



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

1

PAGE NO.

Unit 4: - Semi Supervised Learning & Reinforcement learning

- ① Markov Decision Process (MDP)
- ② Bellman equations
- ③ Policy evaluation using Monte Carlo
- ④ Policy iteration and Value iteration
- ⑤ Q-learning
- ⑥ State-Action Reward-State-Action (SARSA)
- ⑦ Model-based Reinforcement Learning.

Markov Decision Process (MDP) :-

→ It is a feedback-based Machine learning Technique in which an agent learns to behave in an environment by performing the actions and seeing the result of actions.

- for good action - positive feedback
bad action - negative feedback.

- Agents learn automatically using feedback without any labeled data.

- Agent learn by their experience.
- RL used for sequential decision making and for long term goals.
like - Game playing, Robotics.

Definition :- Reinforcement learning is a type of machine learning method where an intelligent agent (computer program) interacts with the environment and learns to act within that.

Eg:- Robotic dog - arm movement learn

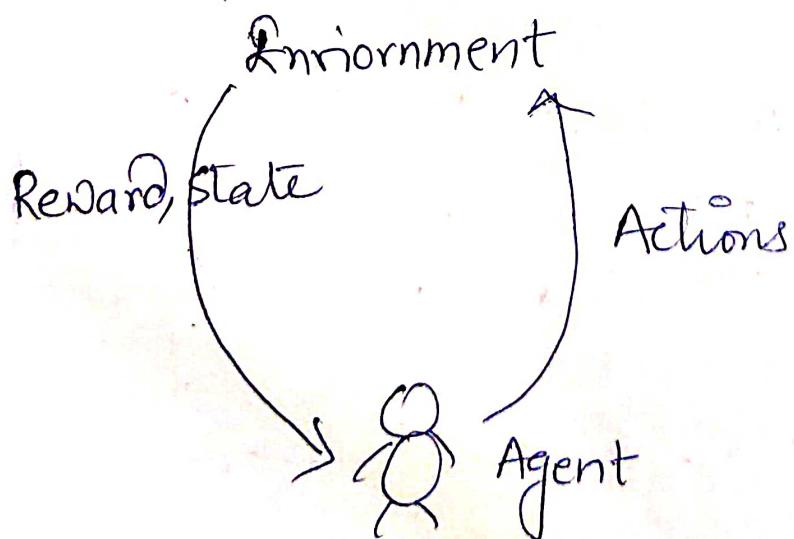
- It's a part of AI & AI agent

- Agent do these three things

① take action

② change state/remain in the same state

③ Get feedback.





POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

(2)

PAGE NO.

Terms used in Reinforcement Learning:

- ① Agent : entity perceive/explore the environment and act upon it.
- ② Environment: situation in which agent is present or surrounded.
- ③ Action : moves by agent in environment.
- ④ State: situation returned by environment.
- ⑤ Reward: feedback returned.
- ⑥ Policy : strategy applied by agent for next action.
- ⑦ Value: long-term return
- ⑧ Q-value: Similar to value, take one additional parameter

Key features of Reinforcement Learning:-

- ① Instructions not given to agent about environment and what action need to be taken.
- ② Based on hit & trial process.
- ③ Agent takes the next action and changes states according to the feedback of the previous action.
- ④ Agent get a delayed reward.
- ⑤ Environment is stochastic, agent explore to get max positive reward.

Approaches to implement Reinforcement Learning

Three way to implement reinforcement in ML

① Value based :-

- It find the optimal value (maximum) at a state under any policy.

② Policy-based :-

It find optimal policy for the maximum future rewards without using the value func.

It use 2 policy

(i) Deterministic :-

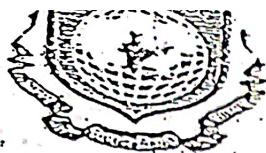
Same action at all state

(ii) Stochastic :-

Probability determines the produced action.

③ Model-based :-

Virtual-model is created and agent learn from the environment by exploring it.



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO. ②

Elements of Reinforcement Learning:-

4 main elements of RL.

- ① Policy
- ② Reward signal
- ③ Value fun
- ④ Model of Environment

→ Policy means how agent is behaving at a given time. It map previous states of environment to the actions taken on those states.

→ Core element of RL.

Deterministic policy $a = \pi(s)$

Stochastic policy $\pi(a|s)$

$$= P[A_t=a | st=s]$$

② Reward Signals:-

- RL use reward signal on each action.
- If good action & bad action is there, award signals are send to agent accordingly.

- agent goal is to maximize the total no. of rewards for good actions.
- If rewards are bad, agent change the policy.

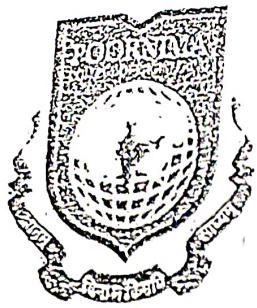
③ Value func: →

Value func gives info about how good the situation and action are.

- A reward indicates the immediate signal for each good and bad action, while value-func specifies the good state and action for the future.
- Value func depend on Reward
if no reward — no value

④ Model: →

- It tell the behavior of the environment.
- model is used for planning,
- Model make inference about the environment will behave.
- If state and action are given [model] can predict the next state and reward.



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES (4)

How RL work?
first understand with the help of example.

RL need two main things

① Environment:-

e.g: Room, maze, football ground.

② Agent

Like AI robot.

s_1	s_2	s_3	s_4	$R = +1$ diamond
s_5	s_6 WALL	s_7	s_8	$R = -1$ fire
s_9	s_{10}	s_{11}	s_{12}	

→ Agent can't cross s_6 (WALL).

if reach $s_4 \rightarrow$ Reward +1

reach $s_8 \rightarrow$ Reward -1

→ Only four actions → move — up, down, left, right

→ He have to use minimum/fewer steps.

Let path $s_9 - s_5 - s_1 - s_2 - s_3 + 1$ reward.

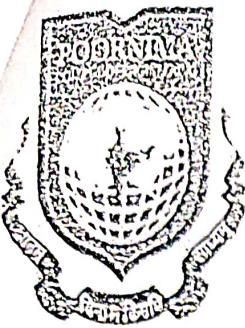
So when he will go he will put value 1 on the previous steps.

$V=1$  s_1	$V=1$  s_2	$V=1$  s_3	 s_4
$V=1$  s_5	 s_6	 s_7	 s_8
$V=1$  s_9	 s_{10}	 s_{11}	 s_{12}

- But what will happen when he start moving from the block, which has 1 value at both side (s_1)
- from $s_1 \rightarrow$ he can move s_2 and s_5
- So this approach is not good for solution
- Bellman equations, concept used behind the reinforcement learning.

The Bellman Equations :-

- This algo is associate with dynamic programming and used to calculate the value of a decision problem at a certain point by including the values of previous states.
- Key elements used in Bellman eq! .
 - ① Action performed by agent, a_t
 - ② State occurred by performing action a_t , s_t



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES P

- Reward/Feedback 'R'
- discount factor is Gamma 'γ'

Bellman eq:-

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

↓

Value /
Calculate at
a particular pointReward at
a particular
state s by
performing an
actionDiscount
factor (0-1)The value at
the previous
state
(futuristic
reward)

(Instantaneous reward)

Let

1st block:

$$V(s_3) = \max [R(s_3, a) + \gamma V(s')]$$

$$V(s') = 0$$

Let $\gamma = 0.9$

$$V(s_3) = 1$$

2nd Block

$$V(s_2) = 0 + 0.9 = 0.9$$

no Reward

3rd Block

$$\begin{aligned} V(s_1) &= (0 + 0.9 \times 0.9) \\ &= 0.81 \end{aligned}$$

4th Block

16

$$V(S_5) = 9 \times 9 = 73$$

$$V(S_9) = 0.66$$

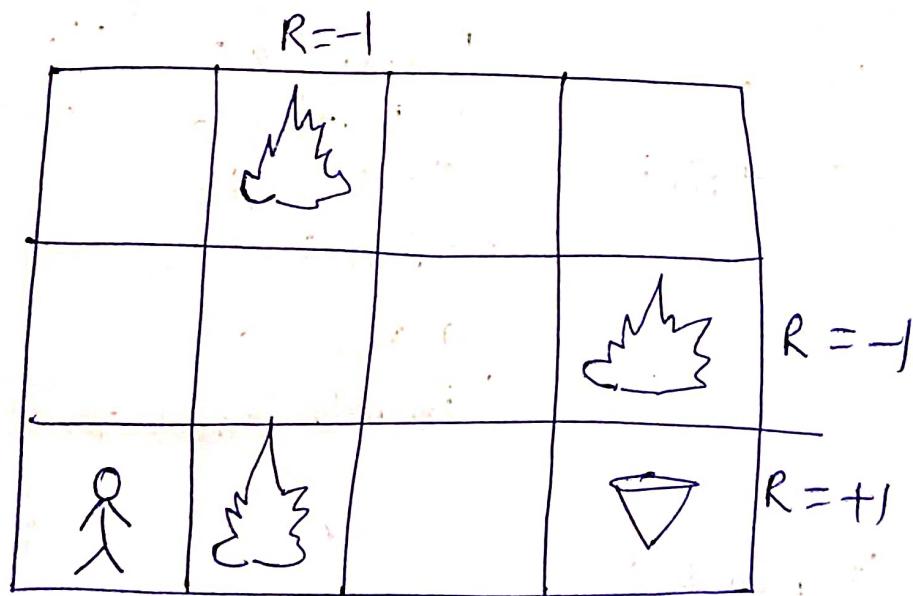
Calculate for rest value

$$S_7 = 0.9$$

$$S_{10} = 0.73$$

$$S_{11} = 0.81$$

$$S_{12} = 0.73$$





Types of Reinforcement Learning: →

2 type of RL

① Positive R

② Negative R

In this strength and freq. of the behavior is increased.

(Add something to increase the tendency that expected behavior would occur again)

Advantage of Positive RL: -

① Maximize Performance

② Sustain change for a long period of time

③ Too much reinforcement of states which can lead to an overload of states which can diminish the results.

It is opposite of PvERL.

It increase the tendency that the specific behavior will occur again by avoiding the negative condition.

Advantages of PvERL: -

① Increase Behavior

② Provide defiance to a min standard of performance

③ It only provides enough to meet up the min behavior

How to represent the agent state?

- We can represent the agent using Markov state that contain all the required info for the history
- State s_t is Markov state if it follows the given condition:

$$P[S_{t+1} | s_t] = P[S_{t+1} | s_1, \dots, s_t]$$

- Markov state follow Markov property which state that "future is independent of the past and can only be defined with the present".
- RL work on fully observed environments, where the agent can observe the environment and act for new state \rightarrow This whole process is known as Markov Decision process.



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

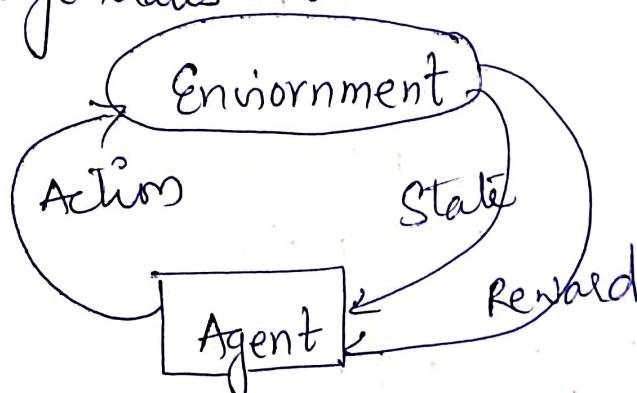
Markov Decision Process: →

01

PAGE NO.

MDP is used to formalize the reinforcement learning problems. If the environment is completely observable, then its dynamics can be modeled as a Markov Process.

- In MDP, agent constantly interacts with the environment and performs actions; at each action, the environment responds and generates a new state.



MDP contains a tuple of four elements (S, A, P_a, R)

(1) State (S)

It is a set of tokens that represent every state that the agent can be in.

(2) Action (A)

- Set of all possible Actions,

- $A(S)$ defines the set of actions that can be taken being in state S .

③ Reward (R_a)

R_a indicates the reward for simply being in the state s .

$R(s, a)$ — Reward for being in a state s and taking an action a .

$R(s, a, s')$ → Indicates the reward for being in state s , taking an action a and ending in a state s' .

④ Policy: →

A policy is a solution to the MDP. A policy is a mapping from s to a ,

e.g.: - States : s

Model : $T(s, a, s') \sim P(s' | s, a)$

Actions : $A(s), A$

Reward ; $R(s), R(s, a)$
 $R(s, a, s')$

→ probability to lead s' state by a action from s state

Policy : $\pi(s) \rightarrow a$
 π^*

MDP uses Markov property

Markov Property :-

If the agent is present in the current state s_1 , performs an action a_1 and move to the state s_2 , then the state transition from $s_1 \rightarrow s_2$ only depends



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

(2)

PAGE NO ...

on the current state and future action and states do not depend on past actions, rewards or states.

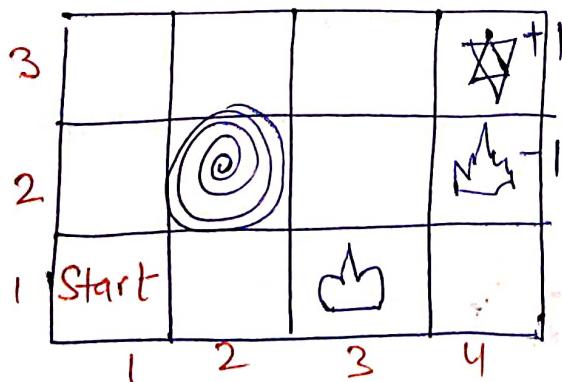
eg:- Chess game - player only focus on the current state and do not need to remember past actions or states.

Markov Process:-

MP is a memoryless process with a seq. of random States $S_1, S_2 \dots S_t$ that uses the Markov property.

Markov process is also known as Markov Chain, which is a tuple (S, P) on state S and func P .

eg:-



Two soln exist: ① R R U U R from (start)

Prob of UP = 0.8

Left = 0.1

Right = 0.1

② U U R R R

$$0.8 + 0.8 + 0.1 + 0.1 + 0.1 =$$

Policy Evaluation using Monte Carlo:-

Any method which solves a problem by generating suitable random numbers, and observing that fraction of no. obeying some property or properties can be classified as Monte-Carlo method.

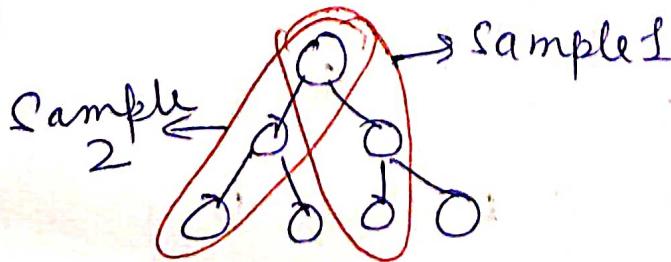
- Monte Carlo method are ways of solving the RL problem based on averaging sample returns.
- We have policy π and accordingly calculate State-value func ($V_{\pi}(s)$)
action-value func ($q_{\pi}(s)$)

So state value is cumulative sum of rewards from the starting state to end state by following a policy π .

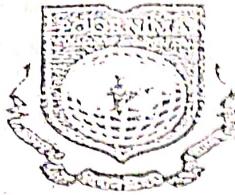
$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

e.g: video game (seq. of all states, frames from start to end game).

- In Monte-Carlo we explore all states from the starting state



→ then avg all the returns.



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO. ③

Two approaches

- ① first visit
- ② Every visit

- Keep counter $N(s)$ for each state s
- $S(s) = \text{sum of diff return (cumulative)}$

First visit MC policy evaluation

Initialize:

$\pi \leftarrow \text{policy}$

$V \rightarrow \text{state value}$

Returns (s) \leftarrow empty list $s \in S$

Repeat forever for (k episodes)

$S_s = 0$ initially.
 $N_s = 0$

Generate episode using π

When reach at that state
for each state s appearing in episode

$N_s = N_s + 1$ $G \leftarrow \text{return first occurrence of } s$

$S_s = S_s + G_t$ Append G to Returns (s)

$V_s = \frac{S_s}{N_s}$ $V(s) \leftarrow \text{avg (Return}(s)\text{)}$

Problems of Monte Carlo: →

- ① Work for episodes only not for continuous or infinite problems
- ② Need to complete an episode first before we can update all values.

Policy iteration and Value iteration

- Bellman expectation eq' for state value under policy π

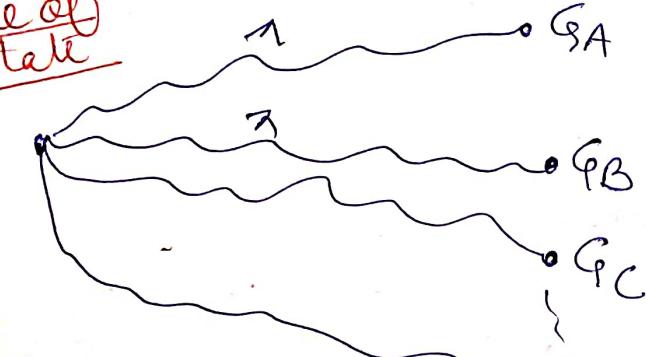
$$V_{\pi}(s) = \sum_a \pi(a|s) \sum p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

Recursive relation b/w current state's value of successive state value of taking action a in state s

In this equation adding all values so eq is updated

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum p(s', r | s, a) [r + \gamma V_k(s')]$$

e.g. Value of state

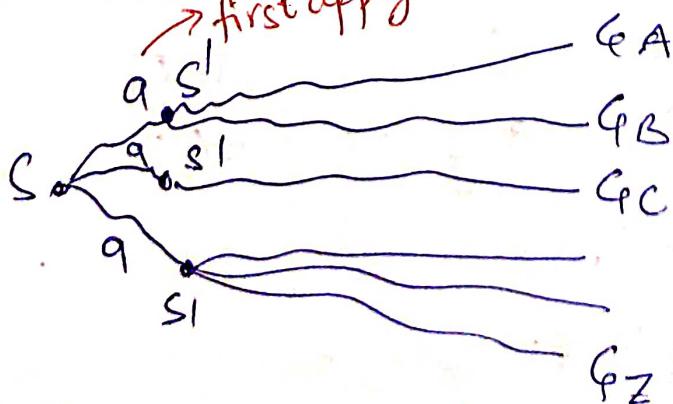


Total 26 discount episodes

$$\therefore V_{\pi}(s) = \frac{1}{26} \sum (g_A + g_B + \dots + g_Z)$$

② Value of action:-

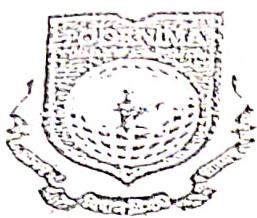
→ first apply action



$$q_{\pi}(s, a) = E_{\pi}[G_t]$$

$$S_t = S, A_t = a$$

$$\hat{q}_{\pi}(s, a) = \frac{1}{26} \sum (g_A + g_B + \dots + g_Z)$$



Poornima COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

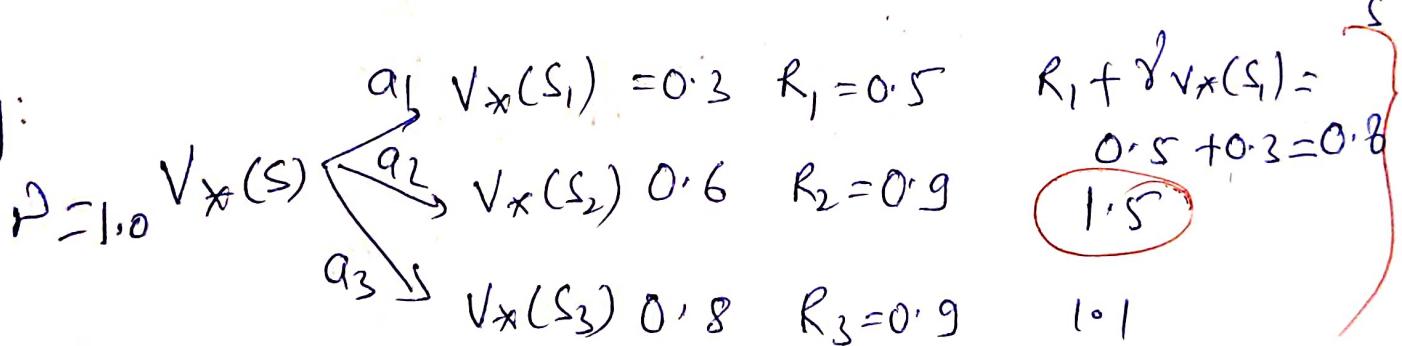
PAGE NO. 9

What is optimal state value and optimal action value.

$$V_{\pi}(s) = \max_{\pi} V_{\pi}(s)$$

$$q_{\pi}(s, a) = \max_{\pi} q_{\pi}(s, a)$$

eg:



Max
1.5

→ How to find optimal policy

Policy → What action to take in a certain state?

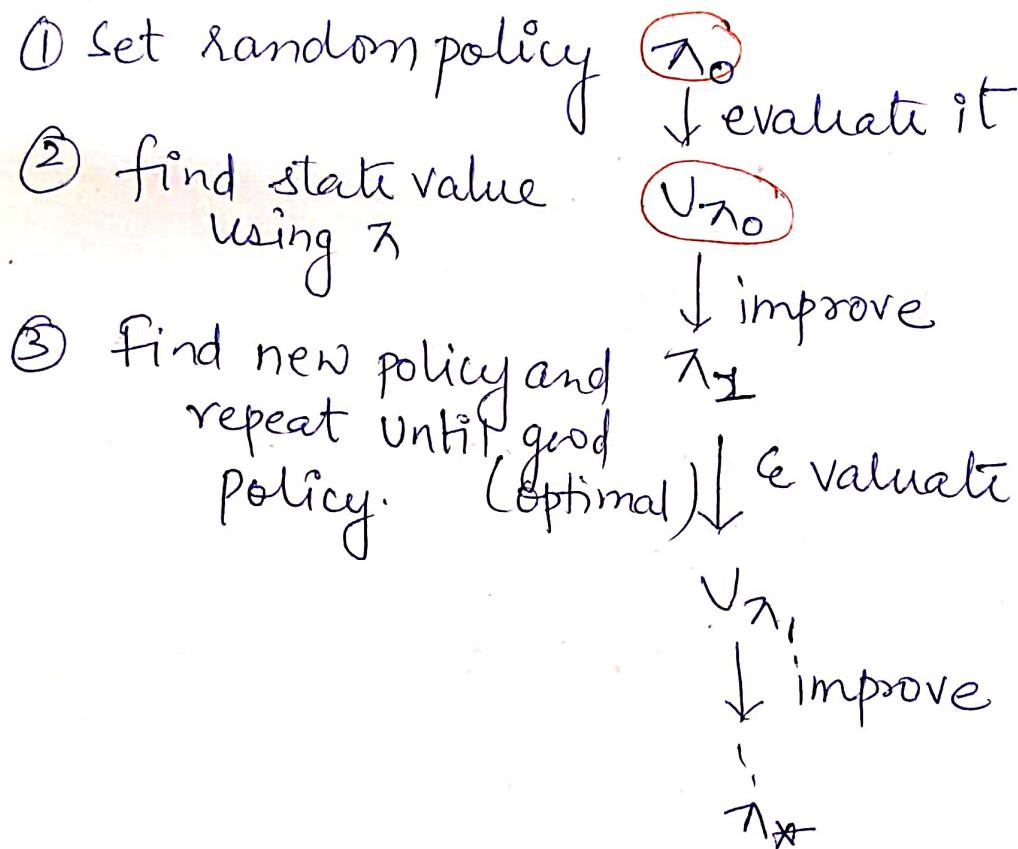
- Ans:
- ① Policy iteration
 - ② Value iteration

Policy iteration

- It has 2 parts

- ① Policy evaluation - compute state value
 $V_{\pi}(s)$ for all states
- ② Policy improvement - Improve the policy

for this



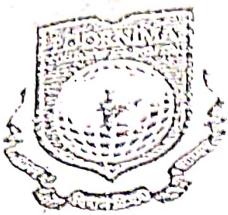
Steps:-

- ① Set π as evaluate policy
- ② Set arbitrary value for all states
- ③ Using π compute new values for the states

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

for all states

- ④ Repeat from 3 until convergence of the state values.



POORNIMA COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

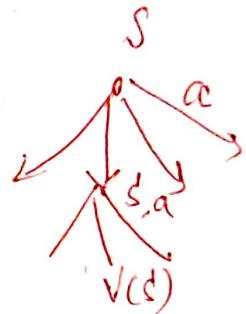
PAGE NO.

- After that calculate the policy improvement.

$$\pi(s) \leftarrow \arg \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V(s')]$$

Value iteration :-

- ① Compute value of state using actions until all value convergence



- ② find optimal value after that

- ③ Return new policy

- In this we are not calculating policy again and again only after optimal value, calculating single time only.

Limitations of 2 algorithm :-

- ① Synchronous update
- ② Need knowledge of the environment dynamics
- ③ Not learning from experience

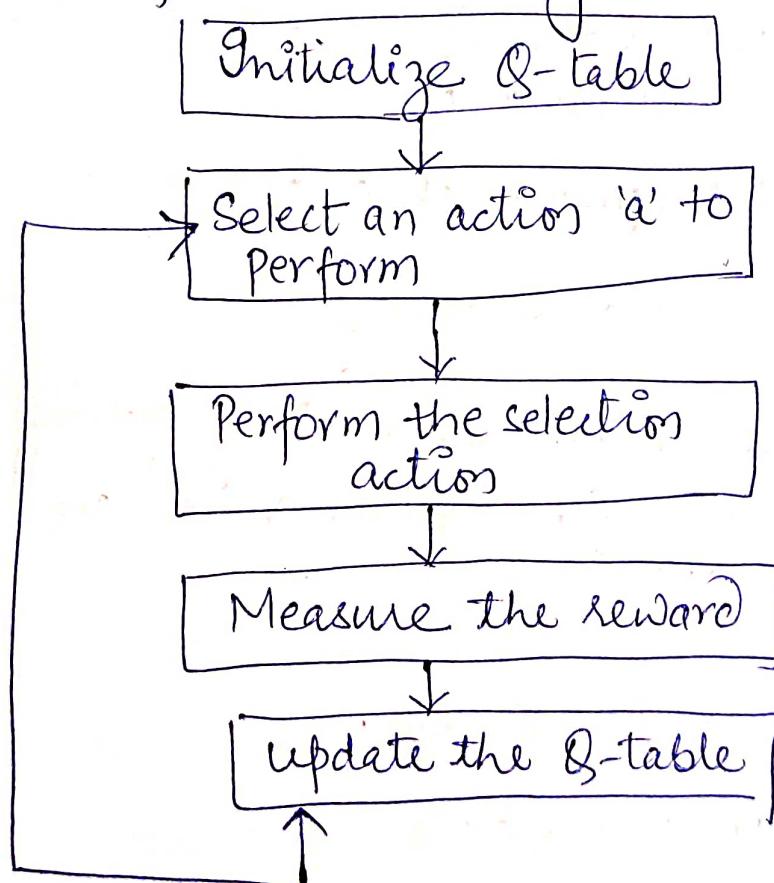
Q: Difference b/w Policy iteration and Value iteration.

Reinforcement Learning Algorithms:-

- RL algo are basically used in AI application and gaming application.
- Main algorithms are
 - ① Q-Learning
 - ② SARSA (State Action Reward State Action)

Q-Learning \rightarrow

- It is a off policy RL algorithm, used for the temporal difference learning.



- The main objective of Q-learning is to learn the policy which can inform the agent that what actions should be taken for maximizing the reward under what circumstances.



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

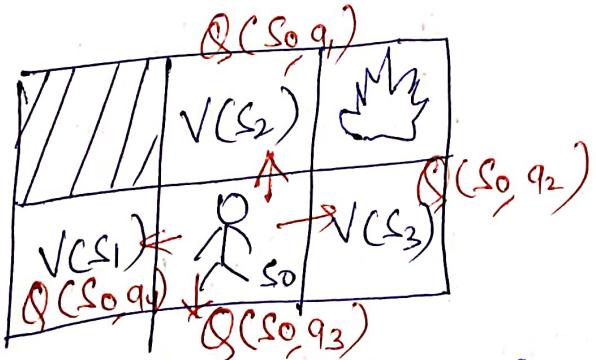
PAGE NO. 6

- find the best action to take at a current state.
- Goal of agent in Q-L is to max the value of Q.
- Value of Q-learning can be derived from the Bellman eq.

$$V(s) = \max \left[R(s,a) + \gamma \sum_{s'} P(s,a,s') V(s') \right]$$

↓ Reward ↓ Discount factor ↓ end state
 ↑ ↑ ↑

Probability



- This MDP so agent only cares for the current state and future state.

→ Q represent quality of action at each state.

- So instead of writing value, we are using pair of state & action $Q(s_i, a_i)$

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

$$V(s) = \max [Q(s,a)]$$

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} P(s,a,s') \max Q(s', a')$$

- A Q-table or matrix is created while performing the Q-Learning. It contains state & action pair and initializes the values to zero.

② SARSA (State Action Reward State Action)

- It's a slight variation of the popular Q-learning algo.
- In RL, it uses 2 policy

① On-policy:- Agent learns the value function according to the current policy

② Off-policy:- Agent learns the value function according to the action derived from another policy.

Q-learning is off-policy technique and uses the greedy approach to learn Q-value.

- SARSA uses on-policy technique and uses the current policy to learn the Q-value.

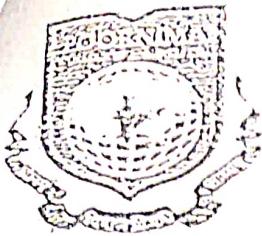
e.g:

$$\text{Q Learning} = Q(s_t, a_t) = Q(s_t, a_t) + \alpha (R_{t+1} + \gamma \max Q(s'_{t+1}, a'_{t+1}) - Q(s_t, a_t))$$

$$\text{SARSA} = Q(s_t, a_t) = Q(s_t, a_t) + \alpha (R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$

(SARSA)(s, a, r, s', a')
next state next action

form
page



Poornima COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO. 2

Model-based Reinforcement Learning :-

- Let you are playing chess.
- U made a move, then ur opponent makes a move.
- After few move u realize that your 1st move was not correct.
- So, next time we think about the opponent next move and create a space-tree in our mind and made a best move.
- Now if we replace usself by AI-agent then you get Model-based Reinforcement learning.

Model

$$MDP(S, A, P, R)$$

In model $M^0 [S, A, P^0, R^0]$

- So model M^0 going from state S to S' after performing action A , is subject to the probability $P^0(S'|S, A)$ have Reward $R^0(S'|S, A)$.

Model is anything the agent can use to predict how the environment will respond to its actions, concretely the state transition $T(s'|s,a)$ or reward func $r(s')$.

We have 2 type of models.

Model-free RL

- We perform sampling and simulation to estimate rewards.
- They perform actions either in the real world or in comp't and reward $\in \frac{\text{+ve}}{-\text{ve}}$

Model-based RL

- If we define a cost func ourselves, we can calculate the optimal actions using the model directly.
- It reduce no. of interaction with the real environment during the learning phase.

- Main loop of Model-based RL

