

At Station: TH:

## Reinforcement Learning -

### Problem Formulation and Terminology:-

Locations 1-9

Agent:  (Robot)

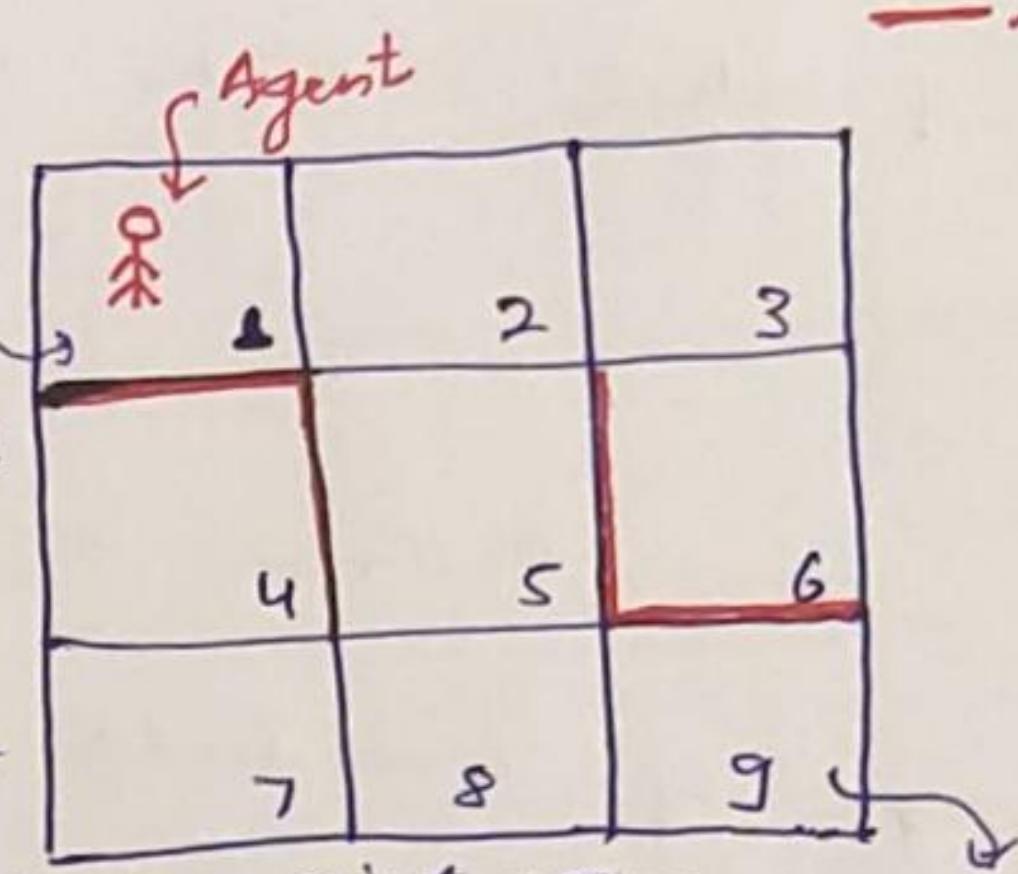
—: walls

Environment: Setup of a problem.

States:  $L_1 - L_9$  (Nine locations) state  
Tells you what is the current state of an agent in the environment.

Actions:

$\rightarrow \leftarrow \downarrow \uparrow$  all actions that an agent can



Take at any point of time.

Reward: Reward is a function of the action taken at any point of time. destination  
on which you perform our action  $g \times 3$  (Grid)  
on state 5 and gets the instantaneous reward.  
eg-  $R(L, \rightarrow) \rightarrow R(2, )$ .

### Bellman Equations:

Value-function:  $v(s)$  :- value function is a function on a given state, you return a numerical value called as value in the

state  $s$ .  $s \xrightarrow{\text{current state}} s' \text{ (next state)}$

eg-  $1 \xrightarrow{\rightarrow} 2$

$2 \xrightarrow{\downarrow} 5$

value of next state of future states

immediate rewards

$v(s) = \max_a (R(s, a) + \gamma \cdot v(s'))$

↓

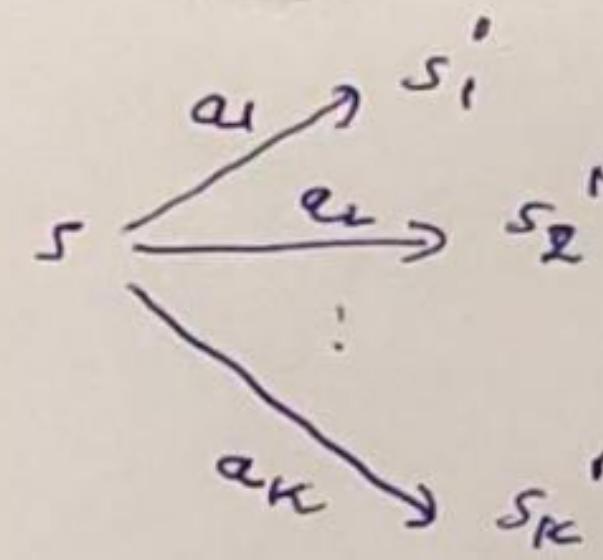
Bellman equations

$a \in a_1, a_2, a_3$  are all actions on state  $s$ .

$0 \leq \text{discounting Factor} \leq 1$

Bellman equations says that the value in any state  $s$  is the maximum of the immediate rewards that you will get by performing all actions on state  $s$ .

↳ Bellman equations can be represented as an recursive equation.



↳ That problem, can be represented as an optimization problem.

↳ So a Bellman equation represent the weighted (sum) average of immediate returns that we are getting and not giving much importance to future rewards by using a discount factor. ( $0 \leq r \leq 1$ )

e.g -

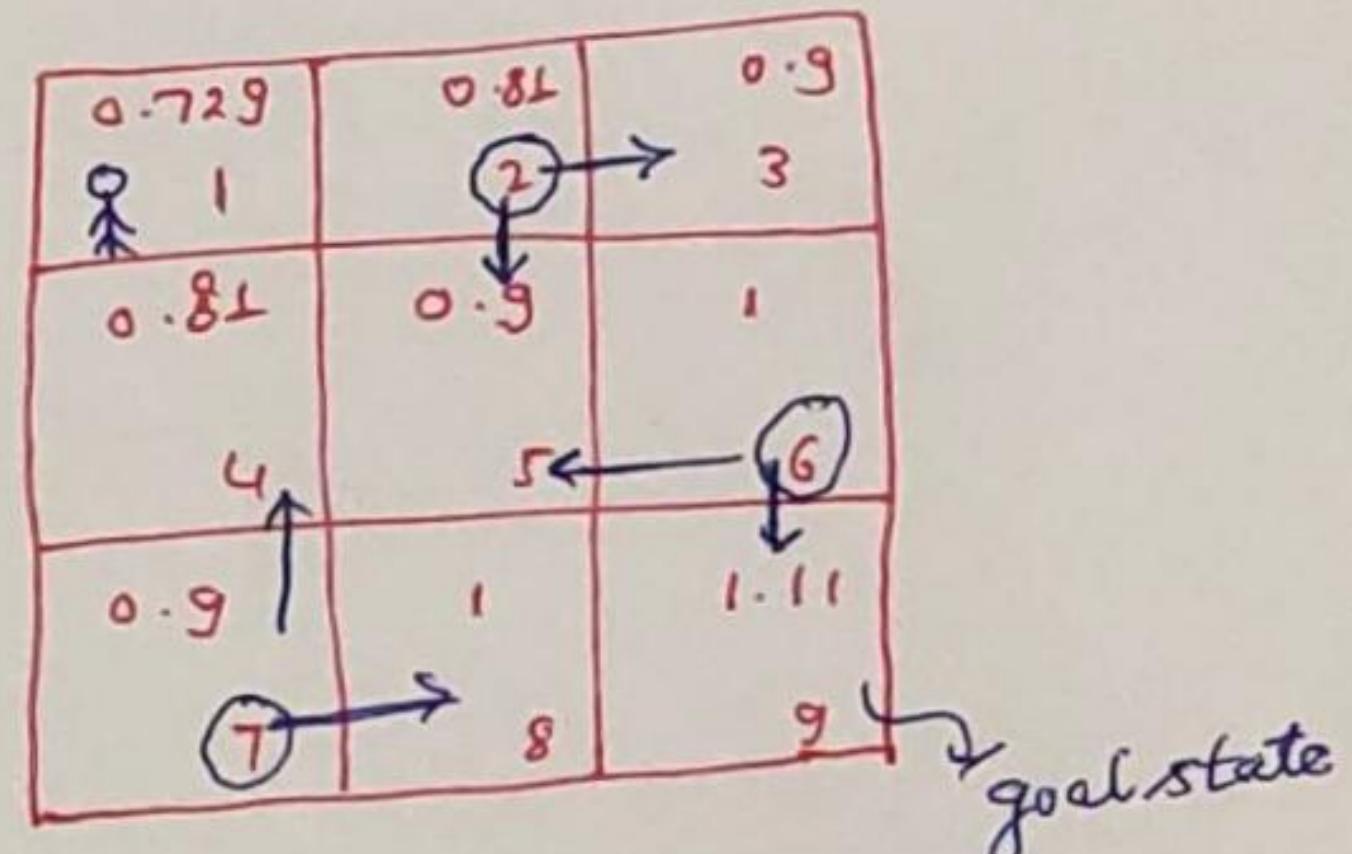
$$v(s) = \max_a (R(s, a) + r v(s'))$$

(constant rewards) future value  
discount factor

Let  $r = 0.9$

Assume  $R(s, a) = 0 \forall s \notin LG$

$$v(LG) = 1.11$$



$$v(LG) = ?$$

$$v(s) = \max_a (\gamma v(s'))$$

$$v(LG) = 1.11$$

$$v(LG) = 1 \quad (0.9 \times 1.11)$$

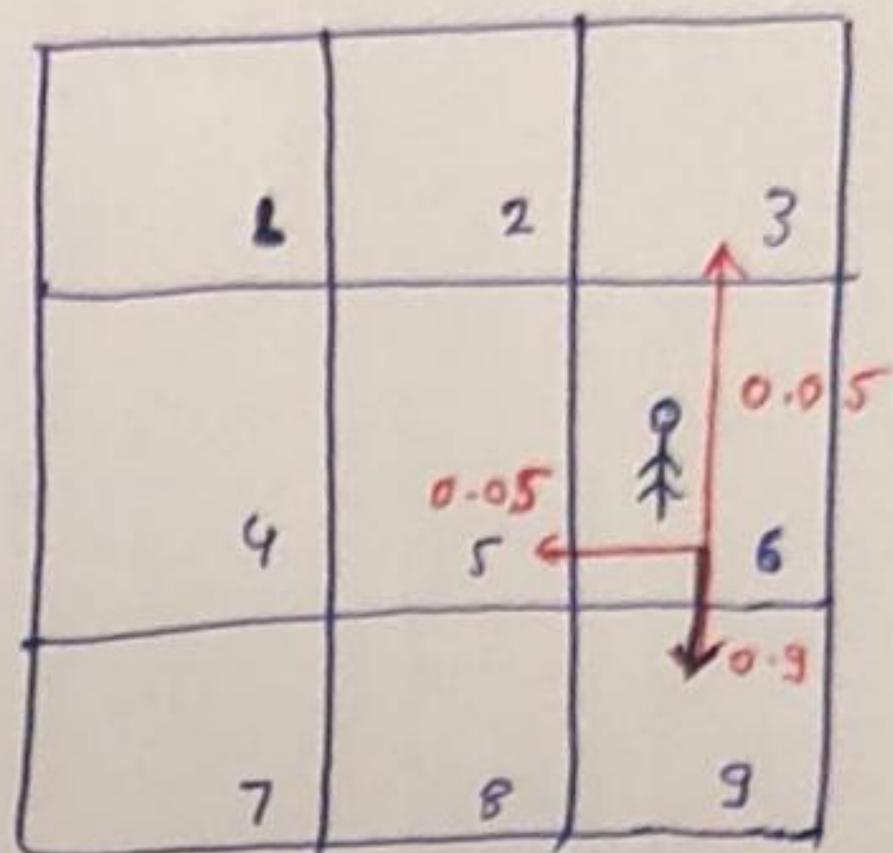
$$v(LB) = 1 \quad (0.9 \times 1.11)$$

$$v(L7) = 0.9$$

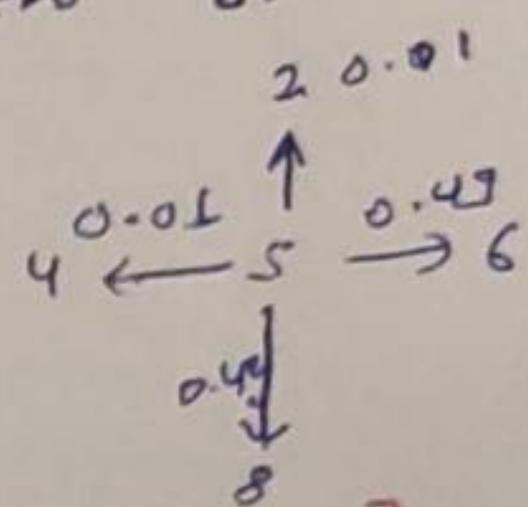
$$v(L2) = 0.81$$

↳ Always go to higher state than the current state if we are able to obtain the value state matrix.

Markov Decision Process:- (stochastic Process):-



means at each state instead of a one deterministic action, in MDP (stochastic Process), probabilities are associated with each action then such type of a system is called as MDP.



$$v(s) = \max_a \left[ r(s, a) + \gamma \sum_{s'} p(s, a, s') v(s') \right]$$

mathematical representation of this eqn.

Expected value =  $0.9 * v(L) + E[v(s')] + 0.05 * v(LB) + 0.05 * v(LS)$

$$v(s) = \max_a \left[ r(s, a) + \gamma \sum_{s'} p(s, a, s') v(s') \right]$$

Expected value of  $v(s')$

eg.  $E(P) = 0.90 * 1 + (0.10 * -1) = 0.8$



# POORNIMA

## FOUNDATION

### LECTURE NOTES

Campus: PCE Course: BTECH in CSE Class/Section: III Yr. Section- A Date: .....

Name of Faculty: Praveen Kumar Yadav Name of Subject: Machine Learning Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*Introduction to Reinforcement Learning*

**IMPORTANT & RELEVANT QUESTIONS:**

**FEED BACK QUESTIONS (AFTER 20 MINUTES):**

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:

SFT 2019/2020: - IV:



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

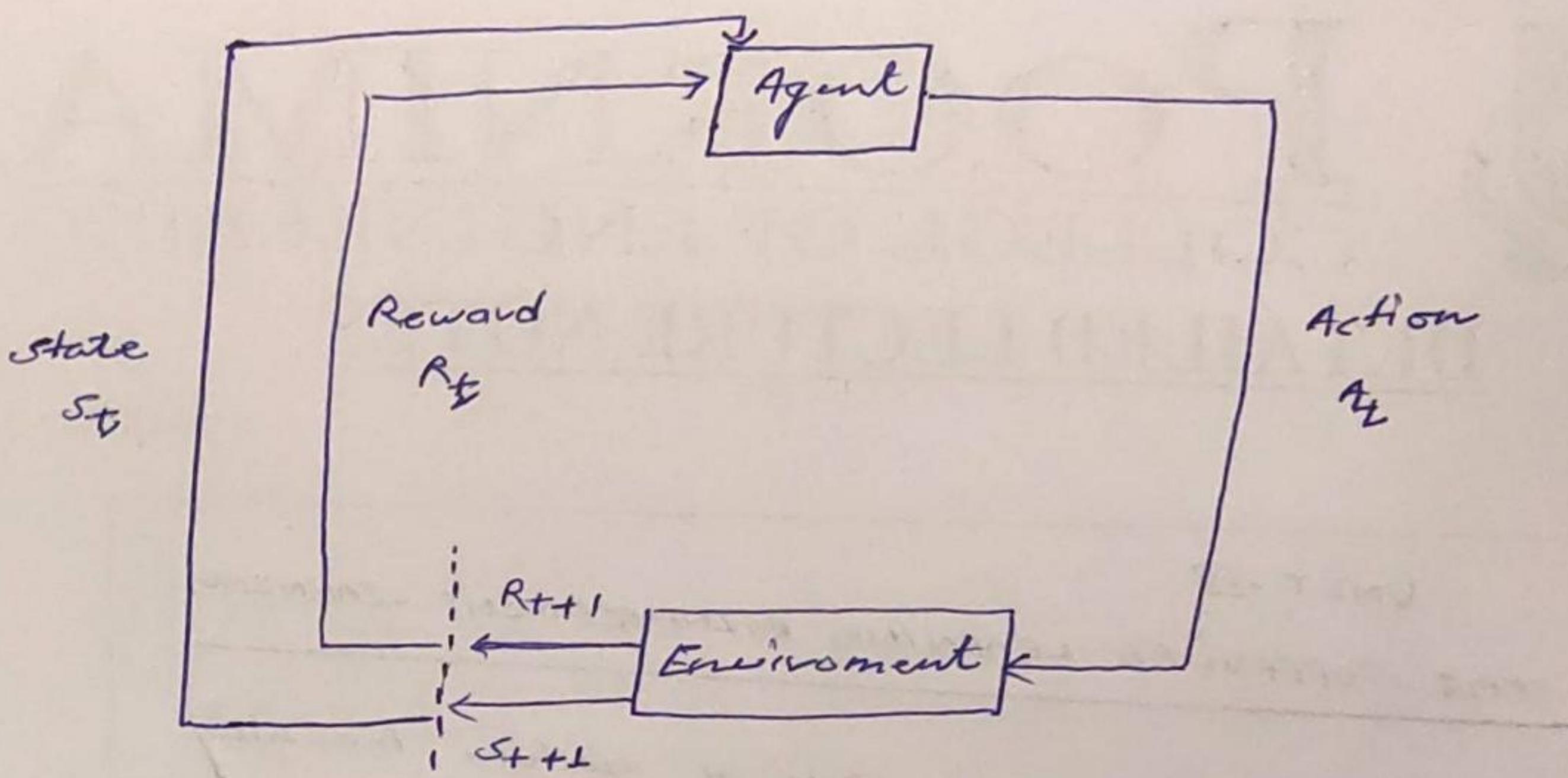
UNIT - IV

SEMI SUPERVISED LEARNING, REINFORCEMENT LEARNING

Reinforcement Learning:- is a type of machine learning technique that enables an agent to learn in an interactive environment by trial and error by using feedback from its own actions and experience.

→ Reinforcement learning uses rewards and punishments as signals for positive and negative behaviour.

→ Objective of reinforcement learning is to find a suitable action model that would maximize the total cumulative reward of the agent



Key Elements of RL:-

Agent :- An agent is a SW program that learn to make intelligent decisions.

eg- chess player  
Mario

Environment :- The environment is the world of the agent. The agent stays within environment.

eg- chessboard. of Environment  
Super Mario Bros

display of agent  
Mario

State and Action :- A state is a position or a moment in the environment that the agent can be in.

→ state is usually denoted by  $s$ .



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

The Agent interacts with the environment and moves from one state to another by performing an action.

→ The action is usually denoted by  $a$ .

Rewards:- The agent interacts with the environment by performing an action and moves from one state to another.

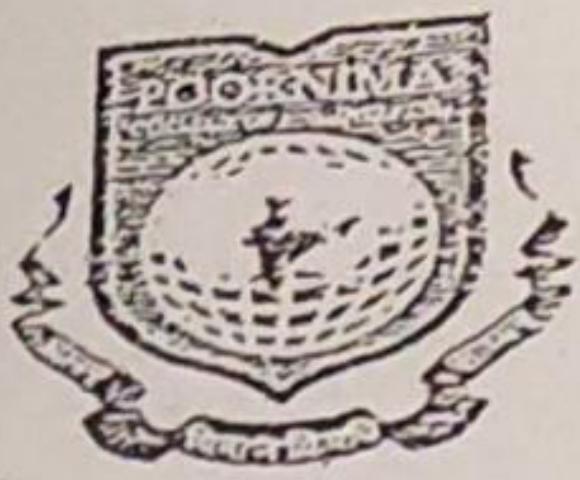
A reward is nothing but a numerical value, say +1, for a good action and -1 for a bad action.

→ Reward is denoted by  $v$ .

In RL setting, we will not teach the agent what to do or how to do it. instead we will give a reward to the agent for every action it

it does. we will give the reward to the agent when it performs a good action and negative award to the agent when it performs a bad action.

→ so RL can be viewed as a trial and error learning process where the agent tries out different actions and learns the good actions, which gives a positive reward.



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

#### The RL ALGORITHM:-

step 1: Firstly, the agent interacts with the environment by performing an action.

step 2: By performing an action, the agent moves from one state to another.

step 3: The agent will receive an reward based on action it performed.

step 4: Based on reward, the agent will understand whether the action is good or bad.

step 5: if the action was good, i.e. if the agent received a +ve reward, then the agent will prefer performing that action, else the

the agent will try performing other actions in  
search of a +ve reward.



# POORNIMA FOUNDATION

## LECTURE NOTES

Campus: PCE

Course: BTECH in CSE

Name of Faculty: Praveen Kumar Yadav

Date (Prep.): .....

Class/Section: III Yr. Section- A

Date: .....

Name of Subject: Machine Learning

Code: 6CS4-02

Date (Del.): ..... Unit No.: ..... Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*Terms Related to Reinforcement Learning*

### IMPORTANT & RELEVANT QUESTIONS:

### FEED BACK QUESTIONS (AFTER 20 MINUTES):

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:



# POORNIMA

---

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

policy :

→ deterministic policy - tells the agent to perform one particular action in a state.

$$a_t = \mu(s_t)$$

Action  $a_t$        $\leftarrow$        $\downarrow$        $\rightarrow$  state  $s$  at Time  $T$   
Time  $T$ .  
(one action).

policy

for  $\forall$   $\mu(a) = \text{Down}$

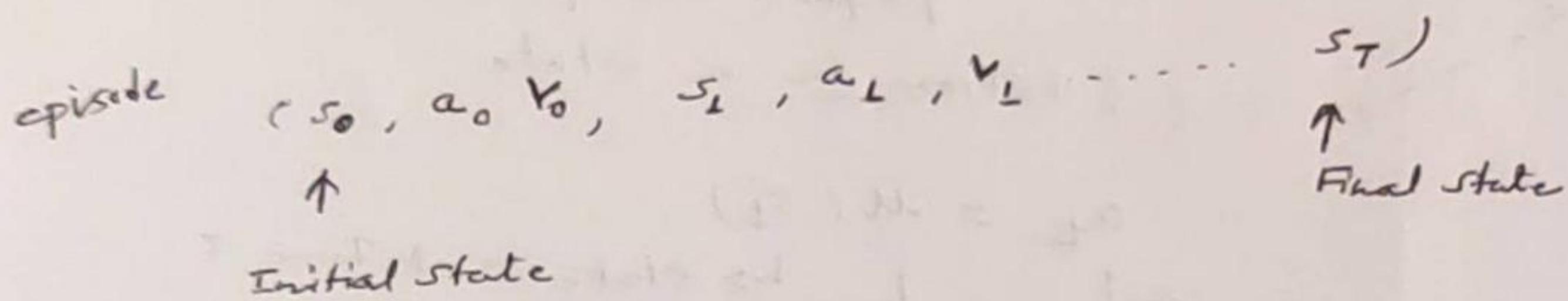
ii) stochastic policy :- Maps the state to a probability distribution over an action space.

Means for a given state  $s$ , the stochastic policy will return a prob. dist over an action space, Here the agent will perform different actions each time based on prob. distribution returned by stochastic policy.

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

Episode:- Agent-environment interaction starting from the initial state until final state is reached is called an episode.

- ↳ Also known as Trajectory. denoted by  $\pi$ .
- ↳ objective.
- ↳ To learn optimal policy.



Finite Horizon- If the agent-environment interaction stops at a particular time step, then horizon is called a finite horizon.

Return and discount Factor:-

return and discount Factor:-  
↳ A return can be defined as the sum of rewards obtained by agent in an episode.

$$R(\tau) = v_0 + v_1 + v_2 + \dots + v_T$$

$$R(z) = \sum_{t=0}^{\tau} v_t$$

The goal of our agent is to maximize the return.



# POORNIMA

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

Discrete Tasks -

continuous Tasks

$$\hookrightarrow R(\tau) = v_0 + v_1 + v_2 + \dots + v_\infty$$

$$R(\tau) = \gamma^0 v_0 + \gamma^1 v_1 + \gamma^2 v_2 + \dots + \gamma^n v_\infty$$

$$R(\tau) = \sum_{t=0}^{\infty} \gamma^t v_t$$

discount factor value < small → more importance to immediate rewards than future rewards.

large - more importance to future rewards compared to immediate rewards.

$\gamma = 0$  - immediate rewards  
→ nothing learns since it concentrates on immediate reward.

$\gamma = 1$  → learning continuous every time

The value Function -

function, also known as state-value function, denotes the value of a state

$$v^\pi(s) = [R(\tau) \mid s_0 = s]$$

↳ value function can be defined as expected return that the agent would obtain starting from state  $s$  by following policy  $\pi$ .

$$v^\pi(s) = E[R(z) \mid s_0 = s]$$

$$z \sim \pi$$



# POORNIMA

## FOUNDATION

### LECTURE NOTES

Campus: PCE Course: BTECH in CSE Class/Section: III Yr. Section- A Date: .....

Name of Faculty: Praveen Kumar Yadav Name of Subject: Machine Learning Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: ..... IV

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*Markov Decision Process*

**IMPORTANT & RELEVANT QUESTIONS:**

**FEED BACK QUESTIONS (AFTER 20 MINUTES):**

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:



# POORNIMA

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

Markov Decision Process:- provides a mathematical framework for solving the RL problem.

The Markov property and Markov chain:-

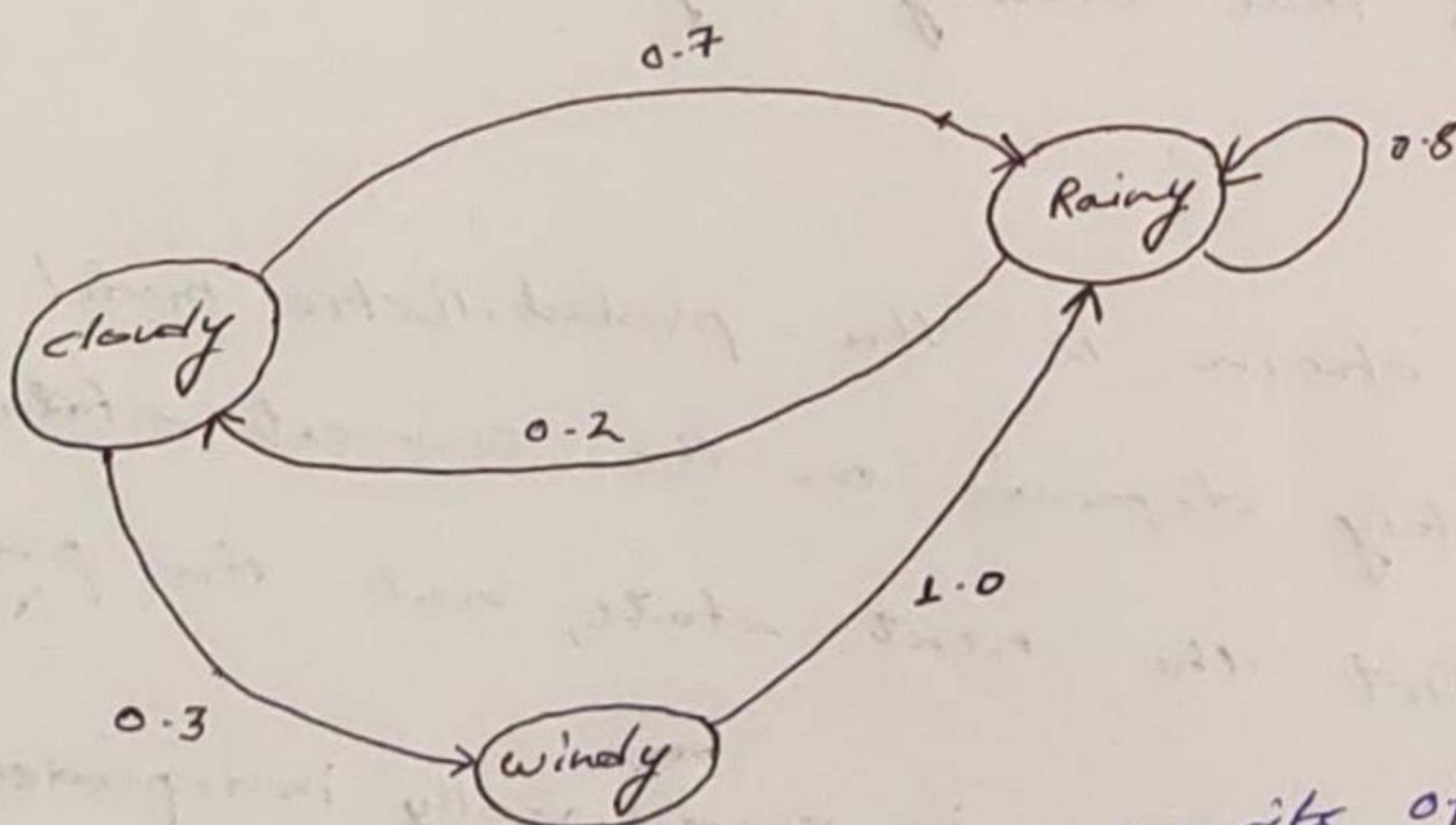
The Markov property states that the future depends only on present and not on the past. The Markov chain, also known as Markov process, consisting a sequence of states that strictly obey the Markov property.

→ Markov chain is the probabilistic model that solely depends on the current state to predict the next state, not the previous states, i.e. future is conditionally independent of the past.

Transition Probability:- moving from one state to another is called a transition and its probability is called as transition probability.

Is denoted by  $P(S'|S)$ .

Current State	Next State	Transition Probability
cloudy	Rainy	0.7
cloudy	windy	0.3
Rainy	Rainy	0.8
Rainy	cloudy	0.2
windy	Rainy	1.0



Markov chain or Markov process consists of a set of states along with their transition probabilities.



# POORNIMA

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

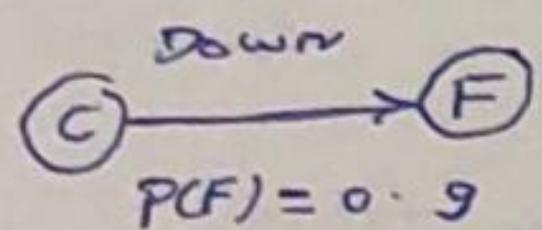
Markov Reward Process:- (MRP):

- ↳ MRP is an extension of Markov chain along with one additional function i.e reward function.
- ⇒ so MRP consists of states, a transition probability and also a reward function.
- ⇒ Reward function denotes rewards that we receive in each state. Reward function is denoted by  $R(s)$ .
- ⇒ MRP consist of states  $s$ , a transition prob.  $p(s'|s)$ , and a reward function  $R(s)$ .

The Markov Decision Process:- MDP is an extension  
of MRP with actions.

so MDP consists of set of states  $s$ , a transition  
probability  $P(s'|s, a)$ , a reward function  $R(s)$   
 $R(s, a, s')$

and also action.



$$R(A, \text{right}, B) = -L$$

$$\hookrightarrow E(x) = \sum_{i=1}^n x_i p(x_i)$$

Action space:- The set of all possible actions  
in the environment is called as

Action space.

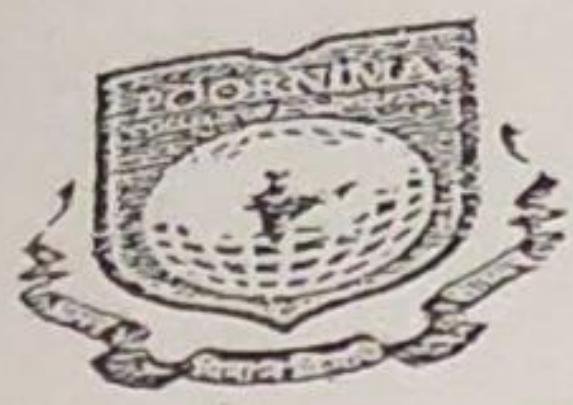
e.g. Grid world - Action - {left, right, up, down}

discrete Action space

continuous Action space.

Policy:- A policy defines the agent behaviour in  
an environment.  
policy tells the agent what action to perform  
in each state.

↑ Firstly we choose random policy.



# POORNIMA

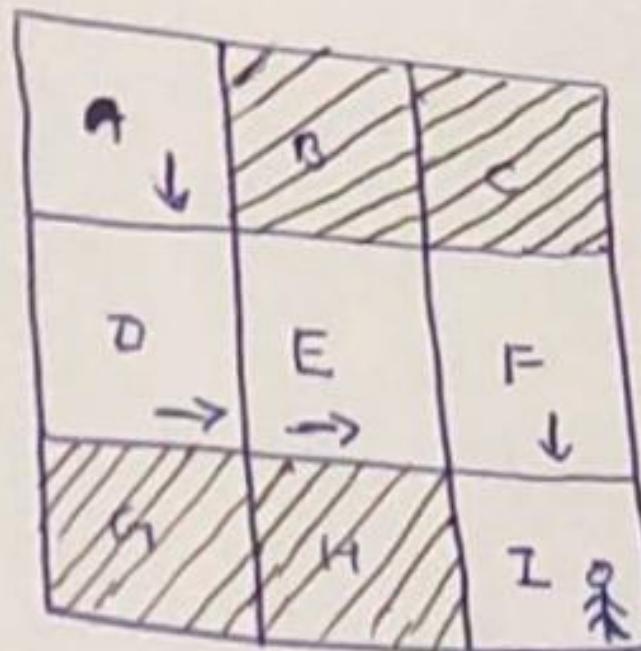
COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

optimum policy

State	Action
A	Down
D	Right
E	Right
F	Down



↳ optimum policy tells the agent to perform the correct action in each state so that the agent can receive good reward.



# POORNIMA

## FOUNDATION

### LECTURE NOTES

Campus: PCE Course: BTECH in CSE Class/Section: III Yr. Section- A Date: .....

Name of Faculty: Praveen Kumar Yadav Name of Subject: Machine Learning Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: 10 .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*BELLMAN EQUATION*

**IMPORTANT & RELEVANT QUESTIONS:**

**FEED BACK QUESTIONS (AFTER 20 MINUTES):**

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:



# POORNIMA

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

BELLMAN EQUATION:- The bellman equation, named after Richard Bellman, helps us to solve Markov decision process. (MDP).

→ BELLMAN Equation widely used for finding the optimal value and Q function recursively. that can be used to find the optimal policy.

BELLMAN EQUATION FOR VALUE FUNCTION:- BELLMAN

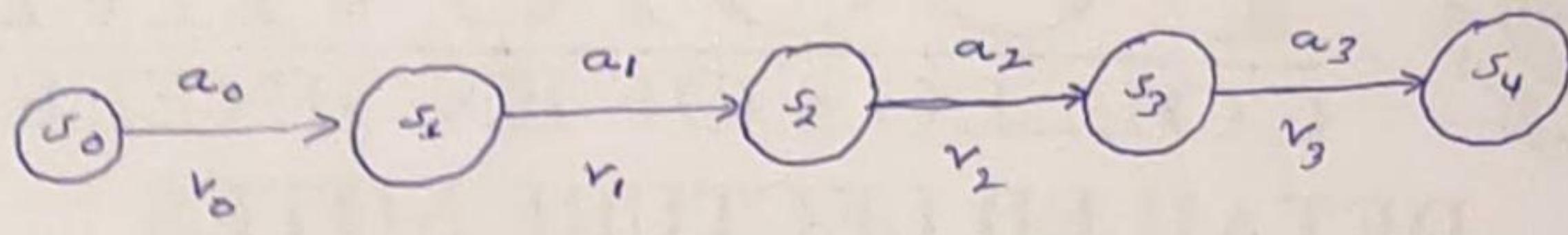
Equation states that the value of a state can be obtained as a sum of immediate reward and discounted value of the next state.

- BELLMAN equation of the value function can be expressed as-

$$V(s) = R(s, a, s') + \gamma V(s')$$

Here  $R(s, a, s')$  - implies the immediate reward obtained while performing an action  $a$  in state  $s$  and moving to new state  $s'$ .

For eg- say we generate a trajectory  $\zeta$  using some policy  $\pi$ :



Let suppose we need to compute the value of state  $s_2$ .  
so according to BELLMAN equation-

$$V(s_2) = \underbrace{R(s_2, a_2, s_3)}_{r_2} + \gamma V(s_3)$$

Discounted value  
of next state.

$$V(s_2) = r_2 + \gamma V(s_3)$$

IN More Generalised Form -

$$V^\pi(s) = R(s, a, s') + \gamma V^\pi(s')$$

where  $\pi$  implies that we are using policy  $\pi$ . RHS term  $R(s, a, s') + \gamma V^\pi(s')$  is often called as - BELLMAN

backups.

Note - Previous equation will work only when we have deterministic environment.

$\varphi$ -Function:-  $\varphi$ -function is a scalar function.



# POORNIMA

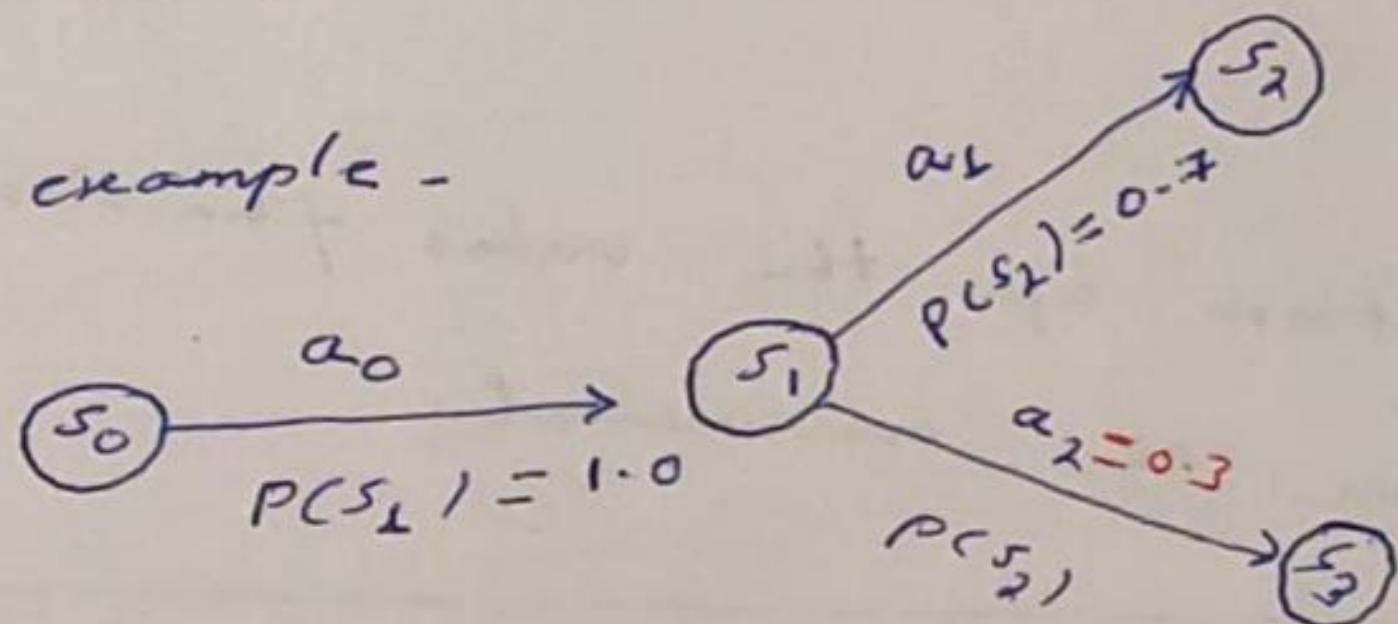
## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

For stochastic Environment-

For example -



Here our Bellman equation with the expectation  
(i.e weighted average), i.e sum of Bellman  
Backup multiplied by the corresponding  
transition probability.

$$v^\pi(s) = \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v^\pi(s')]$$

Here  $p(s'|s, a)$  - denotes the transition probability  
of reaching  $s'$  by performing an action  $a$   
in state  $s$ .

$[r(s, a, s') + \gamma v^\pi(s')]$  - Bellman backup.

These we can write -

$$v(s_1) = P(s_2 | s_1, a_1) [R(s_1, a_1, s_2) + v(s_2)]$$

$$+ P(s_3 | s_1, a_1) [R(s_1, a_1, s_3) + v(s_3)]$$

$$v(s_1) = 0.70 [R(s_1, a_1, s_2) + v(s_2)] + 0.30 [R(s_1, a_1, s_3) + v(s_3)]$$

These the bellman equation of the value function including stochasticity present in the environment.

$$v^\pi(s) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma v^\pi(s')]$$

→ decide the quality of an action.

BELLMAN EQUATION of the Q Function:- states that the Q value of a state-action pair can be obtained as a sum of the immediate reward and the discounted Q value of a next state action pair.

$$q(s, a) = R(s, a, s') + \gamma q(s', a')$$

Here  $R(s, a, s')$  - implies the immediate reward obtained while performing an action  $a$  in state  $s$  and moving to the next state  $s'$ .

$\gamma$  - is the discount factor.

$v(s) = \max_a \sum_{s'} R(s, a, s') + \gamma v(s')$  where state  $s$  is given and we take action  $a$  to get state  $s'$  which is the best reward.

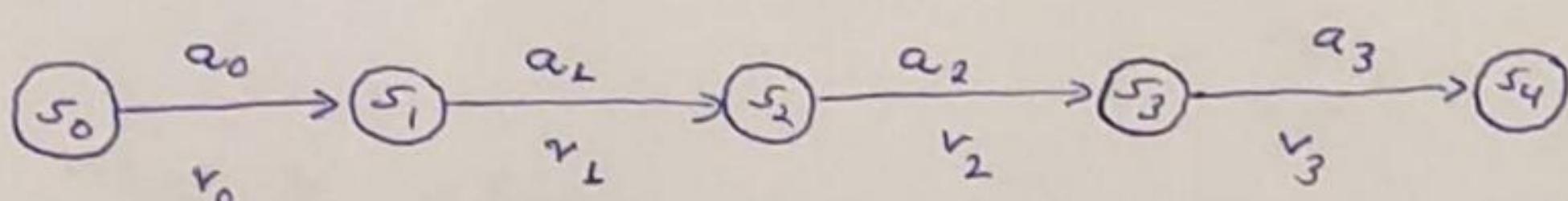


# Poornima COLLEGE OF ENGINEERING DETAILED LECTURE NOTES

PAGE NO. ....

$Q(s', a')$  - is the  $Q$  value of next-state-action pair.

For eg -



Suppose we need to compute the  $Q$ -value of a state-action pair  $(s_2, a_2)$ .

$$Q(s_2, a_2) = \underbrace{R(s_2, a_2, s_3)}_{\text{immediate Reward}} + \underbrace{\gamma Q(s_3, a_3)}_{\text{discounted } Q \text{ value of next-state-action pair.}}$$

$$Q(s_2, a_2) = r_2 + \gamma Q(s_3, a_3)$$

In More Generalized - representation -

$$Q^\pi(s, a) = R(s, a, s') + \gamma Q^\pi(s', a')$$

Here  $\pi$  - implies that we are using the policy  $\pi$  and RHS is Bellman Backup.

Thus -

For stochastic Environment -

$$\phi^{\pi}(s, a) = \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma \pi(s'|a)]$$

**Q-Function:-** Q-function is a simple function of state  $s$  and action  $a$ .  
 Q-function describes at a given state  $s$  which is the best action should I take in order to get the highest reward.

$$Q(s) = \max_a [R(s, a) + \gamma \sum_{s'} P(s, a, s') \cdot V(s')]$$

↓

$Q(s, a)$  :- Q-Function  
 → Goal of Q function is to maximize the total reward  
 so Deep-Q-Learning is all about computing the Q function

$$V(s) = \max_a g(s, a)$$

$$g(s, a) = R(s, a) + \gamma \cdot \sum_{s'} p(s', a, s') \cdot v(s')$$

$$E[u(s')]$$

can be  
expressed as

Max  $\varrho(s'a')$

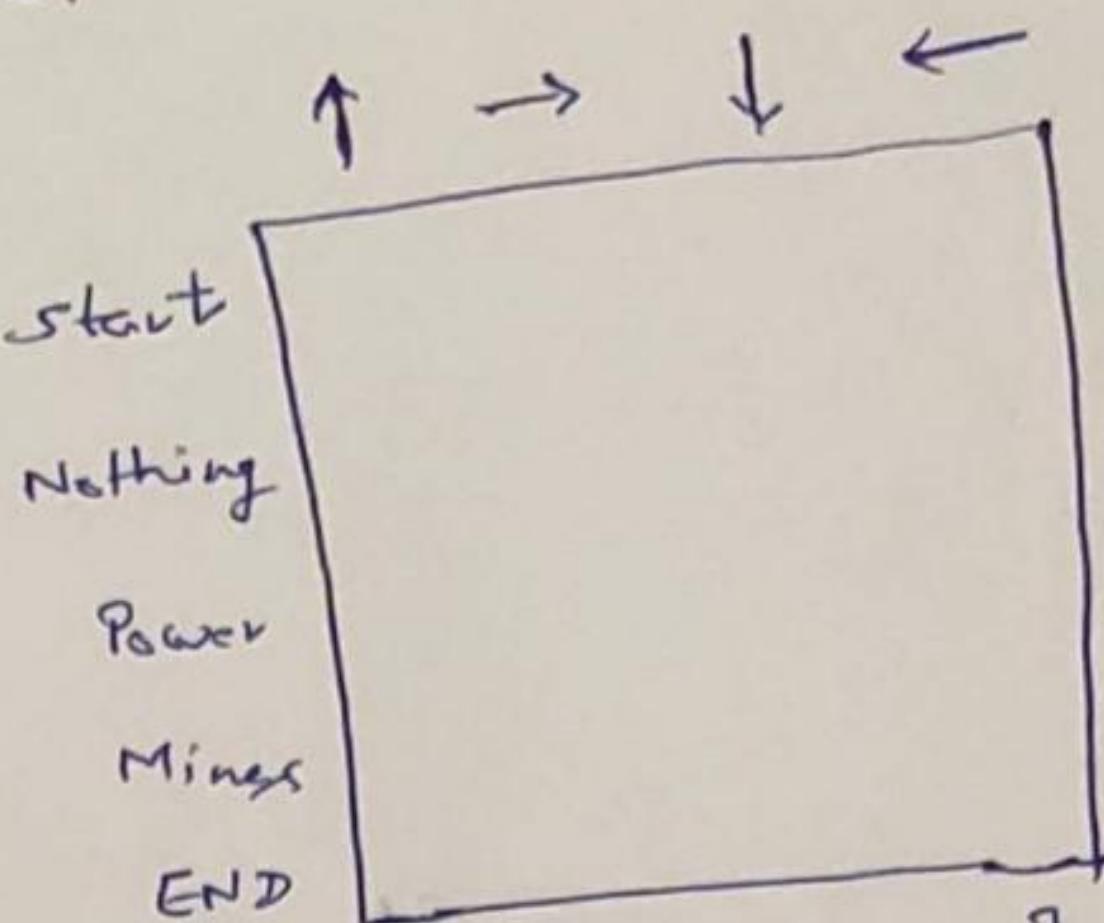
$$Q(s, a) = R(s, a) + \gamma \sum_{s'} [P(s, a, s') \cdot \max_{a'} Q(s', a')]$$

Here  $g(s, a)$  is also a recurrence (recursive equation)

Let  $\rho(s, a, s') = \frac{1}{\gamma} g(s, a) + \gamma \cdot \max_{a'} g(s', a')$

So a Q-table is just a simple look up table where we calculate the maximum expected future rewards for actions at each state.

↳ Q-table, columns are the actions start and rows are states.



$$Q^{\pi}(s_t, a_t) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots] [s_t, a_t]$$

given state  
and action

Q-Learning alg<sup>o</sup>m process-

Initialize Q-table

choose an action

Perform an action

Measure reward

update Q-table

After a lot of iteration  
a good Q-table is ready.

## Q-learning algorithm-

For each  $s, a$  initialize the table entry  $\hat{q}(s, a)$  to zero.  
 observe the current state  $s$ .

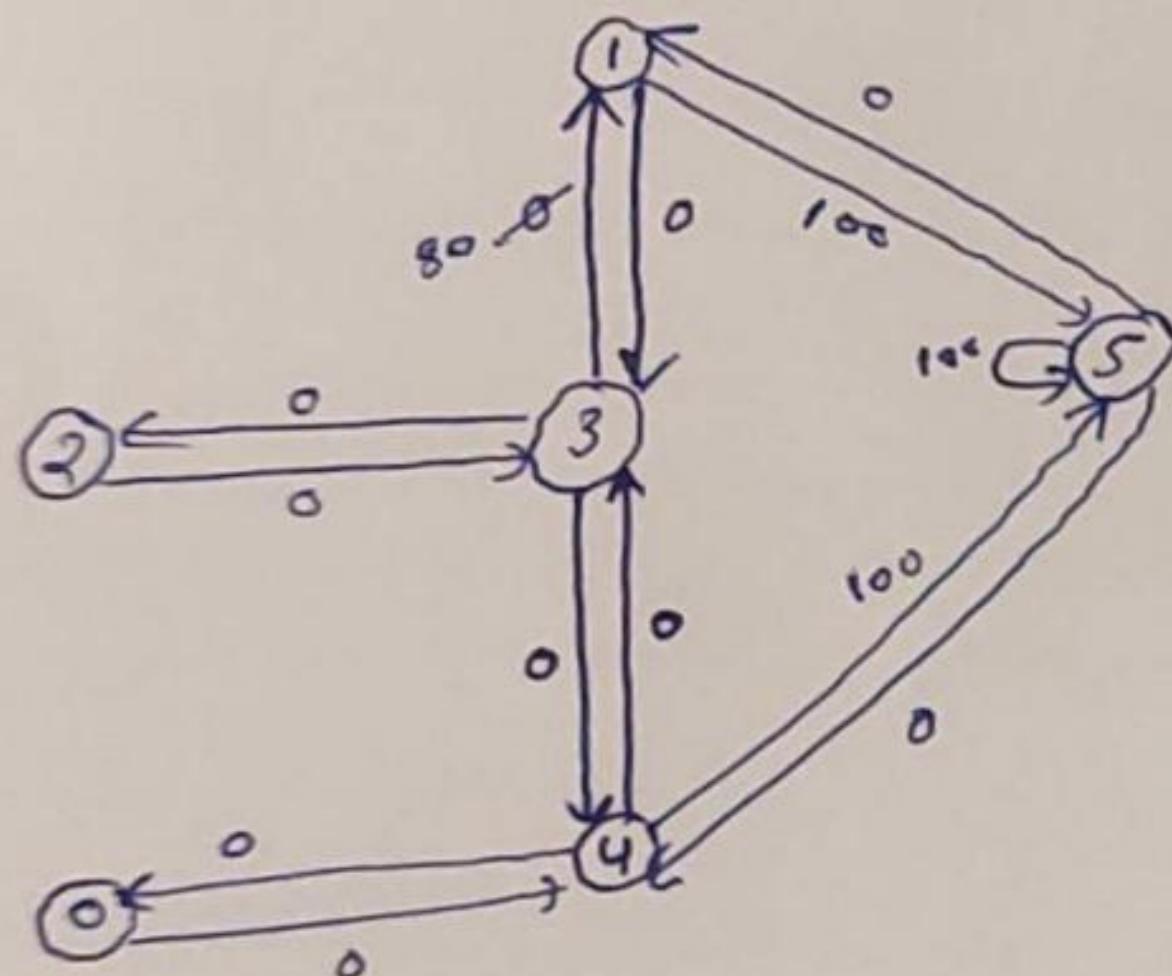
Do forever:

- select an action  $a$  and execute it.
- receive immediate reward  $r$ .
- observe the new state  $s'$
- update the table entry for  $\hat{q}(s, a)$  as follows:

$$\hat{q}(s, a) \leftarrow r + \gamma \max_{a'} q'(s', a')$$

- set  $s \leftarrow s'$

eg-



$$\gamma = 0.8$$

$$s=3 \quad a=2, 1, 4 \\ r=0 \quad s'=1$$

$$\hat{q}(3, 1) = 0 + 0.8 \max \hat{q}(0, 100) \\ = 0 + 0.8 * 100 = 80$$

$$s = s' = 1$$

$$s=1 \quad a=3, 5 \\ r=100, \quad s'=5$$

$$\hat{q}(1, 5) = 100 + 0.8 \max (0, 0) \\ = 100$$

Same process is repeated until all zero replaced with converged values.



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

#### Q-Learning:-

↓ quality of a action.

Q-function - denoted as  $Q(s, a)$

- basically describe at a given state  $s$ , which is the best action should I take in order to get highest reward.

- Quality of Reward.

- Goal - Maximize the total reward.

$$\text{ideal case } R = R_t + R_{t+1} + R_{t+2} + \dots + R_T$$

e.g.  $\textcircled{1} \rightarrow \boxed{100}$  in next 2 days.

$\textcircled{1} \rightarrow 100$  in next 200 days

Here  $r = +2$

(near related to future)

Total reward =  $0 + 2$

(Here is reward is delayed future)

Rewards should come early as possible. For later rewards we penalize the reward if the come

late.. by a discount factor  $r$ .

$$R_t = \frac{r_{t+L} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-L} R_T}{1 - \gamma}$$

discount factor \$0 < \gamma < 1\$  
 uncertainty going to increase as you go in

$$Q(s, a) = \gamma \underset{\text{Immediate reward}}{\uparrow} + \gamma \max Q(s', a')$$

future.

Q-function

BELLMAN EQ

Goal of  $\hat{q}$ -learning - learn  $q$  function  
of  $q$  function

approximation

$s_0$

$s_1$

$s_2$

Left      Right

0.36      0.41

Jump      0.87

q-values <  $B_{st}$   
action for  
a state

$$\varphi(s, a) \xrightarrow{\text{Target}} \hat{\varphi}(s, a)$$

Approximation.

$$mse\cosss = (\hat{y}(s,a) - y(s,a))^2$$



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO.

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import os
```

```
from collections import deque
```

```
from keras.models import Sequential
```

```
from keras.layers import Dense
```

```
from keras.optimizers import Adam
```

```
import random
```

```
%matplotlib inline
```

```
class Agent:
```

```
    def __init__(self, state_size, action_size):
```

```
        self.state_size = state_size
```

```
        self.action_size = action_size
```

```
        self.memory = deque(maxlen=2000)
```

```
        self.gamma = 0.95 # discount factor
```

# Exploration vs Exploitation Tradeoff.

# Exploration: Good in the beginning → helps you to try out various random things.

# Exploitation: Sample good experiences from the past (Memory).

self. epsilon = 1.0 # 100% Random Exploration in the beginning.

self. epsilon-decay = 0.995

self. learning-rate = 0.001

self. epsilon-min = 0.01

1 — 0  
(Random) (Based  
visual  
experience  
learning  
from past)  
C 1% probability that from past  
you can opt the  
random exploration)



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

AIT JOURNALY : - H:

Policy Evaluation using Monte-Carlo Methods:-

Model Free Methods in RL :- In MDP, we are given all components to solve a problem, but what happen if we are not given some of the components. (for eg. transition probabilities and rewards i.e known as dynamics of the environment.)

can we solve that problem?  
(for example - for a given position in chess what will be the next state/next position)  
(eg-2. learning & programming from net)

↳ This type of learning is called model-free learning.  
↳ In Model-free, we just focus on figure out the value functions directly from the interaction with the environment.  
↳ So basically a Monte-Carlo Method is a model free method which is used to solve that kind of problems as discussed previously.

Monte-Carlo Approach:- Monte-Carlo method is a statistical technique used to find an approximate solution through sampling. It approximates the expectation of a random variable by sampling and when the sample is greater, the approximation will be better.

↳ Expectation of a random variable  $x$  is given by.

$$E(x) = \sum_{i=1}^N x_i P(x_i)$$

$$E_{\text{Emp}}[x] \approx \frac{1}{N} \sum_{i=1}^N x_i$$

we can estimate the expected value of  $x$  for some  $N$  times and compute the average value of  $x$  as expected value of  $[x]$

↳ Monte Carlo is a model free method, meaning it doesn't require the model dynamics to compute the value function.

so a value function in R.L is denoted

as -

$$v^\pi(s) = E [R(z) | s_0 = s]$$

↑  
in  $\pi$

expected return of a state  
by following a policy  $\pi$



# POORNIMA

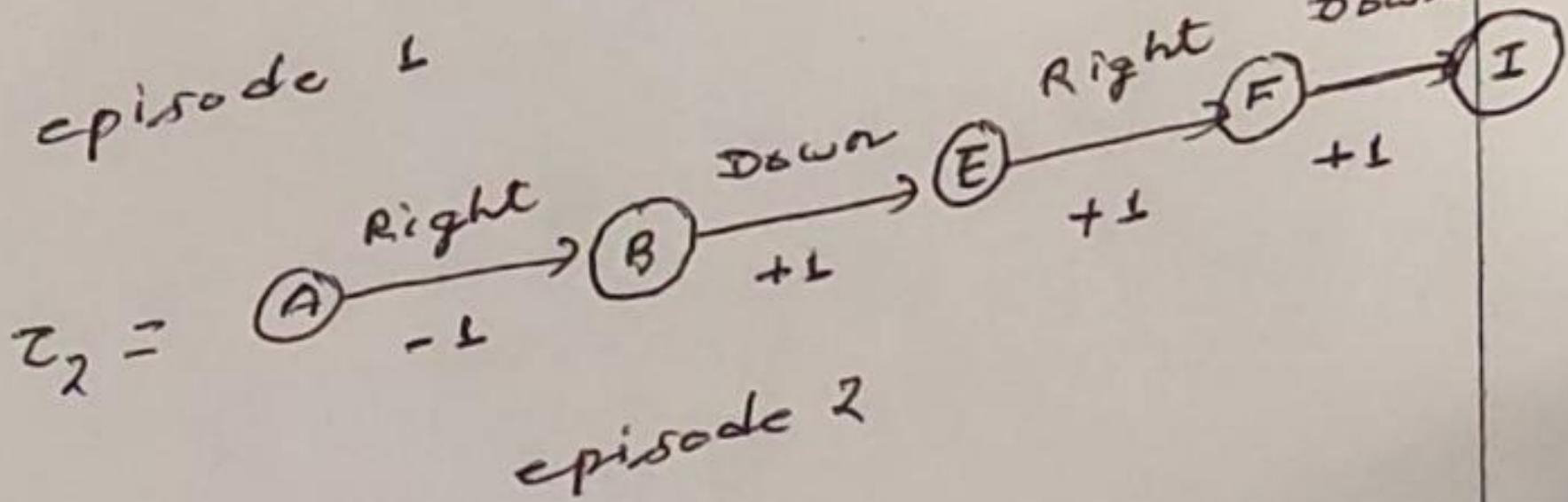
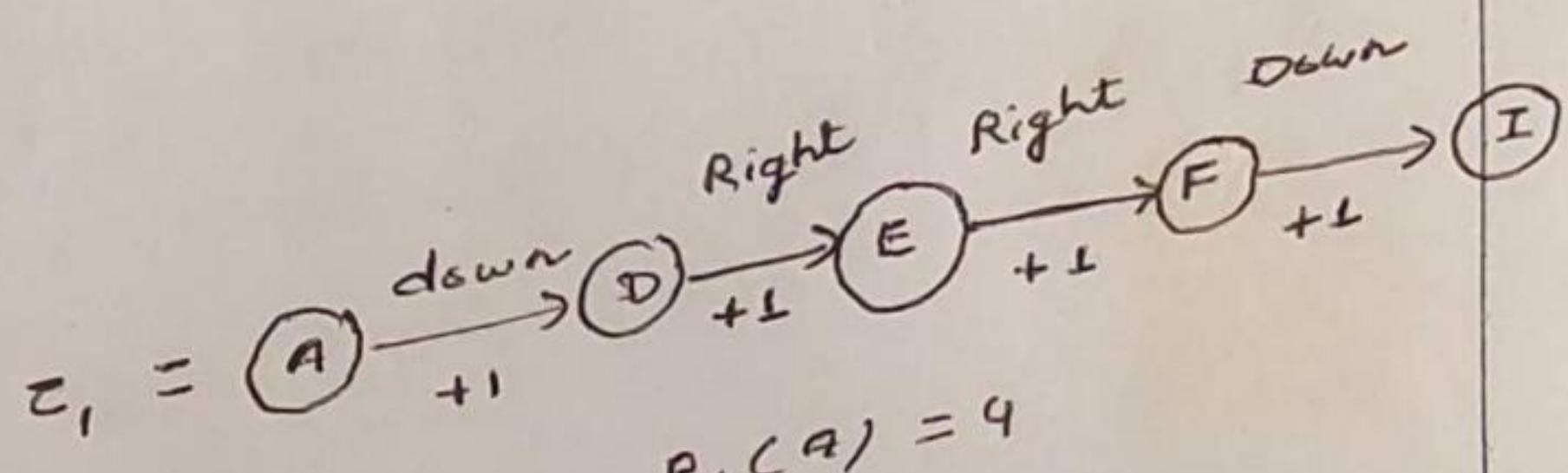
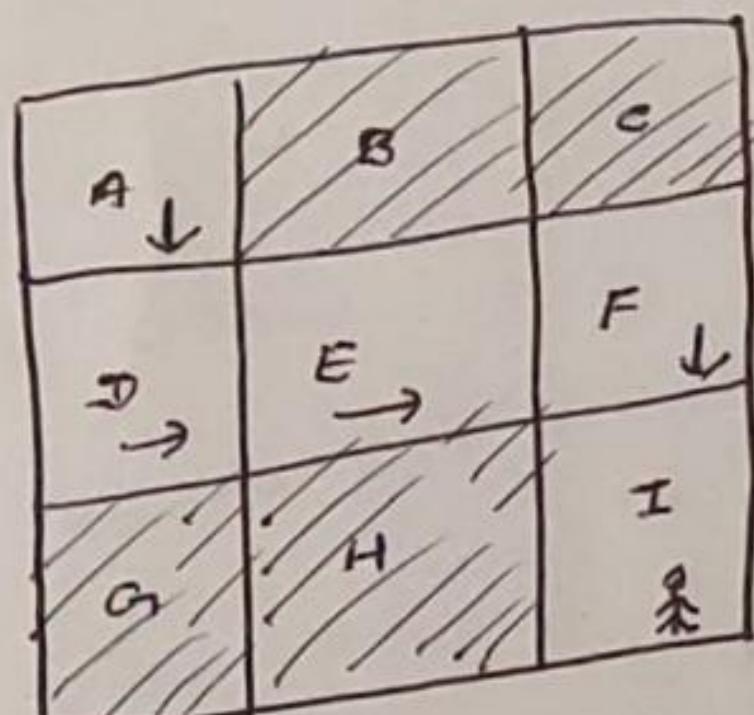
## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

In Monte Carlo prediction method, we approximate the value of a state by taking the average return of a state across  $N$  episodes instead of taking the expected return.

For eg - ① -



$$R_2(\text{A}) = 2$$

$$\tau_3 = \tau_1$$
$$R_3(\text{A}) = 4$$

so according to Monte-Carlo Method, the value of a state can be approximated by computing the average return across  $N$  episode i.e. (trajecotrys)

$$v(s) \approx \frac{1}{N} \sum_{i=1}^N R_i(s)$$

$$v(a) = \frac{1}{3} \sum_{i=1}^3 R_i(a)$$

$$v(a) = \frac{1}{3} (R_1(a) + R_2(a) + R_3(a))$$

$$v(a) = 3.7$$



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. \_\_\_\_\_

Algorithm :- MC Reduction C

$L(\text{total\_rewards}, t) = \text{sum of values of state}$   
across several episodes  
 $N(s) = \text{counter, number of times a state}$   
is visited over several episodes]  
[policy  $\pi$  is given as input]

step 1 - for  $M$  number of iteration:  
1. Generate an episode using the policy  $\pi$ .  
2. Store all rewards in the episode in the list  
called rewards.

3. For each step  $t$  in episode:  
1. Compute the return of state  $s_t$  as  
 $R(s_t) = \text{sum}(\text{rewards}[t :])$

2. Update total return of state  $s_t$   
 $\text{total\_return}(s_t) = \text{total\_return}(s_t) + R(s_t)$

3. Update the counter of  $N(s_t) = N(s_t) + 1$

step 2: Compute the value of a state by just

steps 1

-Taking the average i.e.

$$v(s) = \text{total\_return}(s) / N(s)$$



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

Model-based Reinforcement Learning:- In model-based learning, an agent will have a complete description of the environment.

↳ Here the agent knows the complete dynamics of its environment like transition probability and rewards that the agent will going to earn, then type of learning is called model-based learning.

↳ Here the agent will know the model dynamics to find the optimal policy.

edt



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

Policy iteration and value iteration:-

Reinforcement Learning

Agent - policy:  $P(a|s)$

(policy is a function that maps state to an action)  
or prob. distribution for a action in a given state)

$a \rightarrow$  action

$s \rightarrow$  state

$r \rightarrow$  reward

Environment - Transition dynamics - describe how a states in environment changes.

$P(s', r | s, a)$

↑  
next state

↑ current state

and how the agent gets rewards. so its a joint distribution of next state and reward based on the current state and action applied by the agent.

③ state -

④ action -

⑤ Reward -

state transition prob-

$$P(s'|s, a) = \sum_{r \in R} P(s', r | s, a)$$

Total discounted return -

(weighted sum of rewards).

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{N-1} R_{t+N}$$

$$\gamma \in [0, 1]$$

↑  
Discounted Factor.

value of a state -

(with respect to a policy  $\pi$ )

$$v_\pi(s) = E_\pi [G_t | s_t = s]$$

$$v_\pi(s) = E_\pi [G_t | s_t = s] \\ = E_\pi [R_{t+1} + \gamma R_{t+2} | s_t = s]$$

$$= E_\pi [R_{t+1} + \gamma G_{t+1} | s_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) \left[ v + \gamma E_\pi [G_{t+1} | s_{t+1} = s'] \right]$$

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) \underbrace{\left[ v + \gamma v_\pi(s') \right]}_{\text{next state value}}$$

↑  
current state

{ Bellman - equation



# Poornima

COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

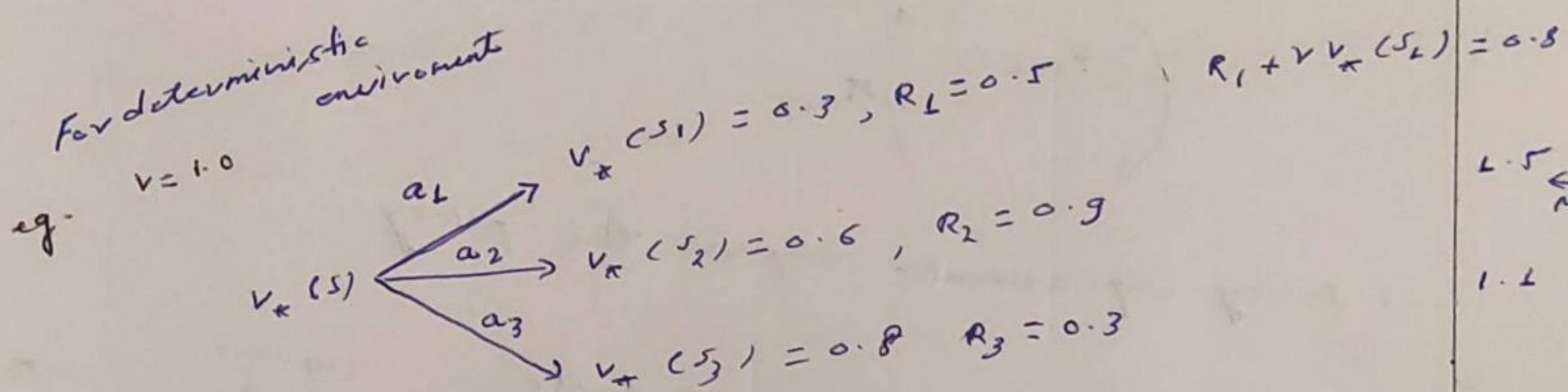
BELLMAN optimality Eq<sup>n</sup>-

$$v_*(s) = \max_a E [R_{t+1} + \gamma v_*(s') \mid s_t = s, a_t = a]$$

$$v_*(s) = \max_a \sum_{s' r} p(s', r | s, a) [r + \gamma v_*(s')]$$

For deterministic environment

$\gamma = 1.0$



$$v_*(s) = R_2 + \gamma v_*(s_2) = 1.5$$

## Finding optimal policy

policy - means what action to take in a certain state to maximize the reward.

① policy iteration algorithm -  $\rightarrow$  to find the optimal policy.

② value "

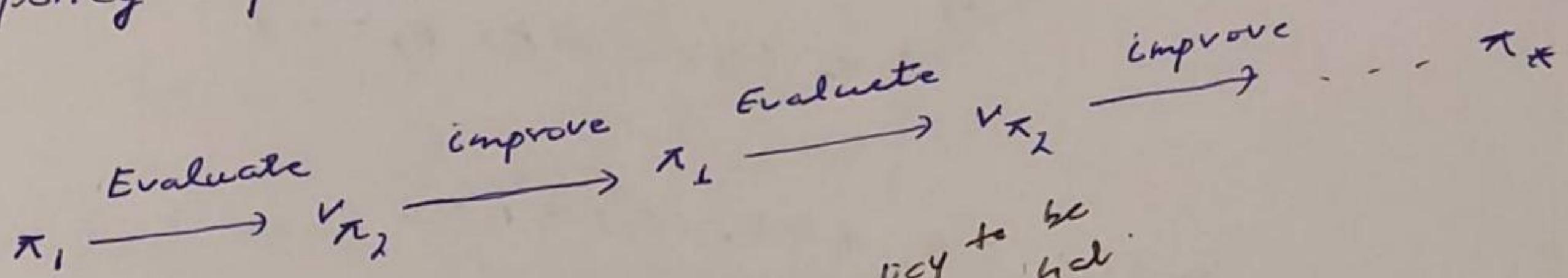
① Policy Iteration Algorithm:-

it has two parts

i) policy Evaluation - compute values for the states, for all state in environment.

$$\left( \begin{array}{c} v_{\pi(s)} \\ \vdots \\ v_{\pi(s)} \end{array} \right)_{\text{new}}$$

ii) policy improvement - improve the policy



set the random policy.

iterative policy evaluation ( $\pi$ )  $\downarrow$  random policy to be evaluated at initial

Evaluate  $\pi_1 \rightarrow v_{\pi_1}$   $\xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{Evaluate}} v_{\pi_2} \xrightarrow{\text{improve}} \dots \pi_k$

1. set the arbitrary value for the states.

2. set the arbitrary value for the states.

3. Using  $\pi$  compute new values for the states.

$$v_{\text{new}(s)} \leftarrow \sum_{s', r} p(s', r | s, \pi(s)) [r + v_{\text{old}}(s')]$$

for all states

4. Repeat from (2) until convergence of the state values.



# Poornima

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

eg:-

$$\pi(s) = \begin{cases} a_1 & \text{if } s = s_1 \\ a_2 & \text{if } s = s_2 \\ a_3 & \text{if } s = s_3 \\ a_4 & \text{if } s = s_4 \end{cases}$$

policy improvement stage :-

For all state  $s \in S$ . update  $\pi(s)$  as

$$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} \sum_{s', r} p(s', r | s, a) [r + \gamma v(s')]$$

Repeat

Repeat

for all  $s \in S$

$$v(s) \leftarrow \sum_{s', r} p(s', r | s, \pi(s)) [r + \gamma v_{old}(s)]$$

Policy  
improvement  
evaluation

until convergence

for all  $s \in S$

$$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} \sum_{s', r} p(s', r | s, a) [r + \gamma v(s')]$$

policy  
improve-

until convergence

Return  $\pi(s)$

value-iteration Algo<sup>m</sup>

Repeat

For all  $s \in S$

$$v(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\text{old}}(s')]$$

until convergence

$$v_*(s) = v(s) \text{ for all } s \in S$$

$$v_* = v(s) \text{ for all } s \in S$$

$$\pi(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]$$

$$\text{return } \pi(s)$$



# Poornima COLLEGE OF ENGINEERING

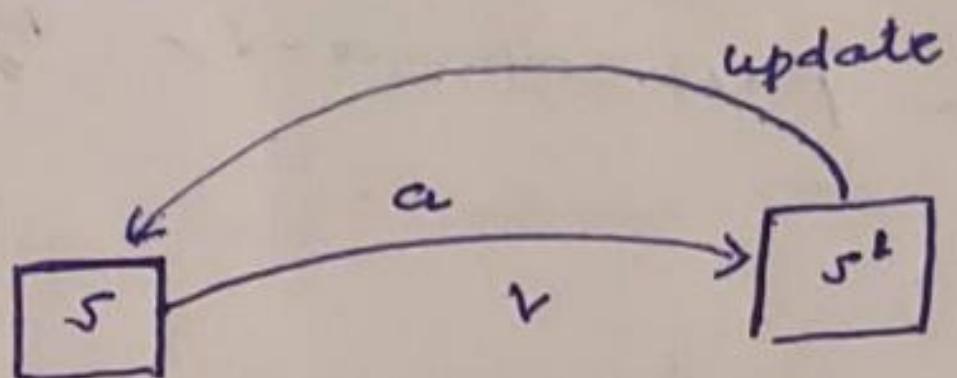
## DETAILED LECTURE NOTES

PAGE NO. ....

27/10/2014: 14:

### Temporal Difference Learning:-

↳ also known as TD(0).



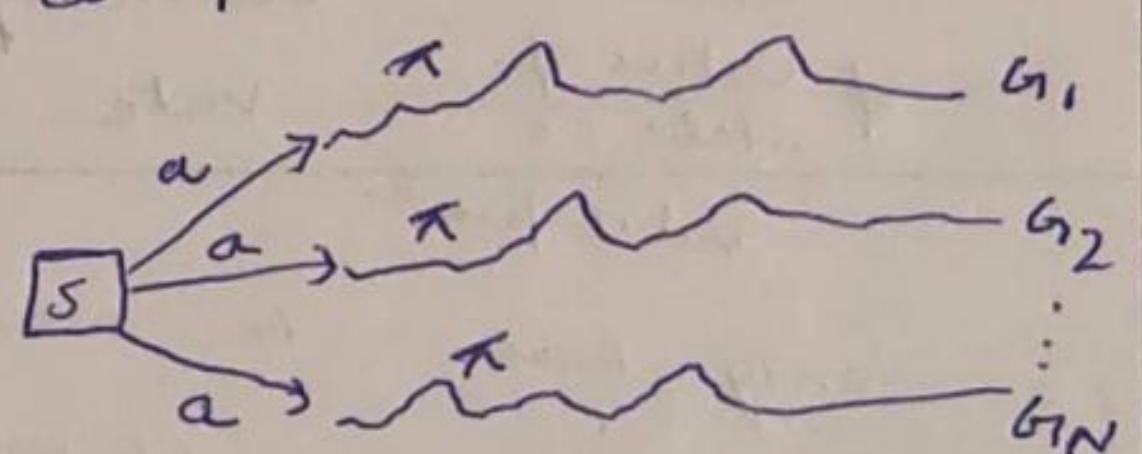
Here we can simply apply an action on state  $s$  and move to next state  $s'$  and then we can obtain the update on previous state by observing the value of  $v$  as well as our prev guess of next state  $s'$ .

$v_{\pi}(s)$   
(Value Function)

value of a state is the expected return that we obtain if we apply policy  $\pi$  from state  $s$ .

$q_{\pi}(s, a)$   
↳ q-value is something

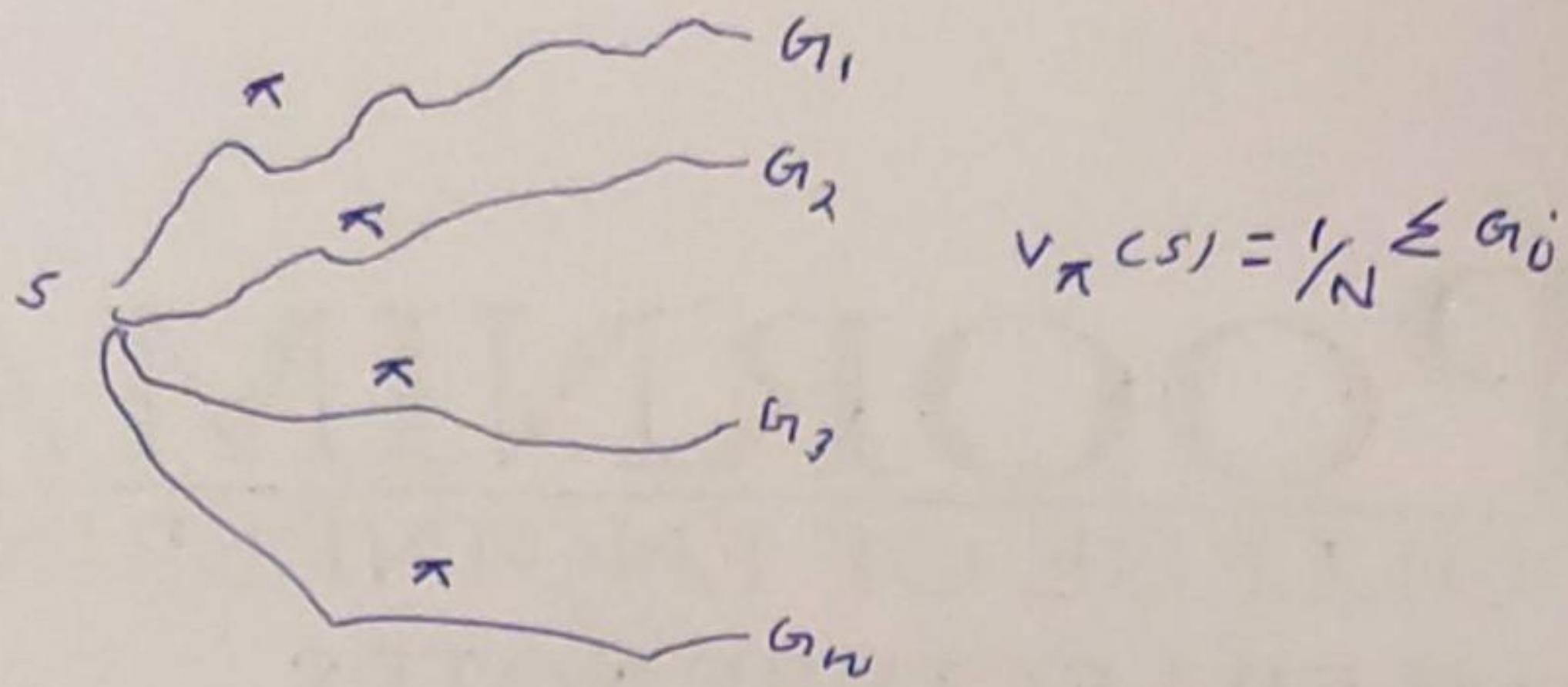
that is associated with the state as well as a particular action.



$$q_{\pi}(s, a) = \frac{1}{N} \sum G_i$$

$v^*(s)$  - optimal value of value function  
(maximum return that we can get)

$v^*(s) = v_{\pi^*}(s)$   
that depends upon the policy  $\pi^*$  that can give maximum rewards.



↳ optimal q-value for state-action pair.

$$q^*(s, a) = q_{\pi^*}(s, a)$$

that is the maximum q-value that you can get from that state and action pair onwards by following policy  $\pi$ .

$$\boxed{\pi(s) = \underset{a}{\operatorname{argmax}} \ q^*(s, a)}$$

TD(0) update rule:-

$$v(s) \xrightarrow[a]{\pi} s'$$

$v(s) \leftarrow v(s) + \alpha [r + \gamma v(s') - v(s)]$ 

new updated value of state  
old value of state i.e value func.

reward  
Learning rate

discount Factor.  
TD Error

old value of state  
New value of next state.

TD(0) only have to wait until the next time step to update the old value.



# POORNIMA

COLLEGE OF ENGINEERING

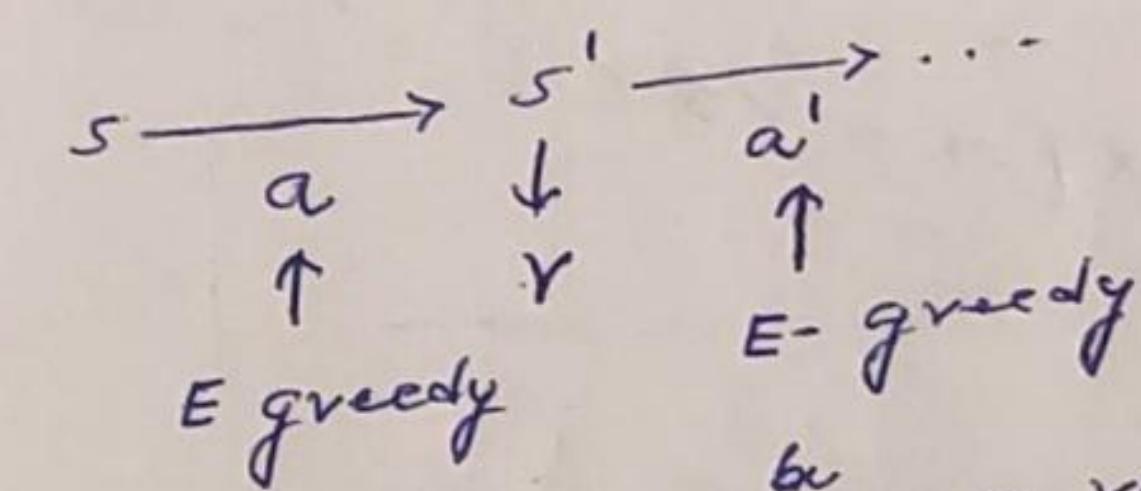
## DETAILED LECTURE NOTES

PAGE NO. ....

SARSA algorithm:- (SARSA is an acronym for state-action Reward-state-Action).

↳ SARSA is an on-policy Temporal Difference learning. [Here the agent behaves using one policy and tries to improve the same policy].

↳ It's only talk about the Q-value not the value of state that is value function.



(epsilon - greedy policy)  $\xrightarrow{b}$  6.  $\xrightarrow{100\% \text{ randomization at start.}}$   $\xrightarrow{\text{Discount Factor}}$   $\pi(s) = \begin{cases} \text{argmax } Q(s, a), \\ a \\ \text{with } p=1-\epsilon \\ \text{or} \\ \text{random } a, \text{ with } p=\epsilon \end{cases}$   $\epsilon = 0.02$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$$

↑  
Learning rate

SARSA update Rule

## SARSA algorithm

initialize the Q-table  
say  $Q(s, a) = 0$  for all  $(s, a)$  pairs.

Repeat for  $N$  episodes

$a = \epsilon\text{-greedy at } s_0$

Repeat for  $k$ -steps.

at any state  $s$ .

① Execute  $a$ , get reward  $r \rightarrow$  state  $s'$

② Next action  $a' = \epsilon\text{-greedy at } s'$

③ Get  $Q(s', a')$  from the table

④ update

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$$

⑤ Set

$$a \leftarrow a', \quad s \leftarrow s'$$

END OF FOR

END OF FOR



# Poornima COLLEGE OF ENGINEERING DETAILED LECTURE NOTES

PAGE NO. ....

cg-

	1	2	3	4
1	S	F	F	F
2	F	H	F	H
3	F	F	F	H
4	H	F	F	G

q - Table with random values.

State	Action	Value
(1,1)	Up	0.5
(1,2)	:	:
(4,1)	Up	0.3
(4,2)	Down	0.5
(4,3)	Left	0.1
(4,3)	Right	0.8
(4,4)	Right	0.5

highest  
(4, 3)  
Episilon -  
value

→ Right      (4, 3) Up      0.6  
 (4, 3) Down      0.3  
 (4, 3) Left      1.0  
 (4, 3) Right      0.9

Let reward  $r = 0$  learning rate  $\alpha = 0.1$  discount factor  $\gamma = 1$

Now according to SARSA update rule

$$Q(s, a) = Q(s, a) + \alpha (r + \gamma Q(s', a') - Q(s, a))$$

$$Q((4, 2), \text{right}) = Q((4, 2), \text{right}) + \alpha (r + \gamma Q((4, 3), \text{right}) - Q((4, 2), \text{right}))$$

$$Q((4, 2), \text{right}) = Q((4, 2), \text{right}) + 0.1 (0 + 1 \times Q((4, 3), \text{right}) - Q((4, 2), \text{right}))$$

$$Q((4, 2), \text{right}) = 0.8 + 0.1 (0 + 1 + Q((4, 3), \text{right}) - 0.8)$$

Now again select  $(Q((4, 3), \text{right}), \text{right})$  with 0.9 randomly

$$Q((4, 2), \text{right}) = 0.8 + 0.1 (0 + 1 + 0.9 - 0.8)$$

$$Q(4, 2, \text{right}) = 0.81$$