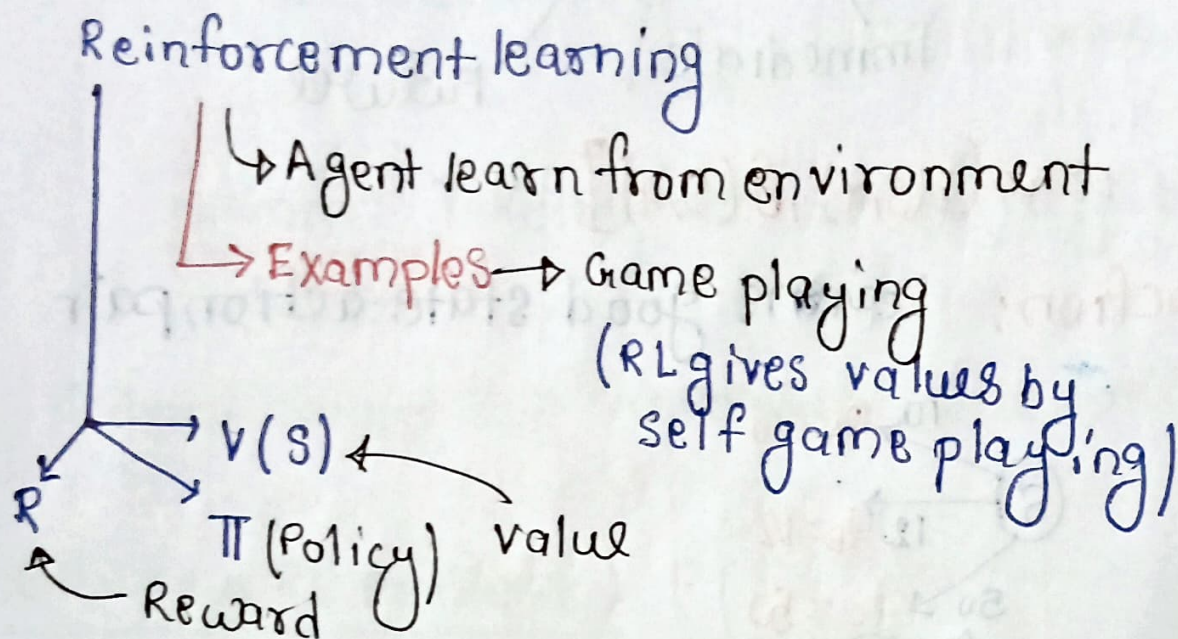


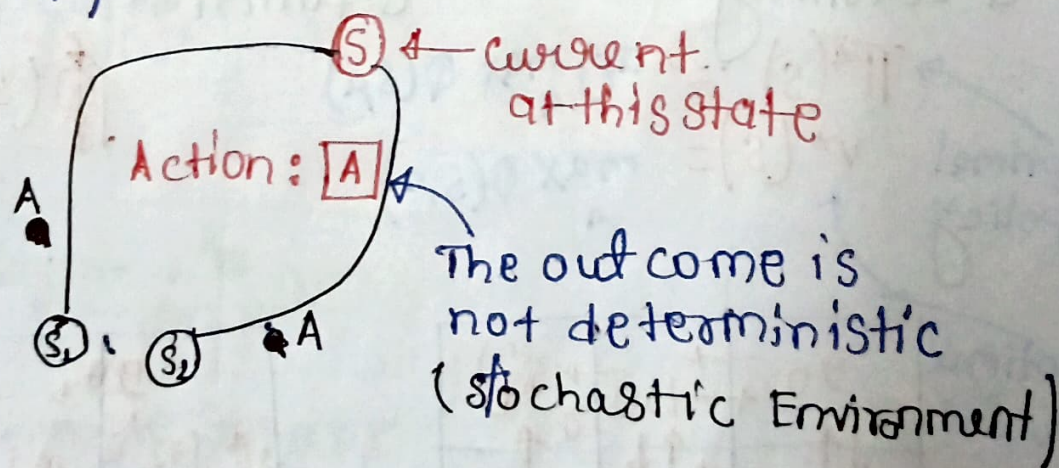
► Reward hacking

► In next class; RL algorithms and policies.

Lecture - 20 (6-November-24) | Part A



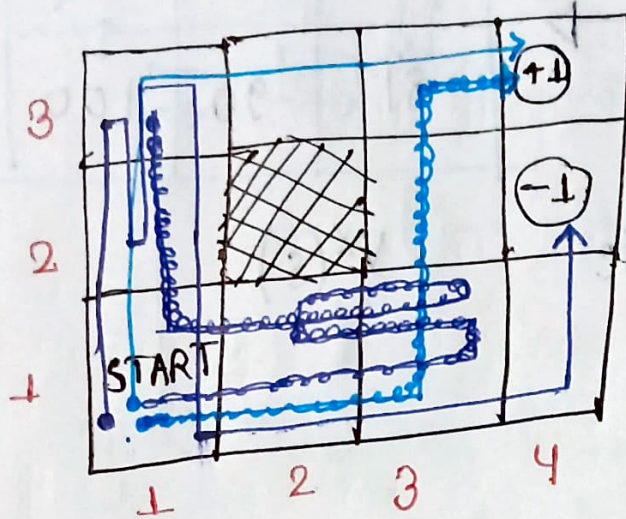
Markov Decision Process (MDP)



Our assumption is that world is deterministic

Passive learning

- Already has some particular policy ' π '
- Agent is trying to evaluate the policy



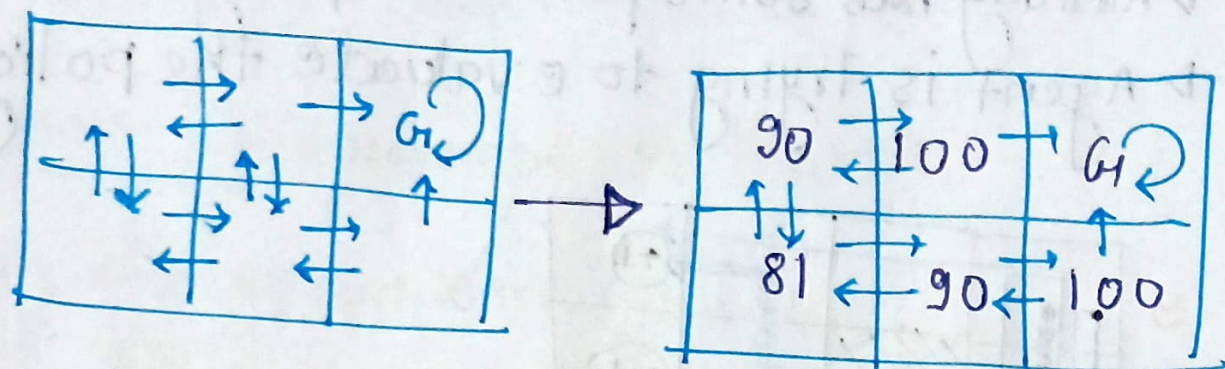
eg. (a) $(1,1) \rightarrow (1,2) \rightarrow (1,3) \rightarrow (2,3) \rightarrow (3,3) \rightarrow (4,3) [+1]$

(b) $(1,1) \rightarrow (1,2) \rightarrow (1,3) \rightarrow (1,2) \rightarrow (1,3) \rightarrow (1,2)$

$\rightarrow (1,1) \rightarrow (2,1) \rightarrow (3,1) \rightarrow (4,1) \rightarrow (4,2) [-1]$

n sequences and try to improve the policy
by evaluating the policy

Epoch



$r(s, a)$: Immediate
reward
values

$V^*(s)$

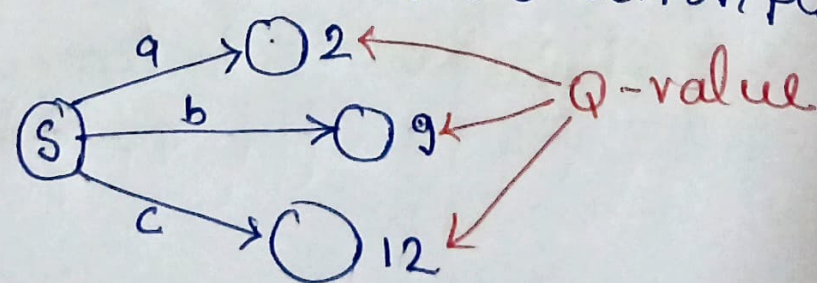
$$V^*(s) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} \dots$$

Assuming $\gamma = 0.9$

▷ Agent does not have idea about complete space state that's why q-learning come into picture

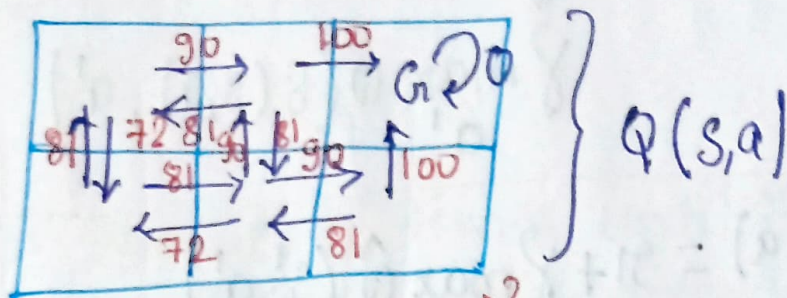
Q-Learning

Q-value : It will give a value/cost for each state-action pair



→ Q-value provides goodness of the state-action pair.

→ Q-learning is how we derive and use Q-value



$$r_t + \gamma r_{t+1} + \gamma^2 r_{t+2}$$

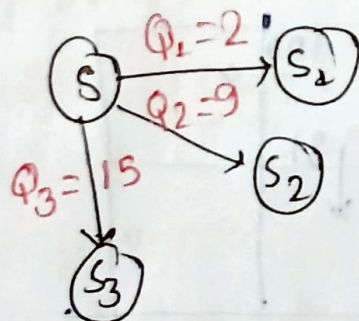
$$\pi^*s = \arg \max_a Q(s,a)$$

$$v^*s = \max_a Q(s,a)$$

Now:

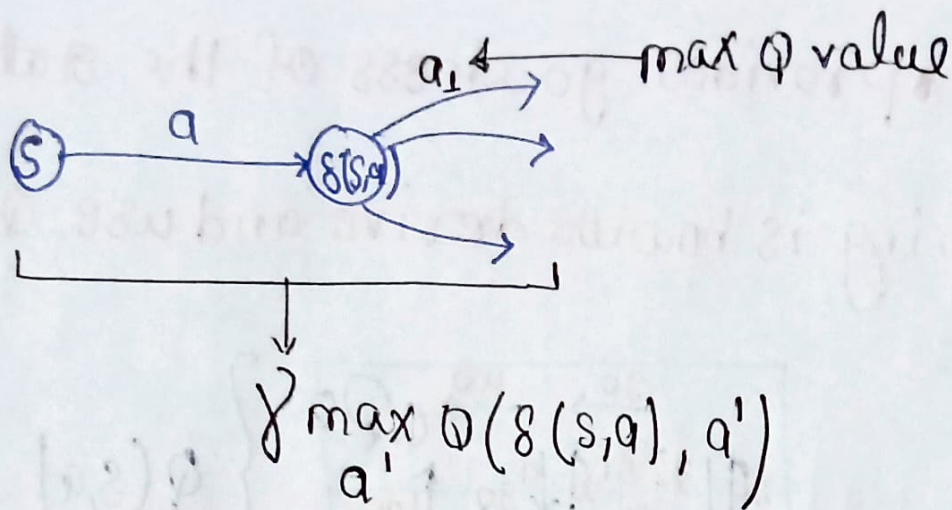
$$Q(s,a) \equiv \boxed{r(s,a)} + \boxed{\gamma v^*(\delta(s,a))} \quad \text{Discount function}$$

$$= \{r(s,a) + \gamma \max_{a'} Q(\delta(s,a), a')\}$$



This still depends on $r(s,a)$ and $\delta(s,a)$

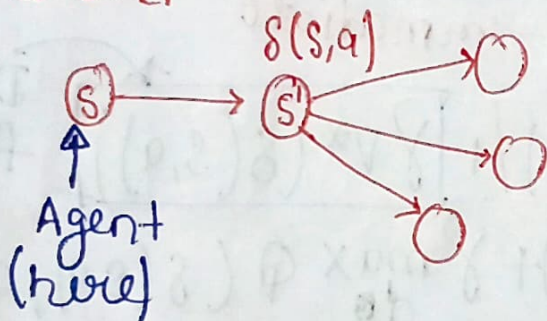
Recursive definition



$$\hat{Q}(s, a) = r + \gamma \max_{a'} \hat{Q}(s', a') \quad \text{where } s' = \delta(s, a)$$

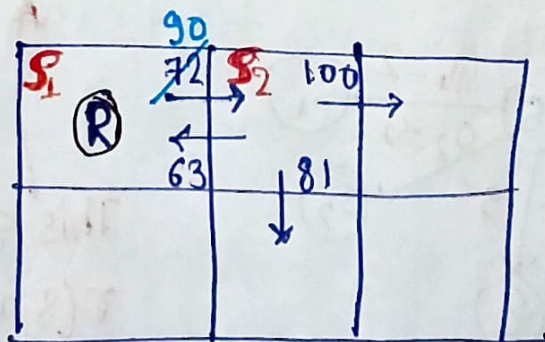
Bellman Equation

► In Q-learning, we are looking 1 step ahead



► Modify $Q(s, a)$ value

► Hyper parameter



initial state: s_1

$$\hat{Q}(s_1, \text{Right}) \leftarrow r + \gamma \max_{a'} \hat{Q}(s_2, a')$$

γ : Gamma

$\gamma = 0.9$ (Discount factor)

$$\begin{aligned} & 0 + 0.9 \times \max(63, 81, 100) \\ &= 0.9 \times 100 \\ &= \boxed{90} \end{aligned}$$

▷ We will do 'n' number of times until the values are converging (reaching some values that does not change)

▷ Start with some initial Q-values

similar to Local Search

↳ ~~this~~ We might get stuck with some local optima (sub-optimal value)

We could use something like this

$p(a|s) \propto$

→ Simulated Annealing

→ Prioritized sweeping

→ Temporal Difference Learning

} Optimizations

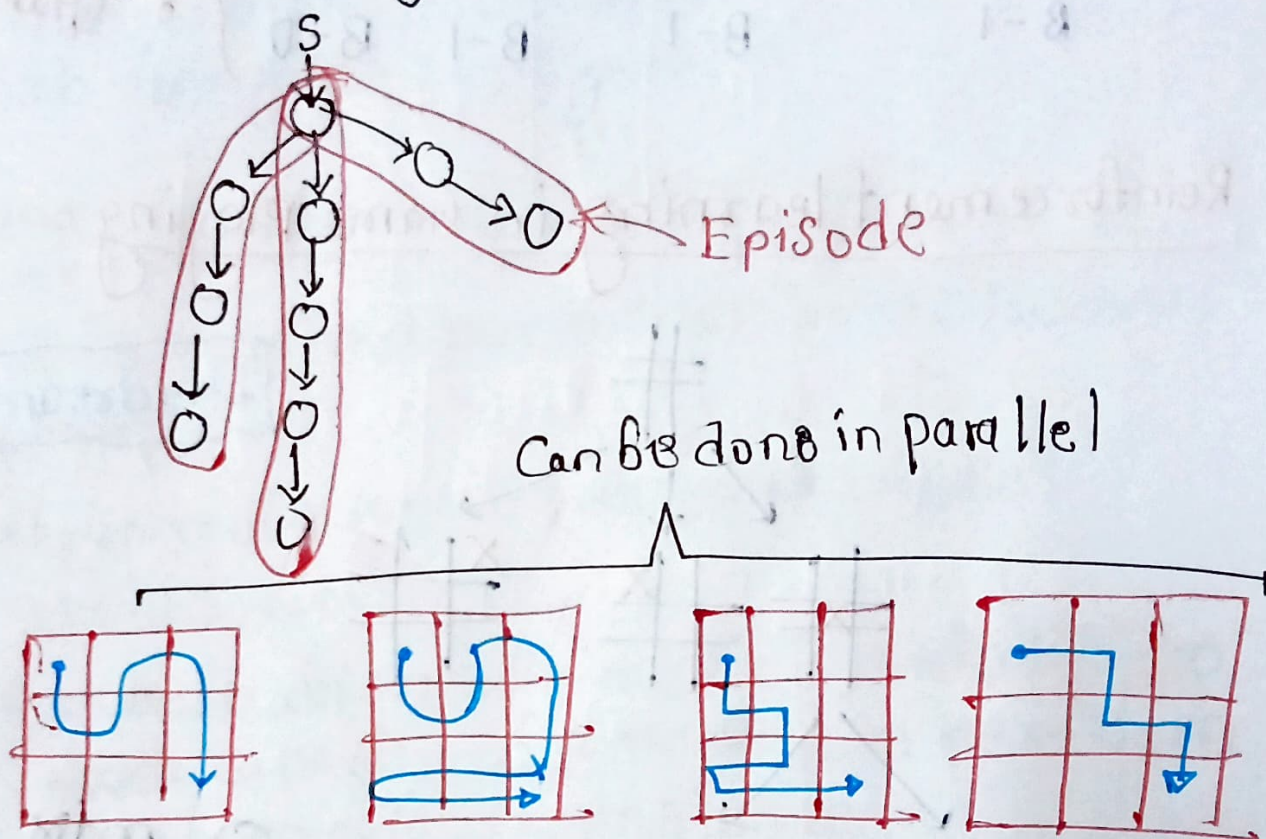
▷ Deep Reinforcement learning

▷ Probabilities for stochastic environment

Lecture-20 (6/11/24) / Part-B

Monte Carlo Policy Evaluation

▷ We devise policy evaluation to different episode.

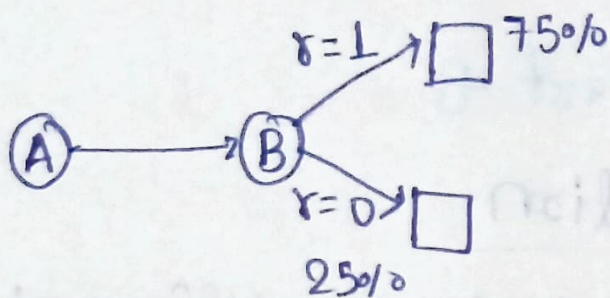


▷ simulated environment is used.

▷ we update Q value after each episode.

Temporal Difference learning

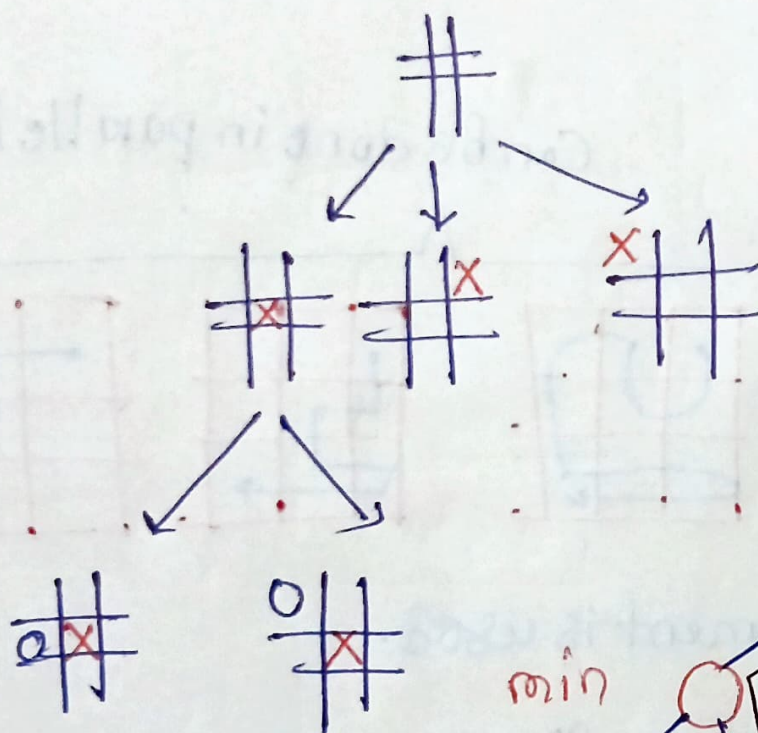
- In Q learning we check only one step ahead
- In Temporal difference learning, agent check upto 1 episode and update Q -value



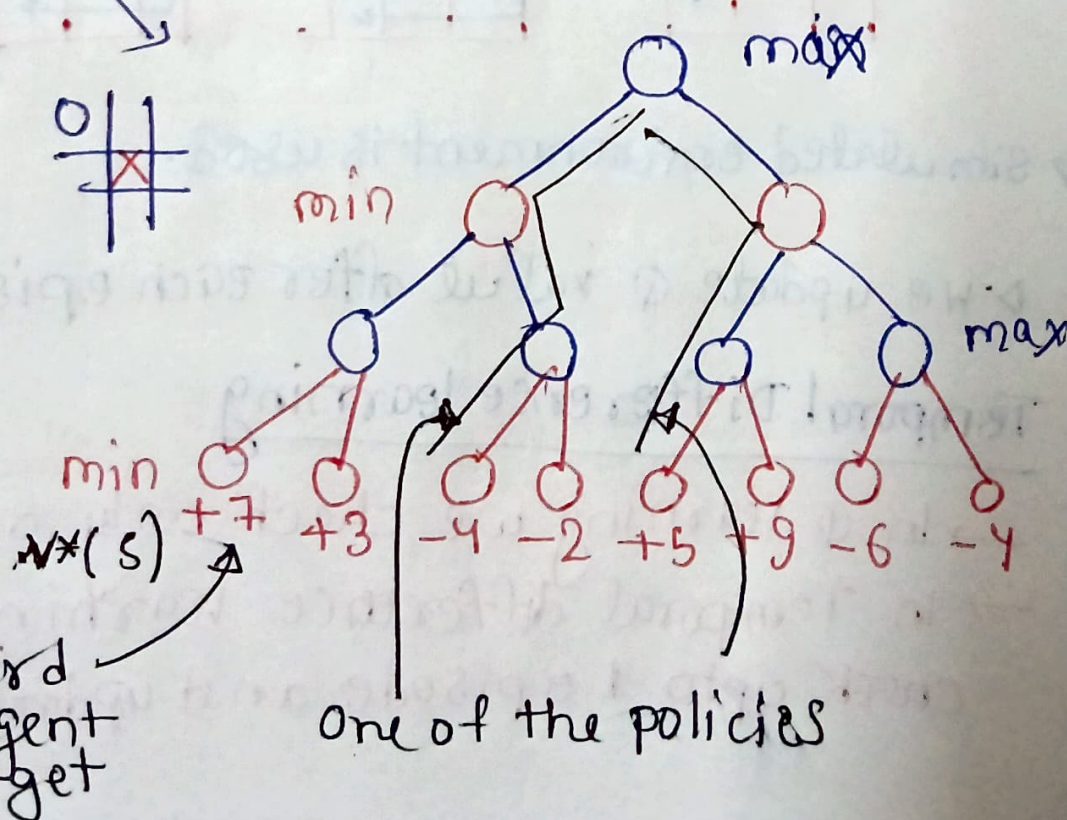
$A=0, B=0 \quad B=1 \quad B=1 \quad B=1$
 $B=1 \quad B=1 \quad B=1 \quad B=0$

} 8-episodes

Reinforcement learning in Game Playing



Daydreaming



▷ How to learn and survive in the environment?

▷ Q learning

▷ Markov Decision process (MDP)

▷ Optimization

▷ Real world example

▷ Expected question: How to update Q-values ~~and~~ given some grid?

Summary

Fundamentals About Neural Network

→ It is inspired by neurons in the human brain

→ Capacity of a model is ability to model varied functions

→ Generalization of a model is ~~ability to~~ ~~model~~ capacity of a model to perform well on previously unseen data

[Underfitting
Overfitting]

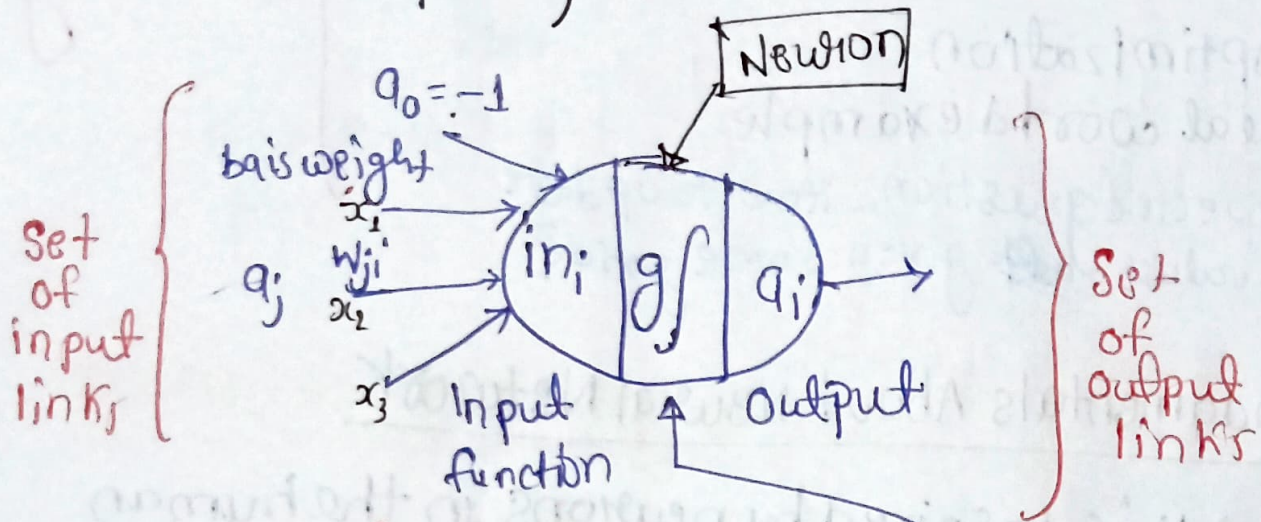
→ Low capacity: Unable to obtain low error on training data

{ Train Data
Test Data }

→ Learn very specific patterns
train ~~data~~ error is very low but test error is quite high comparatively

Simple Neural Network

(Perceptron)



→ Set of nodes (neurons/units)

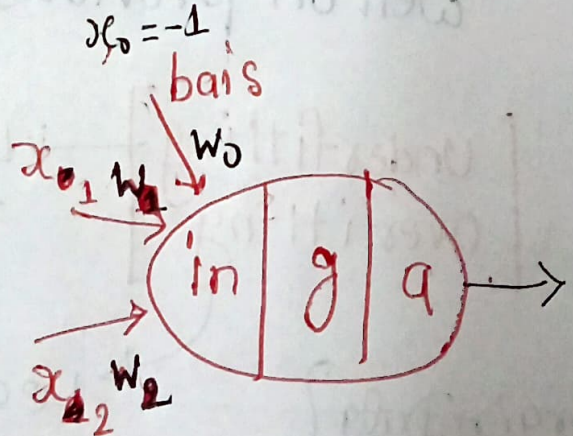
→ ~~computation~~ connection

→ Activation function

Activation function

OR-Gate

x_1	x_2	output
0	0	0
0	1	1
1	0	1
1	1	1



$$in = \sum_{i=0}^2 w_i x_i \quad (\text{weighted function for input value})$$

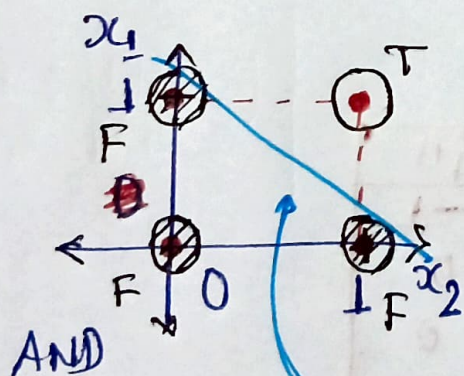
$$a = g(in)$$

$$a = \begin{cases} 0 & \text{if } in \leq 0 \\ 1 & \text{if } in > 0 \end{cases}$$

$$in = -1 + w_1 x_1 + w_2 x_2 \quad (\omega_0 = 1, \omega_1 = 2, \omega_2 = 2)$$

$$= -1 + 2x_1 + 2x_2$$

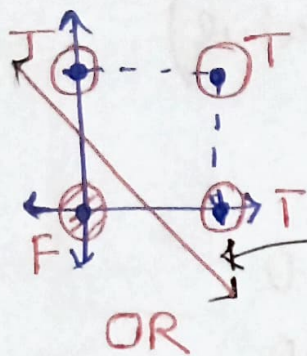
$$in = \begin{cases} -1 & \text{if } (x_1 = 0, x_2 = 0) \longrightarrow 0 \\ 1 & \text{if } (x_1 = 1, x_2 = 0) \longrightarrow 1 \\ 1 & \text{if } (x_1 = 0, x_2 = 1) \longrightarrow 1 \\ 3 & \text{if } (x_1 = 1, x_2 = 1) \longrightarrow 1 \end{cases}$$



○ : True
 ⊗ : False

levels

we want to draw some line so we can discriminate this two levels



can be represented using some linear function

AND - gate

x_1	x_2	output
0	0	0
0	1	0
1	0	0
1	1	1

$$a = \begin{cases} 0 & \text{if } in \leq 0 \\ 1 & \text{if } in > 0 \end{cases}$$

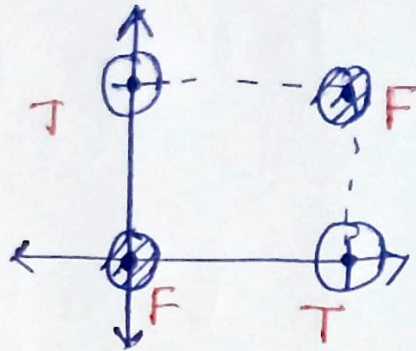
$$\omega_0 = -1, \omega_1 = 1, \omega_2 = 1$$

$$in = -1 + x_1 + x_2$$

x_1	x_2	in
0	0	-1
0	1	0
1	0	0
1	1	1

X-OR Gate

x_1	x_2	output
0	0	0
0	1	1
1	0	1
1	1	0



No line can be drawn for distinguishing between 1 & 0 levels