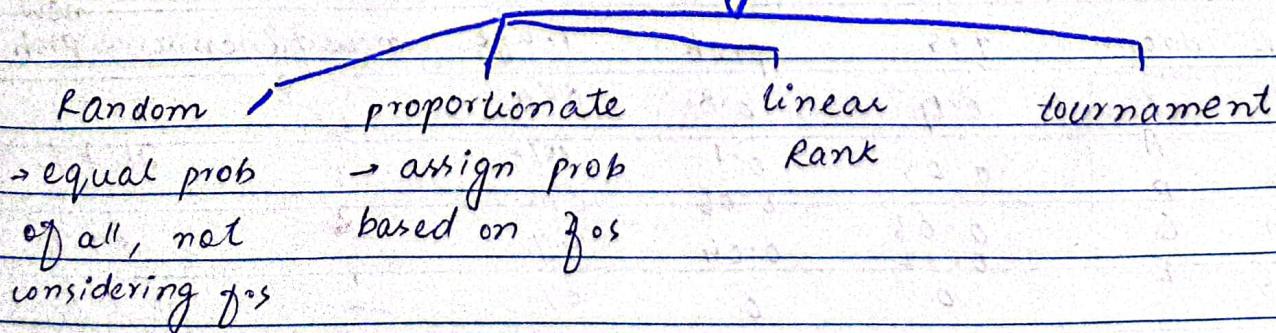
**Mutation**:- when the population is not diverse enough, we use mutation to handle it.

↳ Takes a single candidate & randomly changes some aspect (gene) of it.

• Mutation rate  $\rightarrow$  less  $\Rightarrow$  convergence takes time  
more  $\Rightarrow$  original generation changed

• if generation gap is 1, completely replace all parents. If it is 0.5, half of generation is old, half is new.

### Selection Algos



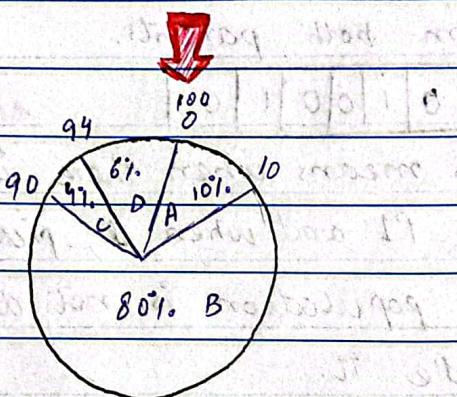


## → Proportionate Selection

chromosome	Fitness Score	Prob	%age
A	0.05	$0.05/0.5 = 0.1$	10%
B	0.4	0.8	80%
C	0.02	0.04	4%
D	0.03	0.06	6%

0.5      011100

- Roulette wheel scenario:



To implement, define ranges for A, B, C, D ... & generate a random num & see where it lies.

If a chromosome say E has 0% prob then it will never go on the roulette wheel  $\Rightarrow$  disadvantage.

E      0      1      0      1      0%

## → linear Rank selection

rank	chrom	f.s	prob	%age	new fitness.sc	new prob
1	B	0.9	0.8	80%	8	33.33%
2	A	0.05	0.1	10%	4	26.7%
3	D	0.03	0.06	6%	3	20%
4	C	0.02	0.04	4%	2	13.33%
5	E	0	0	0%	1	6.67%

15

$$\boxed{\text{New fitness} = (P - r) + 1}$$

$P$  = population size

$r$  = rank

to incorporate user adjusted slope:

$$\boxed{N \cdot F = \frac{(P-r)(\max - \min)}{(P-1)} + \min}$$

i.e.  $\max = 8, \min = 3$

$$N \cdot F = \frac{(5-1)(8-3) + 3}{(5-1)} = 8$$

m-method :-

r	chrom	f.s	prob	$N \cdot \text{prob}$
1	B	0.4	0.8	(0.8)
2	A	0.05	0.1	(0.8)(1-0.8)
3	D	0.03	0.06	(0.8)(1-0.8)^2
4	C	0.02	0.04	(0.8)(1-0.8)^3
5	E	0	0	(0.8)(1-0.8)^4

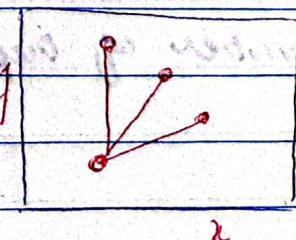
### • Survival of the most diverse

It can be as good to be different as it is to be fit.

→ results in uniformity

$$\uparrow \text{Diversity} = \sum_j \text{distance}_j^2$$

$$\downarrow \text{Div} = \frac{1}{\sum_i d_i^2} \uparrow$$



0.5

12 : 3 : 4

5 : 6



$$0.8 \times 4 = \dots \quad 4 \times 0.5 = 2$$

1/4

function = fitness + Diversity

### Generation gap:

fraction of current population that gets replaced by the offspring. Between 0 - 1

Generation gap \* P = offspring #

P = 6

gap = 0.5

$$6 \times 0.5 = 3$$

C1

C2

C3

C4

C5

C6

} replaces least fittest 3 states  
by offsprings.

Amount of Elitism used: fraction that represents fraction of Best Individuals of a population that will not get replaced.

0 - 1. only used if gen. gap is not used.

Q :-

50

↓

$$\text{generation gap} = \frac{30}{50} = 0.6$$

(5) + 45  
elitism / 15 old      30 new

$$\text{elitism} = \frac{5}{50} = 0.1$$

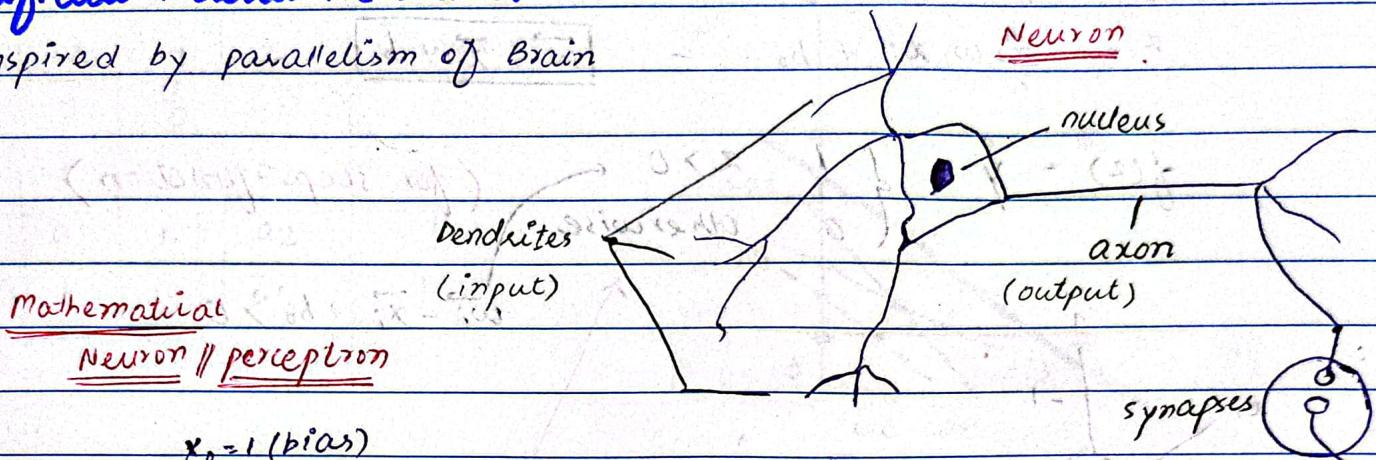
20  
↓  
total old  
30  
↓  
new

Number of duplicates allowed :

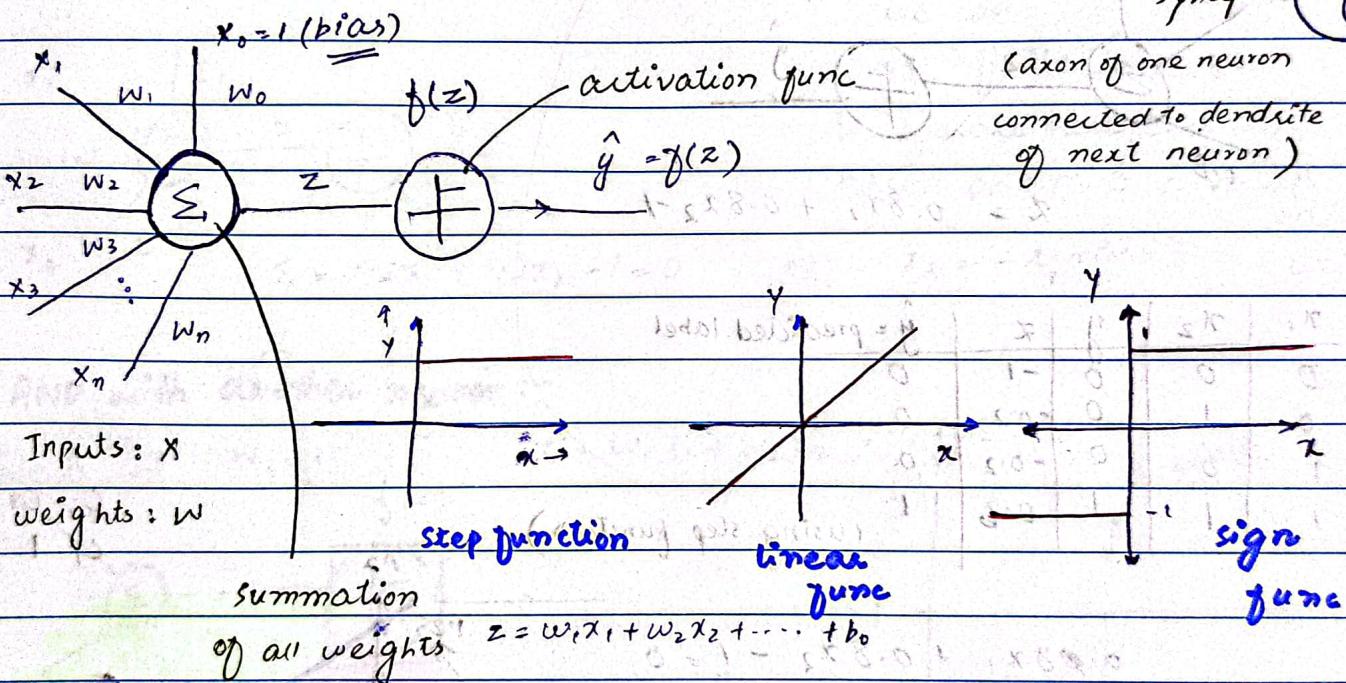
## Artificial Neural Network:

inspired by parallelism of Brain

### Neuron

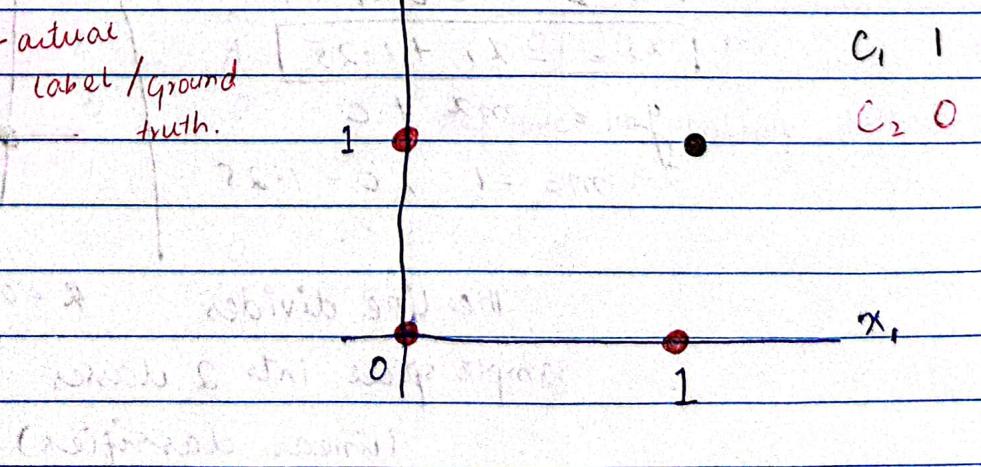


Mathematical Neuron // perception



### AND GATE:

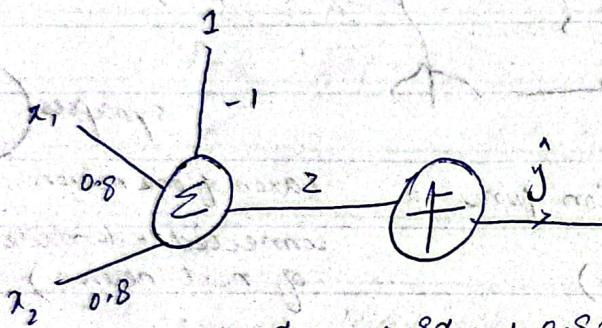
$t^1 x_1$	$t^2 x_2$	$y$ ← actual label / ground truth.
0	0	0
0	1	0
1	0	0
1	1	1



$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b_0$$

$$z = \sum_i w_i x_i + b_0 = \vec{w} \cdot \vec{x} + b_0$$

$$f(z) = \hat{y} = \begin{cases} 1 & z > 0 \\ 0 & \text{otherwise} \end{cases} \quad (\text{for step function})$$



$$z = 0.8x_1 + 0.8x_2 - 1$$

$x_1$	$x_2$	$y$	$z$	$\hat{y}$ = predicted label
0	0	0	-1	0
0	1	0	-0.2	0
1	0	0	-0.2	0
1	1	1	-0.6	1

(using step function)

$$0.8x_1 + 0.8x_2 - 1 = 0$$

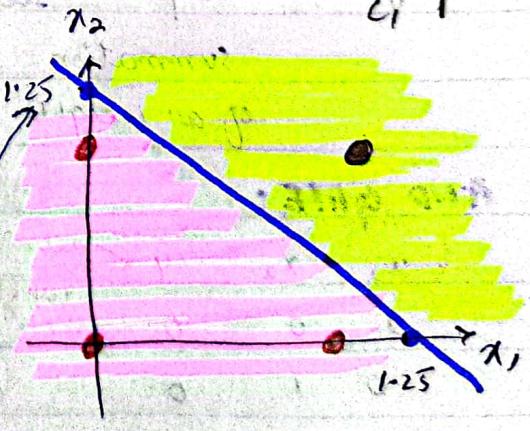
$$0.8x_2 = -0.8x_1 + 1$$

$$x_2 = -x_1 + 1.25$$

$$y = mx + c$$

$$m = -1, c = 1.25$$

The line divides sample space into 2 classes  
(linear classifier)

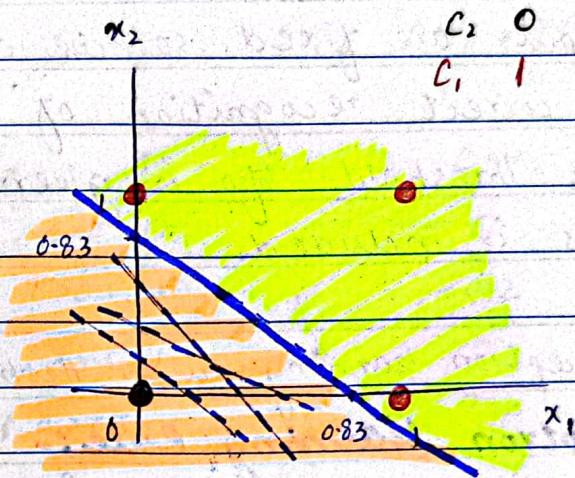


$$x_1 + x_2 - 1.25 > 0$$

C2

OR:

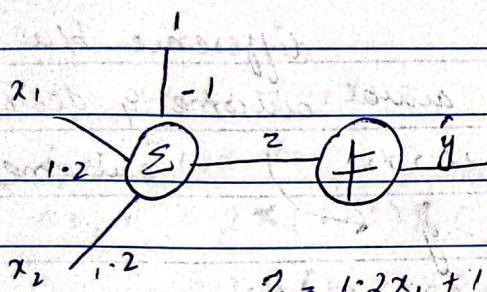
$x_1$	$x_2$	$y$	$z$	$\hat{y}$ = pred. val
0	0	0	-1	0
0	1	1	0.2	1
1	0	1	0.2	1
1	1	1	1.4	1



we can

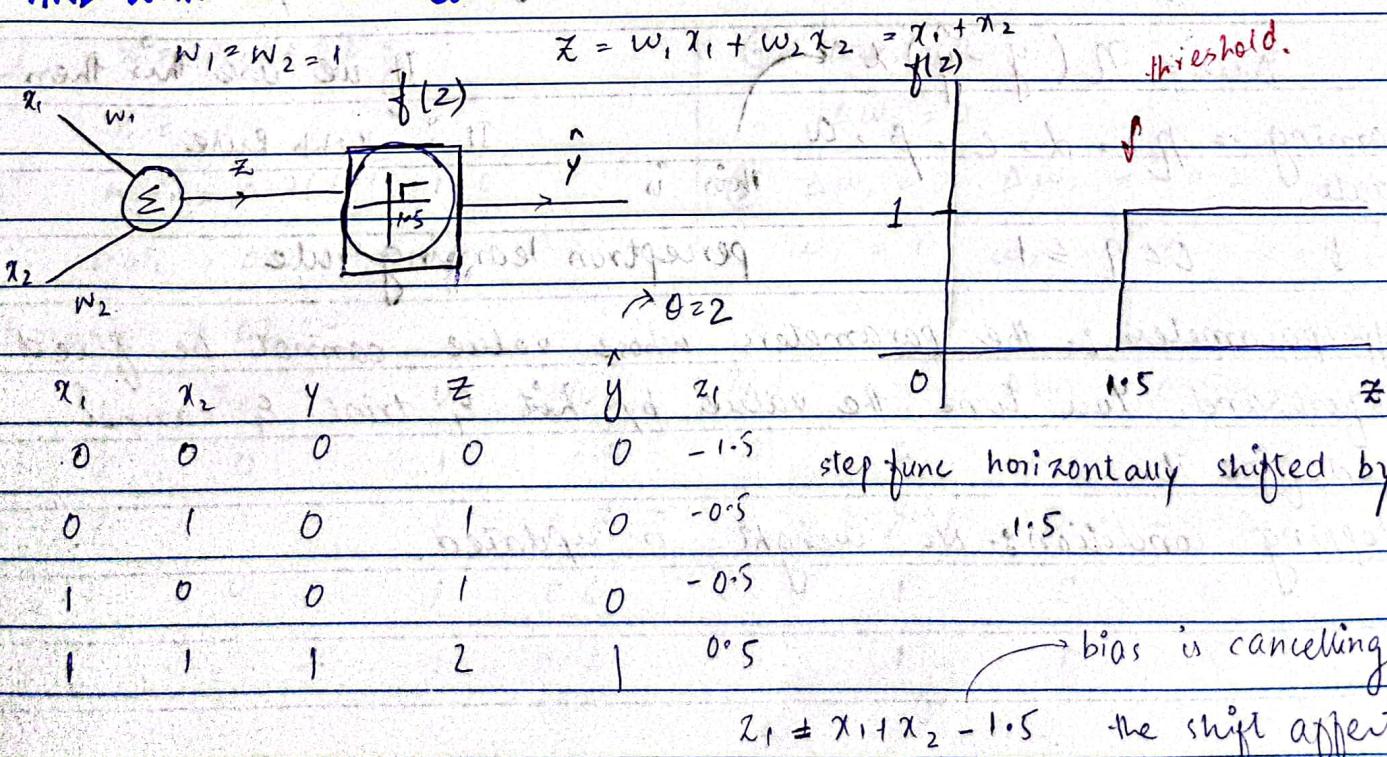
have infinite lines

to separate these 2 classes  
based on m  $\epsilon_i$  c values  
adjusting weights



$$z = 1.2x_1 + 1.2x_2 - 1 = 0 \Rightarrow x_2 = -x_1 + 0.83$$

AND with another neuron:-



- inputs are fixed so we can only adjust the weights for correct recognition of patterns.
- If threshold info is given, then there is no bias. otherwise bias is included.

~~perception learning rule, Hebb rule~~

### Perception learning Rule:-

$x_1$	$x_2$	$y$	$z$	$\hat{y}$	difference b/w actual outcome & desired outcome.
0	0	0	0	0	
0	1	0	0.5	1	$z = 0.5x_1 + 0.5x_2$
1	0	0	0.5	1	$\text{error} = y - \hat{y}$
1	1	1	1	1	$w_{\text{new}} = w_{\text{old}} + \Delta w$

$$\text{threshold} = 0.4$$

$$\Delta w = (y - \hat{y}) x$$

$$\Delta w = \eta (y - \hat{y}) x$$

If we use this then

$$\text{learning rate} = \eta = d = c = \beta = a$$

this is

$$\downarrow \quad 0 < \eta \leq 1 \quad \text{perceptron learning rule}$$

hyperparameters :- the parameters whose value cannot be fixed beforehand. You tune the value by hit & trial & cannot predefine it.

Stopping condition:- No weight is updated.

$$\eta = 0.2, \text{ threshold} = 0.5, w_{\text{new},i} = w_{\text{old},i} + \Delta w_i, \Delta w_i = \eta(y - \hat{y})x_i$$

iteration # 01 :-

$x_1$	$x_2$	$w_1$	$w_2$	$y/d$	$\hat{y}$	$\Delta w_1$	$\Delta w_2$
0	0	1	1	0	$0 > 0.5 : 0$	0	0
0	1	1	1	0	$1 > 0.5 : 1$	0	-0.2
1	0	1	0.8	1	$1 > 0.5 = 1$	0	0
1	1	1	0.8	1	$0.8 > 0.5 = 1$	0	0

$$① z = w_1 x_1 + w_2 x_2$$

$$z = x_1 + x_2$$

$$\Delta w_1 = 0.2 \times (0-0) \times 0 = 0$$

$$\Delta w_2 = 0.2 \times (0-0) \times 0 = 0$$

$$w_1 = 1 + 0 = 1 \quad w_2 = 1 + 0 = 1$$

$$② \Delta w_1 = 0.2(0-1)0 = 0$$

$$\Delta w_2 = 0.2(0-1)1 = -0.2$$

$$w_1 = 1$$

$$w_2 = 1 - 0.2 = 0.8$$

$$③ z = x_1 + 0.8 x_2$$

$$z = 1$$

$$\Delta w_1 = 0.2(1-1) \cdot 1 = 0$$

$$\Delta w_2 = 0.2(1-1) \cdot 0 = 0$$

$$④ \Delta w_1 = 0$$

$$\Delta w_2 = 0$$

$$⑤ \Delta w_1 = 0 \quad \Delta w_2 = -0.2$$

$$w_1 = 1 \quad w_2 = 0.8 - 0.2 = 0.6$$

Iteration # 02 :-

$x_1$	$x_2$	$w_1$	$w_2$	$y/d$	$\hat{y}$	$\Delta w_1$	$\Delta w_2$
0	0	1	0.8	0	0	0	0
0	1	1	0.8	0	$0.8 > 0.5 : 1$	0	-0.2
1	0	1	0.6	1	1	0	0
1	1	1	0.6	1	1	0	0

iteration # 03 :-

$x_1$	$x_2$	$w_1$	$w_2$	$y$	$\hat{y}$	$\Delta w_1$	$\Delta w_2$
0	0	1	0.6	0	0	0	0
0	1	0.1	0.6	0	0.6 > 0.5 = 1	0	-0.2
1	0	0.1	0.4	10	1	0	0
1	1	0.1	0.4	11	1	0	0

3 (2)

$$\Delta w_1 = 0.2(0-1)0 = 0$$

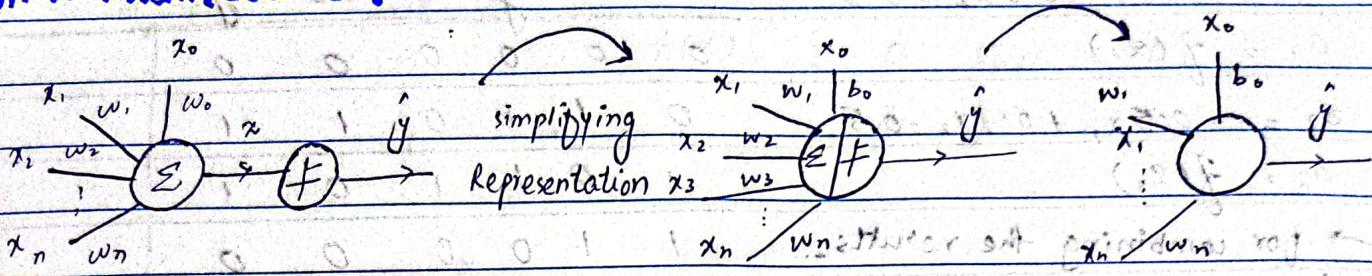
$$\Delta w_2 = 0.2(0-1)1 = -0.2$$

$$w_1 = 1 \quad w_2 = 0.6 - 0.2 = 0.4$$

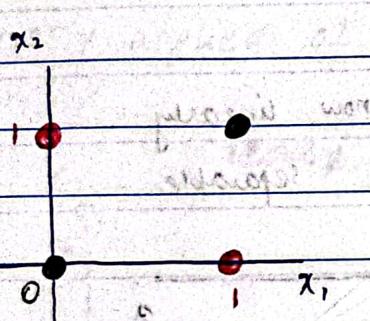
iteration # 04 :-

$x_1$	$x_2$	$w_1$	$w_2$	$y$	$\hat{y}$	$\Delta w_1$	$\Delta w_2$
0	0	1	0.4	0	0	0	0
0	1	1	0.4	0	0	0	0
1	0	10	0.4	1	01	0	0
1	1	1	0.4	1	1	0	0

## ANN Architecture :



$$XOR = x_1 \oplus x_2$$

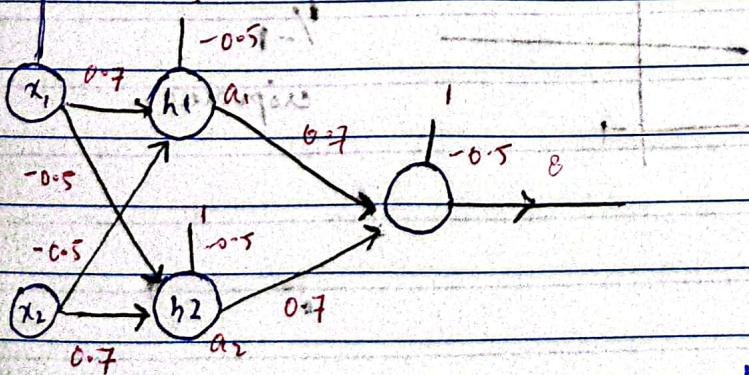


$$C1 = 1 \quad C2 = 0$$

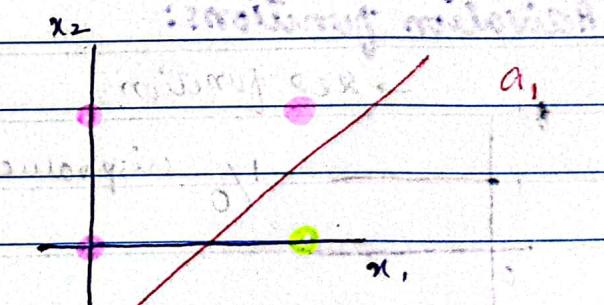
not linearly separable as 1 line

cannot divide the classes so we need

2 perceptrons.  
circle for uniformity (it's not a perceptron)



level 0      level 1      level 2  
input layer    hidden layer    output layer  
→ Multilayer perceptron (MLPs)



## Feed Forward NN:

The inputs only propagate forward

## Recurrent NN:

inputs go backwards through some loop.

$$z_1 = 0.7x_1 - 0.5x_2 - 0.5$$

$$a_1 = f(z_1)$$

$$z_2 = -0.5x_1 + 0.7x_2 - 0.5$$

$$a_2 = f(z_2)$$

$$x_1 \quad x_2 \quad y \quad a_1 : a_2 \quad \hat{y}$$

$$0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

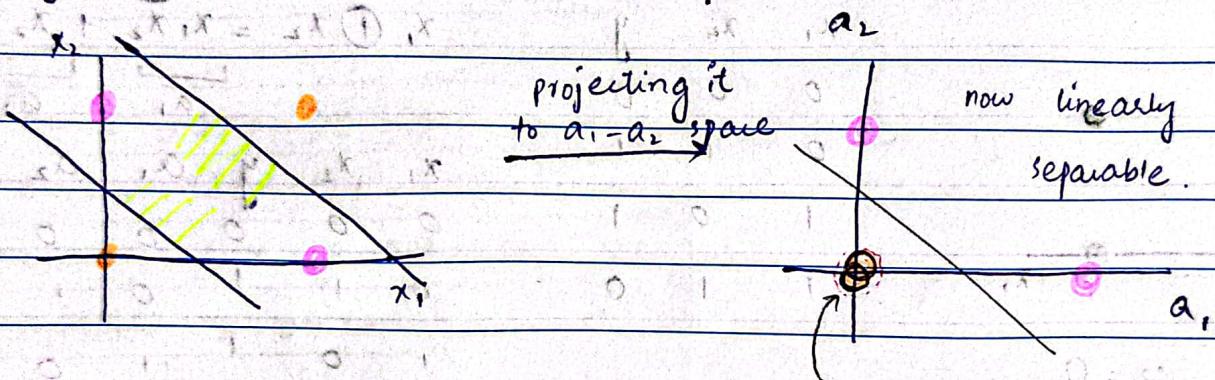
$$0 \quad 1 \quad 1 \quad 0 \quad 1 \quad 1$$

$$1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 1$$

$$1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

for combining the results

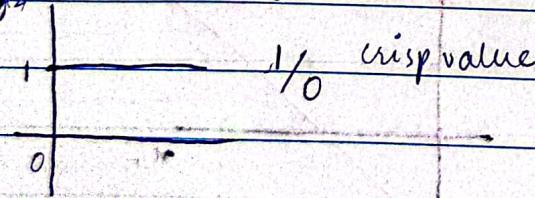
$$\hat{y} = 0.7a_1 + 0.7a_2 - 0.5$$



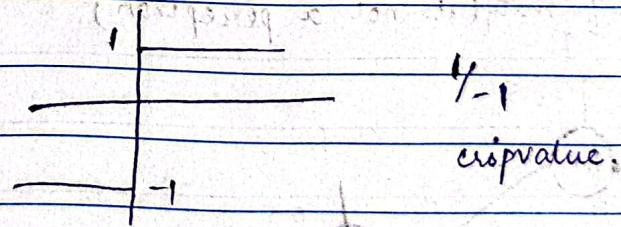
2 points here  
at single point.

### Activation functions:

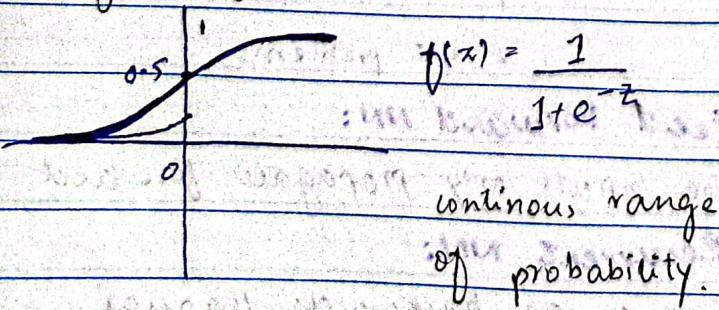
→ step function



→ sine func (Bipolar step)

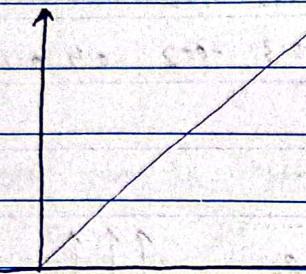


→ sigmoid function





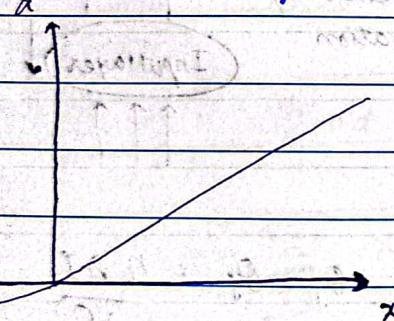
### Rectified linear (ReLU)



$$f(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ x & \text{for } x > 0 \end{cases}$$

→ only applied at the hidden layers. Not at output/input layers.

### Leaky ReLU



$$f(x) = \begin{cases} 0.01x & \text{for } x \leq 0 \\ x & \text{for } x > 0 \end{cases}$$

→ only applied at the hidden layers. Not at output/input layers.

QUESTION

ANSWER

QUESTION

ANSWER

QUESTION

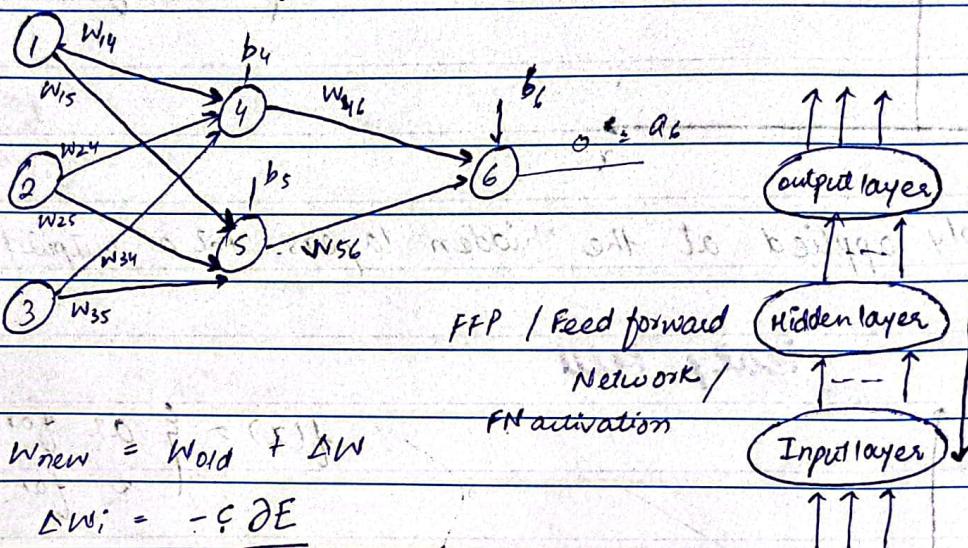
ANSWER

$x_1$	$x_2$	$x_3$	$w_{14}$	$w_{15}$	$w_{24}$	$w_{25}$	$w_{34}$	$w_{35}$	$w_{46}$	$w_{56}$	$b_4$	$b_5$	$b_6$	$d$
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1	0

i

j

k



Delta Rule :-

$$w_{new} = w_{old} + \Delta w$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$$E = \frac{1}{2} \sum_k (d_k - o_k)^2 = \frac{1}{2} (d - o)^2$$

$$o_j = o_j + \eta \frac{\partial E}{\partial o_j}$$

$$z_4 = w_{14}x_1 + w_{24}x_2 + w_{34}x_3 + b_4 = -0.7$$

$$a_4 = \sigma(z_4) = 0.331$$

$$z_5 = w_{15}x_1 + w_{25}x_2 + w_{35}x_3 + b_5 = 0.1$$

$$a_5 = \sigma(z_5) = 0.524$$

$$z_6 = w_{46}a_4 + w_{56}a_5 + b_6 = -0.104$$

$$o = o_6 = \sigma(z_6) = 0.473$$

$$w_{jk} = w_{46}$$

$$j=4 \quad k=6$$

$$\frac{\partial E}{\partial w_{46}} = \frac{\partial E}{\partial o} \times \frac{\partial o}{\partial z_6} \times \frac{\partial z_6}{\partial w_{46}}$$

$$\frac{\partial o}{\partial z} / \partial z$$

$$\sigma'(z) = \sigma(z)[1 - \sigma(z)]$$

$$\frac{\partial E}{\partial w_{46}} = -(d - o) \times o(1 - o) \times a_4$$

$$= o[1 - o]$$

$$= -(0 - 0.473) \times 0.473[1 - 0.473] \times 0.331$$

$$\boxed{\theta = \text{zero}}$$

$$= 0.839$$

$c = 0.1$   
learning rate

$$\begin{aligned} w_{46\text{-new}} &= w_{46\text{-old}} + \Delta w_{46} \\ &= -0.3 - (0.1 \times 0.039) \\ &= -0.3039 \end{aligned}$$

$w_1^1, w_1^2$

① ①  
② ②  
③ ③

$$w_{56\text{-new}} = w_{56\text{-old}} + \Delta w_{56} \quad \text{updated weights}$$

$$\begin{aligned} \frac{\partial E}{\partial w_{56}} &= \frac{\partial E}{\partial o} \times \frac{\partial o}{\partial z_6} \times \frac{\partial z_6}{\partial w_{56}} \\ &= -(d - o) \times o(1-o) \times a_5 \times \frac{\partial z_6}{\partial w_{56}} \\ &= 0.0617 \end{aligned}$$

$w_{46} = -0.3039$

$w_{56} = -0.20617$

$b_6 = 0.088$

$$w_{56\text{-new}} = -0.2 - (0.1 \times 0.0617) = -0.20617$$

$$b_6\text{-new} = b_6\text{-old} + \Delta b_6 \quad \Delta b_6 = -c \frac{\partial E}{\partial b_6}$$

$$\begin{aligned} \frac{\partial E}{\partial b_6} &\approx \frac{\partial E}{\partial o} \times \frac{\partial o}{\partial z_6} \times \frac{\partial z_6}{\partial b_6} \times (o - 1)o \times (o - b) \\ &= -(d - o) \times o(1-o) \times 1 \\ &= 0.117 \end{aligned}$$

$$\begin{aligned} b_6\text{-new} &= 0.1 - (0.1 \times 0.117) \\ &= 0.0883 \end{aligned}$$

use old values  
for one whole  
iteration.

hidden layer :-

$$w_{14\text{-new}} = w_{14\text{-old}} + \Delta w_{14}$$

$$\Delta w_{14} = -c \frac{\partial E}{\partial w_{14}}$$

Backward Error  
propagation

$$\frac{\partial E}{\partial w_{14}} = \frac{\partial E}{\partial o} \times \frac{\partial o}{\partial z_6} \times \frac{\partial z_6}{\partial a_4} \times \frac{\partial a_4}{\partial z_4} \times \frac{\partial z_4}{\partial w_{14}}$$

$$= -(d - o) \times o(1-o) \times w_{46} \times a_4 (1-a_4) \times x_1$$

$$= -(0 - 0.473) \times (0.473) [1 - 0.473] \times -0.3 \times 0.331 [1 - 0.331] \times 1$$

$$\frac{\partial E}{\partial w_{15}} = -0.0078$$

$$(0.3 \times 0.1 \times 0.0078) = -0.0078$$

$$E = -0.0078$$

$$w_{14-010} = 0.2 - (0.1 \times -0.0078) = 0.20078$$

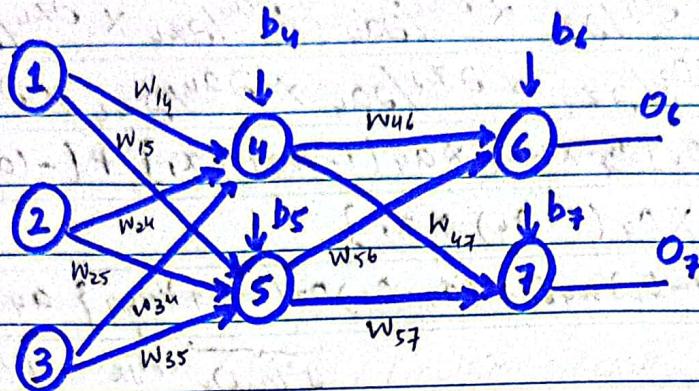
$$b_4-\text{new} = b_{4,010} + \Delta b_4$$

$$\frac{\partial E}{\partial b_4} = -(d-0) \times (0)(1-0) \times w_{46} \times a_4 \times (1-a_4) \times 1 = -0.0078$$

$$\Delta b_4-\text{new} = -0.4 \times (0.1 \times -0.0078) = -0.39922$$

$$\begin{aligned} \frac{\partial E}{\partial w_{15}} &= \frac{\partial E}{\partial 20} \times \frac{\partial 20}{\partial z_6} \times \frac{\partial z_6}{\partial a_5} \times \frac{\partial a_5}{\partial z_5} \times \frac{\partial z_5}{\partial w_{15}} \\ &= -(d-0) \times 0(1-0) \times w_{56} \times a_5(1-a_5) \times x_1 \\ &= -(0-0.473) \times 0.473(1-0.473) \times 0.2 \times 0.1 \times (-0.1) \times 0.523 \times [(-0.523) \times \\ &\quad 5.882 \times 10^{-5}] \times [0.150 \times (0.85)] \\ &= -0.3 - 0.1(-5.88 \times 10^{-5}) \times 1.0 = 1.0 \end{aligned}$$

$$w_{15,\text{new}} = -0.299$$



$$Z_4 = w_{14}x_1 + w_{24}x_2 + w_{34}x_3 + b_4$$

$$a_4 = \sigma(Z_4) = f(x_4)$$

$$Z_5 = w_{15}x_1 + w_{25}x_2 + w_{35}x_3 + b_5$$

$$a_5 = \sigma(Z_5)$$

$$Z_6 = w_{46}a_4 + w_{56}a_5 + b_6$$

$$o_6 = \sigma(Z_6)$$

$$Z_7 = w_{47}a_4 + w_{57}a_5 + b_7$$

$$o_7 = \sigma(Z_7)$$

### Output layer weights

Error in this case

$$W_{46\text{-new}} = W_{46\text{-old}} + \Delta W_{46}$$

$$E = \frac{1}{2} \sum_k (d_k - o_k)^2$$

$$\Delta W_{46} = -C \frac{\partial E}{\partial W_{46}}$$

$$= \frac{1}{2} (d_6 - o_6)^2 + \frac{1}{2} (d_7 - o_7)^2$$

$$\frac{\partial E}{\partial W_{46}} = \frac{\partial E}{\partial o_6} \times \frac{\partial o_6}{\partial Z_6} \times \frac{\partial Z_6}{\partial W_{46}} = E_6 + E_7$$

$$\frac{\partial E}{\partial W_{46}} = \frac{\partial E}{\partial o_6} / \partial W_{46} = \frac{\partial E_6}{\partial W_{46}} + \frac{\partial E_7}{\partial W_{46}}$$

$$\frac{\partial E}{\partial W_{46}} = \frac{\partial E}{\partial o_7} / \partial W_{46} = \frac{\partial E_7}{\partial W_{46}} = 0$$

$$\frac{\partial E}{\partial W_{46}} = \frac{\partial E_6}{\partial W_{46}} = -(d_6 - o_6) \times o_6 (1 - o_6) \times a_4$$

(For Output layer)

Generalized Equations

$$\Delta W_{jk} = -C \frac{\partial E}{\partial W_{jk}}$$

$$\frac{\partial E}{\partial W_{57}} = -(d_7 - o_7) \times o_7 (1 - o_7) \times a_5$$

$$= +C [-(d_k - o_k) \times f'(act)_k \times x_j]$$

input to

that node

### Hidden layer weights

$$W_{14\text{-new}} = W_{14\text{-old}} + \Delta W_{14}$$

$$\Delta W_{14} = -C \frac{\partial E}{\partial W_{14}}$$

$$\frac{\partial E}{\partial W_{14}} = \frac{\partial E}{\partial o_6} / \partial W_{14} + \frac{\partial E}{\partial o_7} / \partial W_{14}$$

$$\begin{aligned}\frac{\partial E_6}{\partial w_{14}} &= \frac{\partial E_6}{\partial o_6} \times \frac{\partial o_6}{\partial z_6} \times \frac{\partial z_6}{\partial a_4} \times \frac{\partial a_4}{\partial z_4} \times \frac{\partial z_4}{\partial w_{14}} \\ \frac{\partial E_7}{\partial w_{14}} &= \frac{\partial E_7}{\partial o_7} \times \frac{\partial o_7}{\partial z_7} \times \frac{\partial z_7}{\partial a_4} \times \frac{\partial a_4}{\partial z_4} \times \frac{\partial z_4}{\partial w_{14}} \\ \frac{\partial E}{\partial w_{14}} &= [-(d_6 - o_6) o_6 (1-o_6) w_{46} \times a_4 (1-a_4) \times x_1] + [-(d_7 - o_7) \\ &\quad o_7 (1-o_7) \times w_{47} \times a_4 (1-a_4) \times x_1] \\ &= [-\{(d_6 - o_6) o_6 (1-o_6) w_{46} + (d_7 - o_7) o_7 (1-o_7) w_{47}\} a_4 (1-a_4) x_1]\end{aligned}$$

$\underbrace{w_{jk}}_{f'(act)_k} \underbrace{f'(act)_j}_{f'(act)_i} \underbrace{x_i}_{z_i}$

**Generalized Equation for weights of hidden layer (single)**

$$\Delta w_{ij} = -c \frac{\partial E}{\partial w_{ij}} = -c \left[ -\sum_k \delta(d_k - o_k) f'(act)_k w_{kj} \right] f'(act)_i x_i$$

- $\Delta b_6 = -c \left[ -(d_6 - o_6) o_6 (1-o_6) \times 1 \right]$

- $\Delta b_4 = -c \left[ -\delta(d_6 - o_6) o_6 (1-o_6) w_{46} + (d_7 - o_7) o_7 (1-o_7) w_{47} \right] a_4 (1-a_4)$

### Gradient Descent Algorithm

(Used during backward propagation to update weights)

Batch Gradient

Stochastic

Mini-batch Gradient

Descent

Gradient Descent (SGD)

Descent

(whole batch)

(only 1 updated)

(a specific size of batch)

$E \rightarrow E$   $\rightarrow$  Stopping condition  
threshold

or for a specific # of times.

$2+2+(3\times 2)+2(2\times 2)=14$  times weight updated in prev example  
biases

**Epoch:** entire training set is passed through the model, forward & backward propagation are performed & parameters are updated.

**Overfitting:** gives accurate prediction for training data but not for new data

**Underfitting:** model cannot determine meaningful relationship b/w input & output data. (performs bad on both training & testing data).

**F1-score:** harmonic mean of precision & Recall

$$\text{F1-score} = 2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall})$$

0 - 1. ↑ better

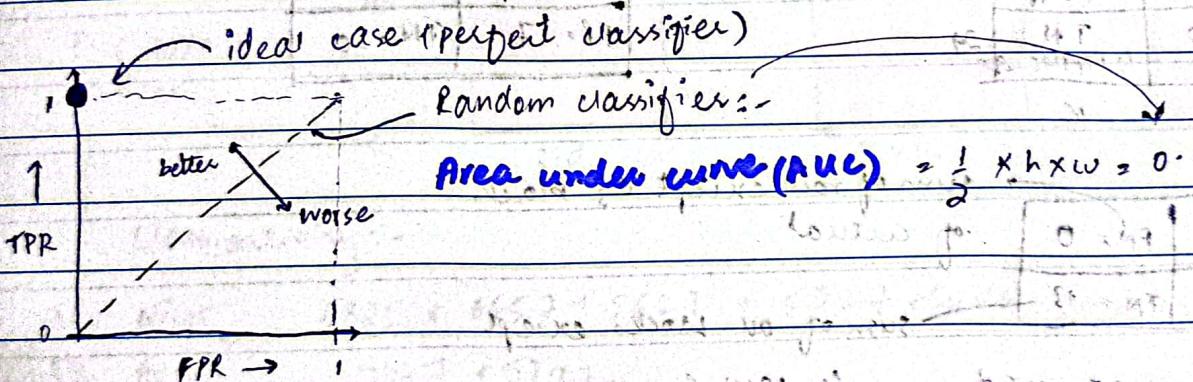
**Recall:** True positive Rate =  $TP / (TP + FN)$

**Specificity:** True negative Rate =  $TN / (TN + FP)$

**False positive Rate:**  $1 - Sp = FP / (FP + TN)$

0 - 1. ↓ better

**ROC curve:** Receiver operating characteristic curve



**Accuracy:** correctly classified

**Misclassification Rate (error):**  $1 - \text{Accuracy}$ .

In case of multi-classes:  $N \times N$

pre	$C_1$	$C_2$	$C_3$
$C_1$	1	1	1
$C_2$	1	2	0
$C_3$	0	1	1

sample	y	$\hat{y}$
1	2	2
2	1	2
3	1	3
4	3	3
5	2	1
6	3	2
7	1	1
8	2	2

Actual      True positives.

$$\text{Accuracy} = \frac{4}{8} = 0.5$$

pred

	3	7	4
actual	0	1	17
	0	0	11

$$Acc = \frac{5+11+1}{5+7+4+1+17+11} = \frac{17}{45}$$

False Negative = // First we need class wise values.

$$\text{Recall} = \frac{TP}{TP + FN}$$

Class-wise calculation:-

C1 :- seosa vs all.

	pred	
actual	TP = 5	FN = 11
-ve	FP = 0	TN = 29
	$\frac{5}{5+11+0+29}$	

C2 :- versicolor vs all

	pred	
actual	TP = 10	FN = 17
-ve	FP = 7	TN = 20
	$\frac{10}{10+17+7+20}$	

C3 :- sum of row except C3 block

	pred	
actual	TP = 11	FN = 0
-ve	FP = 21	TN = 13
	$\frac{11}{11+0+21+13}$	

sum of col of pred  
except C3 block

pre1, recall1, pre2, ...

then average of all this "for overall model's values":-

**Macro. average** :-  $Pre_1 + Pre_2 + Pre_3 \dots Pre_n$

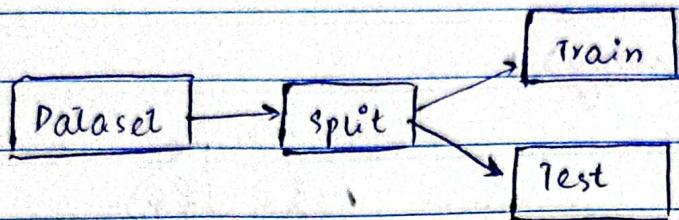
(class imbalance case).

Some weight to all classes, (for balanced data)

Some weight to all samples = **micro. average** :-

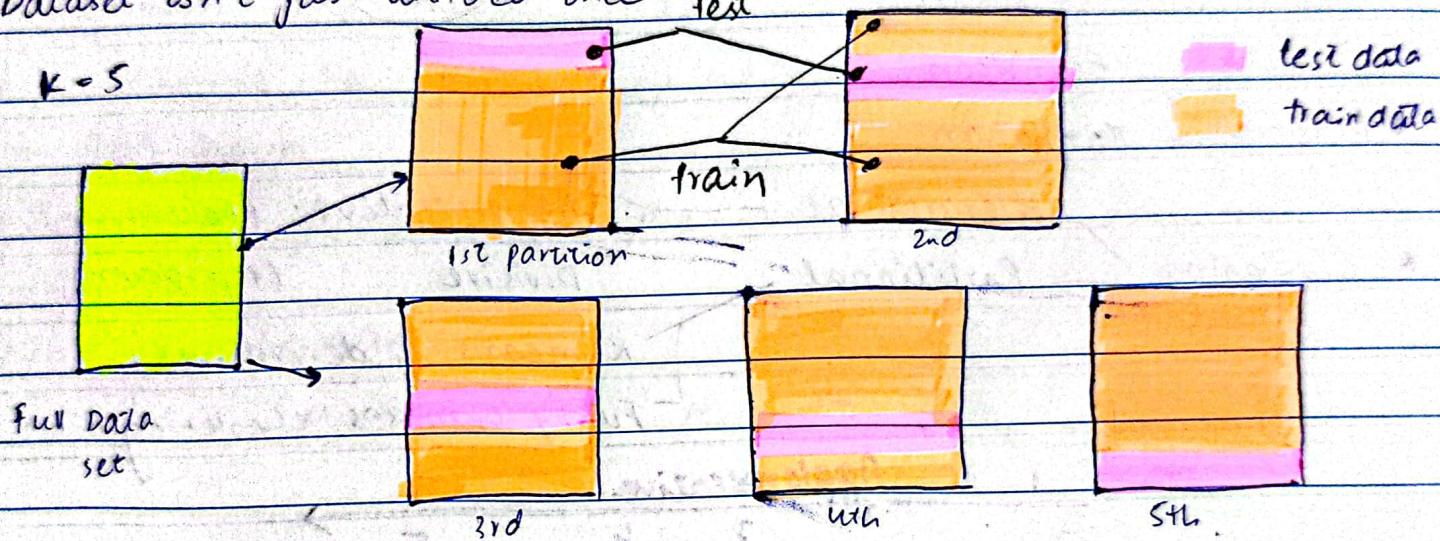
$$Pre = \frac{TP_1 + TP_2 + TP_3}{TP_1 + TP_2 + TP_3 + FP_1 + FP_2 + FP_3}$$

Train test  
split :-



Cross-validation:- K-fold cross-validation.

Dataset isn't just divided once. test



Class activity:-

$$TP_1 \text{ Micro} = 5825 + 6657 + 5627 + 5027 + 5480 + 5095 + 5789 + 5872 +$$

$$TP_2 \text{ Avg} = 5447 + 5727$$

$$TP_3 \text{ Precision} = 57346 + 98485 + 371 + 304 + 362 + 326 + \\ 129 + 393 + 409 + 222$$

$$= \frac{57346}{60000}$$

$$= 0.9557$$

## Clustering (unsupervised learning) :

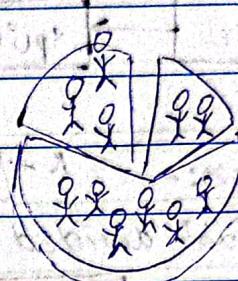
- ① Market Segmentation
- ② Social Network Analysis
- ③ Astronomical Data Analysis
- ④ Organize computer clusters

Euclidean

Distance Manhattan

Hierarchical

Clustering Partitional



Agglomerative (bottom-up)

Divisive (top-down)

K-means & derivatives

Fuzzy c-means clustering

Agglomerative

1 2 3 4 5

a

ab

b

abcde

c

cde

d

de

e

pop + spp + psl

4

Divisive

3

2

1

stop

pop + spp + psl

**K-means algorithm:** cluster  $n$  objects based on attributes

into  $k$ -partitions, where  $k \leq n$ .

↓  
integers.

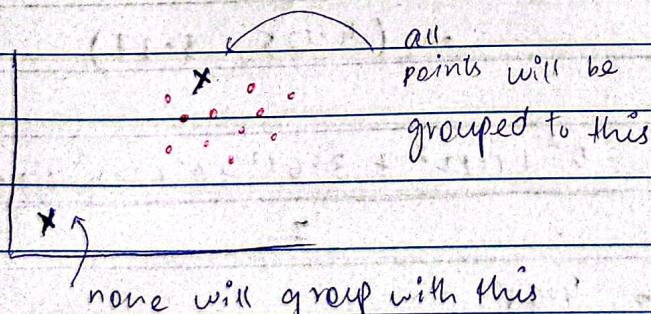
- minimize the sum of square of distances b/w data & corresponding cluster centroid to do grouping.
- First we randomly pick 2 centroids. (will be predefined) # of centroids. calculate distance of each data point from the centroids & group them to the lesser distance centroid point.
- re-calculate the centroids by taking average of data points in that group.
- Repeat this process until the centroid point does not change anymore.

Optimization Objective:-

$$J(c^{(1)}, \dots, c^{(m)}, \mu_{c(i)}) = \frac{1}{m} \sum \|x^{(i)} - \mu_{c(i)}\|^2$$

When doing Random selection:

$$\frac{\sum_i \|dist\|^2}{m} = \frac{1}{m} (\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2})$$



To overcome this issue, we randomly select centroids from the data points themselves.

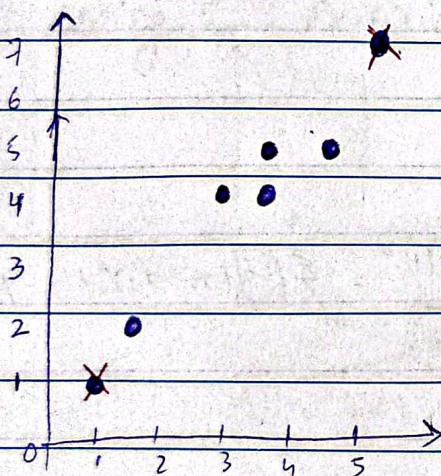
$$\mu_1 = (1, 1)$$

$$\mu_1 = (1.83, 2.33)$$

$$\mu_2 = (5, 7)$$

$$\mu_2 = (4.12, 5.37)$$

x	y	$d_1$	$d_2$	$d_1$	$d_2$	$d_1$	$d_2$
1	1	1	7.2	1.56	2.36	0.55	5.02
2	1.5	2	1.18	6.103	0.46	4.26	0.55
3	3	4	3.605	3.605	2.03	1.76	3.05
4	5	7	7.21	0	5.64	1.85	6.6
5	3.5	5	4.72	2.5	3.14	0.72	4.16
6	4.5	5	5.03	2.06	3.77	0.5302	4.77
7	3.5	4.5	4.3	2.92	2.73	1.068	3.75



$$\mu_{1\text{new}} = \left( \frac{1+1.5+3}{3}, \frac{1+2+4}{3} \right)$$

$$= (1.83, 2.33)$$

$$\mu_{2\text{new}} = \left( \frac{5+3.5+4.5+3.5}{4}, \frac{7+5+5+4.5}{4} \right)$$

$$= (4.125, 5.37)$$

$$\text{Cost/Error of 1st iteration} = \frac{0^2 + 1.012^2 + 3.61^2 + 0^2 + 2.5^2 + 2.06^2 + 2.92^2}{7}$$

$$J = 4.75$$

$$\mu_{1\text{new}} = \left( \frac{1+1.5}{2}, \frac{1+2}{2} \right) = (1.25, 1.5)$$

$$\mu_{2\text{new}} = \left( \frac{3+5+3.5+3.5+3.5}{5}, \frac{4+7+5+5+4.5}{5} \right)$$

$$= (4.9, 5.1) \quad J = 5.58$$

K-mean algorithm  $\rightarrow$  sub-optimal results  $\rightarrow$  because of random initialization of centroids.

k-means++ : picks initial centroids intelligently.

$$u_1 = (1, 1) \quad u_2 = (5, 7)$$

	$x$	$y$	$d_1$	$d_2$
1	1	1	0	7.21
2	1.5	2	1.12	6.1
3	3.0	4	3.61	3.61
4	5	7	7.21	0
5	3.5	5	4.72	2.5
6	4.5	5	5.31	2.06
7	3.5	4.5	4.3	2.92

next centroid will be the one having max distance.