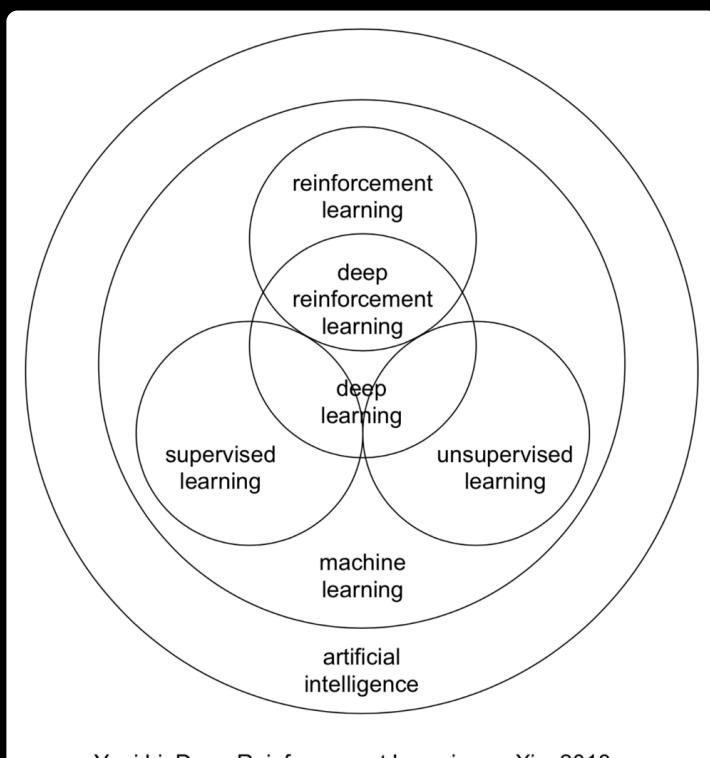
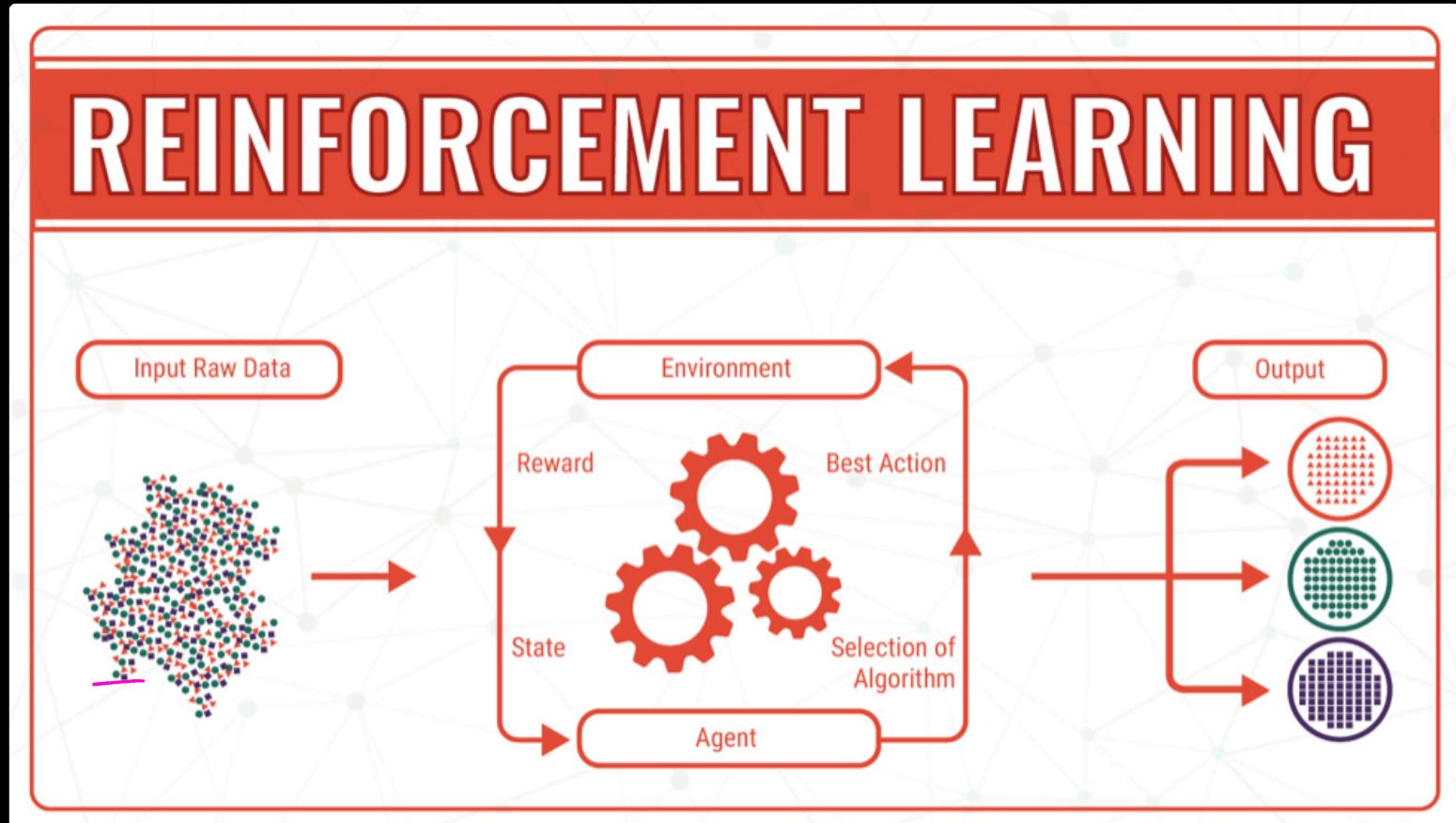


This lecture is completely adapted from Stanford lecture series : CS-231n by Serena Yeung



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 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

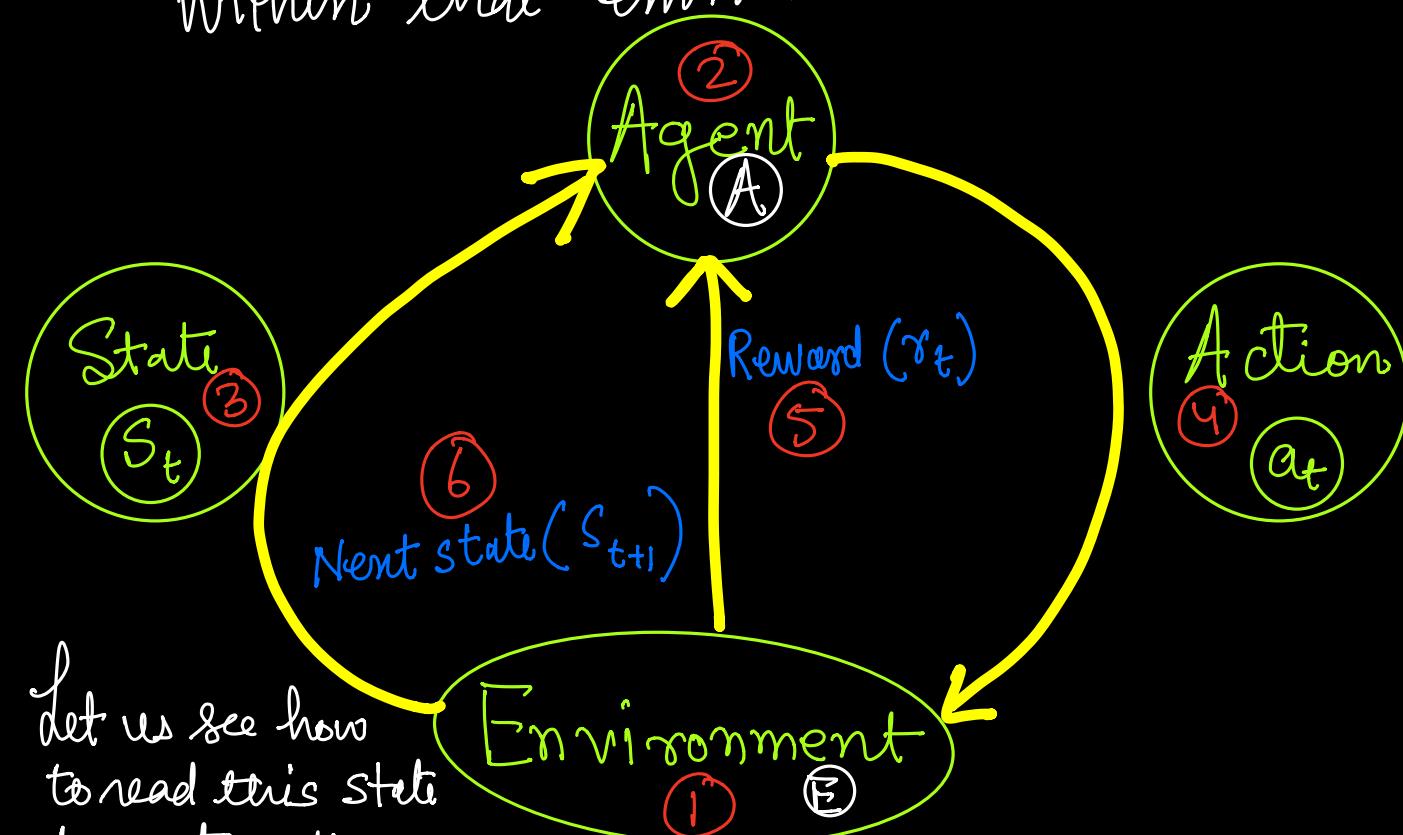
II) Unsupervised Learning

- 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.

III) Reinforcement learning : It is a learning paradigm in which learning happens by explorations. So agent interact with the environment repeatedly w/o any prior knowledge and labels wrt the environt. — 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 .

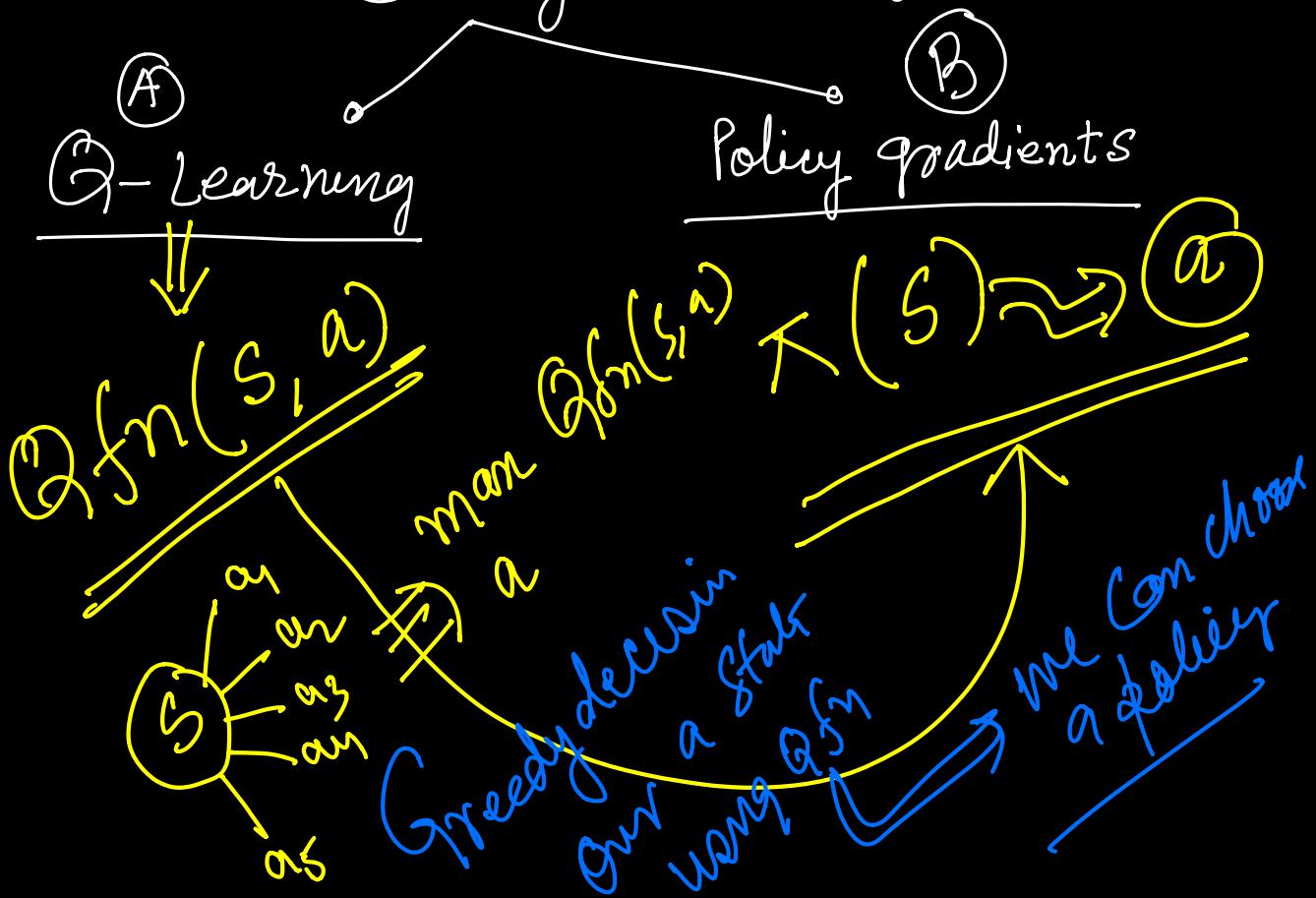
* The final Goal of the agent is :

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

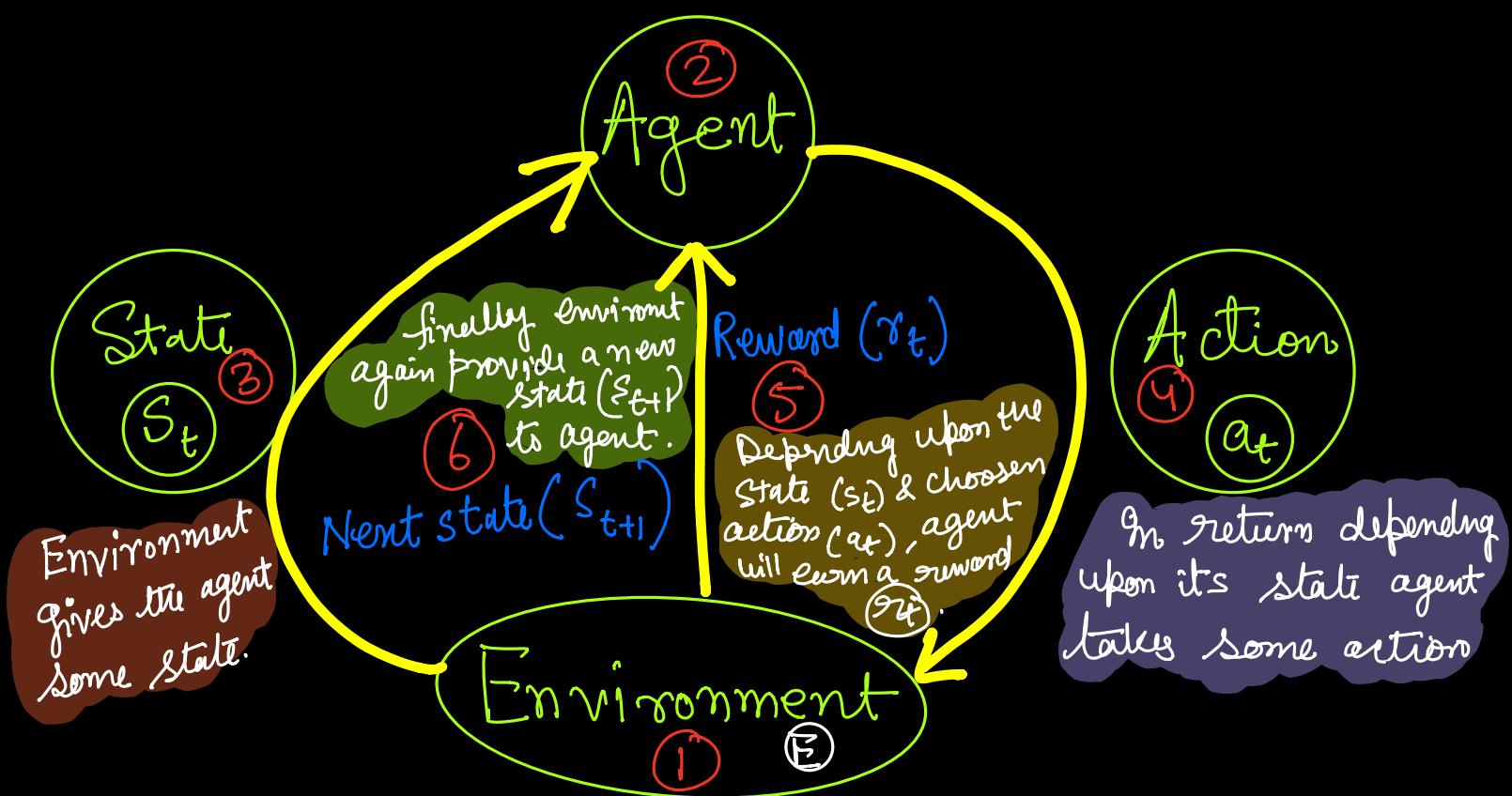
(Future discounted reward maximization)

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

* There are 2 major classes of RL algorithms



B → Basic RL-Setup/framework

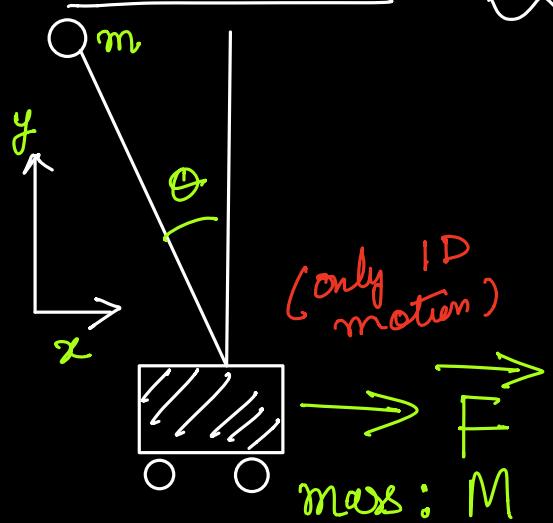


Episodic training :

- 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 (θ), angular speed (ω), 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:

$\langle S, A, P, R, \gamma \rangle$

Set of all possible states:
 s_1, s_2, \dots, s_n

Set of all possible actions:
 a_1, a_2, \dots, a_k

State transition Matrix:

$\langle s_t, a_t \rangle$

s_{t+1}

Reward fn.:
 $\langle s_t, a_t \rangle \rightarrow r_t$

$\langle s_t, a_t \rangle$



(How much current rewards are important as compared to future.)

Discount factor
(0-1)

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.

$$\text{P}^a_{S S'} = P[S_{t+1} = S' | S_t = S, a_t = a]$$

C. 2.2

Immediate reward

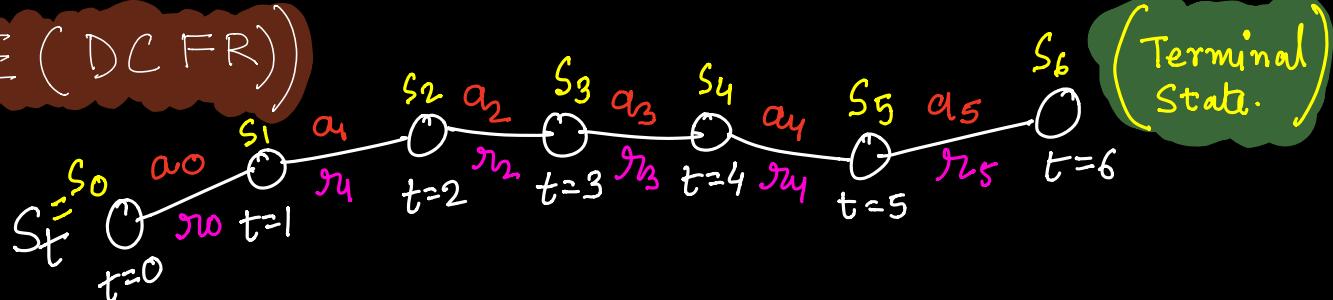
Amount of reward an agent is going to get from the environment when at state s_t agent takes an action a_t .

👉 $r_{s_t}^{a_t} = r(s_t, a_t)$

Greedy approach

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

$$\max(\mathbb{E}(\text{DCFR}))$$



How MDP Operates and Can be Used to Obtain Episodes.

MDP- Execution

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

$$S_0 \sim p(S_0)$$

Initial State

Probability distribution (ISPD)

- ② from $t=0$ (until terminates)

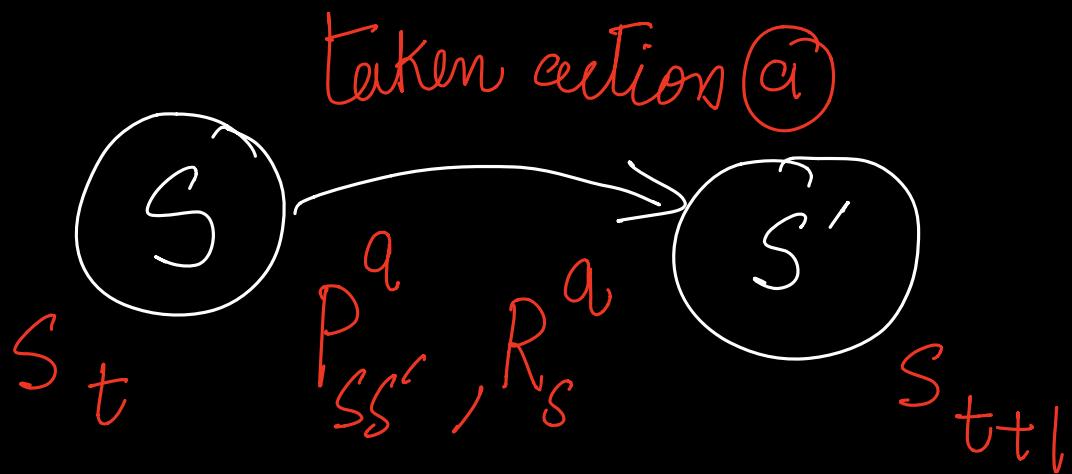
- ②.1 An agent can select an action (a_t) [How!!!]
using some policy that we need to figure out so as to $\max(\mathbb{E}(\text{DCF}))$

Q.2) Depending upon the agent's state s_t and the action taken a_t environment will return the immediate reward r_t as well as next state s_{t+1} using P & R .
 ∴ for a MDP P & R governs.

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

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

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



C. 2.3 But how to take an action \Rightarrow Policy π

A policy π 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

$$as \quad \pi(a_t | s_t) = P[a_t | s_t]$$

C. Q. 4

Goal of MDP

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

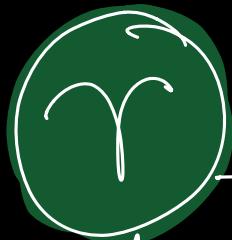
* The actual returns are always in the form of DCFR.

G_t = Discounted Cumulative future Reward from time Step t . [future rewards are geometrically weighted]
 $= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$

for all upcoming future time step

$$= \sum_{K=0}^{\infty} \gamma^k \mathcal{R}_{t+K} \quad (\gamma \in (0,1))$$

{Gamma}



→ close to 0

→ close to 1

"far-sighted"

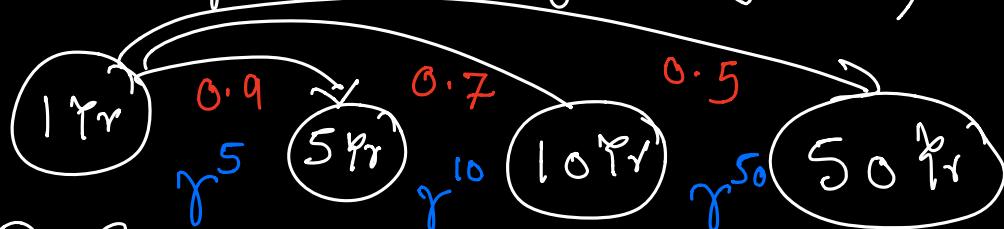
"Myopic" evaluation.

More important to immediate returns.

Care about all future rewards

② Requirement of discount :

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



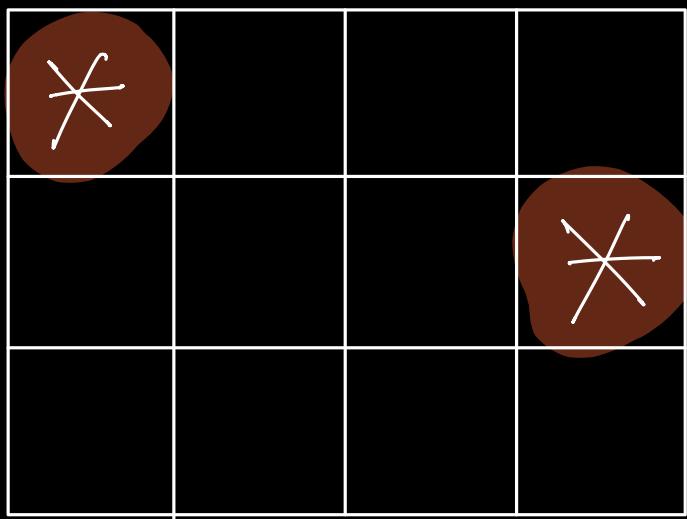
$$\gamma \in \{0 - 1\}$$

b) 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

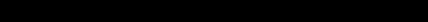
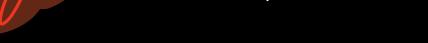
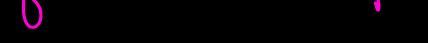
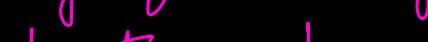
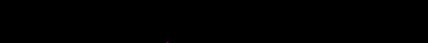
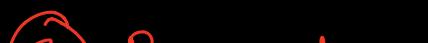
A) States (2D-locations)



B) Actions = {

right, left, up, down}.

$\rightarrow \leftarrow \uparrow \downarrow$



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 (π^*)

Policy \Rightarrow Planning

It donot matter that in whatever state (say s_t) agent is at time t , π^* is going to suggest that "what action (a_t) agent needs to take in order to get the DCFR Maximized." $\forall a_t \in \pi(a_t | s_t)$

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

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

Annotations:

- π^* is highlighted with a blue circle.
- A yellow box contains "Returning the best policy".
- A green box contains "maximizing DCFR".
- A brown box contains "Given a policy".
- A green box contains "where $S_0 \sim p(S_0)$ ".
- A green box contains "Enforcing overall Eta policies".

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

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

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

* There are few attributes associated to a policy:

$\pi(a_t | s_t)$ → Value fn - wrt a given policy (π) , How good is a given state (s) , i.e if agent just follow π

Q-value fn → How much DCFR is expected is the value of that state (s) for a given policy (π) .

wrt a given policy (π) , How good is a state action pair. i.e

Instead of directly following the policy (π) at state (s) , agent first takes an action (a) and then follow the policy (π) . Qvalue of (s, a) is the expected DCFR.

D.I

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

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

Given a policy and a initial state $s_0 \sim p(s_0)$; 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

$\leftarrow (s_1, a_1, r_1) \leftarrow (s_2, a_2, r_2)$

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

Value of state s w.r.t a given policy π

Expected DCFR from state s following the policy π 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 $\pi(s_2)$

s_1, a_1	s_2, a_1	s_3, a_1	$\frac{1}{\pi(s_2)}$
s_1, a_2	s_2, a_2	s_3, a_2	$\frac{2}{\pi(a_1 s_2) \& \pi(a_2 s_2)}$

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

action ① at state ⑤ & then following ⑧.

The expected DCFR, if at state ⑤, ① action has been chosen.

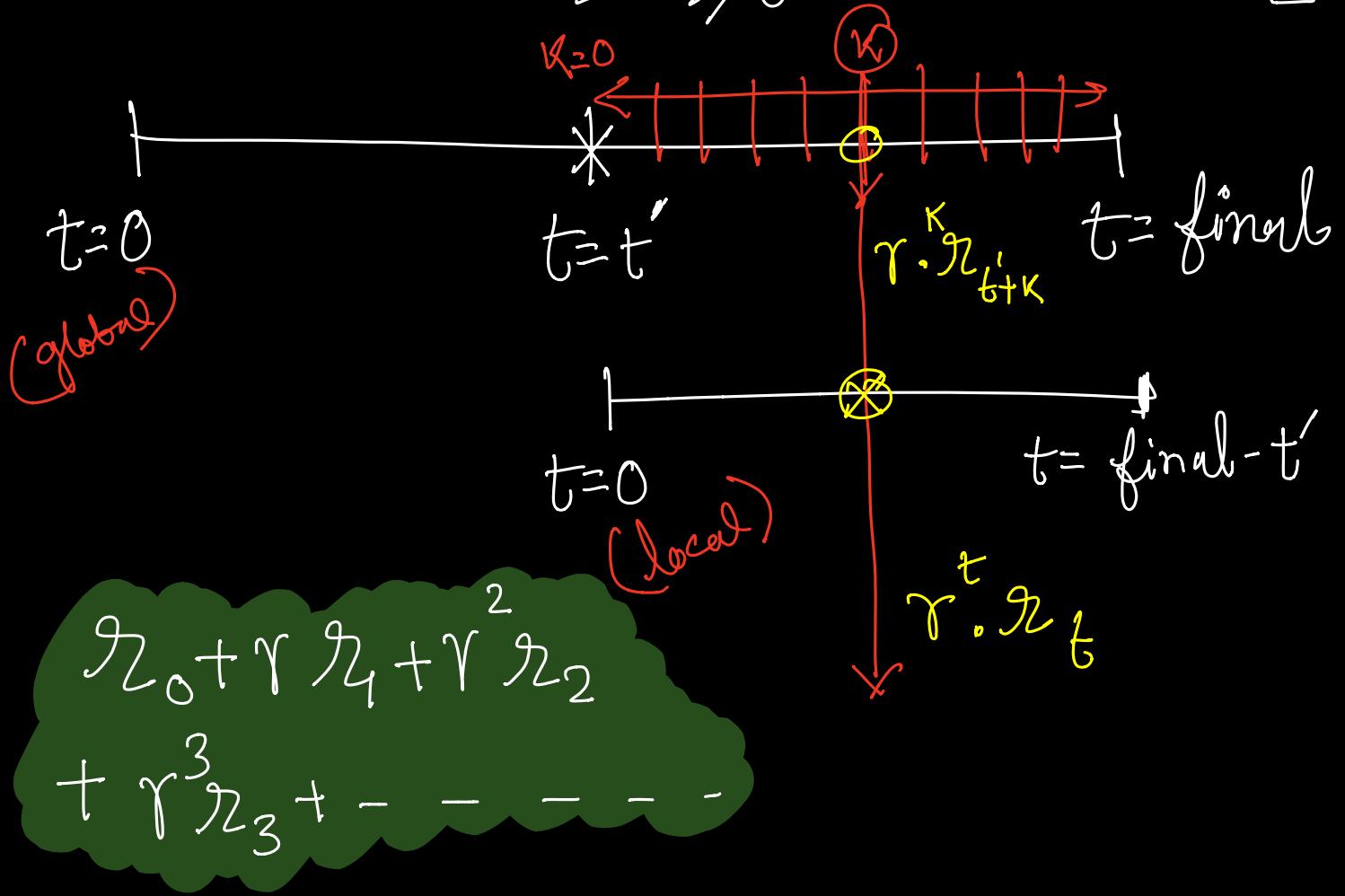
$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid \begin{array}{l} s_t = s \\ a_t = a \\ \pi \end{array} \right]$$

Used for policy optimization

at State ⑤
agent took
action ① then follow ⑧

 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 $V^*(s)$

↳ maximum value function
Over all the policies

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

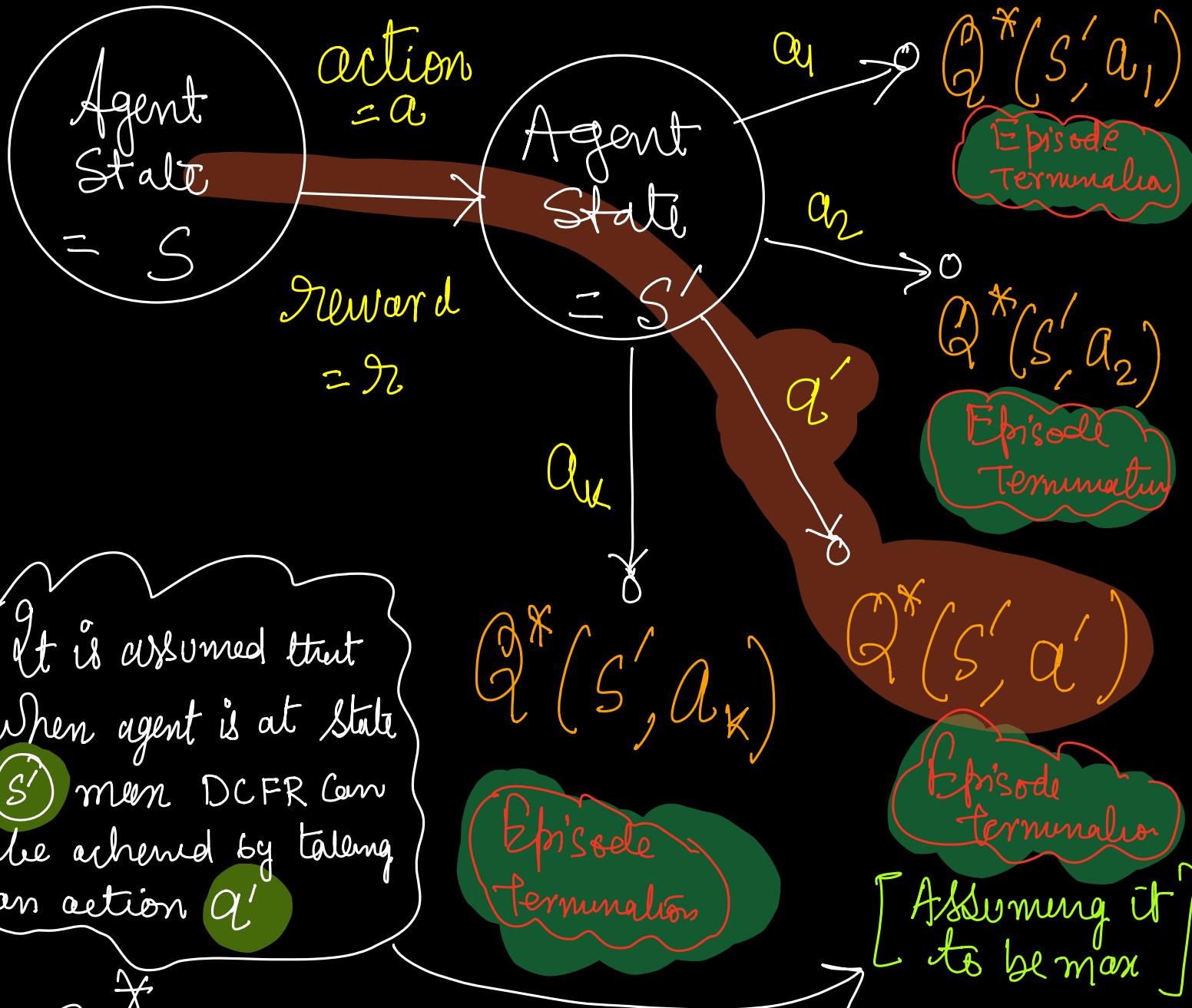
* Optimal action-value of Q-value fn

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

Maximizing over all policies

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

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



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

$Q^*(s', a_1), Q^*(s', a_2)$

$\dots \dots \dots Q^*(s', a_k)$

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

Best Q-value over all actions

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

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

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

it may reach
 to a set of
 states (E)

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

Q^* as well as V^* both functions
 satisfies Bellman equations and are
 defined recursively.

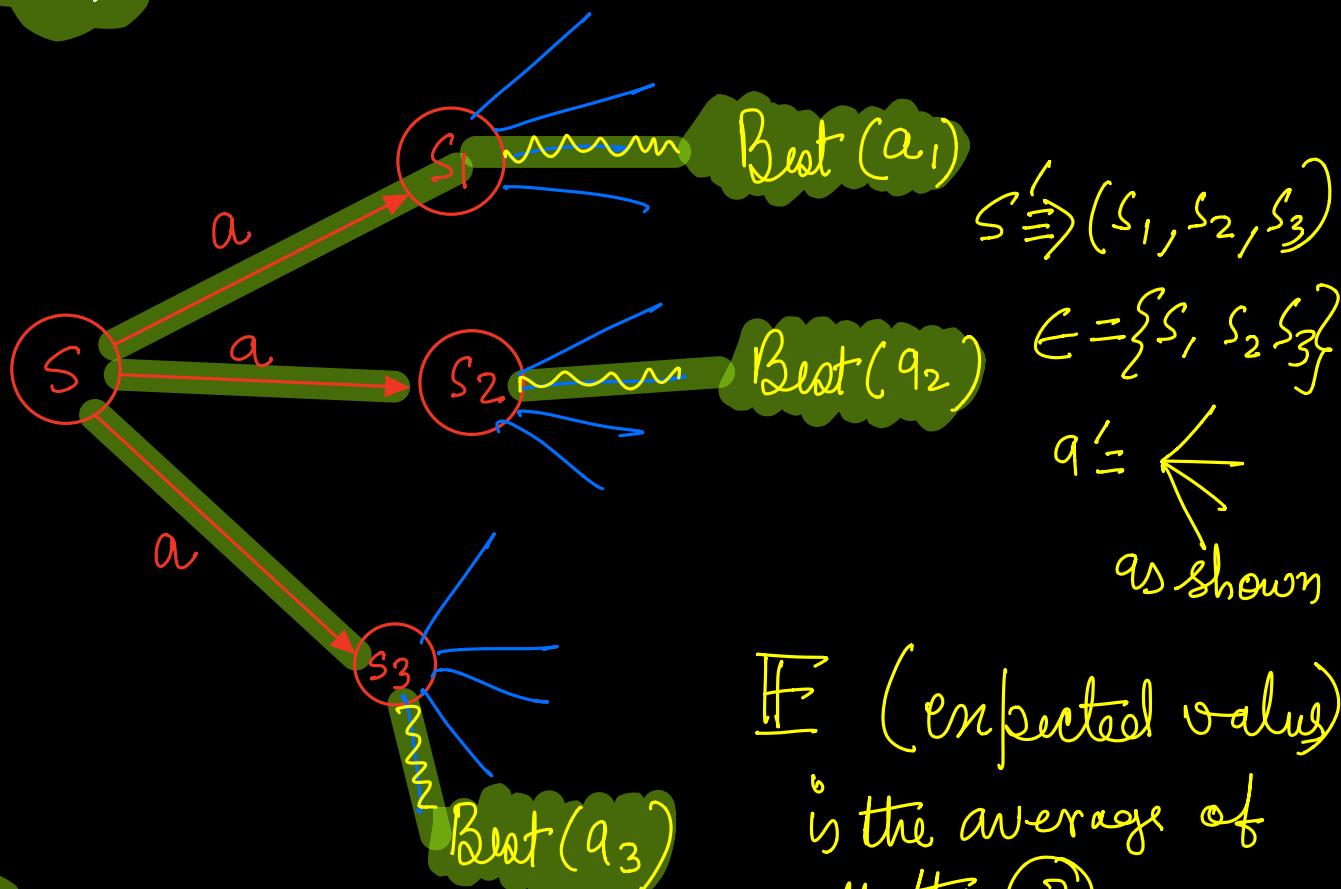
* Expected values are considered in order
 to address the randomness in the process.

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

$$Q^*(s, a) = \mathbb{E}_{s' \in \mathcal{S}' | s, a} [r + \gamma \mathbb{E}_{s'' \in \mathcal{S}'' | s', a'} [Q^*(s'', a'')]$$

it may reach to a set of states (\mathcal{E})

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



\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 (π^*)

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
enables suitable
state-action space
exploration

These p_i 's we are deciding under policy π^* assuming to be given. Later we need to compute them 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)

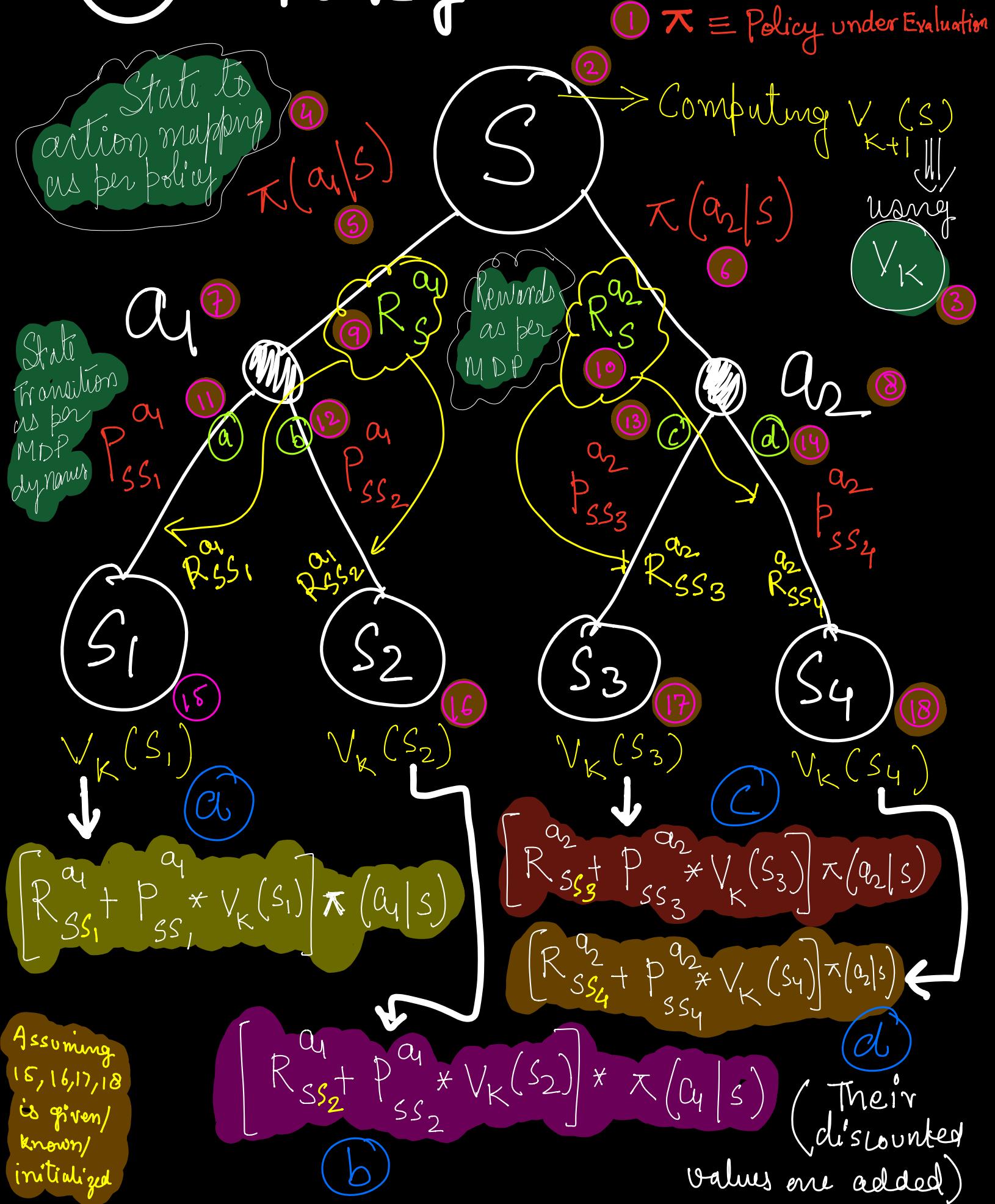
* firstly we need to figure out how to evaluate a policy (Policy Evaluation)

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^{th} step are optimal and will be used to compute the policy/value for $i+1^{th}$ step. [Using Dynamic Programming]

E.I → Policy Evaluation



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

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

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^{th} iteration

Let us evaluate a random policy in small grid world.

E.2

Policy Evaluation: Example

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 $\{Q(s, a) + \gamma q\}$

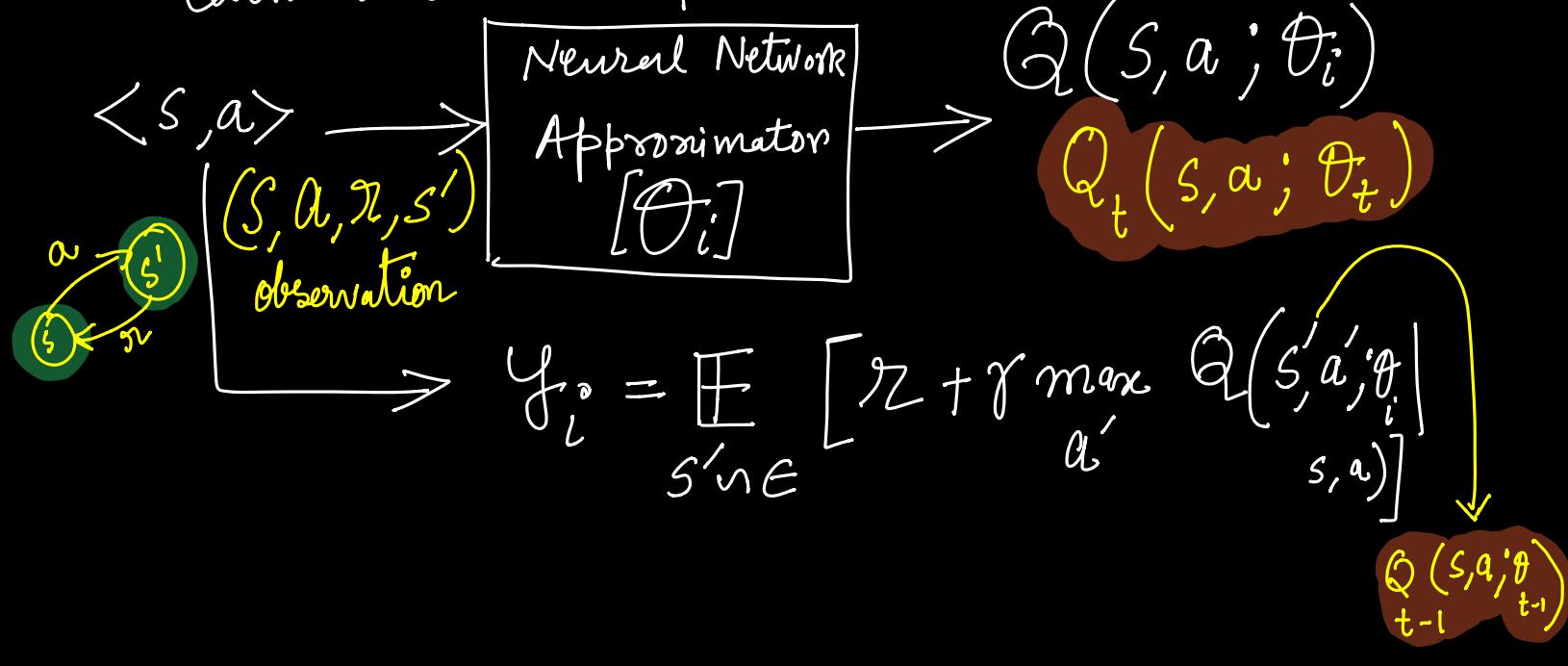
$$Q(s, a; \theta) \approx Q(s, a)$$

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

Deep Q-Learning

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

(Intuition)



G Solving for Optimal Policy (Directly)

Q.1 Introduction to policy gradient

* Instead of using DNN to approximate Q-Value function, why can't directly learn the suitable policy (π), parameterized by θ .

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

for any policy (π) . define its value as:

$J(\theta) = \mathbb{E} \left[\sum_{t=0}^T r_t \mid \pi_\theta \right]$

↓
Expected DCFR
while following policy π_θ

↓
Value of the policy π
Parameterized by θ

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

GOAL: DNN need to optimize for the parameter θ that can realize the optimal policy π^* .

$$\theta^* = \arg \max_{\theta} J(\theta) \text{ Gradient Ascent over } \theta \text{ in REINFORCE algo.}$$

A

Deep Q-Learning

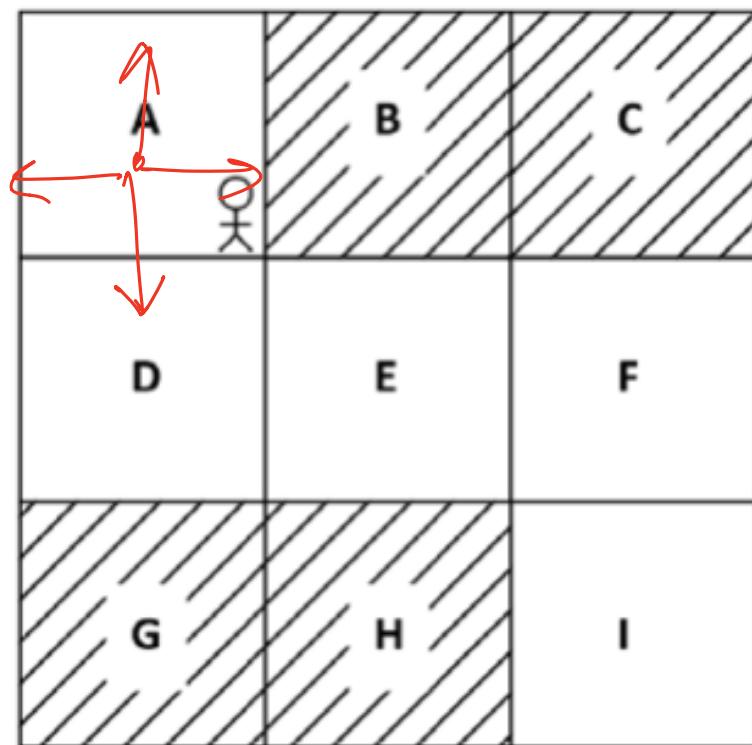


Figure 9.1: Grid world environment

State	Action	Value
A	up	17
A	down	10
B	up	11
B	down	20

Table 9.1: Q-value of state-action pairs

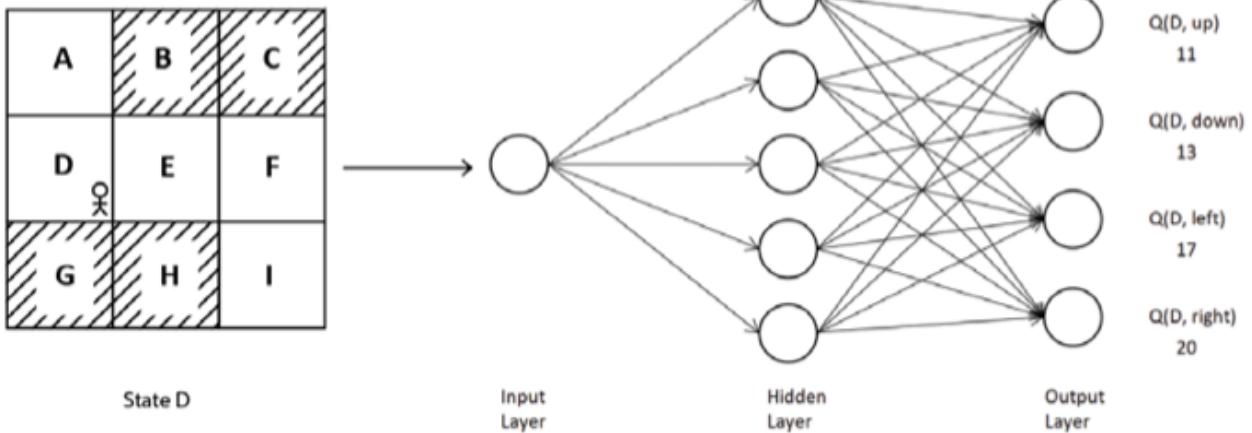


Figure 9.2: Deep Q network

DQN - Learning Q fn.

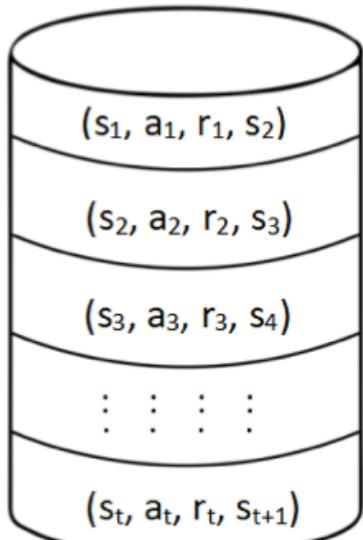


Figure 9.3: Replay buffer

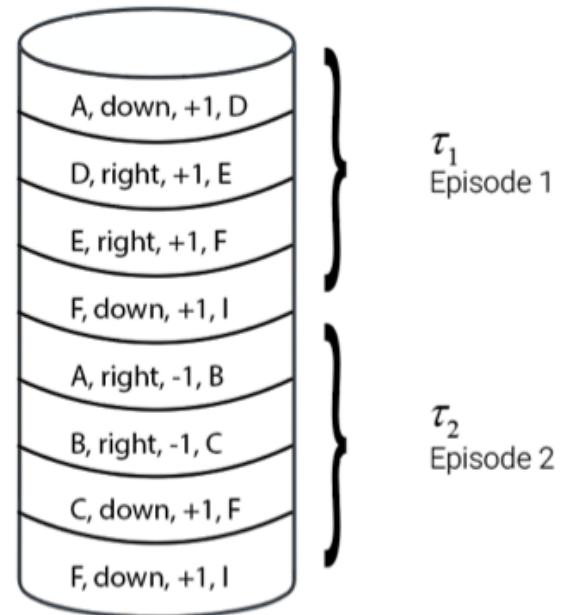


Figure 9.6: Replay buffer

the replay buffer \mathcal{D} .

1. Initialize the replay buffer \mathcal{D} .
2. For each episode perform step 3.
3. For each step in the episode:
 1. Make a transition, that is, perform an action a in the state s , move to the next state s' , and receive the reward r .
 2. Store the transition information (s, a, r, s') in the replay buffer \mathcal{D} .

1

$$\text{MSE} = \frac{1}{K} \sum_{i=1}^K (y_i - \hat{y}_i)^2$$

2

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

3

$$Q^*(s, a) = \underset{s' \sim P}{\mathbb{E}} [r + \gamma \max_{a'} Q^*(s', a')]$$

4

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

5

$$L(\theta) = Q^*(s, a) - Q_\theta(s, a)$$

6

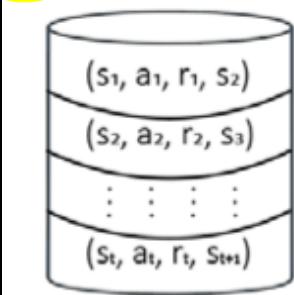
$$L(\theta) = r + \gamma \max_{a'} Q(s', a') - Q_\theta(s, a)$$

7

$$L(\theta) = r + \gamma \max_{a'} Q(s', a') - Q_\theta(s, a)$$

How do we compute this?

8



Randomly Sample
Minibatch of k transitions

$$L(\theta) = \frac{1}{K} \sum_{i=1}^K (r_i + \gamma \max_{a'} Q_\theta(s'_i, a') - Q_\theta(s_i, a_i))^2$$

using same n/w

Figure 9.7: Loss function of DQN

9

Target N/W

(Same n/w may introduce instability)

$$L(\theta) = \frac{1}{K} \sum_{i=1}^K (r_i + \gamma \max_{a'} Q_{\theta'}(s'_i, a') - Q_\theta(s_i, a_i))^2$$

Compute
using θ'

Compute
using θ

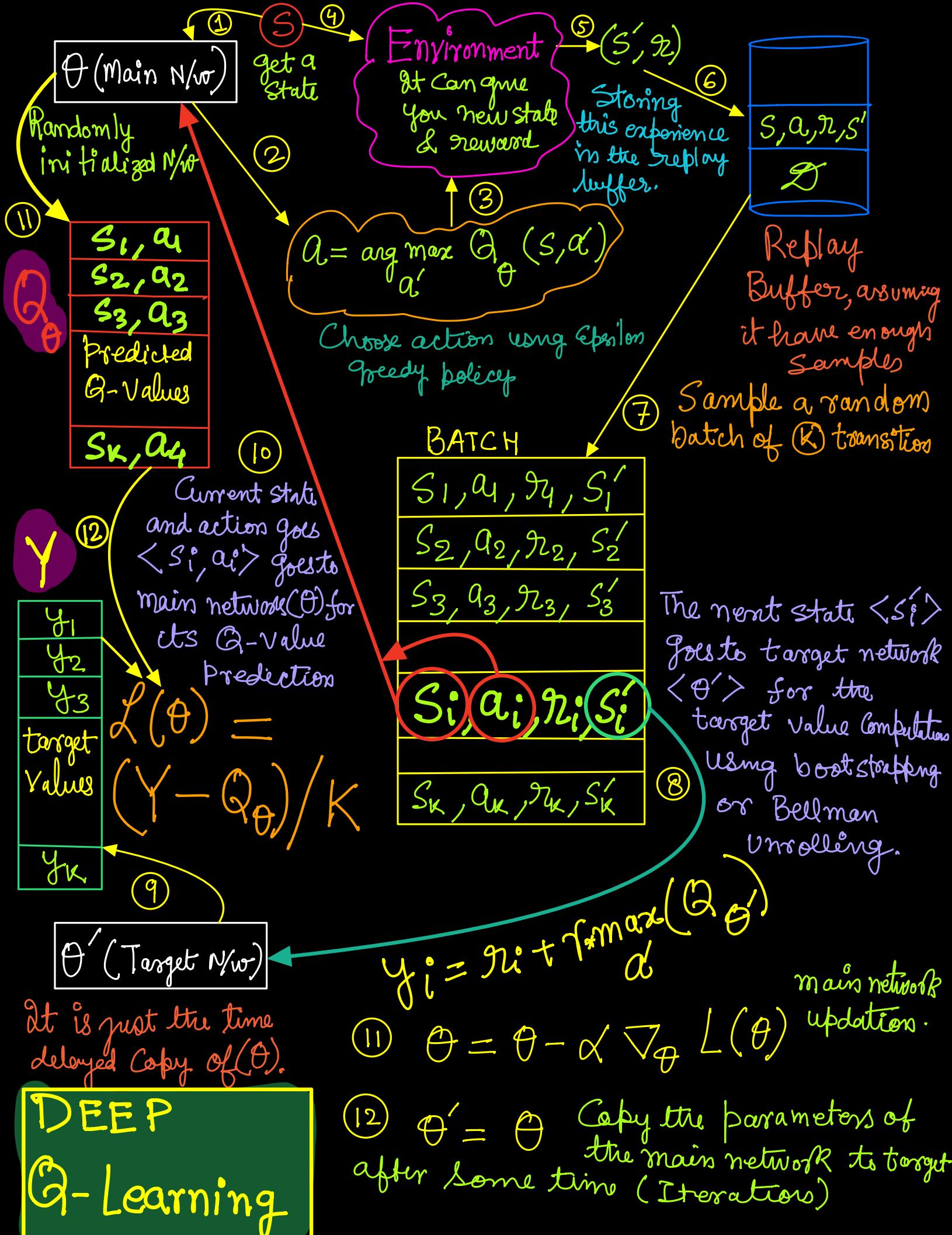
in the next section.

The DQN algorithm

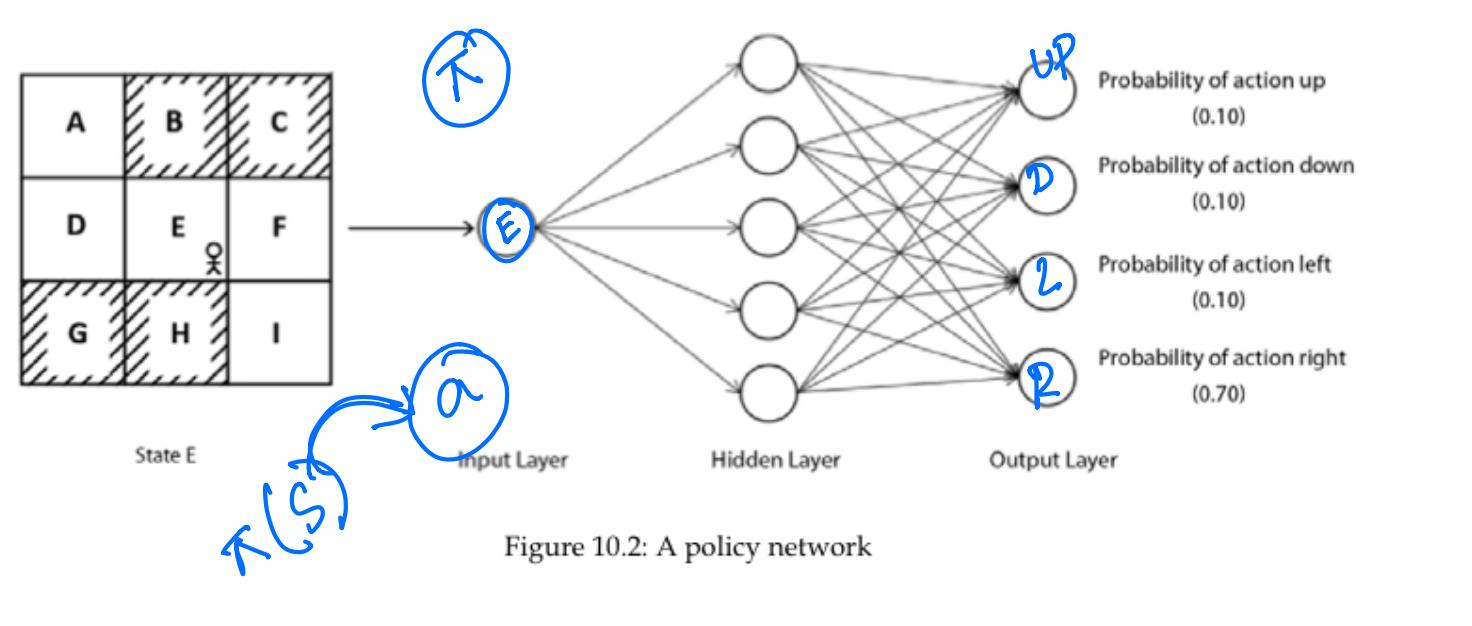
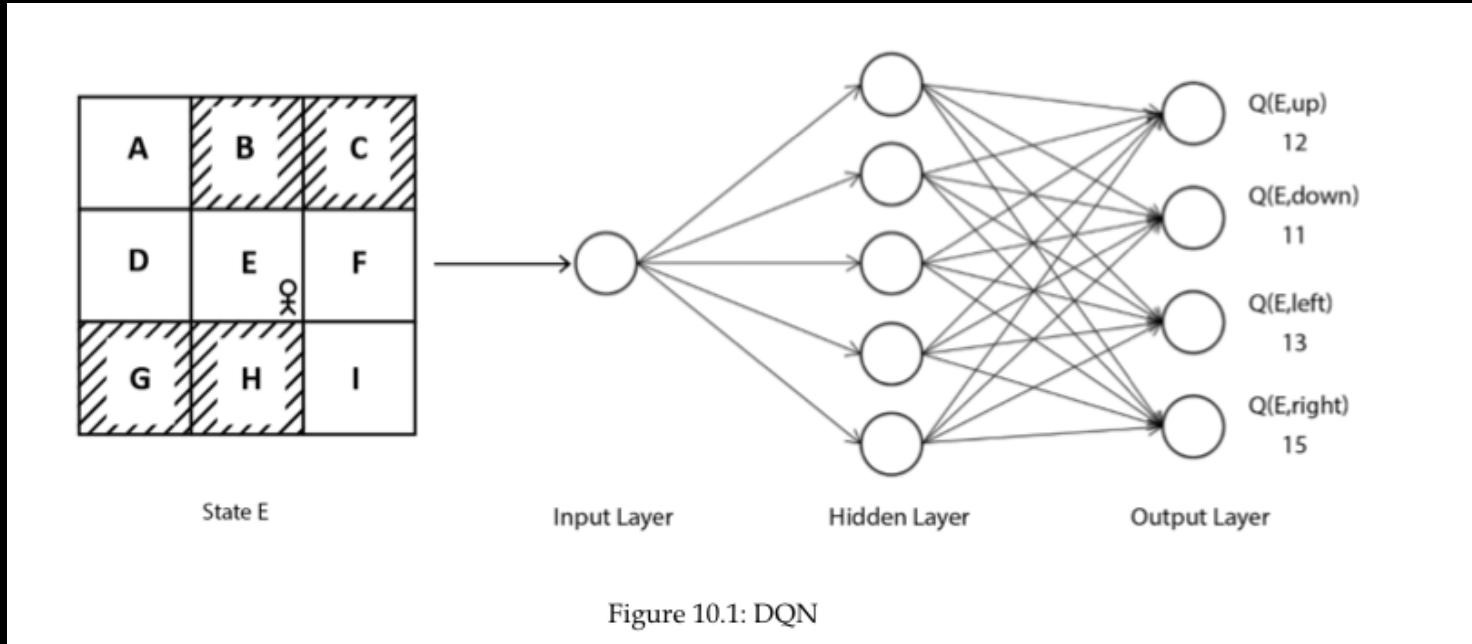
Final Algorithm:

The DQN algorithm is given in the following steps:

1. Initialize the main network parameter θ with random values
2. Initialize the target network parameter θ' by copying the main network parameter θ
3. Initialize the replay buffer \mathcal{D}
4. For N number of episodes, perform step 5
5. For each step in the episode, that is, for $t = 0, \dots, T-1$:
 1. Observe the state s and select an action using the epsilon-greedy policy, that is, with probability epsilon, select random action a and with probability 1-epsilon, select the action $a = \arg \max_a Q_\theta(s, a)$
 2. Perform the selected action and move to the next state s' and obtain the reward r
3. Store the transition information in the replay buffer \mathcal{D}
4. Randomly sample a minibatch of K transitions from the replay buffer \mathcal{D}
5. Compute the target value, that is, $y_i = r_i + \gamma \max_{a'} Q_{\theta'}(s'_i, a')$
6. Compute the loss, $L(\theta) = \frac{1}{K} \sum_{i=1}^K (y_i - Q_\theta(s_i, a_i))^2$
7. Compute the gradients of the loss and update the main network parameter θ using gradient descent: $\theta = \theta - \alpha \nabla_\theta L(\theta)$
8. Freeze the target network parameter θ' for several time steps and then update it by just copying the main network parameter θ



B Policy Gradient



In RL - we need to learn
the optimal policy

* In Continuous action space
(like Car driving, Robot movement)
DQN faces issues.

* θ for parameterizations (θ)

* Policy parameterizations (θ)

⇒ we directly learn
the policy.

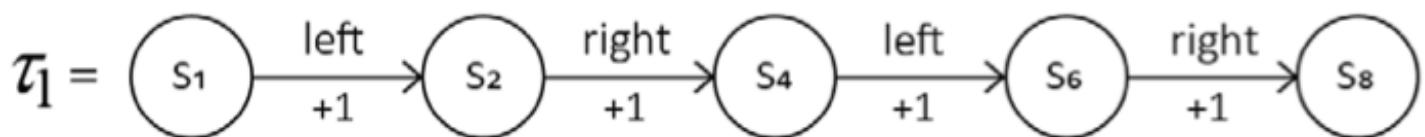


Figure 10.3: Trajectory τ_1

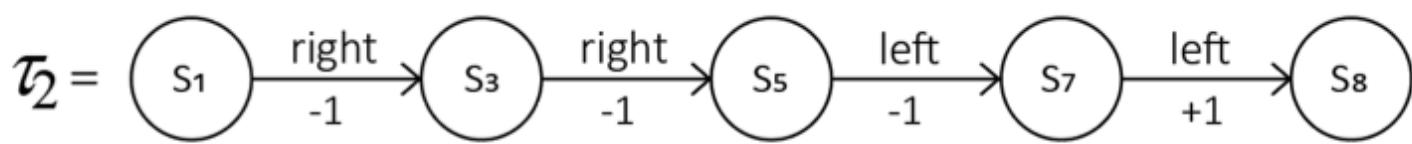


Figure 10.4: Trajectory τ_2

Intuition :

- ① In PG, we use NN to approximate the optimal policy π^* .
- ② Initialize θ randomly.
- ③ We feed state s as an input and it will return the action probabilities.
- ④ Initially as N/W is not trained Random O/P action probabilities we get.

- ⑤ Still we select the action based on the O/P action probability distribution.
- ⑥ Also store (S, A, R, S') (until the end of the episode).
- ⑦ This will become our training data
- ⑧ If agent WIN that episode \Rightarrow Assign high probabilities to all the actions of the episode for that- ⑤ Else reduce.

REINFORCE.

1. Initialize the network parameter θ with random values
2. Generate N trajectories $\{\tau^i\}_{i=1}^N$ following the policy π_θ
3. Compute the return of the trajectory $R(\tau)$
4. Compute the gradients

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

5. Update the network parameter as $\theta = \theta + \alpha \nabla_\theta J(\theta)$
6. Repeat steps 2 to 5 for several iterations

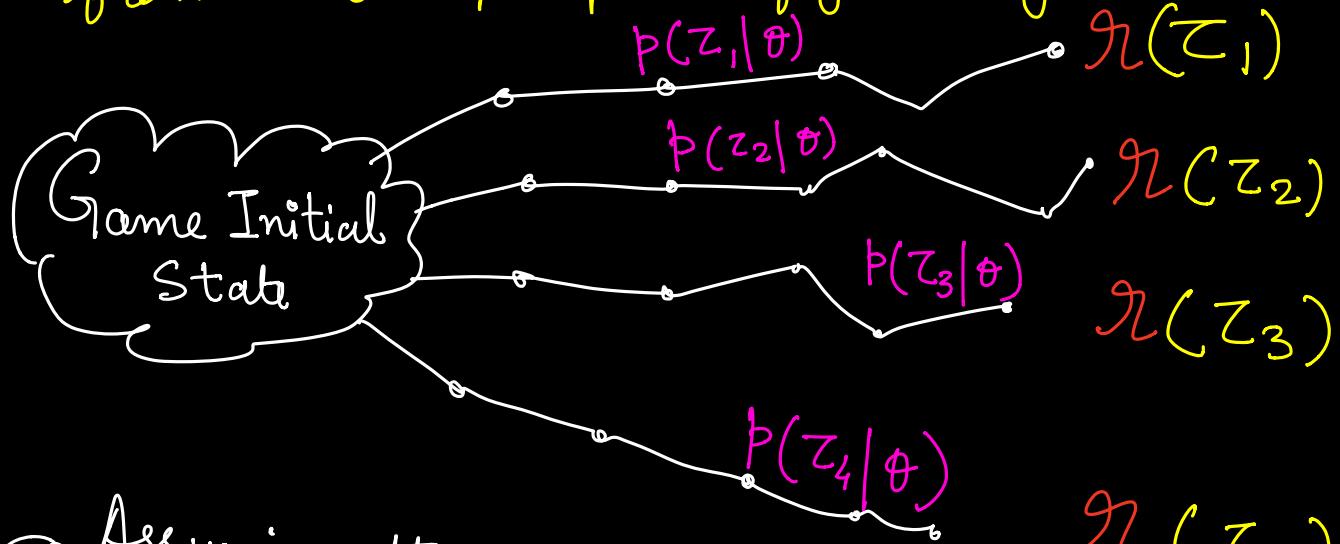
How to proof that

$\nabla_\theta J(\theta) =$ as stated in
step ④,

G.2

REINFORCE algorithm Loss function

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



- * Assuming the agent to be starting from some initial game state and following (π_θ) policy parameter

→ a) Assuming all trajectories to be equally probable : $J(\theta) = \frac{1}{4} \sum_{i=1}^4 \mathcal{R}(\tau_i)$

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

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

Value associated to a policy (π_θ) ,

$$\text{Parameterized over } [\theta] = \mathbb{E}_{\tau_i} [\mathcal{R}(\tau_i)]$$

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

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

$$\therefore \theta^* = \arg \max_{\theta} [J(\theta)]$$

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

G.3

REINFORCE algo Backward Pass

* Forward Pass : $J(\theta) = \int_{\mathcal{Z}} \mathcal{R}(z) * p(z|\theta) dz$ (1)

* Backward pass : $\nabla_{\theta} J(\theta) = \int_{\mathcal{Z}} \mathcal{R}(z) * \nabla_{\theta} p(z|\theta) dz$ (2)

*** This is intractable, Gradient computation of an expectation is problematic when p depends upon θ . [Gradient of expectation]

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

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

$$\therefore \nabla_{\theta} J(\theta) = \int_{\mathcal{Z}} \mathcal{R}(z) * \nabla_{\theta} \log[p(z|\theta)] * p(z|\theta) dz$$

$\therefore \nabla_{\theta} J(\theta) = \mathbb{E}_{\mathcal{Z} \sim p(z|\theta)} [\mathcal{R}(z) * \nabla_{\theta} \log[p(z|\theta)]]$

Expectation of Gradient

G.4 Issue with Gradient & its Computation

- * There can be infinitely many trajectories possible for any parameterized policy (π_θ) 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 (π_θ) are not equiprobable, instead following a distribution. Say $p(z|\theta)$
- * Hence, for a some trajectory say $\tau_1 = \langle (s_0, a_0, r_0), (s_1, a_1, r_1), \dots, (s_n, a_n, r_n) \rangle$
let us compute the probability of (τ_1) given $p(z|\theta)$

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

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

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

Given the policy (π_θ) parameterized by (θ) & learned by some Deep policy N/o maximizing the reward fn $[J(\theta)]$.

Transition
Probabilities

As per Parameterized π_θ Policies.

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

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

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

$$\pi_\theta(a_t | S_t)$$

Probability of taking an action (a_t) at state (S_t) under the policy (π_θ)

Eq-A

(*) 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 (π_θ) 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

$$\therefore \nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim P(z|\theta)} \left[R(z) * \nabla_{\theta} \log [P(z|\theta)] \right] \quad (\text{Eq-B})$$

Expectation of Gradient

The log likelihood of such a sampled trajectory (τ) is using (Eq-A)

$$\log [P(z|\theta)] = \sum_{t \geq 0} \log T(s_{t+1}|s_t, a_t) + \text{Independent of } [\theta]$$

Just differentiating it wrt $[\theta]$.

$$\log \pi_{\theta}(a_t|s_t) \quad (\text{Eq-C})$$

Dependent over $[\theta]$

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

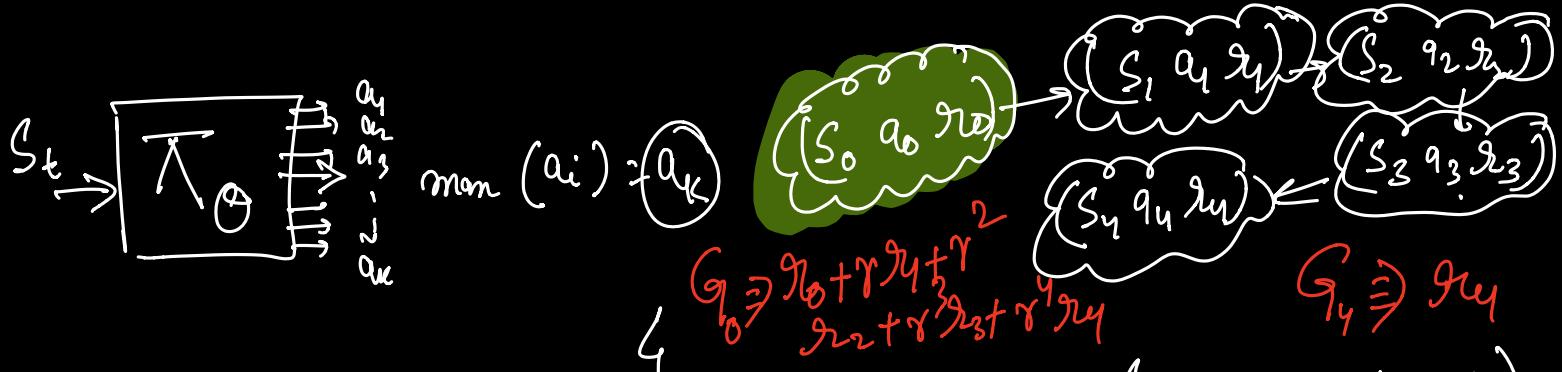
Independent of transition probability
(Eq-D)

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

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

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

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{z \sim p(z|\theta)} \left[\mathcal{I}_2(z) * \sum_{t \geq 0} \nabla_{\theta} \log [\pi_{\theta}(a_t | s_t)] \right]$$



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

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

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

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

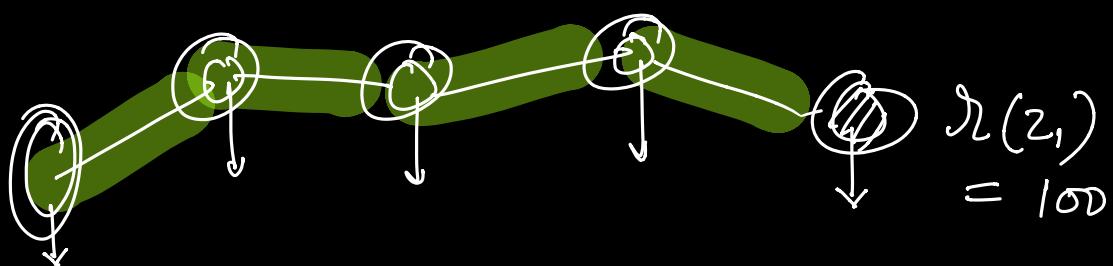
$$G_3 * \nabla_{\theta} \log (\pi_{\theta}(a_3 | s_3)) +$$

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

$\partial L(z)$ Good

Credit Assignment

$\partial L(z)$ (Gradient Estimation)
Unbiased -



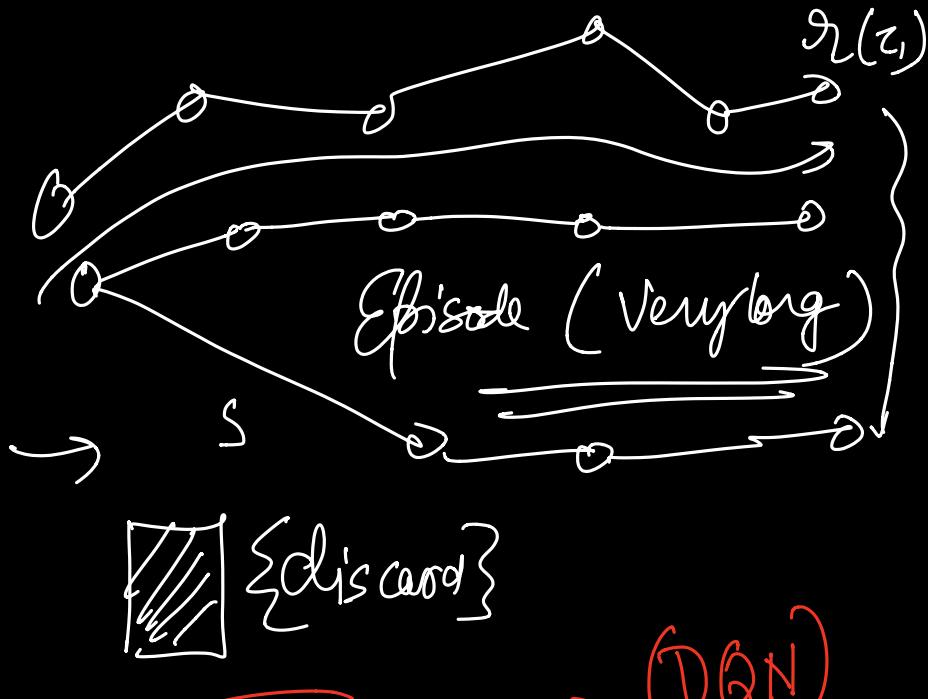
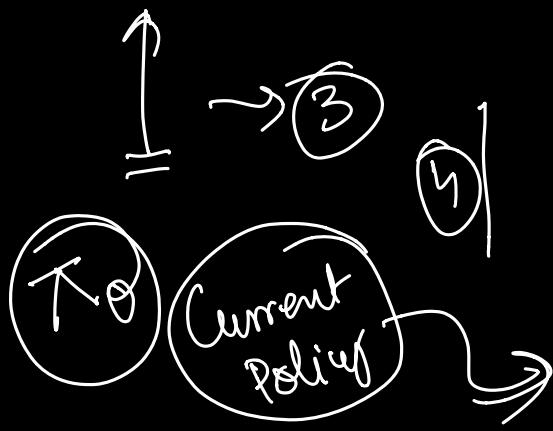
$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \gamma(z_t) \nabla_{\theta} \log(\pi_{\theta}(a|s_t))$$

$$\sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t-t'} \right) \nabla_{\theta} \log(\pi_{\theta}(a|s_t))$$

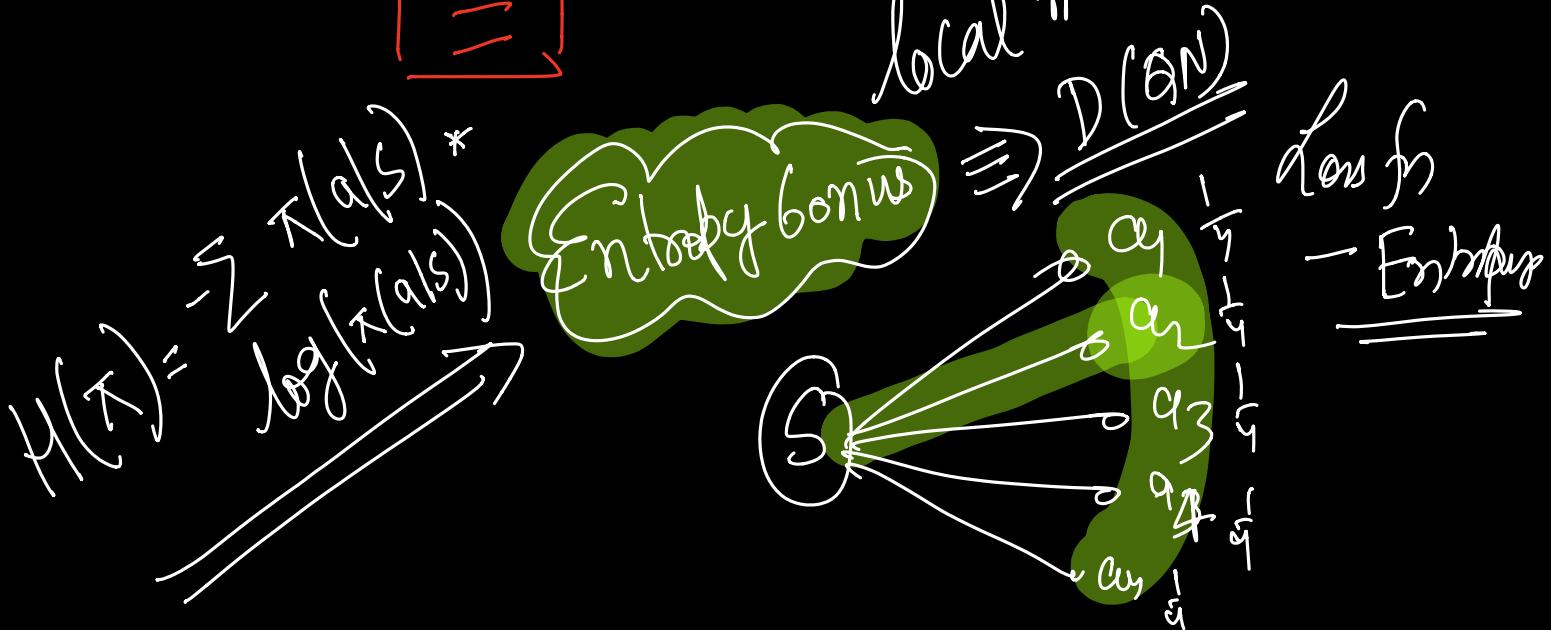
$$\sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t-t'} \gamma_{t'} \right) * (\text{Same})$$

$$\begin{matrix} \gamma(z_1) \\ \gamma(z_2) \\ \gamma(z_3) \end{matrix}$$

Playing game



*Exploration is limited and it may stuck in local minima.



- ① Full episode
 - ② High gradient variance
 - ③ Exploration
 - ④ Confounded samples
- Average
Actor/Critic N/W

Finish properly all the
limitations & issues of
REINFORCE algorithm.

i) Algorithm for Policy Gradient (REINFORCE)

① Initialize the network parameters θ .

② Generate (N) trajectories $\{z^i\}_{i=1}^N \left(z_1^i, z_2^i, \dots, z_N^i \right)$ following the policy π_θ .

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.

(θ) 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~~

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

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.)

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

Baseline (b) : It is the value that can give us the expected return from the state the agent is in. Simplest baseline: $b = \frac{1}{N} \sum_{i=1}^N R(z)$

Can be $\{ \text{OR even moving average} \}$

baseline can be any function but should not dependent upon (θ) . (Does not affect N/W func.)

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 (π) .

↳ Other options are Q-fn & Advantage fn.

How can we learn the baseline functions.

↳ Just like approximating policy using (θ) network

↳ One can use value network (ϕ)

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

Now since the value of a state is floating point number, ϕ 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 θ and value network parameter ϕ
2. Generate N number of trajectories $\{\tau^i\}_{i=1}^N$ following the policy π_θ
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 θ 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 ϕ 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.)

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

$$\nabla_{\theta} J(\theta) = \sum_{t \geq 0} [Q_{\theta}^{\pi_{\theta}}(s_t, a_t) - V_{\theta}^{\pi_{\theta}}(s_t)]$$

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

✗ 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)

R_t

(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.

(II)

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

This we can get from ϕ

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

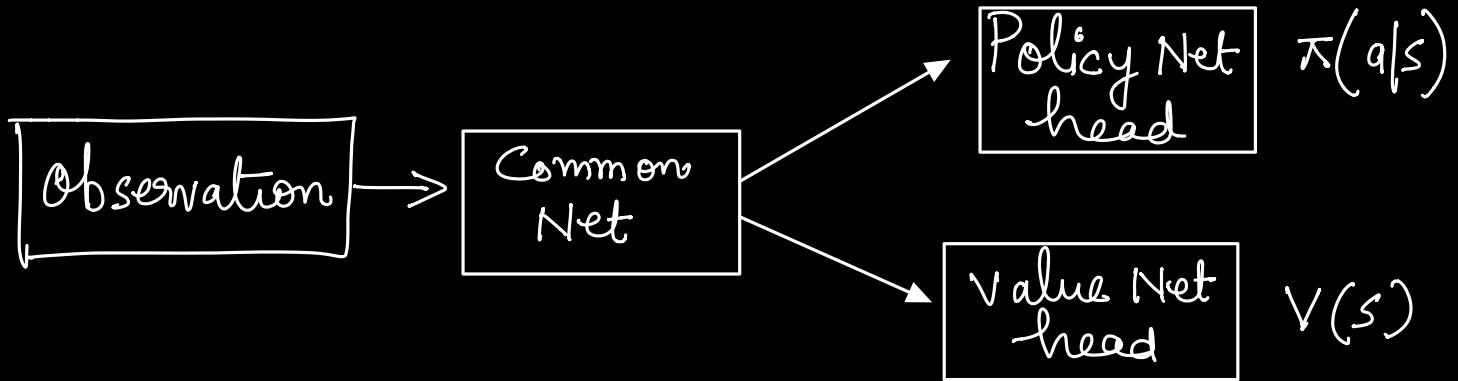
Now Q fn can be Computed for any step, instantly using, immediate reward

ϕ 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 (π).

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

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

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

⋮

⋮

⋮

⋮

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

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

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

	m_1	m_2	m_3	\vdots	m_m

	m_1	m_2	m_3	\vdots	m_m

	m_1	m_2	m_3	\vdots	m_m

	m_1	m_2	m_3	\vdots	m_m

Action

a

Reward

r_t

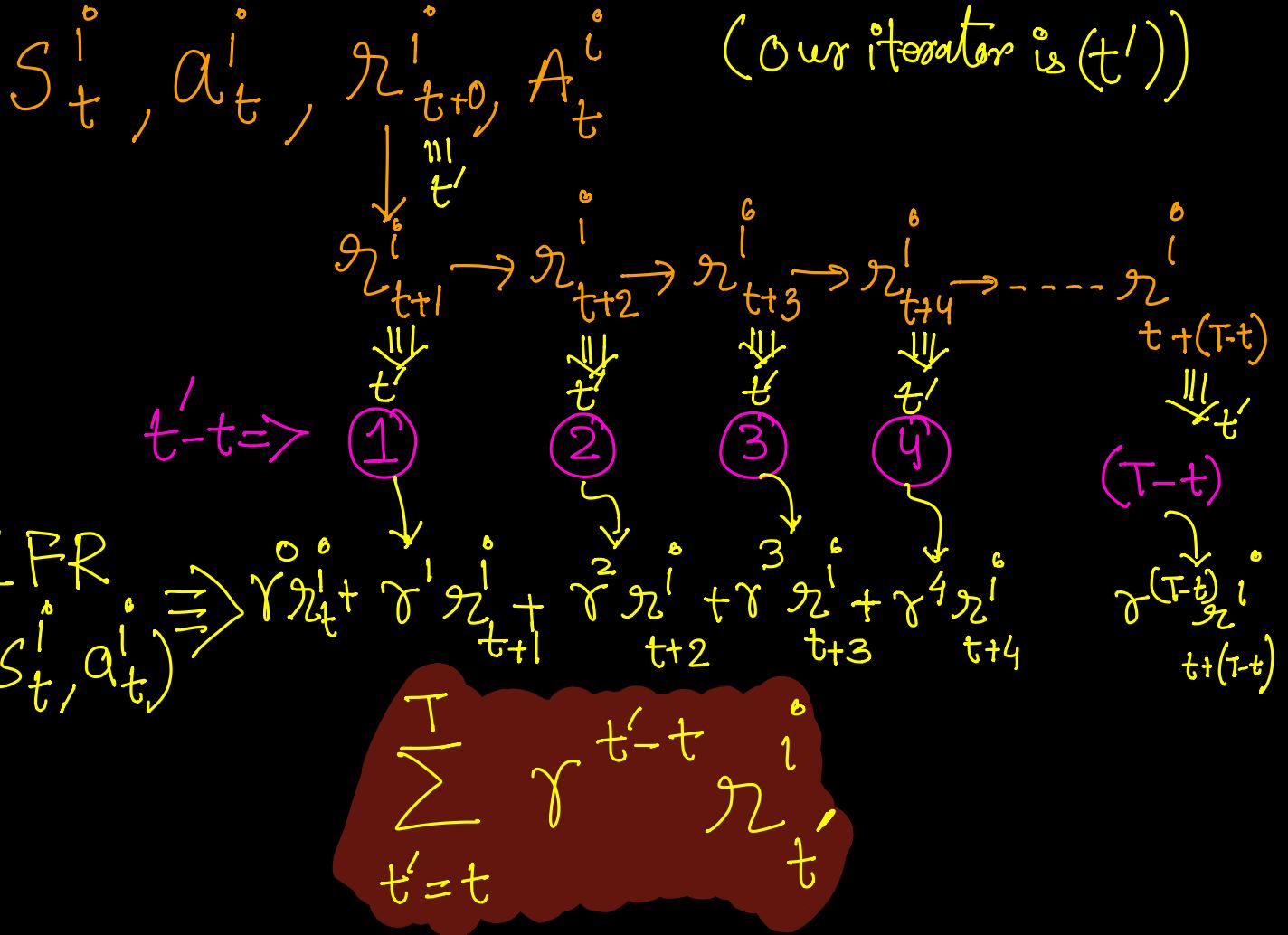
Advantage

A

s_t

i

\Rightarrow State of the agent at t^{th} step in the i^{th} episode.



ALGORITHM (Actor/Critic)

- ① Initialize the
 - Policy N/w $[\theta]$ (Actor)
 - Q learning N/c $[\phi]$ (Critic)
- ② For each training iteration $= (1, 2, \dots)$ do

After every iteration policy will be updated. New data need to be generated.

 - ③ Sample (m) trajectories under the current policy π . (S, a, r, t)

Q.2

$$\Delta \theta = 0$$

Gradient accumulator for Actor N/W.
for a given policy (π), 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 $A_t^{(i)}$

It doesn't
depends upon
the Actor N/W (θ)

$$A_t^{(i)} = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}^{(i)} - V(S_{t'}^{(i)})$$

all future rewards

(R_t)

Critic
network

Minimization of A_t , 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

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

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

Gradient
ascent

Gradient estimator
step.

Q.4

$$\Delta \phi =$$

$$\sum_i \sum_t$$

$$\nabla_{\phi} \| A_t^{(i)} \|_2^2$$

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
Updation

A_t minimization
enforces Bellman
Constraint.

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

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

③ End for

Also incorporates
standard
experience replay
trick too.

old problem of
deep Q-learning:
(*) Learn Q-values
for all $\langle \text{State}, \text{action} \rangle$
pairs. \Downarrow
Critic only learns the
Q-values for generated/
observed $\langle S, a \rangle$ pairs by
the policy under consideration.

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.

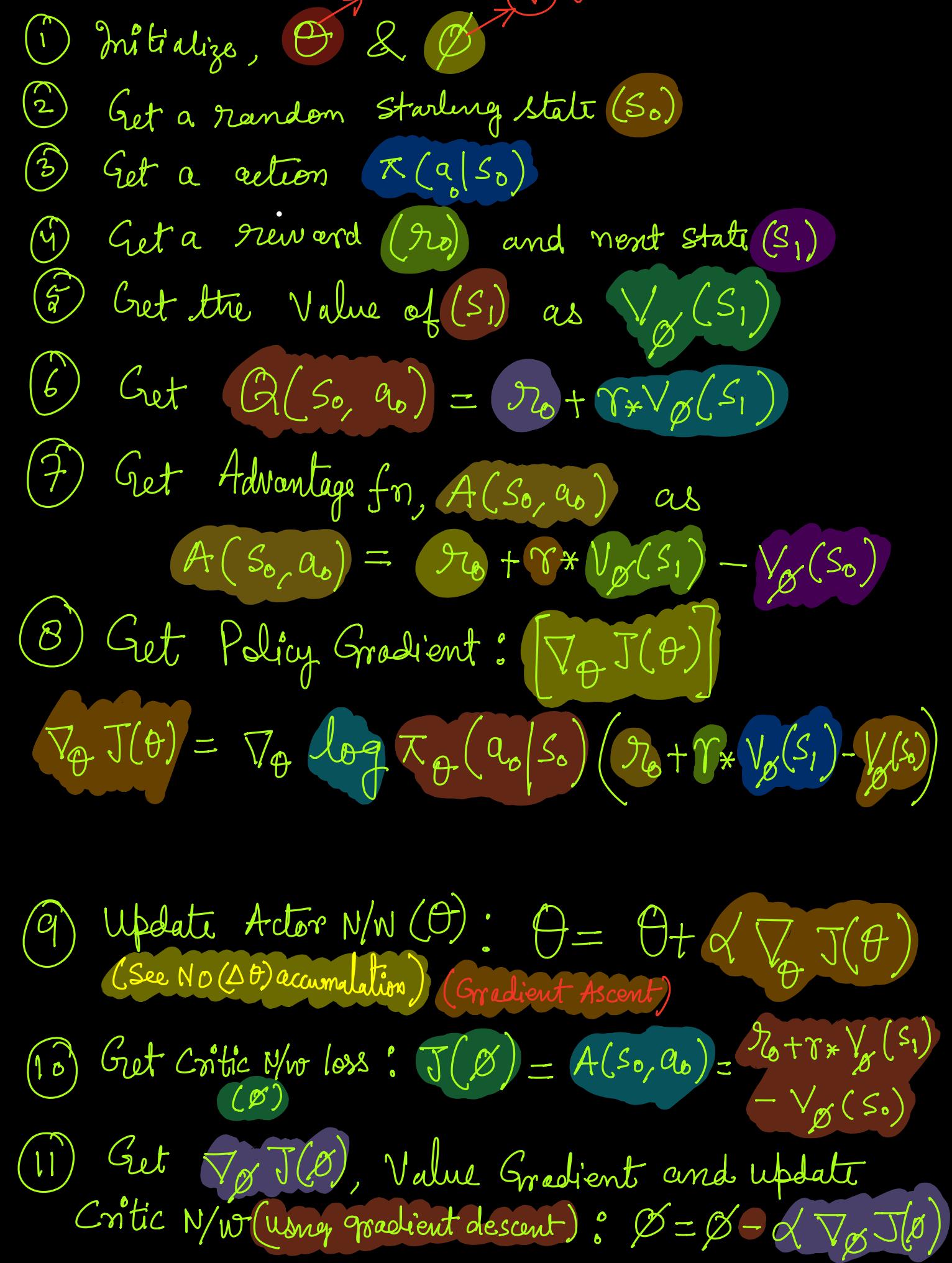
*) We can approximate (R_t) or $Q(S, a)$ as:

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

*) None, 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, θ & ϕ
 - ② Get a random starting state (s_0)
 - ③ Get an 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 $Q(s_0, a_0) = r_0 + \gamma * V_\phi(s_1)$
 - 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 (θ) : $\theta = \theta + \alpha \nabla_\theta J(\theta)$
(See NO($\Delta\theta$) 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 θ and the critic network parameter ϕ
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 θ using gradient ascent:

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

5. Compute the loss of the critic network:

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

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

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

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

↓
Sample
efficient.