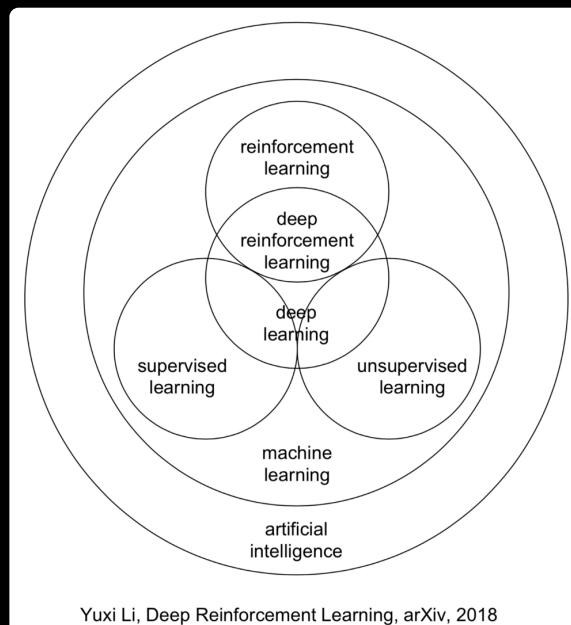
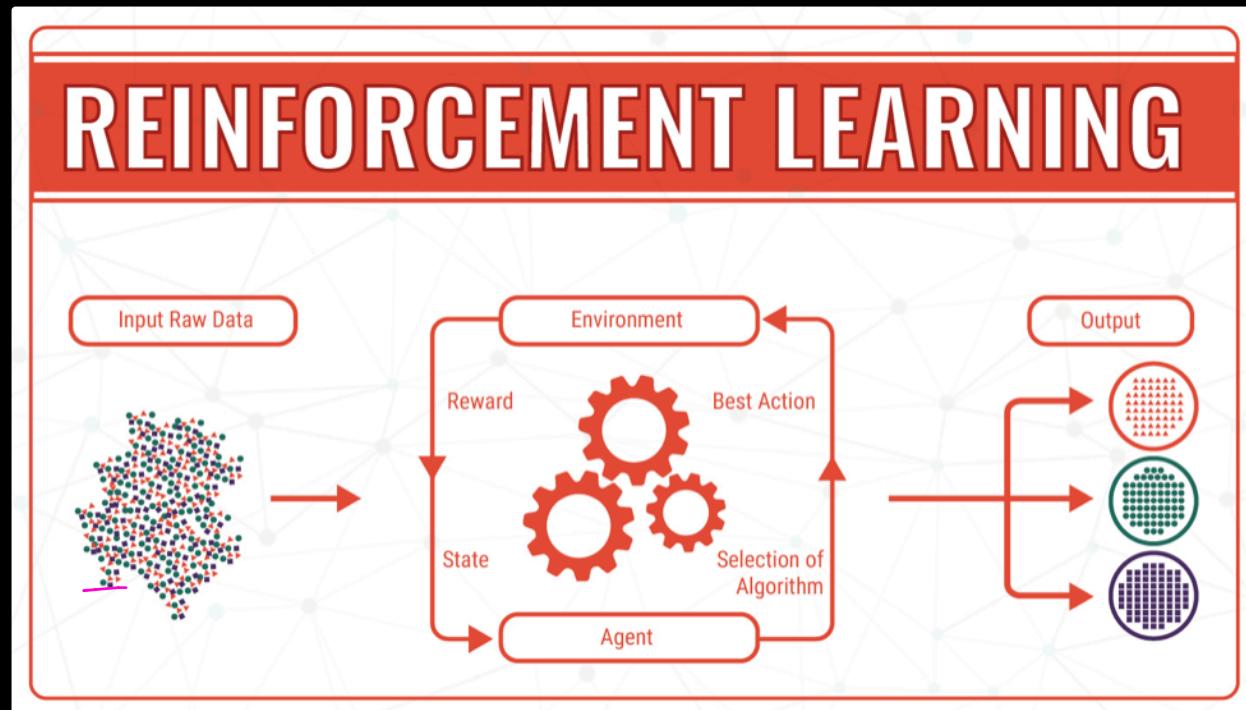


This lecture is completely adapted from Stanford lecture series : CS-231n by Serena Yeung



Yuxi Li, Deep Reinforcement Learning, arXiv, 2018

# A → Introduction to learning paradigms

These are ③ major paradigms of m/c learning

## I) Supervised Learning

→ i) Data:  $(x, y)$   
  data      label

→ ii) Goal: Learn a functional  
  Can map  
 $(x \rightarrow y)$

Examples: Classification (Cat vs Dog),  
Object detection, Semantic segmentation,  
Regression  $[f(x) : y]$ , image  
Caption.

Anything in which  
target is known/  
given at least  
for training

## II) Unsupervised Learning

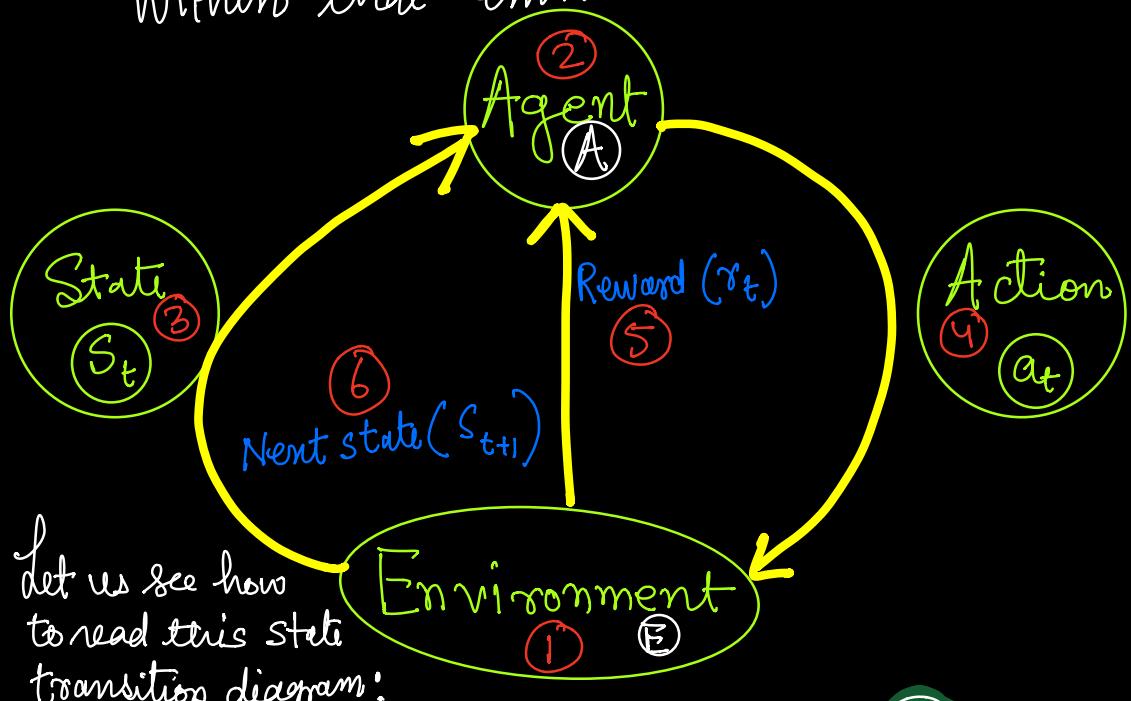
→ i) Data:  $(x)$  w/o  
  labels

→ ii) Goal: learn some  
  underlying hidden  
  relationship between  
  the data that can  
  be used to explore  
  the structure of data

Examples: Clustering,  
dimensionality reduction,  
feature learning, density  
estimation.

In the absence of  
target underlying  
data probability  
distributions can be  
estimated and used for  
various purposes. Such  
Generative models.

③ Reinforcement learning : It is a learning paradigm in which learning happens by exploration. So agent interact with the environment repeatedly w/o any prior knowledge and labels w.r.t the environment. — Completely relying on hit and trial strategy to learn how to behave optimally within that environment.



Let us see how to read this state transition diagram:

\* The agent  $A$  is at some state  $S_t$  at time  $t$ .

→ It is interacting with Environment  $E$ , via. few defined set of actions  $A$ .

Action at any state is chosen only by utilizing the state  $S_t$ . It assumes that  $S_t$  completely characterizes the environmental state. (Req. No history) Markov

\* At any state  $(S_t)$ , taking a particular action  
say  $(a_t)$  is going to fetch a reward  $(R_t)$ .

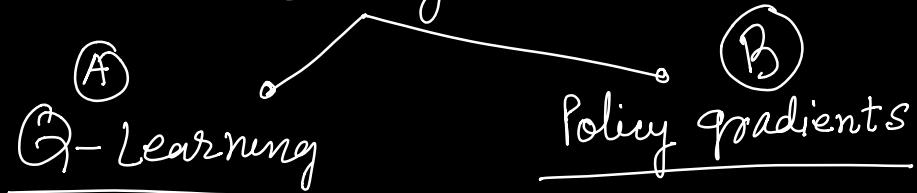
\* The final Goal of the agent is :

To learn how to take an appropriate (optimal)  
action so as to maximize the reward.

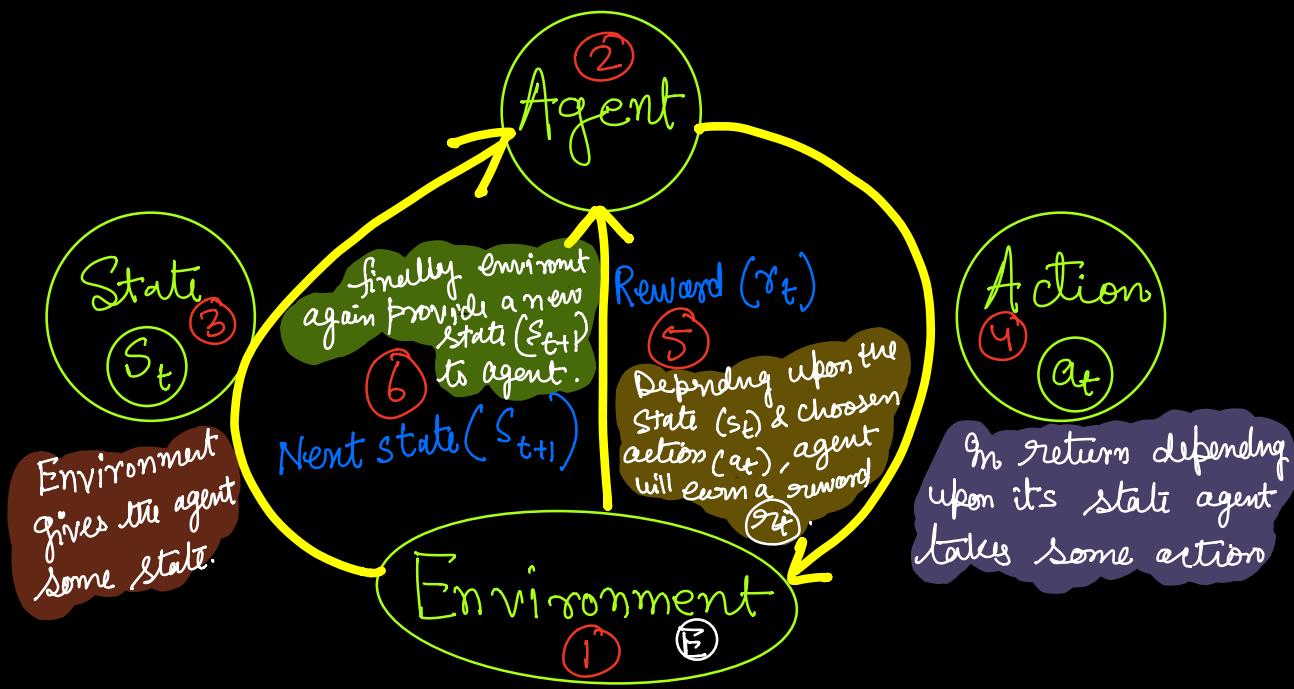
(Future discounted reward  
maximization)

\* Reinforcement learning can be formulated as  
Markov decision process  $\Rightarrow$  MDP.

\* There are 2 major classes of RL algorithms



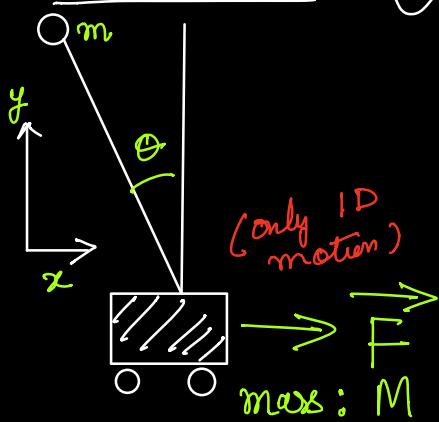
## B → Basic RL-Setup/framework



Episodic training : These six steps keeps on repeating in a loop and considered as an episode. The episode continues until the environment provides a terminal state (final) to the RL-Agent, after which the current episode get terminated.

\* Hence one need to generate data (Episodes) by allowing agent to explore the environment in episodes. later use these episodic data for training the agent to learn how to behave optimally wrt the environment.

### Example-01 : Cart Pole Problem.



(A) Objective: To balance the pole on the top of a movable Cart.

(B) State: The physical state can be encoded as: angle ( $\theta$ ), angular speed ( $\omega$ ), Position ( $x, y$ ), horizontal velocity ( $v_x$ )  
(may be more)

(C) Action: Horizontal force vector applied on the Cart.

(D) Reward:  $+1$  each time, if the pole is upright.

### Example-02 : Atari game playing agent.

(A) Objective: Finish the game with the highest score.

(B) State: Game state can be raw pixel input of the game screens. (Images of game screen)

(C) Action: Game controlling actions, (L/R/U/D).

(D) Reward: At the end of the episode, if win then all  $\langle \text{state}, \text{action} \rangle$  pair will get a  $+1$  reward Else they will get  $-1$  reward.

## C → Mathematical formulation of RL.

- \* RL Can be mathematically formulated using Markov Decision processes (MDP).

### C.1 Markov Property:

The future is independent of the past given the present.

Def: A state  $(S_t)$  is Markov

iff  $P[S_{t+1} | S_t] = P[S_{t+1} | S_1, S_2, \dots, S_t]$

- \* State Captures all relevant information from the history. State is a sufficient statistic of the future.

C.2

## Markov Decision Process

MDP is a decision process, choosing action at each state. It is an environment in which all states are Markov. The MDP can be characterized by tuple:



C. 2.1

## State transition Matrix ( $P$ )

Markov process/chain is a memory less random process of sampling iteratively as per the given STM ( $P$ ), starting from some given (seed) random starting state and following Markov property.

Full game dynamics can be encoded within ( $P$ ). Robot can be interacting with with the gaming environment following  $\textcircled{P} \Rightarrow$  The game rules.

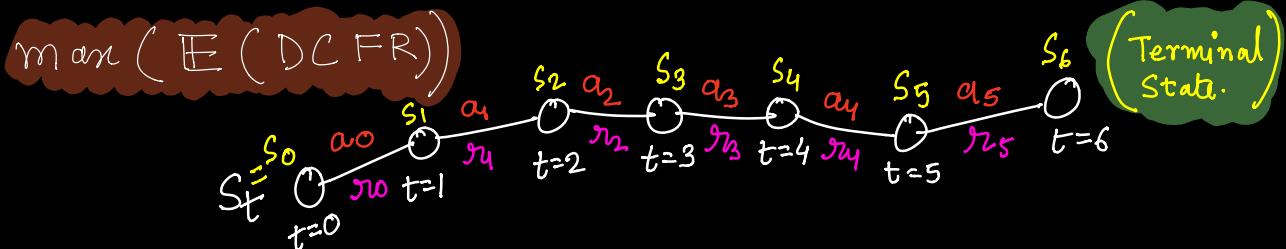
$$\underset{\text{SS'}}{\overset{\text{P}^a}{\leftarrow}} = P[S_{t+1} = S' \mid S_t = S, a_t = a]$$

## C. 2.2) Immediate reward

Amount of reward an agent is going to get from the environment when at state  $s_t$  agent takes an action  $a_t$ .

$$\text{👉 } r_{s_t}^{a_t} = r(s_t, a_t)$$

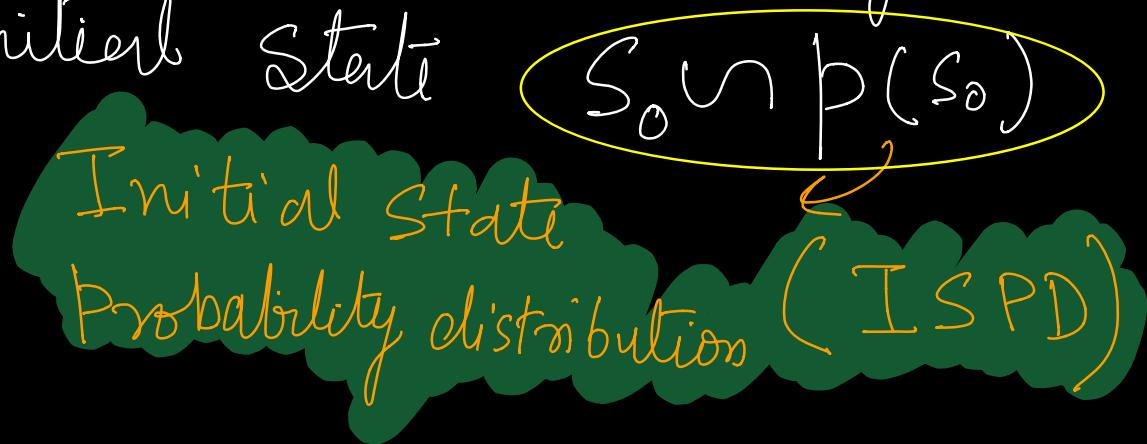
✳️ But in MDP, we are not interested in immediate rewards. An optimal behavior must need to maximize the Average Discounted Cumulative Future Reward (DCFR)



How MDP operates and can be used to obtain Episodes.

MDP-Execution

- ① At time ( $t=0$ ), environment may sample randomly an initial state  $s_0 \sim p(s_0)$



- ② from  $t=0$  (until terminates)

- ②.1 An agent can select an action ( $a_t$ ) [How!!!]  
Using some policy that we need to figure out so as to  $\max(\mathbb{E}(\text{DCF}))$

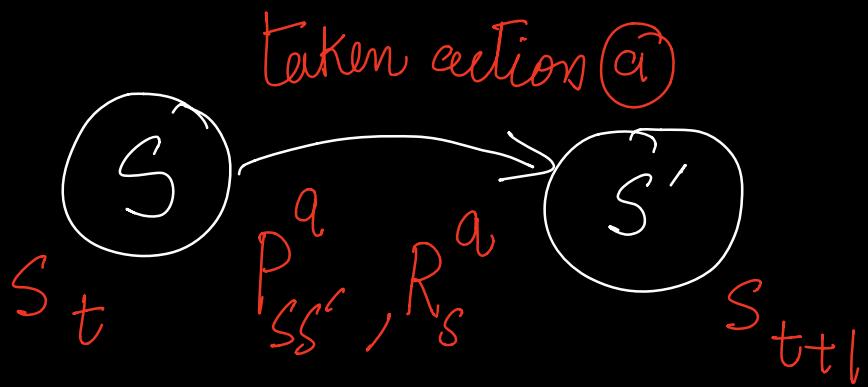
Q.2) Depending upon the agent's state  $S_t$  and the action taken  $a_t$ , environment will return the immediate reward as well as next state  $S_{t+1}$  using  $R$  &  $P$ .

∴ for a MDP  $P$  &  $R$  governs.

$$\Rightarrow P_{ss'}^a = P[S_{t+1} = s' \mid S_t = s, A_t = a]$$

$$\Rightarrow R_s^a = \mathbb{E}[r_t \mid S_t = s, A_t = a]$$

Given a discount factor  $\gamma \in [0, 1]$



(C. 2.3) But how to take an action  $\Rightarrow$  Policy  $\pi$ )

A Policy  $\pi$  is just a mapping of states  $S_t$  to actions  $a_t$  enabling any agent to take action

$$R^{\text{State}} \rightarrow R^{\text{action}}$$

\* It can be deterministic or stochastic (Random). Defined

$$\text{as } \pi(a_t | s_t) = P[a_t | s_t]$$

### C. 2.4 Goal of MDP

The objective of MDP is to find an optimal policy ( $\pi^*$ ), that can maximizes the DCFR.

\* The actual returns are always in the form of DCFR.

$$G_t = \text{Discounted Cumulative future Reward from time Step } t. \quad [ \begin{matrix} \text{future rewards are} \\ \text{geometrically weighted} \end{matrix}]$$

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

for all upcoming future time step

$$= \sum_{K=0}^{\infty} \gamma^K R_{t+K} (\gamma \in (0,1))$$

{Gamma}



→ close to ① "Myopic" evaluation.

→ close to ② "far-sighted"

More important to immediate returns.

Care about all future rewards

### Requirement of discount :

ⓐ future is uncertain, and we care about  $\gamma$   
 If I invest 1K-INR (Today)  $\Rightarrow$  then how do I need to care  
 about my return after (1yr), (5yr), (10yr) ...

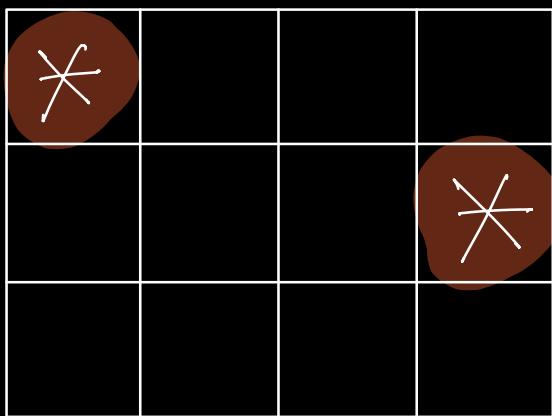


ⓑ Since we are predicting using an inaccurate model of environment hence the trust over model predictions can decrease exponentially. [Reason]

- (c) Mathematically Convenient to discount rewards.
- (d) Avoids infinite returns in Cyclic Markov processes.

## C.2.5 Example Grid World - MDP

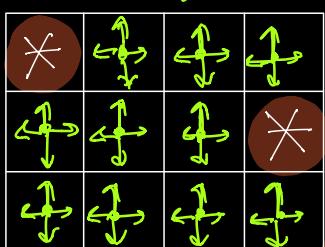
(A) States (2D-locations)      (B) Actions = {right, left, up, down}



(C) Reward = { $r_2 = -1$ }  
negative reward for each transition. It ensures to learn how to reach target fast.

(D) Objective: Starting from any randomly chosen block  $\Rightarrow$  Reach to a terminal state (\*) in the least number of actions.

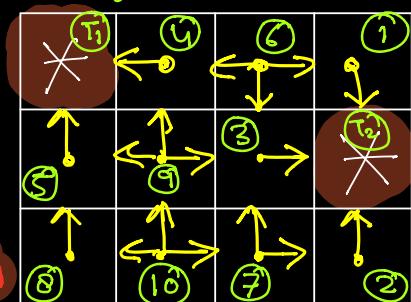
(E) Given the above 2D-Grid world MDP we need to figure out some optimal policy that can minimize the DCF.R. (Out of all the policies  $\pi$ )



Random policy  
( $\pi_1$ )

Randomly choosing any direction.

Visual Policy realization



(B)  
( $\pi_2$ )

Policy philosophy: (a) For immediate neighbours:

One should have to move in a direction that can take you to the terminal state.

(b) Other states: Move in the direction that

can take you closest to one of the immediate neighbouring state.  $\pi(a_t | s_t)$  [Mapping from  $s_t \rightarrow a_t$ ]

## D → Solution of MDP: Optimal Policy ( $\pi^*$ )

Policy  $\approx$  Planning

It does not matter that in whatever state (say  $s_t$ ) agent is at time  $t$ ,  $\pi^*$  is going to suggest that "what action  $(a_t)$  agent needs to take in order to get the DCFR maximized."  $\pi(a_t | s_t)$

\* Such MDP's are realized as stochastic framework with initial probability distribution or transition probability distributions. Hence in order to handle the randomness we always talk about Expected value of the reward maximization.

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} \right]$$

where  $S_0 \sim P(S_0)$

*Enforcing overall best policy*

*Given a policy.*

*The initial state, action & next state got sampled from their respective probability distribution.*

$a_t \sim \pi(\cdot | s_t)$

$s_{t+1} \sim P(\cdot | s_t, a_t)$

④ There are few attributes associated to a policy:

$\pi(a_t | s_t)$  → Value fn - wrt a given policy  $(\pi)$ , How good is a given state  $(s)$ , i.e if agent just follow  $(\pi)$  How much DCFR is expected wrt a given policy  $(\pi)$ , is the value of that state  $(s)$  How good is a state action pair. i.e

Instead of directly following the policy  $(\pi)$  at state  $(s)$ , agent first takes an action  $(a)$  and then follows the policy  $(\pi)$ . Q-value of  $(s, a)$  is the expected DCFR.

D.I

# Value function $[V^\pi(s)]$

Defined as the expected DCFR if agent follows  $\pi$  at state  $s$  until the episode terminates (reaching terminal state)

Given a policy and a initial state  $s_0 \in p(s_0)$ ; Episode is  $(s_0, a_0, r_0) \rightarrow (s_1, a_1, r_1) \rightarrow \dots \rightarrow (s_n, a_n, r_n)$

*Episode carry on till terminal state is reached*

$$V^\pi(s) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = s, \pi \right]$$

Value of state  $s$  wrt a given policy  $\pi$

Expected DCFR from state  $s$  following the policy  $\pi$  to choose future actions.

D.2

## $Q$ -Value function $[Q^\pi(s, a)]$

Let us assume that we have

$$S = \{s_1, s_2, s_3\} \quad A = \{a_1, a_2\} \quad \text{for some } \pi$$

we can have ⑥ state action pairs  $\pi^\star(s_2)$  and state  $s_2$

$s_1, a_1$	$s_2, a_1$	$s_3, a_1$	$\frac{Q^\pi(s_2, a_1)}{Q^\pi(s_2, a_2)}$
$s_1, a_2$	$s_2, a_2$	$s_3, a_2$	$\pi(a_1   s_2) \& \pi(a_2   s_2)$

$Q^\pi(s, a) \Rightarrow$  How good it is to take an

action  $a$  at state  $s$  & then following  $\pi$ .

the expected DCFR, if at state  $s$ ,  $a$  action has been chosen.

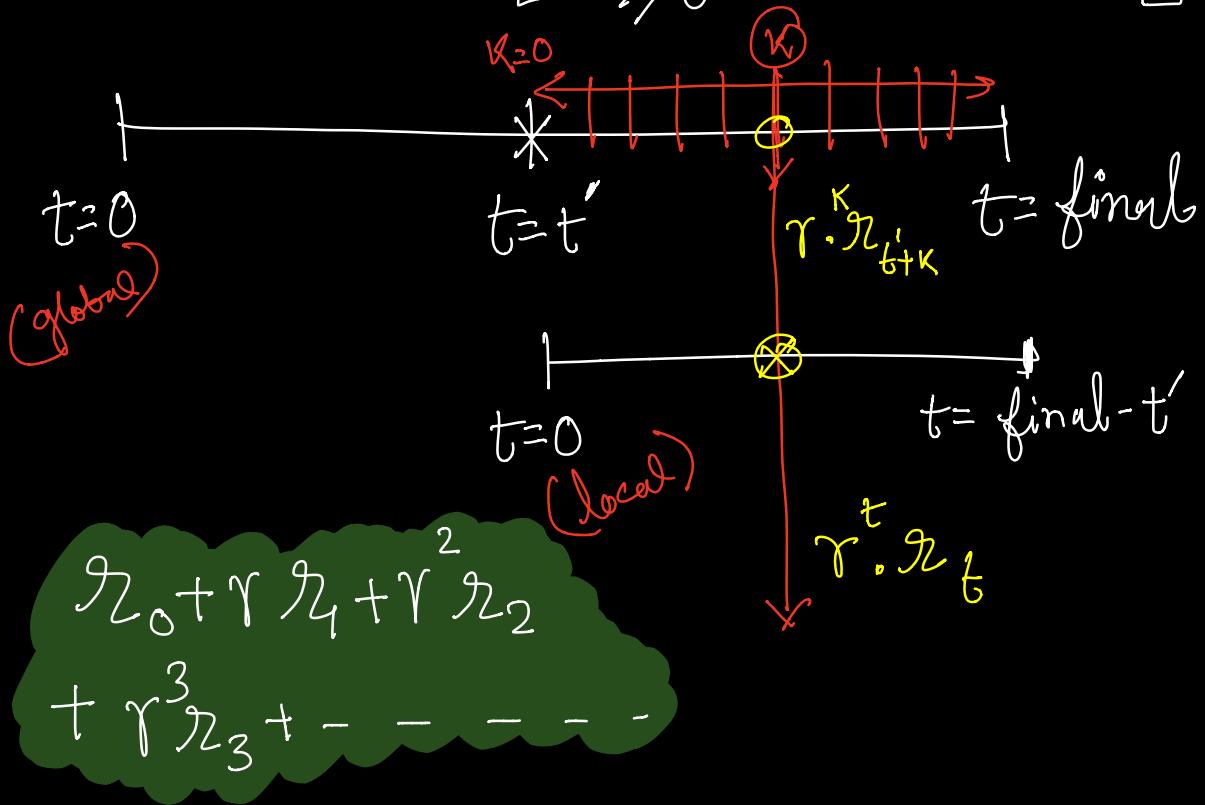
$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = s, a_t = a, \pi \right]$$

Used for  
policy optimization

at state  $s$   
agent took  
action  $a$  then follow  $\pi$

Equivalently

$$Q^{\pi}(s, a) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid \begin{array}{l} s_0 = s, \\ a_0 = a \end{array} \right]$$



## D.3 → Optimal value functions and Q-value fn.

\* Optimal state value fn  $\mathbb{V}^*(s)$

↳ maximum value function

Over all the policies

$$\mathbb{V}^*(s) = \max_{\pi} \mathbb{V}_{\pi}(s)$$

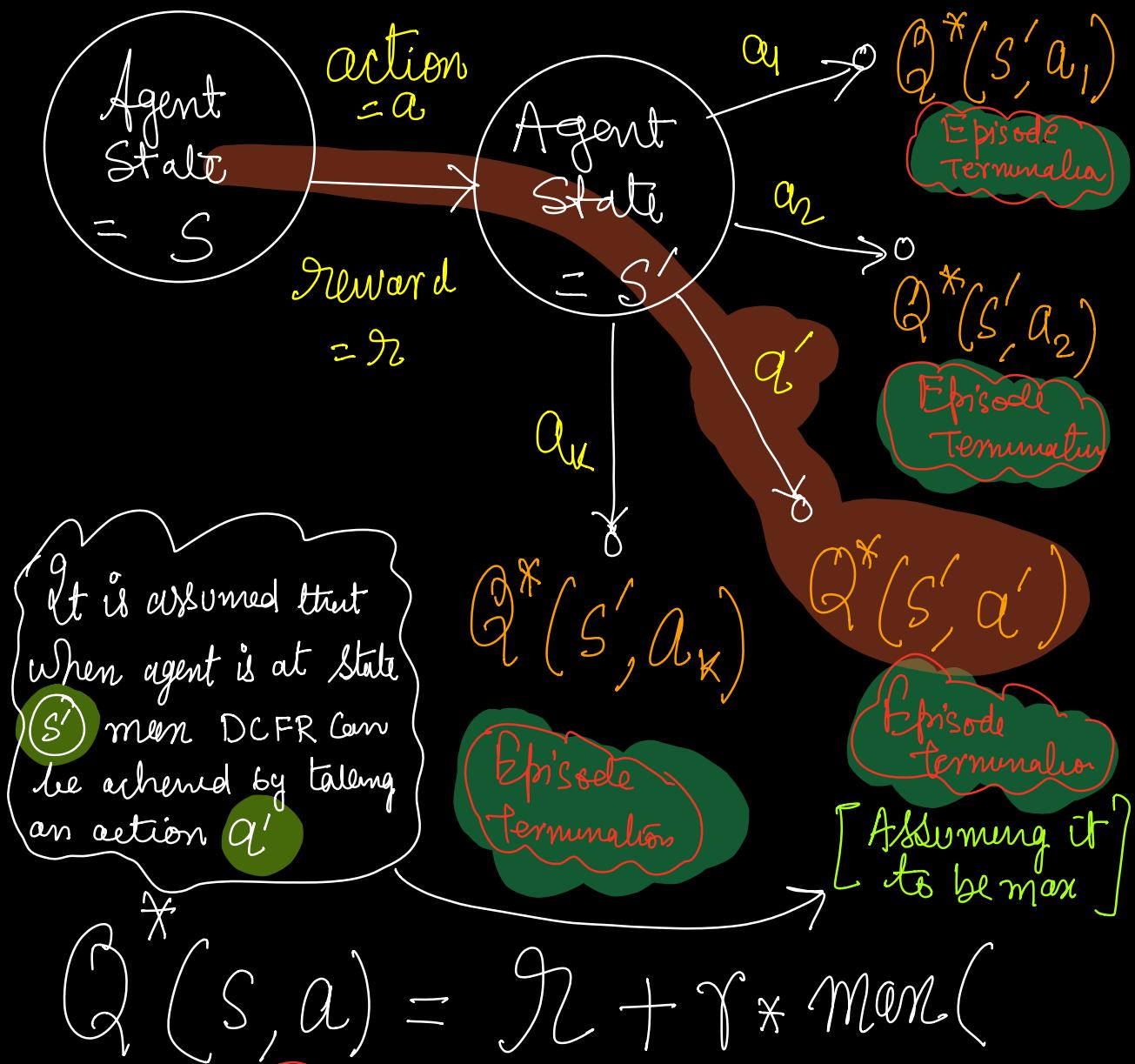
\* Optimal action-value of Q-value fn

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) \quad Q^*(s, a)$$

$$= \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \middle| \begin{array}{l} s_0 = s \\ a_0 = a \end{array} \right]$$

(minimizing over all policies)

Maximum DCFR one can achieve  
from any given (state, action) pair.



$$Q(s, a) = r + \gamma * \max \left( \begin{array}{l} Q^*(s', a_1), Q^*(s', a_2), \\ \dots, Q^*(s', a_k) \end{array} \right)$$

$$= r + \gamma * Q^*(s', a')$$

Since at state  $s$  after taken an action  
 (a) agent will land into state  $s'$  with  
 some probability  $p$ .

$$\Rightarrow s, a, s_1 = 0.3 ; s, a, s_2 = 0.4 ; s, a, s_3 = 0.3$$

$$Q^*(s, a) = \mathbb{E} \left[ r + \gamma \max_{s' \in S} Q^*(s', a') \right]$$

it may reach  
 to a set of  
 states ( $\epsilon$ )

$$\max_{a'} Q^*(s', a') \mid s, a \Big]$$

$Q^*$  as well as  $V^*$  both functions

Satisfies Bellman equations and are  
 defined recursively.

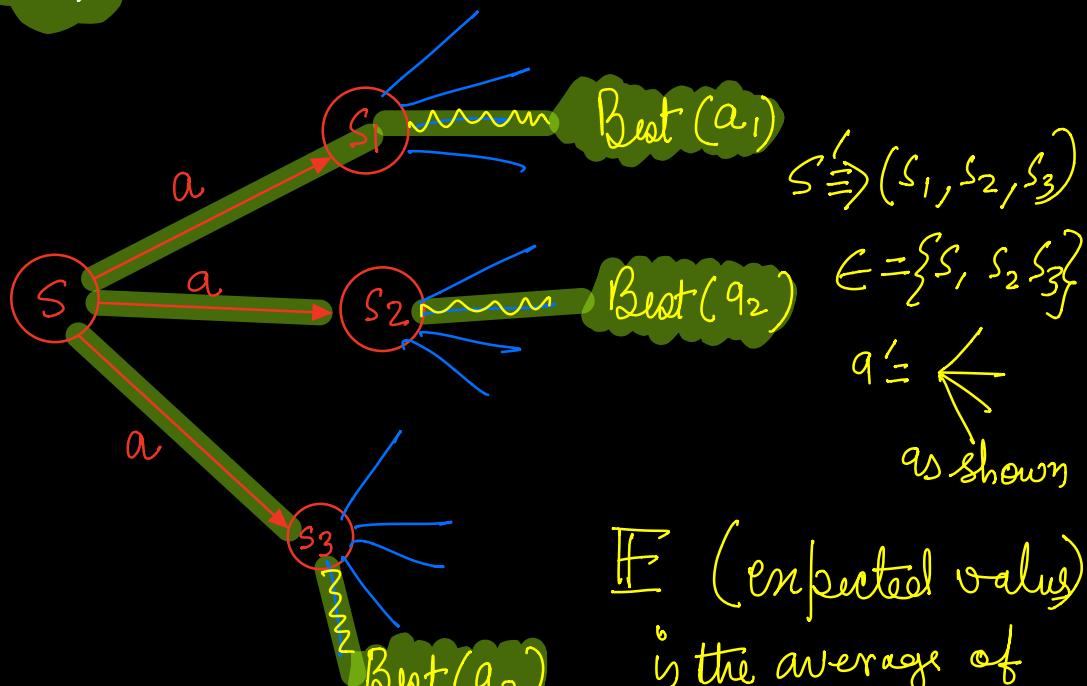
(\*) Expected values are considered in order  
 to address the randomness in the process.

Just see it once more as we are going to use it for Deep Q Learning ( $Q\text{-fn} = \text{Quality of } \langle s, a \rangle \text{ wrt } \pi$ )

$$Q^*(s, a) = \mathbb{E}_{s' \sim \infty} [r + \gamma \times$$

it may reach to a set of states ( $\infty$ )

$$\max_{a'} Q^*(s', a') \mid s, a]$$



$\mathbb{E}$  (expected value)  
is the average of all the 3

basically whatever we observe after taking a decision  
at state  $s$

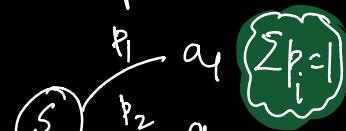
## D.4

# Optimal Policy ( $\pi^*$ )

Policy Def:

- (\*) MDP policies depends only over current state.  
(Not on history)
- (\*) Hence policies are stationary (time-independent).

$$\pi(a|s) = P[A_t=a | S_t=s]$$



Stochastic policies  
ensures suitable  
state, action space  
exploration

These  $p_i$ 's we are deciding under policy  $\pi$ , assuming to be given. Later we need to compute others optimally using DP  $\Rightarrow$  Not Scalable hence approximate it using Policy Network.

- (\*) Policy governs the dynamics of agent.
- (\*) Best/optimal policy need to ensure maximum future reward.
- (\*) In Markov Random Processes, we don't care about the historical reward accumulated so far. We only consider DCFR (Discounted Cumulative Future reward.)  
*This is how we behave optimally*

## (E) How to solve for optimal policy (using DP)

- \* firstly we need to figure out how to evaluate a policy (Policy Evaluation)  
This will tell us how good is that policy.

\* Later using Bellman equation policy iteration or value iterations can be done in order to update the current policy. (Iterative update)

\* Basically assuming policy or value functions at  $i^{\text{th}}$  step are optimal and will be used to compute the policy/value for  $(i+1)^{\text{th}}$  step. [Using Dynamic Programming]

# D.P (Review)

## Bellman Ford Shortest Path Algorithm

\* -ve edges are allowed.

-ve cycles are NOT allowed.

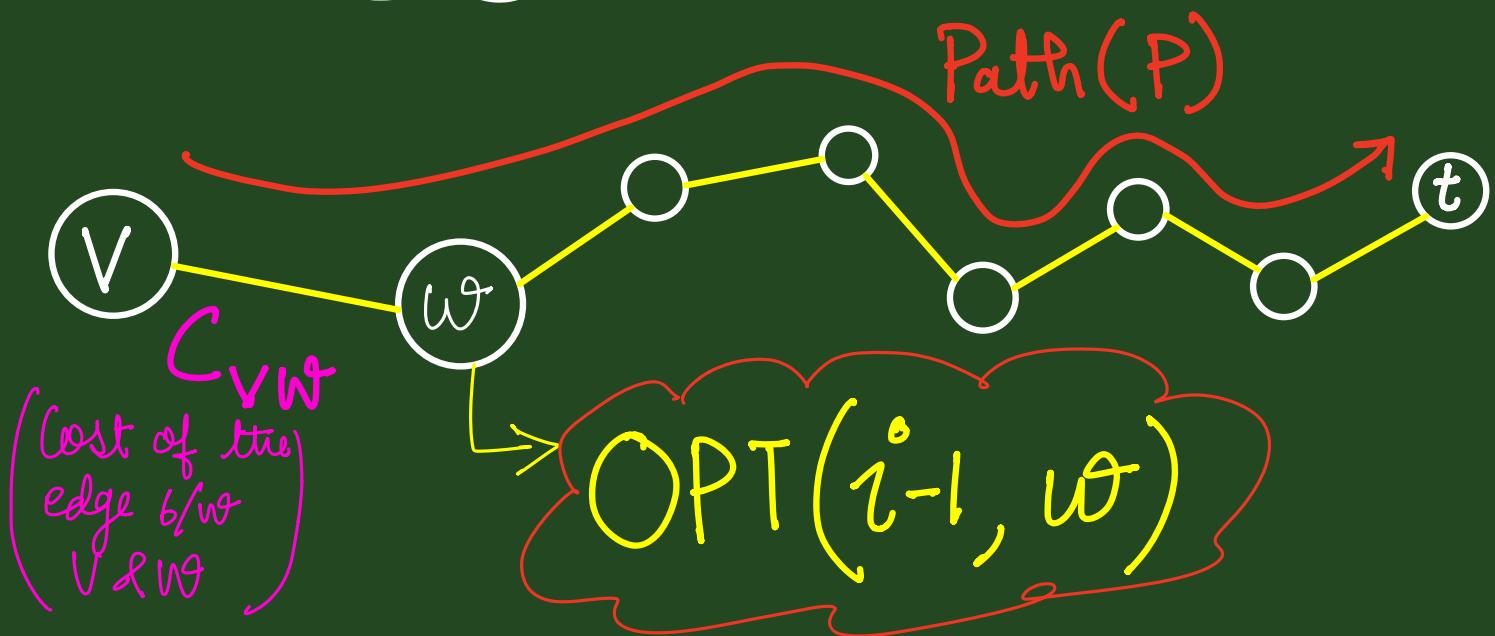
\* Observation: for a given graph  $G = (V, E)$

If  $G = (V, E)$  has NO -ve cycle, then there exist a shortest path from say node  $s$  to  $t$  that is

(a) Simple [No repetition of nodes]

(b) with at most  $|V| - 1 = (n - 1)$  edges.

# \* Setting up the recurrence :



\*  $OPT(i, v) \Rightarrow$  Represent the minimum Cost of say path ( $P$ ) from  $V$  to  $t$  using atmost  $i$  edges.

Since -ve edges are present allowing more edges may reduce the cost.

$$OPT(5, v) \geq OPT(6, v)$$

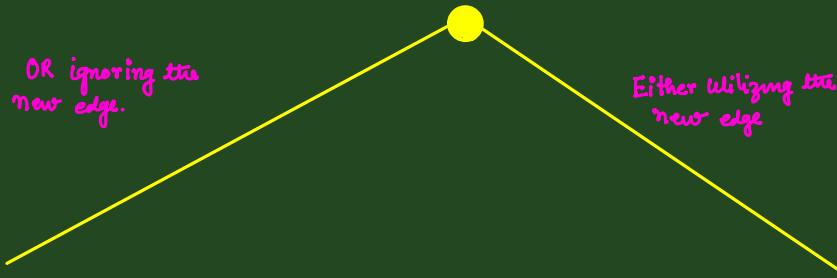
Hence under such setting our originally defined problems solution will be :

$OPT(n-1, s) \Rightarrow$  Minimum Cost of the path ( $P$ ) that can take you from  $s$  to  $t$ , using atmost  $(n-1)$  edges.

How to derive a recurrence relation in terms of smaller subproblems.

$\text{OPT}(i, \mathcal{V})$

Let us assume an optimal path ( $P$ ) from  $\mathcal{V}$  to  $t$  using atmost  $(n-1)$  edges.



$(P)$  either uses atmost  $(i-1)$  edges

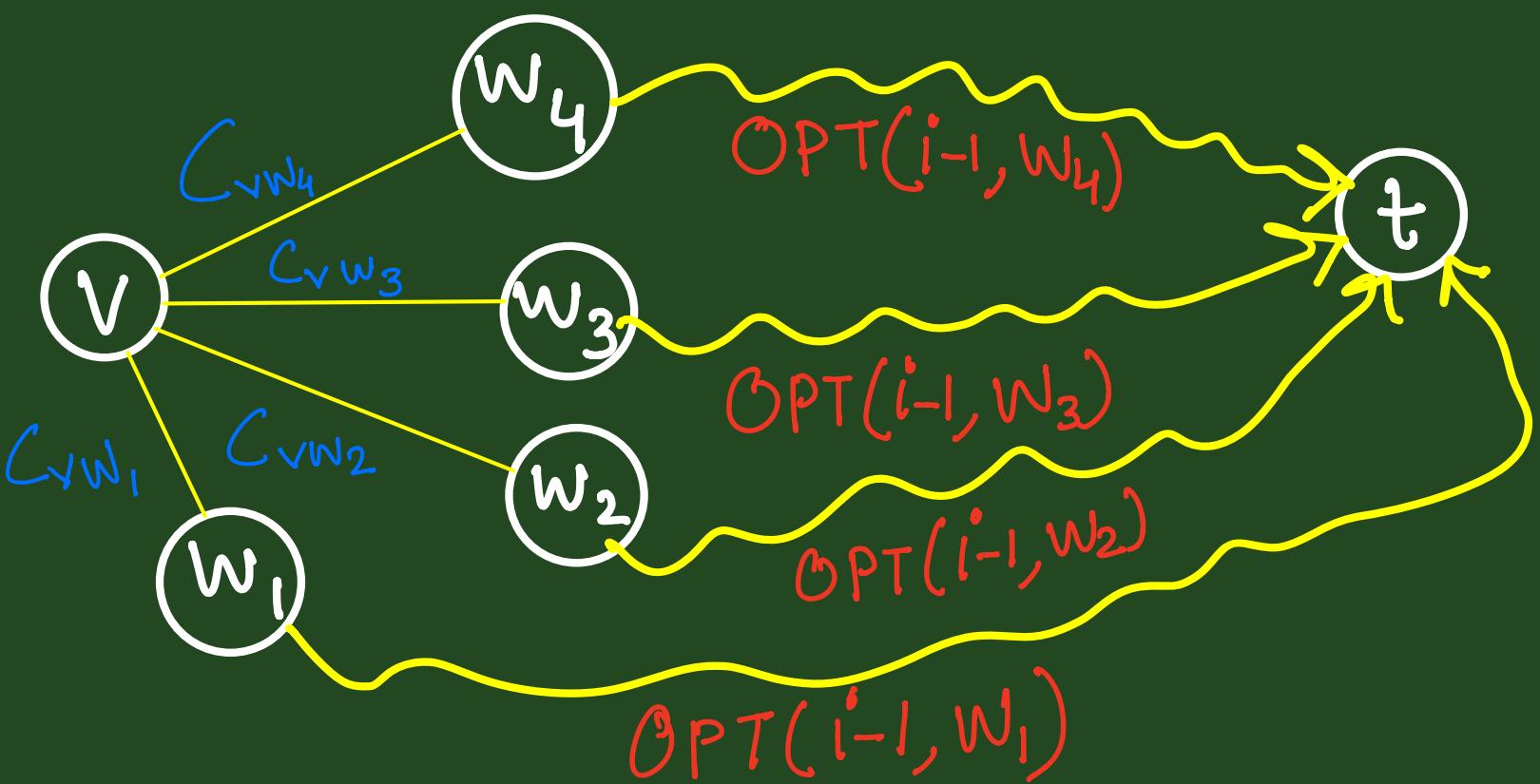


$$\text{OPT}(i, \mathcal{V}) = \text{OPT}(i-1, \mathcal{V})$$

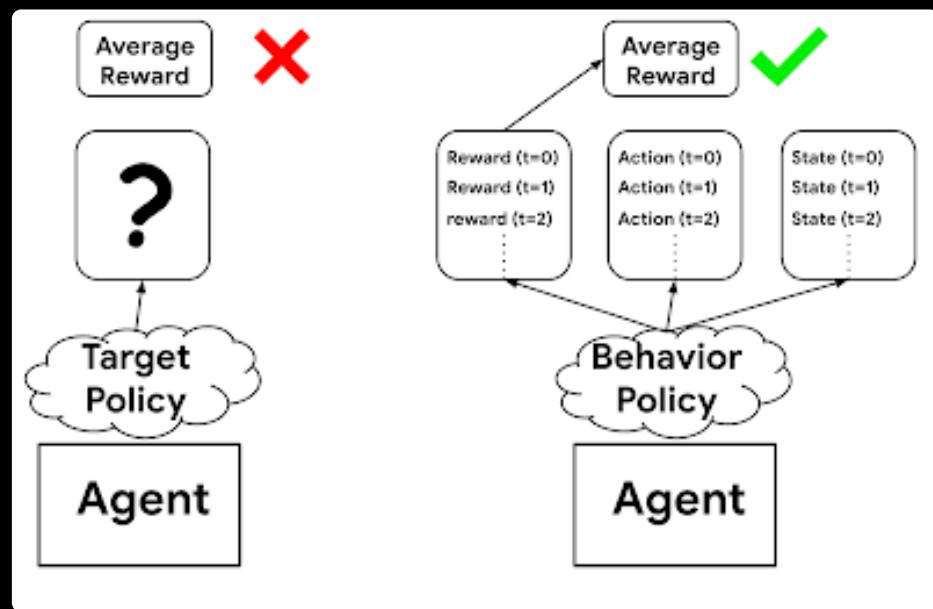
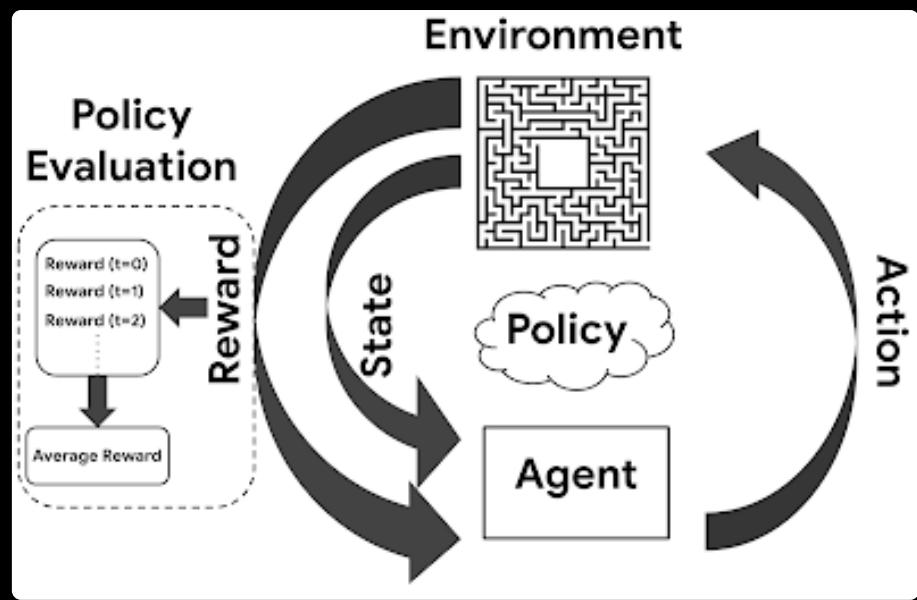
OR  $(P)$  uses atmost  $(i)$  edges



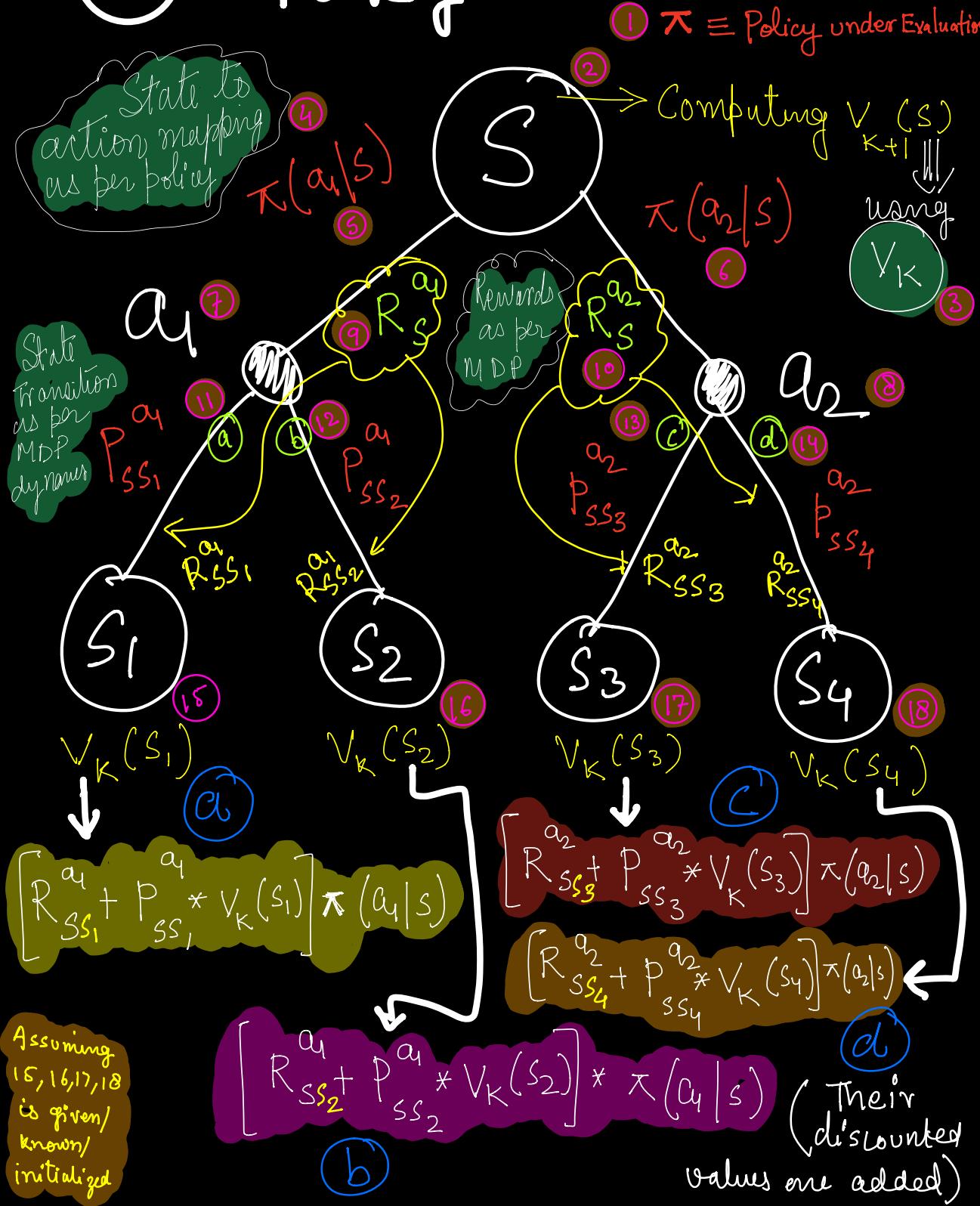
$$\text{OPT}(i, \mathcal{V}) = \text{OPT}(i-1, \mathcal{W}) + C_{\mathcal{V} \mathcal{W}}$$



$$\text{OPT}(i, v) = \min \left[ \begin{array}{l} \text{OPT}(i-1, v) , \\ \min_{w \in V} \left[ \begin{array}{l} \text{OPT}(i-1, w) \\ + c_{vw} \end{array} \right] \end{array} \right]$$



# E.I → Policy Evaluation



$$V_{K+1}(s) = a + b + c + d \quad \left( \text{Expected DCFR of } s \text{ following and using } \right)$$

$$= \sum_{a \in A} \pi(a|s) * \left[ R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a * V_K(s') \right]$$

This probability is doing the expectation.

In closed form:

$$V_{K+1} = R + \gamma P V_K$$

Value vector for all states at  $(K+1)^{\text{th}}$  iteration

Reward function as per the MDP

State Transition Probability

Value fn at  $K^{\text{th}}$  iteration

Let us evaluate a random policy in small grid world.

## E.2 Policy Evaluation: Example

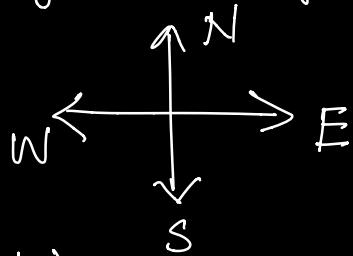
Example : Policy evaluation for smell grid world having four actions :

moving in N | S | E | W

with equal probability :  $\pi(N|.) = \pi(E|.) =$

② Terminal states.

✓/X/X/	1	2	3
4	5	6	7
8	9	10	11
12	13	14	/X/X/



$$\pi(S|.) = \pi(W|.) = 0.25$$

Uniformly random policy.

\*) Reward  $R_2 = -1$

for any transitions  
to ensure step  
minimization

\*) Undiscounted  
episodic MDP ( $\gamma=1$ )

\*) How long does it take to reach to some terminal state from any grid point. [Trajectory length] as per some policy  $\pi$ .

\*) Compute the optimal policy  $\pi^*$  : How to reach

\*) At terminal state episode is finished and looping with zero reward.

\*) Non-terminal state  $(R_2 = -1)$  loop.

$V_K$  for Random Policy ( $K=0$ )

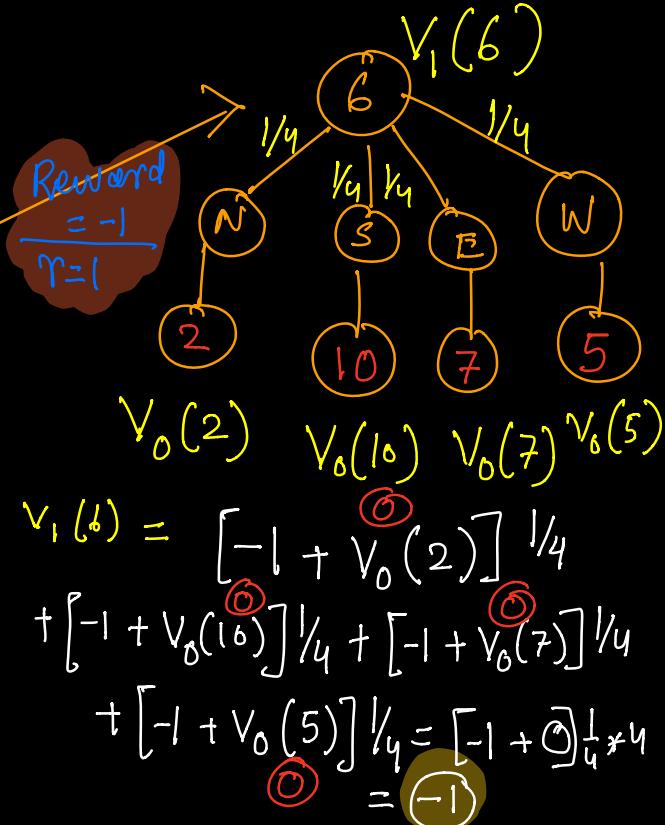
$V_0$	1	2	3
4	0.0	0.0	0.0
0.0	0.0	0.0	0.0
8	9	10	11
0.0	0.0	0.0	0.0
12	13	14	/\
0.0	0.0	0.0	/\

$\Rightarrow V_0(i)$  denotes the no. of steps required to move from state/qmid.  $i$  to a Terminal state  $T$  &  $V_K(T) = 0$  for all  $K's$

For random policy  $V_0$  don't know anything. Later after exploration  $(V_1)$  is computed using  $(V_0)$  which slowly converge to optimality.

$K=1 (V_1)$

$V_0$	1	2	3
4	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
8	9	10	11
-1.0	-1.0	-1.0	-1.0
12	13	14	/\
-1.0	-1.0	-1.0	/\



$K=2 (\vee_2)$

	1	2	3
4	-1.7	-2.0	-2.0
8	-2.0	-2.0	-2.0
12	-2.0	-2.0	-1.7
13	-2.0	-1.7	-1.7
14	-1.7	-1.7	-1.7

why values  
are not reaching  
to the correct  
values.

looking  
back into the  
same state.

$$\begin{aligned} V_2(6) &= \left[ -1 + V_1(2) \right] \\ &\quad -1 + V_1(10) -1 + V_1(7) -1 \\ &\quad + V_1(5) \left[ \frac{1}{4} \right] = (-8) \times \frac{1}{4} \\ &= (-2) \end{aligned}$$

$$\begin{aligned} V_2(14) &= \frac{1}{4} \left[ -1 + V_1(13) -1 + V_1(16) \right. \\ &\quad \left. -1 + V_1(1) -1 + V_1(14) \right] \\ &= \frac{1}{4} \times 7 = 1.75 \end{aligned}$$

$K=3 (\vee_3)$

	1	2	3
4	-2.4	-2.9	-3.0
8	-2.9	-3.0	-2.9
12	-3.0	-2.9	-2.4
13	-2.9	-2.4	-2.4
14	-2.4	-2.4	-2.4

$$V_3(2) = \frac{1}{4} \left[ -1 + V_2(2) \right]$$

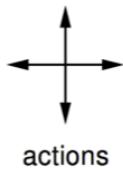
$$\begin{aligned} &-1 + V_2(1) -1 + V_2(3) -1 + V_2(6) \\ &(-1.7) (-2) (-2) \end{aligned}$$

$$V_3(2) = \frac{1}{4} [-11.7] = -2.925$$

$$\begin{aligned} V_3(14) &= \frac{1}{4} \left[ -1 + V_2(14) -1 + \right. \\ &\quad \left. V_2(13) -1 + V_2(1) \right. \\ &\quad \left. -1 + V_2(10) \right] \\ &(-2) (-2) (-1.7) (3) \end{aligned}$$

$$V_3(14) = \frac{1}{4} [-9.7] = -2.425$$

## E.3 Policy Improvement



	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$r = -1$   
on all transitions

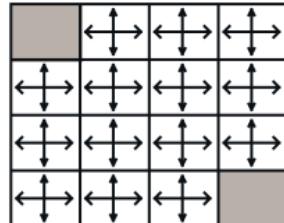
For this MDP, we have seen how to evaluate any policy (say  $\pi = \text{Random}$ ). Given any  $(v_k)$  one take action greedily in order to get refined policy (say  $\pi'$ ).

$v_k$  for the Random Policy

$k = 0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

Just choose greedy actions w.r.t value fn.  
Greedy Policy  
w.r.t.  $v_k$

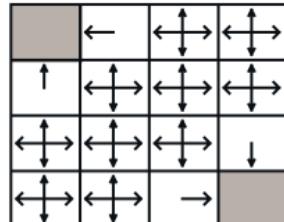


$\pi_{\text{rand}}$   
random policy

$k = 1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

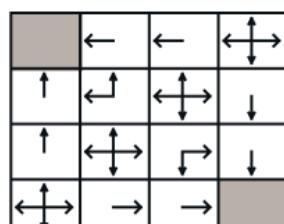
Choosing the direction greedily using value fn.



$\pi_1$

$k = 2$

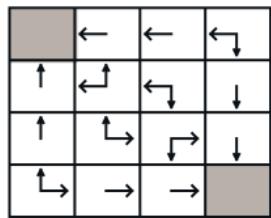
0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0



$\pi_2$

$k = 3$

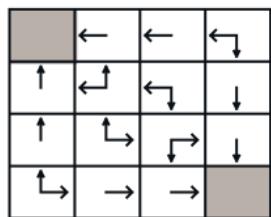
0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0



$\pi_3$

$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

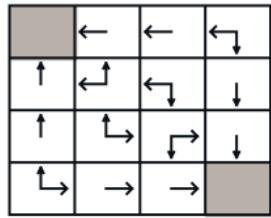


$\pi_3$

optimal policy

$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



$\pi_3$

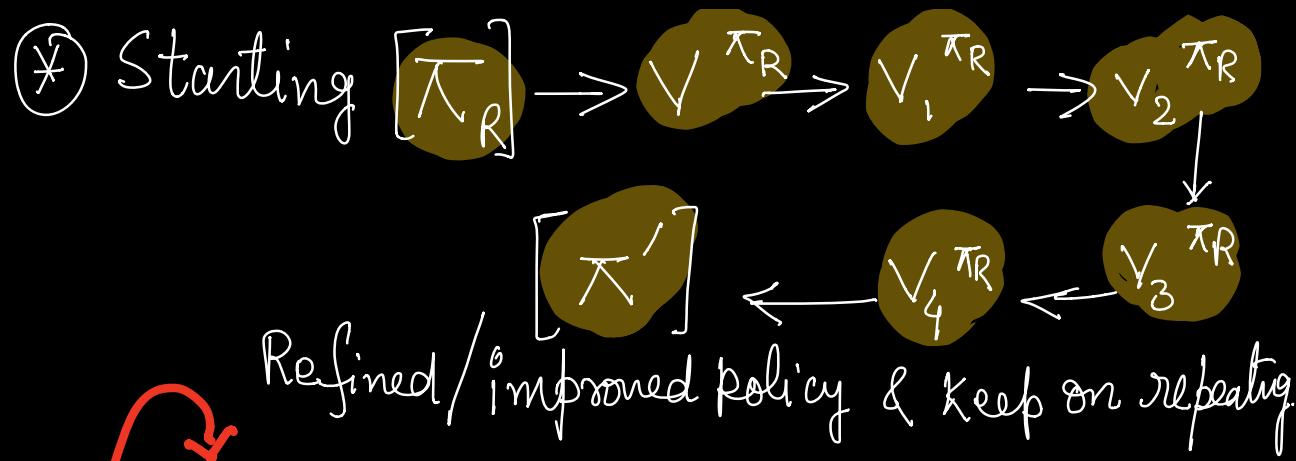
[one can see  
after few  
iterations]

although value fn is changing by policy got saturated

④ Hence we understand how to evaluate a policy.  $\boxed{\text{Policy} \Rightarrow \text{Value}_\pi}$

⑤ Later given a value fn act greedily as per it and get a refined policy

$\boxed{V^\pi \rightarrow \pi' = \text{Greedy}(V_\pi)}$



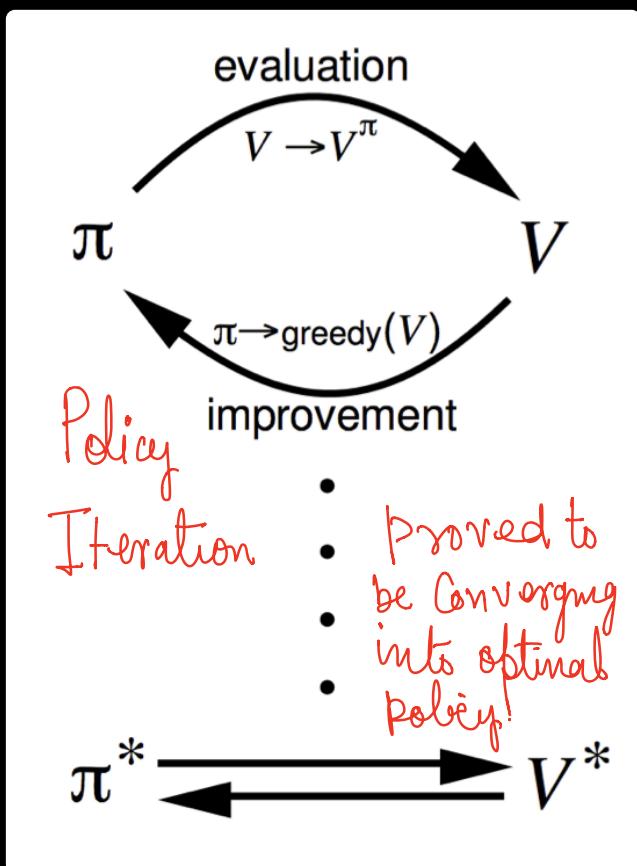
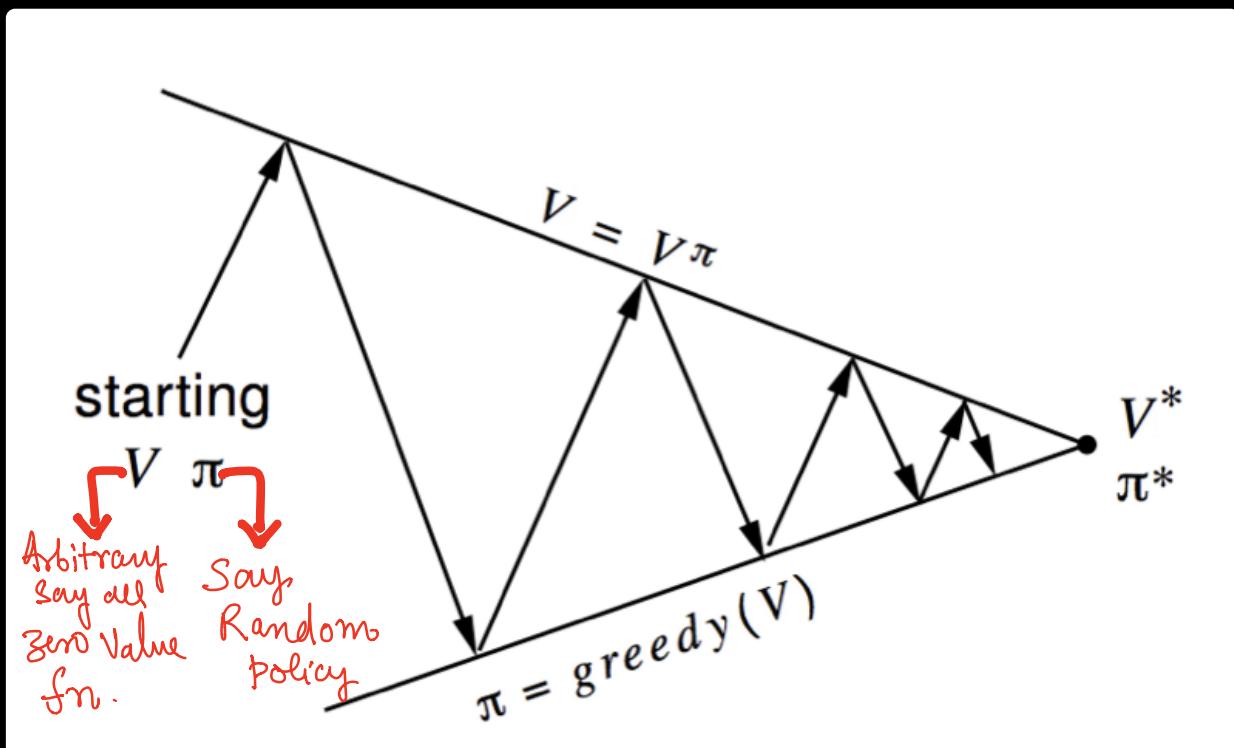
- Iterating**
- ① Given any policy  $(\pi)$
  - ② Evaluate the policy  $(\pi)$  by computing value fn.
  - ③ Keep on iterating over value fn So that its corresponding Greedy policy  $(\pi)$  got Saturated.
  - ④ Improve  $(\pi) \xrightarrow{\text{to}} (\pi')$

**Policy evaluation** Estimate  $v_\pi$

Iterative policy evaluation

**Policy improvement** Generate  $\pi' \geq \pi$

Greedy policy improvement



This has been proved that this process will always converge into the optimal policy  $\pi^*$ .

\* No matter from where one is starting.  
(Initialization just effect the rate of convergence).

# F Solving for Optimal policy (Deep Q-learning)

Coming from above

$$Q^*(s, a) = \mathbb{E}_{S' \in \cup} [R + \gamma \max_{a'} Q^*(s', a')]$$

(already discussed)

Basically learning the Q fn

it may reach to a set of states ( $G$ )

$Q^*$  as well as  $V^*$  both functions  
Satisfies Bellman equation and are defined recursively.

## F.1 Problems with DP Solution

\* As shown above Bellman equation can be used to solve for the optimal Q-value fn. (Evaluate & Improve).

- \* This optimal solution is proved to be converging and easy to compute
- \* But is simply not scalable. for any typical modern day problem say Atari gaming agent  $S \equiv$  State (Raw pixels of a frame)
  - (a) Can be limited but  $S \equiv$  can be exponentially many. Hence one need to Compute Q-value for full  $\langle s, a \rangle$
  - i.e.  $\langle \text{State}, \text{action} \rangle$  space  $\Rightarrow$  Simply Computationally in feasible.

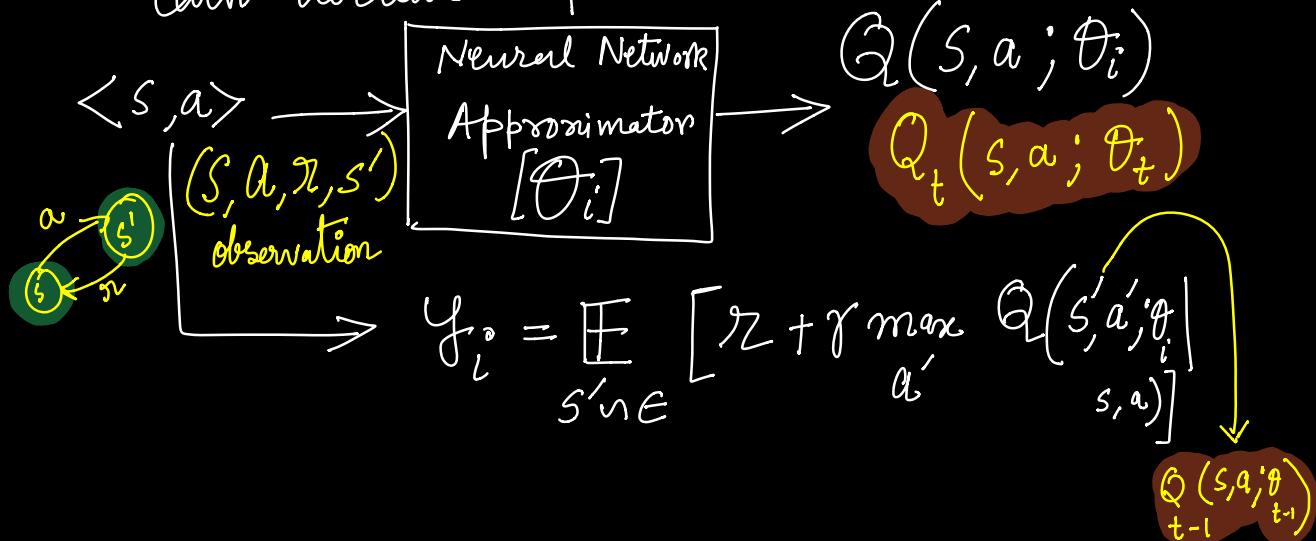
## F.2 Deep Q-learning for Qfn approximation

- \* Instead of computing  $Q(s, a)$ , we need to use some function approximator (Say neural network) to approximate/estimate  $\{Q(s, a) + \gamma, a\}$

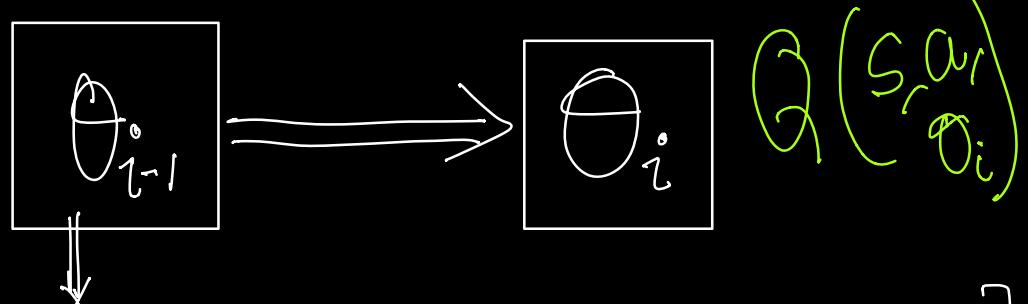
$$Q(s, a; \theta) \approx \hat{Q}(s, a)$$

Using DNN parameters  
to estimate the Q value  
for  $(s, a)$  pair.

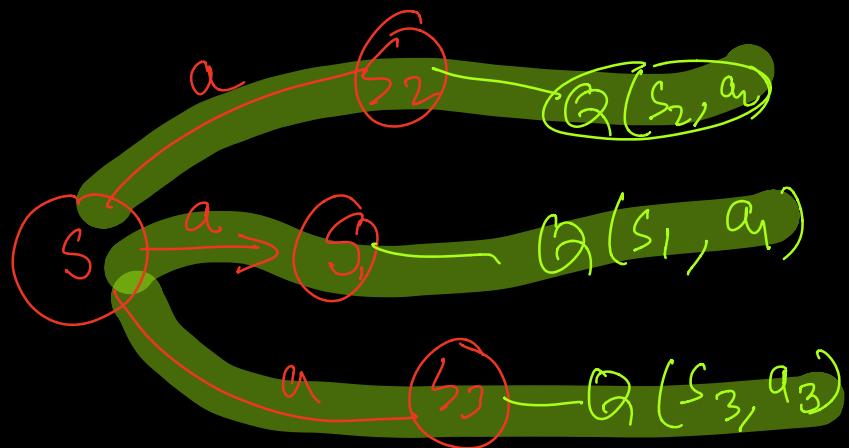
- \* We need a Qfn approximator that satisfies Bellman equation by enforcing Bellman optimality at each iterative step. (Intuition)



$$(S, a, r, S') \quad S \in \{s_1, s_2, s_3\}$$



$$y_i \Rightarrow \underset{S \in \{s_1, s_2, s_3\}}{\mathbb{E}} \left[ r + \gamma \max_{a'} Q(s', a'; \theta_{t-1}) \right]$$



### F.3 Forward Pass: Loss function

$$\therefore L_i(\theta_i) = \mathbb{E}_{s, a \sim P(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$$

↳ Current network o/p.  
(n/w parameters at  $i^{\text{th}}$  step)

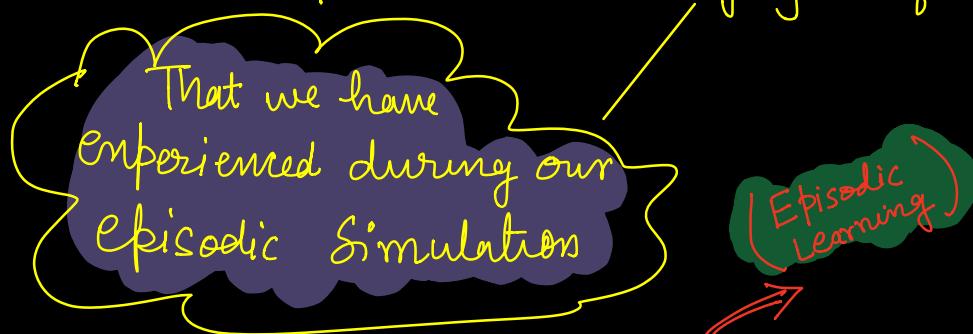
$$y_i = \mathbb{E}_{s' \in E} [r + \gamma * \max_{a'} Q(s', a'; \theta_{i-1}|s, a)]$$

↳ Previous n/w o/p (old parameters at  $(i-1)^{\text{th}}$  iteration)

for addressing the scalability issue, we don't need to compute Q-fn for each & every  $\langle \text{state, action} \rangle$  pair.

\* Using the environment & agent machinery along with the full MDP dynamics, play few game episodes.  
[Play game]

\* Instead of Computing the Q-fn for full  $\langle \text{state, action} \rangle$  pair space, explore only few of them.



Therefore play games  $\Rightarrow$  Generate Episodes

Only consider those  $\langle \text{state, action} \rangle$  pairs that we observe in home state instead of full space

We these episodes for training

## F.4 Game Play (given MDP dynamics OR N/W parameters $\Theta$ )

① Saving the game states into  $\mathcal{T}$  [Training data] Generating training data.  $[Q]$  approximating  $N/W(\Theta)$

$\langle S_t \text{ (or } \emptyset_t), a_t, r_t, S_{t+1} \rangle$

Raw state  $\in$   $S_t$  or  $\emptyset_t$  Or may be processed Raw state

Action performed  $\downarrow$

Reward  $\downarrow$

Next state  $\downarrow$

$\emptyset_t = \text{Processed}(S_t)$

② for  $t = 1$  to  $(M)$  do Generating full game data : M episode training data

②.1 ① Initialize Starting State  $S_1 = \{x_1\}$  May be some starting game screen.

②.2 ① Pre-processed Sequence  $\emptyset_1 = \emptyset(S_1)$  Some processing function (say edges of differences)

②.3 for  $t = 1$  to  $(T)$  do for each time step  $t$  of the game episode

2.3.1 Select a random action (Say  $a_t$ )  $\Rightarrow$  for sufficient state action space exploration  
with some small probability ( $\epsilon$ )  $\Rightarrow$  Small probability to explore new things.  
 select random action or select greedy action as per the current policy

②.3.2 Otherwise select  $a_t = \max_a Q^*(\emptyset(S_t), a; \Theta)$

$(\emptyset(S_t), a_t) \rightarrow \boxed{\Theta \text{ N/W parameter}} \rightarrow a_t \text{ (next actions need to be taken)}$   
 Approximating  $Q^*$

2.3.3 Execute action  $a_t$  in the game emulator and observe the reward  $r_t$  and next image  $(x_{t+1}) \sim \{ \begin{matrix} \text{new state} \\ \text{next state} \end{matrix} (s_{t+1}) \}$

2.3.4 Store the transition

$\langle \phi_t, a_t, r_t, \phi_{t+1} \rangle$  into  $T$  as training data.

3) Return the training data  $T$ .

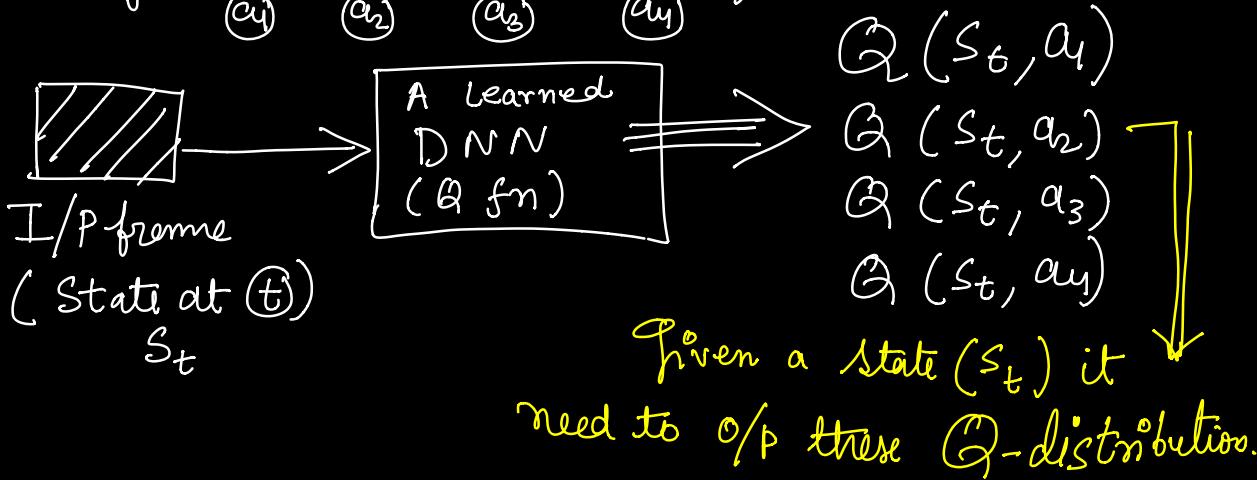
☞ This is just Gameplay algorithm in case actual  $Q^*$  function is known or an approximating N/W is given.

☞ It can be used to generate more training data in  $T$  and incrementally train the N/W more and get better and better N/W parameters ( $\theta^t \Rightarrow \theta^*$ ) at Convergence.

\* What does it mean when we say an approximator deep neural network is learning / learned the  required  $\{Q^*\text{ functions}\}$ . (Intuition)

In most of the games (for example Atan Games) actions are always finite in number and fixed.

(Say UP | DOWN | LEFT | RIGHT)



\* Atan games have around (4-18) actions

Hence  $\nabla_{a_i} Q(s_t, a_i) \Rightarrow$  one scalar value per action and per state.

[Efficiently computing full  $Q(s_t)$  distribution for all actions in one single forward pass]

## F.5 Training issues and Experience Replay

Problem: Learning from consecutive samples

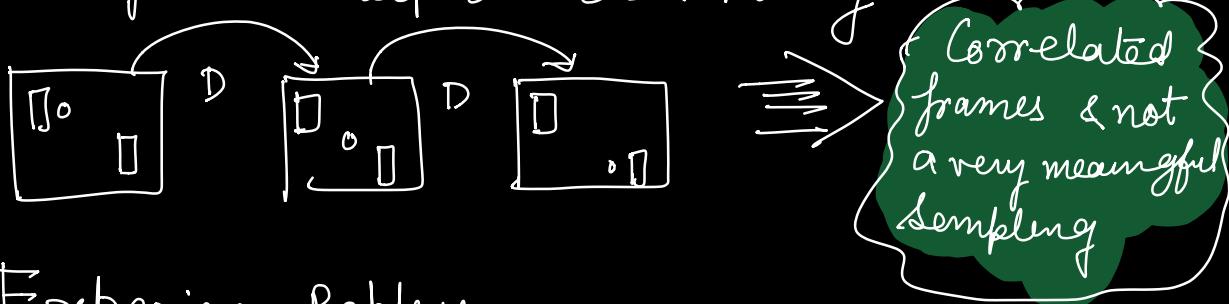
Coming from same episode ( $T$ ) is problematic.

↳ Samples will be highly correlated hence easily learn some trivial dynamics but fails to generalize.

↳ May become biased towards the current o/p.

$\Rightarrow$  If the current action maximizing Q-value  
is (left)  $\xrightarrow{\text{Down}}$  Training samples are dominated  
by left hand side samples  
(Down Side)  
Due to this  $\Leftarrow$

the current Q-Network parameters ( $\theta_{t-1}$ )  
may lead to bad bias towards (DOWN) &  
May even keep on oscillating.



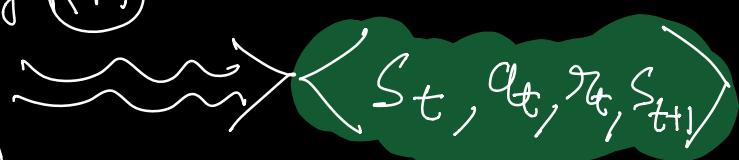
### Experience Replay :

Better idea is to make random minibatches  
out of training data ( $T$ )  $\rightsquigarrow$  Experience &  
train n/w over that  $\rightsquigarrow$  Replay

Replay Memory  $Rm$

OR

Training data  $(T)$



## F.6 Algorithm: Deep Q-learning with Experience Replay.

- ① Initialize the Deep Q-learning N/w parameters ( $\theta_0$ ) randomly. (N/w learning the Q-function)
- ②  $T = \text{Gameplay}(\theta_0)$
- ③ Till Convergence (Working for some  $\theta_i$ )
  - 3.1 Sampling some transitions from  $T$ . Sample a random minibatch of transitions  $\langle \phi_j, a_j, r_j, \phi_{j+1} \rangle$  from  $T$ .
  - 3.2 Set  $y_j = [r_j]$  for terminal  $(\phi_{j+1})$   
for each  $j^{th}$  transition
 
$$y_j = r_j + \gamma \max_{a'} Q[\phi_{j+1}, a' | \theta]$$

for non-terminal  $(\phi_{j+1})$

$$③.3) L_j(\theta_j) = (y_j - Q(\theta_j, a_j | \theta))^2$$

$$③.4) L^{\text{batch}} = L + L_j(\theta_j)$$

New  $T = \text{GamePlay}(\theta_i)$

Forward Pass :

$$L_i(\theta_i) = \underset{s, a \sim P(\cdot)}{\mathbb{E}} \left[ (y_i - Q(s, a | \theta_i))^2 \right]$$

$$y_i = \underset{s' \sim \mathcal{S}}{\mathbb{E}} \left[ r + \gamma * \max_{a'} Q(s', a' | \theta_{i-1}) \mid s, a \right]$$

Backward Pass : Differentiating loss w.r.t N/w parameters  
(Gradient update)

$$\nabla_{\theta_i} L_i(\theta_i) = -2 * \underset{s, a \sim P(\cdot)}{\mathbb{E}} \left[ y_i - Q(s, a | \theta_i) \right]$$

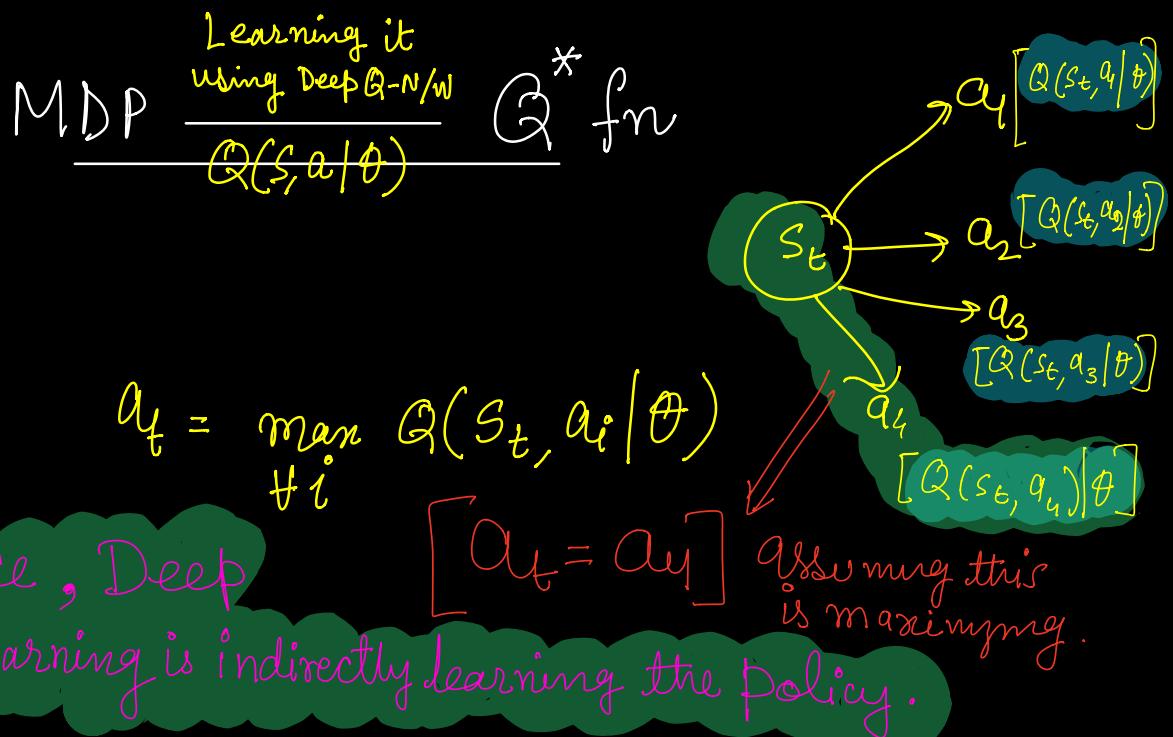
$$* \nabla_{\theta_i} [Q(s, a | \theta_i)]$$

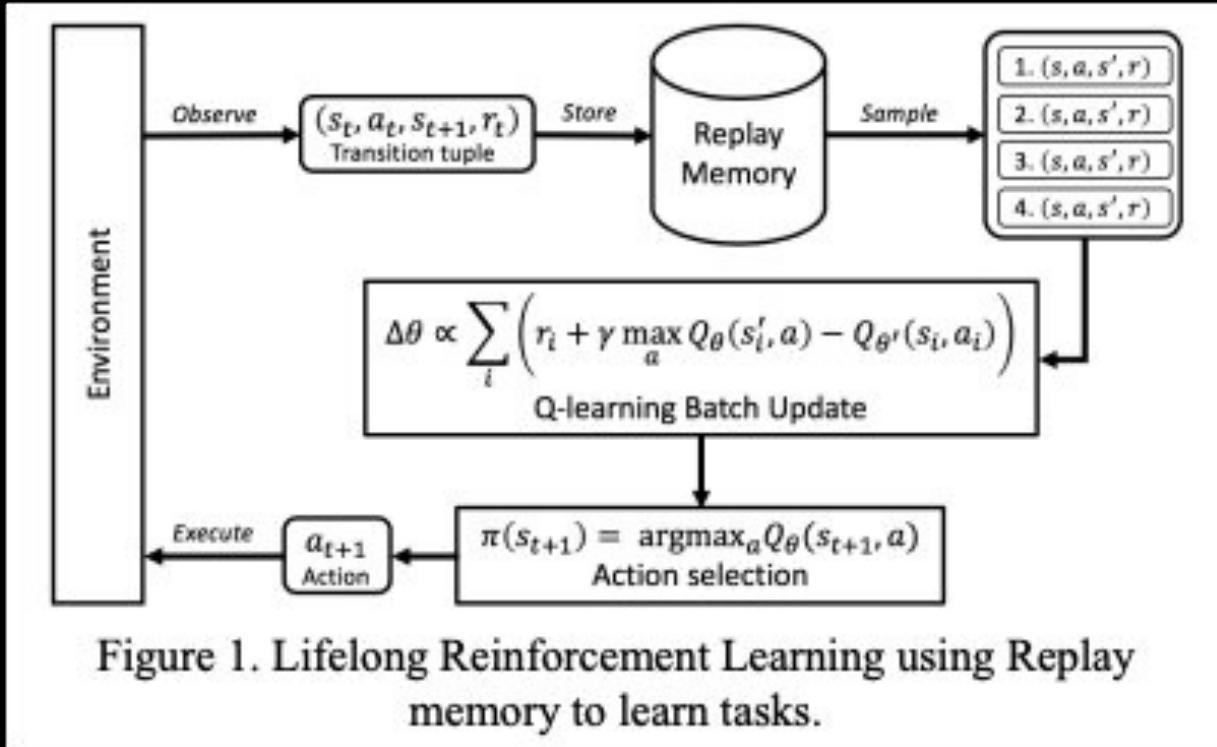
Bellman's estimate

$Q$ -fn estimate using  
Current N/w.

\* Hence, we have seen that Deep Q-Learning Network tries to estimate the optimal Q-function starting from some random estimate and keep on refining it iteratively by imposing Bellman Optimality. That is ensured by the loss function.

\* Given any MDP its corresponding optimal Q fn is its solution. Given  $Q^*$  fn for any MDP, one can plan to define the policy by moving greedily using  $[Q^*]$ .

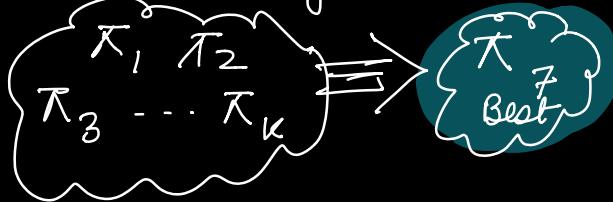




\* Deep Q-learning N/o learns the optimal  $Q$ -function and using it greedily provides us the suitable policy  $\pi$   $\Rightarrow$  Hence it gives us an implicit handle over the policy ( $\pi$ ).

\* Can we find out the best policy out of the collection of policies?

How can we learn the policy explicitly w/o learning an Q-value estimate



# G Solving for Optimal Policy (Directly)

## G.1 Introduction to Policy Gradient

\* Instead of using DNN to approximate Q-Value function, why can't directly learn the Suitable Policy ( $\pi$ ), parameterized by  $\theta$ .

 Class of parameterized policies:  $\pi = \{\pi_\theta \mid \theta \in \mathbb{R}^m\}$

for any policy ( $\pi$ ):  $J(\theta) = \mathbb{E} \left[ \sum_{t=0}^T r_t \mid \pi_\theta \right]$   
 define its value as:  
 ↓  
 Value of the policy  $\pi$  parameterized by  $\theta$   
 ↓  
 Expected DCFR while following policy  $\pi_\theta$

\* Given any policy  $(\pi_\theta)$ , how much DCFR on an average one can extract  $\Rightarrow$  value of  $\pi_\theta \sim J(\theta)$

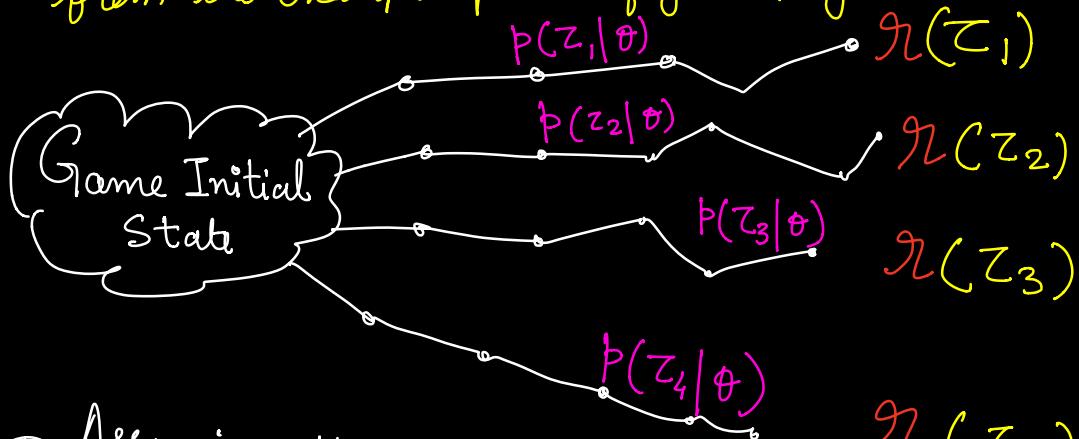
GOAL: DNN need to optimize for the parameter  $\theta$  that can realize the optimal policy  $\pi^*$ .

$\theta^* = \arg \max_{\theta} J(\theta)$  Gradient Ascent over  $\theta \in \mathbb{R}^m$  REINFORCE also.

## G.2

# REINFORCE algorithm Loss function

\* Trajectory ( $\tau$ ) :  $\tau = (\langle s_0, a_0, r_0 \rangle, \langle s_1, a_1, r_1 \rangle, \langle s_2, a_2, r_2 \rangle, \dots)$   
 Sample trajectory from the example episode of game play. (Sampled randomly).



\* Assuming the agent to be starting from some initial game state and following  $(\pi_\theta)$  policy parameter over  $(\theta)$ .  
 → a) Assuming all trajectories to be equally probable :  $J(\theta) = \frac{1}{4} \sum_{i=1}^4 R(z_i)$

→ b) Assuming trajectories follow  $P(z | \theta)$  distribution

$$J(\theta) = \sum_{i=1}^4 P(z_i | \theta) * R(z_i)$$

Value associated to a policy  $(\pi_\theta)$ , Parameterized over  $[\theta]$ .

$$= \mathbb{E}_{\pi_\theta} [R(z_i)]$$

$$\begin{aligned}
 \text{Expected Reward} &= \mathbb{E} [\mathcal{R}(z)] \\
 &\quad z \sim p(z|\theta) \\
 &= \int (\mathcal{R}(z) * p(z|\theta)) dz
 \end{aligned}$$

Considering all  
 trajectories that one  
 can encounter while  
 agent starts from some initial state & following  $(\pi_\theta)$  policy.

$$\therefore \theta^* = \underset{\theta}{\operatorname{argmax}} [J(\theta)]$$

(\*) REINFORCE algorithm uses  $J(\theta)$  for its  
 forward pass as a Gain fn [Need to be  
 maximized]

# G.3

## REINFORCE algo Backward pass

$$(*) \text{ Forward Pass : } J(\theta) = \int_z r(z) * p(z|\theta) dz \quad (1)$$

$$(*) \text{ Backward pass : } \nabla_{\theta} J(\theta) = \int_z r(z) * \nabla_{\theta} p(z|\theta) dz \quad (2)$$

\*\*\* This is intractable. Gradient computation of an expectation is problematic when  $p$  depends upon  $\theta$ . [Gradient of expectation]

Trick :  $\nabla_{\theta} p(z|\theta) = p(z|\theta) * \left[ \frac{\nabla_{\theta} p(z|\theta)}{p(z|\theta)} \right] \quad (3)$

Putting eq.(3) into eq.(2)  $= p(z|\theta) * \nabla_{\theta} \log[p(z|\theta)]$

$$\therefore \nabla_{\theta} J(\theta) = \int_z r(z) * \nabla_{\theta} \log[p(z|\theta)] * p(z|\theta) dz$$

Expectation of Gradient  $\therefore \nabla_{\theta} J(\theta) = \mathbb{E}_{z \sim p(z|\theta)} [r(z) * \nabla_{\theta} \log[p(z|\theta)]]$

## G.4

# Issue with Gradient & its Computation

- (\*) There can be infinitely many trajectories possible for any parameterized policy  $\pi_\theta$  estimated by DNN.
- (\*) The reward function  $[J(\theta)]$  has been defined as an Expectation over all such trajectories. Hence  $J(\theta)$  can be computed essentially by Monte-Carlo Sampling of trajectories.
- (\*) All trajectories  $(\tau_1, \tau_2, \dots, \tau_\infty)$  for some given  $\pi_\theta$  are not equiprobable, instead following a distribution. Say  $p(\tau|\theta)$
- (\*) Hence, for a some trajectory say  $\tau_1 = \langle (s_0, a_0, r_{00}), (s_1, a_1, r_{10}), \dots, (s_n, a_n, r_{n0}) \rangle$   
let us compute the probability of  $(\tau_1)$  given  $p(\tau|\theta)$

$$p(\tau_1 | \theta) = [T(s_1 | s_0, a_0) * \pi_\theta(a_0 | s_0)]$$

$$* [T(s_2 | s_1, a_1) * \pi_\theta(a_1 | s_1)]$$

$$* [T(s_2 | s_2, a_2) * \pi_\theta(a_2 | s_2)]$$

Given the policy  $\pi_\theta$  parameterized by  $(\theta)$  & learned by some Deep policy N/W maximizing the reward for  $[J(\theta)]$ .

Transition  
Probabilities

As per  
Parameterized  $\pi_\theta$   
Policy.

$$(*) \quad P(\mathcal{Z} | \theta) = \prod_{t \geq 0} T(S_{t+1} | S_t, a_t) * \pi_\theta(a_t | S_t) \quad \text{Eq-A}$$

Estimating the probability of randomly selected trajectory ( $\mathcal{Z}$ ), while following policy ( $\pi_\theta$ )

The transition probability of going to ( $S_{t+1}$ ) state after taken an action ( $a_t$ ) at state ( $S_t$ )

Probability of taking an action ( $a_t$ ) at state ( $S_t$ ) under the policy ( $\pi_\theta$ )

- \*) The major problem is the  $J(\theta)$  requires  $P(\mathcal{Z} | \theta)$  which in turn requires transition probability ( $T$ ) which is unknown. Our policy network is estimating ( $\pi_\theta$ ) by maximizing  $J(\theta)$ .
- \*) Now the question is : Can we compute  $J(\theta)$  without transition probabilities ( $T$ ).
- \*) But for backpropagation  $\nabla_\theta [J(\theta)]$  only gradient of  $J(\theta)$  is required and it is not depending upon ( $T$ ).

As already shown

$$\therefore \nabla_{\theta} J(\theta) = \mathbb{E}_{\substack{\text{Expectation} \\ \text{of Gradient}}} \left[ R(z) * \nabla_{\theta} \log [P(z | \theta)] \right] \quad (\text{Eq-3})$$

The log likelihood of such a sampled trajectory ( $z$ ) is using (Eq-4)

$$\log [P(z | \theta)] = \sum_{t \geq 0} \log T(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t) \quad (\text{Eq-4})$$

Independent of  $[\theta]$

Dependent over  $[\theta]$

$$\nabla_{\theta} \log [P(z | \theta)] = \sum_{t \geq 0} \nabla_{\theta} \log [\pi_{\theta}(a_t | s_t)]$$

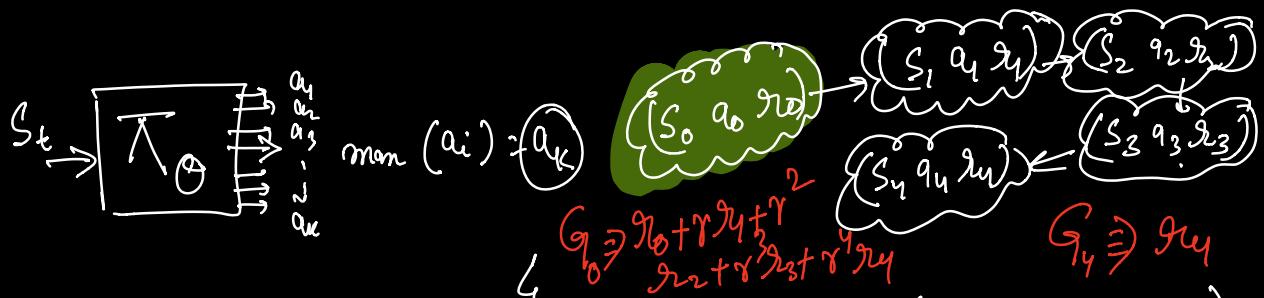
(\*) Hence the derivative of log likelihood of the trajectory ( $z$ ), which is required for the computation of the Gradient of the Reward function i.e  $\nabla_{\theta}(J(\theta))$  is independent of Transition probability  $[T]$ .

Independent of transition probability  
(Eq-5)

This also implies that  $\nabla_{\theta}(J(\theta))$  is also independent of  $[T]$

Hence putting (Eq-D) into (Eq-B)

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{z \sim p(z|\theta)} \left[ r(z) * \sum_{t \geq 0} \nabla_{\theta} \log [\pi_{\theta}(a_t | s_t)] \right]$$



$$\nabla_{\theta} (J(\theta)) = \sum_{t=0} \nabla_{\theta} \log (\pi_{\theta}(a_t | s_t)) * G_t$$

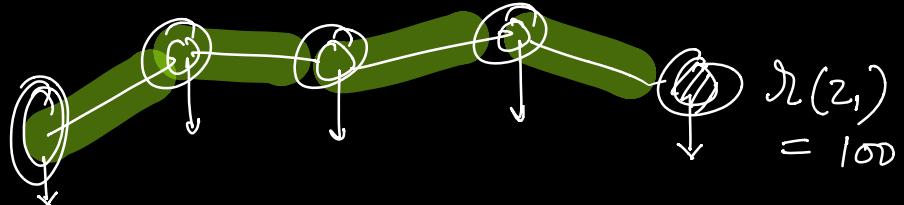
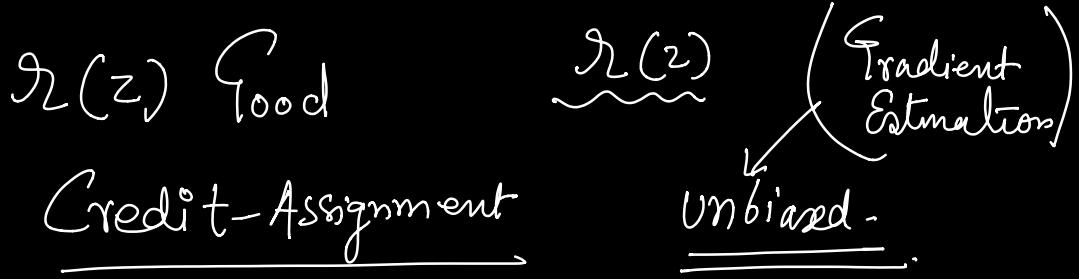
$$\hat{=} G_0 * \nabla_{\theta} \log (\underline{\pi_{\theta}(a_0 | s_0)}) +$$

$$G_1 * \nabla_{\theta} \log (\pi_{\theta}(a_1 | s_1)) +$$

$$G_2 * \nabla_{\theta} \log (\pi_{\theta}(a_2 | s_2)) +$$

$$G_3 * \nabla_{\theta} \log (\pi_{\theta}(a_3 | s_3)) +$$

$$G_4 * \nabla_{\theta} \log (\pi_{\theta}(a_4 | s_4))$$



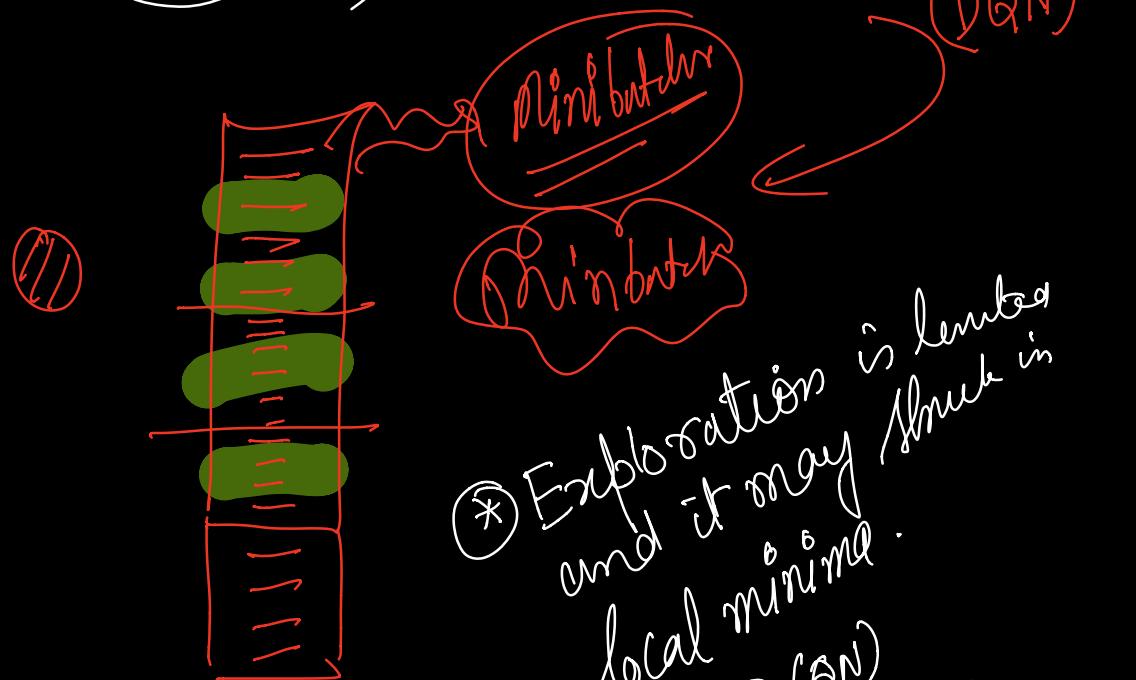
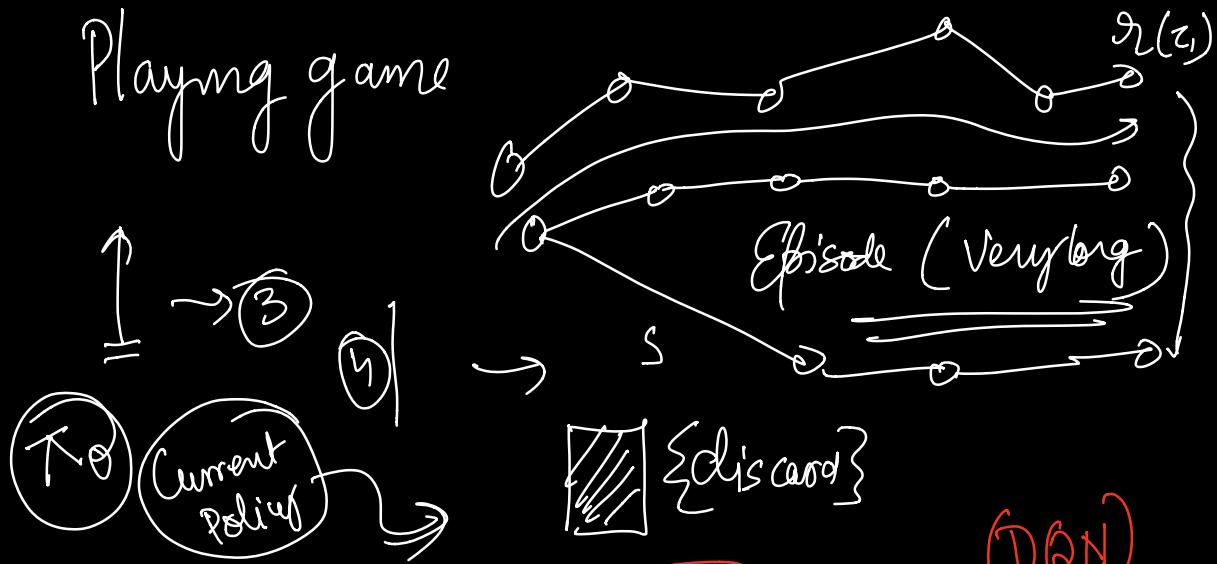
$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \gamma(z_t) \nabla_{\theta} \log(\pi_{\theta}(a_t | s_t))$$

$$\sum_{t \geq 0} \left( \sum_{t' \geq t} \gamma_{t'} \right) \nabla_{\theta} \log(\pi_{\theta}(a_t | s_t))$$

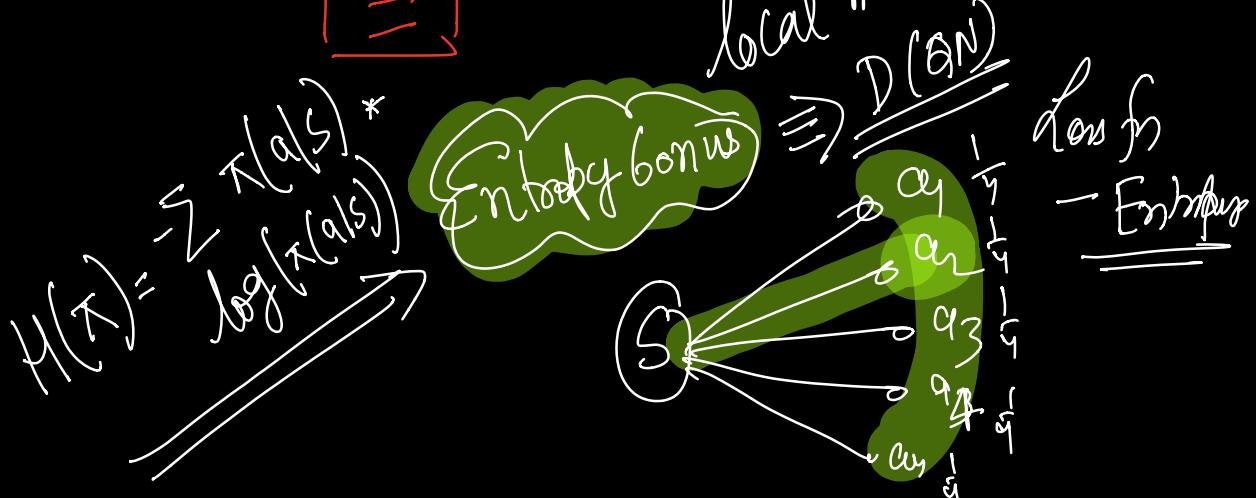
$$\sum_{t \geq 0} \left( \sum_{t' \geq t} \gamma^{t-t'} \gamma_{t'} \right) * (\text{Same})$$

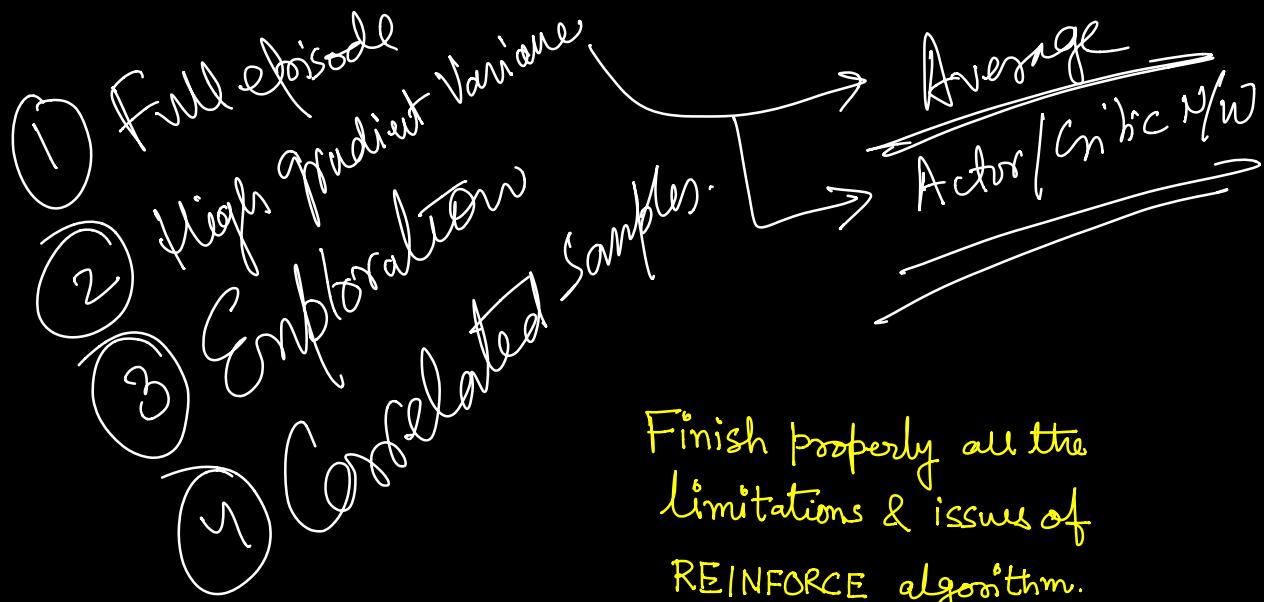
$\gamma(z_1)$   
 $\gamma(z_2)$   
 $\gamma(z_3)$

Playing game



\* Exploitation is limited and it may stuck in local minima.





Finish properly all the  
limitations & issues of  
REINFORCE algorithm.

## i) Algorithm for Policy Gradient (REINFORCE)

① Initialize the network parameters  $\theta$ .

② Generate ( $N$ ) trajectories  $\{z^i\}_{i=1}^N \left( z^1, z^2, \dots, z^N \right)$   
 following the policy  $\pi_\theta$ .  
looping

③ Compute the return of the trajectory  $R(z)$

④ Compute the gradient:

$$\nabla_\theta J(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) * R(z) \right]$$

⑤ Update the network parameter

$$\theta = \theta + \alpha \nabla_\theta J(\theta)$$

⑥ Repeat steps ② to ⑤ for several iterations.

( $\theta$ ) getting updated after each iteration. Hence  
 policy got changed (become better) and then  
 generate new data / episode / trajectories for  
 training.  
 It is an On-policy method.

## How to manage high Variance in the Gradient

\* Actually we are using policy to generate the trajectory and then computing  $\nabla_{\theta} J(\theta)$  to update the policy itself.

↳ This will in turn improve the policy after each iteration.

↳ Hence returns vary greatly introducing high variance in the gradient updates.

### ii) Policy gradient with reward-to-Go

for Vanilla PG

$$\nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) * R(z) \right]$$

Proper Credit Assignment

$$\text{Where } R(z) = \sum_{t=0}^{T-1} r_t$$

let us say Reward-to-Go defined as ( $R_t$ )

Sum of the rewards of the trajectory starting from the state ( $s_t$ ).

$$\therefore R_t = \sum_{t'=t}^{T-1} r(s_{t'}, a_{t'})$$

Instead of  $R(z)$  use  $R_t$

An action is only  
as good as the DCFR  
it can generate

### (iii) Policy gradient with the baseline

→ Trajectory lengths are different.

→ Favoring early rewards

Can we Normalize the Reward to Go

↳ In order to reduce the variance  
one can subtract some baseline  $(b)$   
from the  $(R_t)$  (Just like DC-  
Component/Average)

Can we develop  
the notion of  
"Expected"

$$\nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \times (R_t - b) \right]$$

Baseline  $(b)$ : It is the value that can give  
us the expected return from the state the agent  
is in.

Simplest baseline:  $b = \frac{1}{N} \sum_{i=1}^N R(z)$  {OR  
even moving average}

Baseline can be any function but should not be dependent upon  $(\theta)$ . (Does not affect N/W form)

Other obvious choices can be

↳ Value function ( $V(s_t)$ )

Value of a state is the expected return an agent would obtain starting from the state following  $(\pi)$ .

↳ Other options are Q fn & Advantage fn.

How Can we learn the Baseline functions.

↳ Just like approximating policy using  $(\theta)$  network

↳ One can use value network  $(\phi)$

$$\nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (R_t - V(s_t))$$

Now since the value of a state is floating point number,  $\phi$  can be trained by minimizing MSE.

$(R_t)$   $\Rightarrow$  Actual Return

$V(s_t)$   $\Rightarrow$  Predicted Return

$$\therefore J(\phi) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} (R_t - V_\phi(s_t))^2$$

We need to minimize  $J(\phi)$  hence

$$\phi = \phi - \beta \nabla_\phi J(\phi)$$



### Algorithm – REINFORCE with baseline

The algorithm of the policy gradient method with the baseline function (REINFORCE with baseline) is shown here:

1. Initialize the policy network parameter  $\theta$  and value network parameter  $\phi$
2. Generate  $N$  number of trajectories  $\{\tau^i\}_{i=1}^N$  following the policy  $\pi_\theta$
3. Compute the return (reward-to-go)  $R_t$
4. Compute the policy gradient:

*looking*

$$\nabla_\theta J(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) (R_t - V_\phi(s_t)) \right]$$

5. Update the policy network parameter  $\theta$  using gradient ascent as  $\theta = \theta + \alpha \nabla_\theta J(\theta)$
6. Compute the MSE of the value network:

$$J(\phi) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} (R_t - V_\phi(s_t))^2$$

7. Compute gradients  $\nabla_\phi J(\phi)$  and update the value network parameter  $\phi$  using gradient descent as  $\phi = \phi - \alpha \nabla_\phi J(\phi)$
8. Repeat steps 2 to 7 for several iterations

Ⓐ Advantage fn ( $A_t^i$ )  $\doteq$   $t^{th}$  step of  $i^{th}$  episode.

$\rightarrow$  How "good" it is to take an action "a", at state "s".

when  $Q(s, a) \geqslant \text{Value}(s)$

(Action  $a$  is better than expected/average.)

$$A_t^i \Rightarrow Q(s_t, a_t) - V(s_t) \quad \text{for some policy } (\pi).$$

$$\nabla_{\theta} J(\theta) = \sum_{t \geq 0} \left[ Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t) \right]$$

$$\text{Scaling fn for log likelihood} \xleftarrow{*} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

\* Minimization of the advantage function will also enforces Bellman optimality.

\* Computation of Advantage fn needs ② neural networks (a) for Q fn (b) for value fn [This is not optimal]

This is computationally inefficient & expensive.

$Q$  fn, Can be

→ (a) Approximated via Monte Carlo based Random Sampling or [Reward-as-Go].

(I)  $Q(S, a) \Rightarrow$  All future reward policy itself says take action  $a$  at  $S$  (Doing it over the Experience)

(Major Concern is that, we need to have the full trajectory or trajectories)

→ (b) It can be approximated by using the definition of  $Q$ -fn and using Bellman optimality.

$$Q(S, a) \xrightarrow{(a)} R + \gamma V(S') \Rightarrow Q(S, a) = R + \gamma V(S')$$

This we can get from  $\phi$

$Q$  Value can be Computed using Value fn N/w.

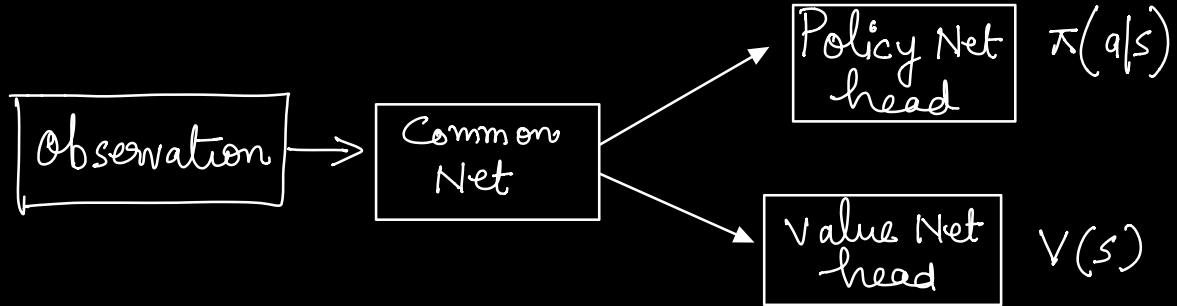
Now  $Q$  fn can be Computed for any step, instantly using, immediate reward

$\phi$  O/P of value fn.

Hence

$$A_t^i = A(S, a) = R + \gamma V(S') - V(S)$$

## Actor Critic Algorithm



Observation

Sample  $m$  trajectories under the current policy ( $\pi$ ).

$$m_1 \langle s_0^1, a_0^1, r_0^1 \rangle \langle s_1^1, a_1^1, r_1^1 \rangle \dots \langle s_t^1, a_t^1, r_t^1 \rangle$$

$$m_2 \langle s_0^2, a_0^2, r_0^2 \rangle$$

$$m_3 \langle s_0^3, a_0^3, r_0^3 \rangle$$

:

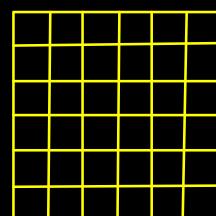
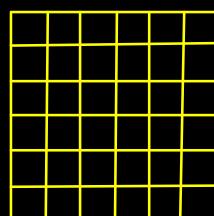
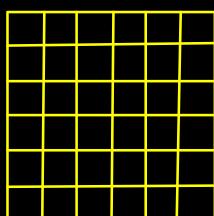
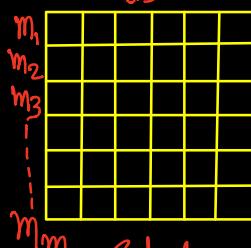
:

:

$$m_m \langle s_0^m, a_0^m, r_0^m \rangle$$

$$t=1 \ t=2 \ t=3 \ \dots \ t=t$$

$$\langle s_t^m, a_t^m, r_t^m \rangle$$



State

Action

Reward

Advantage

(s)

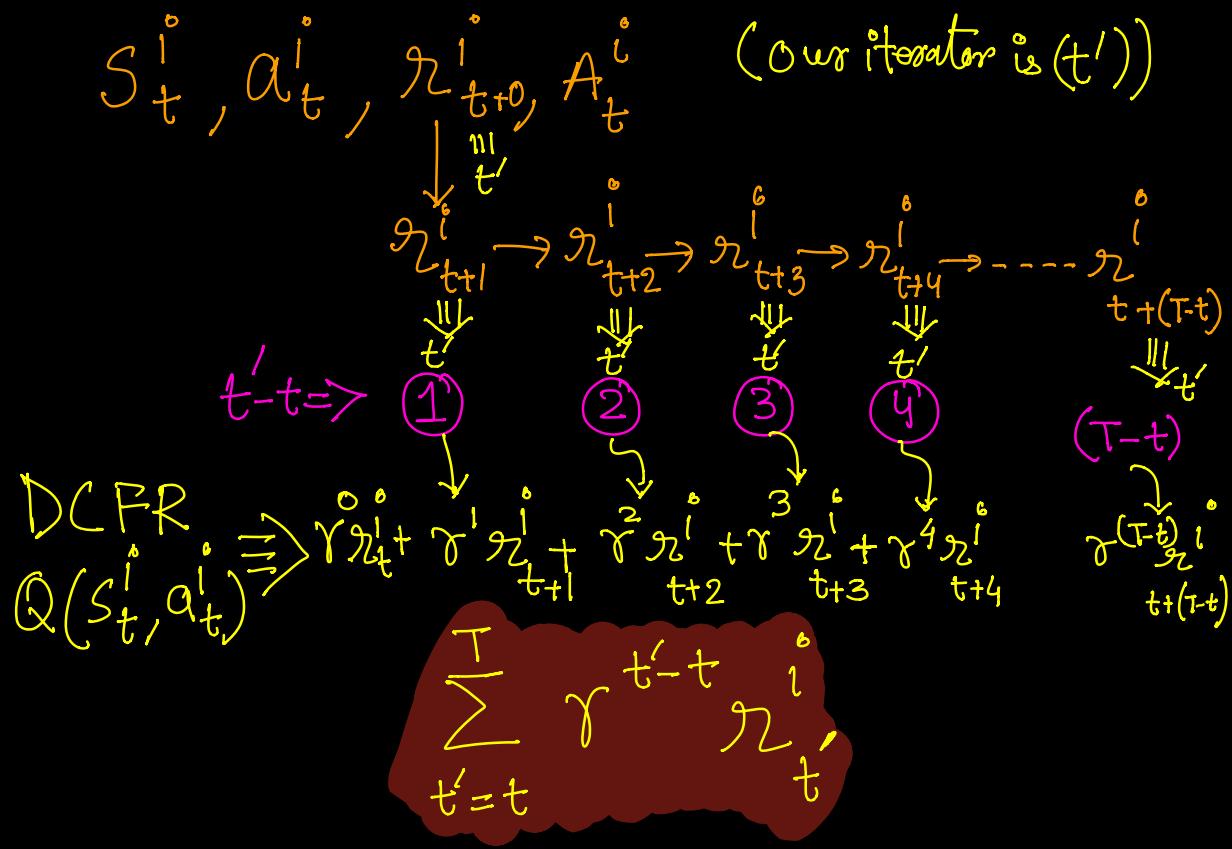
(a)

(r)

(A)

$s_t$

⇒ State of the agent at  $t^{th}$  Step  
in the  $i^{th}$  episode.



## ALGORITHM (Actor/Critic)

- ① Initialize the
  - Policy w/  $[\theta]$  (Actor)
  - Q learning w/  $[\phi]$  (Critic)
- ② For each training iteration =  $(1, 2, \dots)$  do
  - After every iteration policy will be updated. New data need to be generated.
  - ②.1 Sample  $m$  trajectories under the current policy  $\pi$ .  $(S, a, r, t)$

Q.2

$$\Delta \theta = 0$$

Gradient accumulator for Actor N/W.  
for a given policy ( $\pi$ ), before training reset.

Q.3

for  $i = 1, 2, \dots, m$  do

for each episode / trajectory

Q.3.1

for  $t = 1, 2, \dots, T$  do

for each step ( $t$ )

Advantage fn  
for the  
( $t^m$ ) step of  
( $i^m$ ) episode

$$Q.3.1.1 A_t^{(i)} = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}^{(i)} - V(s_t^{(i)})$$

$\underbrace{t' = t}_{\text{all future rewards}}$   $\underbrace{(R_t)}$   $\underbrace{\emptyset}_{\text{Critic network}}$

Minimization of  $A_t^{(i)}$ , enforces Bellman Eq.

Getting the  
updates from

the  $i^m$  episode's  
 $t^m$  time step.  
scaled by ( $A_t^{(i)}$ )

Q.3.1.2

$$\Delta \theta = \Delta \theta + A_t^{(i)} \nabla_\theta \log(\pi_t^{(i)} | s_t^{(i)})$$

Accumulating  
all the policy  
updates for (Actor N/W).

Gradient  
ascent

Gradient estimator  
step.

Q.4

$$\Delta \phi =$$

$$\sum_i \sum_t \nabla_\phi \|A_t^{(i)}\|^2$$

Minimizing the advantage  
fn, Discounted accumulated  
future rewards getting  
close to the predicted  
value function.

for all episodes  
and for all  
time steps, just  
accumulate them  
all.

Critic N/W  $A_t$  minimization  
Update enforces Bellman  
Constraint.

$$2.5 \quad \theta = \theta + \alpha \Delta \theta$$

$$2.6 \quad \phi = \phi - \beta \Delta \phi$$

③ End for

old problem of  
deep Q-learning:  
\* Learn Q-values  
for all  $\langle \text{State}, \text{action} \rangle$   
pairs.

Critic only learns the  
Q-values for generated/  
observed  $\langle S, a \rangle$  pairs by  
the policy under consideration.

Also incorporate  
standard  
experience replay  
trick too.

Are these algorithms Sample efficient (as they are)  
on-policy

\* Instead of generating the full trajectory  
and then only Compute the Returns/Reward,  
Can we do something efficient.

\* We can approximate  $(R_t)$  or  $Q(S, a)$  also

$$R \approx r + \gamma(V(S'))$$

\* Now, we don't need to wait till the end of  
episode, to Compute reward.

\* One can Compute gradient at each step &  
update the policy network parameters. [At each  
Step]

Let us see the flow : 

- ① Initialize,  $\Theta$  &  $\phi$
- ② Get a random starting state ( $s_0$ )
- ③ Get a action  $\pi(a_0|s_0)$
- ④ Get a reward ( $r_0$ ) and next state ( $s_1$ )
- ⑤ Get the Value of ( $s_1$ ) as  $V_\phi(s_1)$
- ⑥ Get  $G(s_0, a_0) = r_0 + \gamma * V_\phi(s_1)$
- ⑦ Get Advantage fn,  $A(s_0, a_0)$  as  

$$A(s_0, a_0) = r_0 + \gamma * V_\phi(s_1) - V_\phi(s_0)$$
- ⑧ Get Policy Gradient :  $\left[ \nabla_\theta J(\theta) \right]$   

$$\nabla_\theta J(\theta) = \nabla_\theta \log \pi_\theta(a_0|s_0) (r_0 + \gamma * V_\phi(s_1) - V_\phi(s_0))$$

- ⑨ Update Actor N/W ( $\theta$ ) :  $\theta = \theta + \alpha \nabla_\theta J(\theta)$   
 (See  $N\theta(\Delta\theta)$  accumulation) (Gradient Ascent)
- ⑩ Get Critic N/W loss :  $J(\phi) = A(s_0, a_0) = r_0 + \gamma * V_\phi(s_1) - V_\phi(s_0)$
- ⑪ Get  $\nabla_\phi J(\phi)$ , Value Gradient and update  
 Critic N/W (using gradient descent) :  $\phi = \phi - \alpha \nabla_\phi J(\phi)$

# The actor-critic algorithm

The steps for the actor-critic algorithm are:

1. Initialize the actor network parameter  $\theta$  and the critic network parameter  $\phi$
2. For  $N$  number of episodes, repeat step 3
3. For each step in the episode, that is, for  $t = 0, \dots, T-1$ :
  1. Select an action using the policy,  $a_t \sim \pi_\theta(s_t)$
  2. Take the action  $a_t$  in the state  $s_t$ , observe the reward  $r$ , and move to the next state  $s'_t$
  3. Compute the policy gradients:

$$\nabla_\theta J(\theta) = \nabla_\theta \log \pi_\theta(a_t | s_t) (r + \gamma V_\phi(s'_t) - V_\phi(s_t))$$

4. Update the actor network parameter  $\theta$  using gradient ascent:

$$\theta = \theta + \alpha \nabla_\theta J(\theta)$$

Both updates  
are done at  
each &  
every step of  
the episode

5. Compute the loss of the critic network:

$$J(\phi) = r + \gamma V_\phi(s'_t) - V_\phi(s_t)$$

Sample  
efficient.

6. Compute gradients  $\nabla_\phi J(\phi)$  and update the critic network parameter  $\phi$  using gradient descent:

$$\phi = \phi - \alpha \nabla_\phi J(\phi)$$