THIS PIECE OF STUDY MATERIAL HAS BEEN BROUGHT TO YOU BY

MUSTUDENTS UNITED

contributed by - Iram Mohd Ahmed Shaikh of college - M.H Saboo Siddik College



FOR REMOVAL OF CONTENT OR CREDITS CONTACT US AT-INSTAGRAM ID-MUSTUDENTSUNITED

OR

MAIL US AT -MUSTUDENTSUNITED@GMAIL.COM



Anjuman - I -Islam's

M.H SABOO SIDDIK COLLEGE OF ENGINEERING

8, Saboo Siddik Polytechnic Rd., Byculla, Mumbai – 400 008. Department of Computer Science & Engineering (AI & ML) (2024-2025)

For Assignment [5 Marks]

Class: CSE(AIML) - BE

Subject Name : RL

Subject Code: csporsois

Assignment No:	1		
Title:	Module 1,2,3		
Date of Performance:	15/13/25		
Date of Submission:	22/3/25		
Roll No:	211714		
Name of the Student:	Iram mohd Ahmed Shaikh.		

Evaluation:

Rubric	Marks	
On time Submission & Completion (2)	0	
Technical Content (2)	n	
Presentation & Organization (1)	0)	
Total (5)	US	
	On time Submission & Completion (2) Technical Content (2) Presentation & Organization (1)	

1

Signature of the Teacher:

Date:

	RL O
1	Explain the concept of agent, environment, state, action, seward and policy in reinforcement learning
1	Agent: The agent is the decision-making or learner in the RL System. It interact with the environment takes actions, and learns from the feedback (rewards or penalities) to improve its performance.
	Example - A Self-driving car, a chess-playing AI or a kobotic
2	Environment: The environment is everything that the agent interact with. It respond to the agent's action & determine the new State and reward.
	Example - For a self driving car, the environment includes the road, traffic and pedestrains.
	In video game, the game world acts as the environment.
3	State (s) : A state is a representation of the current situation in the environment. It contains all the necessary information that the agent needs to make a decision.
	In a chess game, the board configuration is the State
undaram	In self driving car, the state include the car's position, speed and nearby obstacles. FOR EDUCATIONAL USE

Action (a) . An action is a move or decision taken by the agent that affects the environment. The set of possible actions depends on the problems. Example - In chess, an action is moving a piece In a self-driving car, an action could be turning left, right or accelerating Reward (x) . A reward is a numerical value that provides feedback on the agent's action. The goal of the agent is to maximize the cumulative reward over time. Example - In a chess game, winning a match gives a high reward, while losing gives a negative reward. In a self-driving car, avoiding a collision may give a reward, while crashing result in a penalty. Policy (TL) : A policy defines the agent's behavior by mapping state to actions. It guides the agent in selecting actions to maximize future rewards. Policies can be deterministic & Stochastic. Example - In chess, a policy might suggest moving the queen to maximize board control. In a robotic arm, a policy determines how to grip an object: FOR EDUCATIONAL USE undaram

		H1000
	Application of Q-learning	
	Approacion	
	Game AI	
	Robotics	
	Autonomous	
	Finance:	
**	CHARLES OF THE STATE OF THE STA	
5445 TO 101		

Describe different action-value estimation methods

1 Greedy Action Selection Method

always chooses the action that currently has the highest estimated value, meaning the action that is believed to give the best reward based on previous experience

E-Greedy Action Selection Method

It is an extension of the greedy action selection strategy. It introduce an element of random exploration by selecting a random action with a small probability & and selecting the action with the highest estimated value (greedy choice) with a probability of 1-8.

How it works :

Most of the time " you pick the best option you know (with probability 1-81).

Occassionally, you pick a random option to check if there's something better (with probability E).

E is a small number, like 0.1 for 10%). So, 90% of the time you pick the best option, & 10% of the time you pick randomly.

Upper Confidence Bound Action Selection Method

It is more sophisticated approach than E-greedy
for balancing exploration & exploitation in decision
making problems, like the multi-armed bandit problem:

How it works: The key idea behind UCB is that uncertainity plays a crucial role in choosing which option (or arm) to explore next. When you're uncertain about how good a particular option is you should explore it more; As you gather more information, the uncertainity decreases, and you can start exploiting the options that perform well.

What is the Upper-Confidence-Bound (UCB) action selection in RL? UCB is a Strategy action selection in Reinforcement Learning used in the multi-armed bandit problem to execute balance exploration & exploitation It select action based on: $Q(a) + c \int Int N(a)$ where (D(a) = estimated action value c = exploration factor N(a) = number of times action a has been chosen t = total number of actions taken. VCB ensures Exploration of less-sampled actions while favouring high-reward actions. UCB ensures it reduced randomness compared to Elepsilon E-greedy methods Advantage: No need to set a learning rate. Naturally balances exploration vs exploitation. Theoretical guarantees.

	Limitations :
2	Assumes rewards are bounded and stochastic. Require maintaining counts and averages for each action
3	Not easily Scalable to complex environment.
	Tolls + Coro
	(DEC) = Estimated action value and a couplex of times action a loss here choses that a total number of actions token.
	socidos heigens of less-sameial ections in farmes se sendand actions barrens se sendander beautas 18 carrens ac
	shootso, these souled by
	desertage of to set of learning vale.

1	d] Robot vaccum cleaner goal is to clean an
	entire room.
	Justin to the score of the
	e] Game playing agent is to win the game.
I	the summer to surjection to the surjection of th
1	f] waxehouse Robot goal is to pick & deliver
	packages efficiently
	arrivada so produce hancah ad dosergue loop a storie
0] Healthcare Treatment Agent igoal is to improve
	patients health outcomes.
	the section of greaters and the agent to
	indesitional subject actions care beneficial and nitich
-	the time agent learns to chanse action that waringe
-	the moulative remords, which exercise a long-term success.
-	2 synthe house Jupon oils lader earlist loop you said in
-	sand software the agent's behavior by sanfarcing ections
	hat tend to accomplishment of these goals.
	no agreement relation below the agent develop an
-	tend out panaidon vol y silve un protecte Jomita
+	transactions there a mention aldiners
	Signax
	the state of the good as the gome not just
	ellure pieces (sub-goals)
	The senote over Should pricitize precision over just
-	
-	
-	
Second Persons	self dying for goal is to work the desiration
	the private cas goal to be work the desiration
-	

in RL, and how they are different from regular value Function.

Optimal Value Function are crucial in reinforcement learning because they represent best possible performance an agent can achieve from any given state or state-action pair, assuming it follows the optimal policy thereafter. These function serve as a benchmark for evaluating other policies:

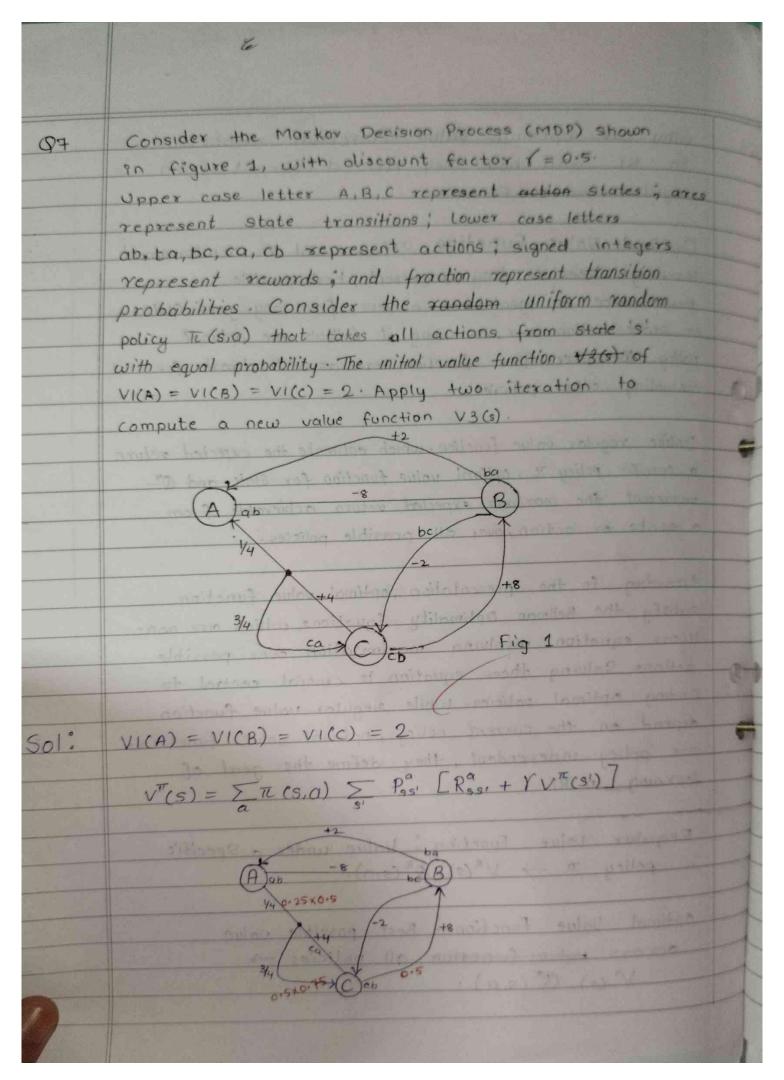
Unlike regular value function, which estimate the expected return a specific policy it, optimal value function for state- and or represent the maximum expected return achievable from a state or action, over all possible policies.

According to the presentation, optimal value function satisfy the Bellman Optimality Equations, which are non-linear equation involving maximization over possible actions. Solving these equation is crucial central to finding optimal policies. While segular value function depend on the current policy, optimal value functions are policy-independent, they define the goal of learning in Ri.

Regular Value Function: Value under a Specific policy $\pi \to V^{\pi}(s)$, $Q^{\pi}(s,\alpha)$.

optimal value function: Best possible value across value function all policies ->

V" (s), Q* (s,a):



```
Assign 1
  CONVERTED MINISTER STATES OF STATES OF
Iteration 1: . )
V2(A) = 1 [-8+(0.5) V(B)]
     = [-8 + 0.5 \times 2]
      = -8+0.5 x 2
V2(B) = 0.5 [2+(0.5) V(A)] + 0.5 [-2+(0.5) V(C)]
    = 0.5 [2+1]+0.5 [-2+1]
     = 0.5 \times 3 + 0.5 (-1)
     = 1.5 + (+0.5)
     = 1 [2810] 278.0
    e-F1 2F8-0+12-81 2-01 + 12-01 20-0
V2 (C) = 0.125 [4+0.5 V1(A)] + 0.5 [8+0.5 V1(B)] +
           0.375 [4+0.5 V.(C)]
     = 0.125 [4+1] + 0.5 [8+1] + 0.375 [4+1]
     = 0.125 [5] + 0.5 [9] + 0.375 [5]
     = 0.625 + 4.5 + 1.875
Iteration 2:
 V(S) = PSSP [Rt+1 + YV(S')]
V3(A) = 1[-8+0.5 V2(B)]
      = 1 [-8 +0.5(1)]
      = -7.5
```

```
V_3(B) = 0.5[2 + (0.5)V_2(A)] + 0.5[-2 + (0.5)V_2(c)]
     = 0.5 [2+(0.5)(-7)]+0.5 [-2+(0.5)(7)]
     = 0.5 [2+(3.5)] + 0.5 [-2+3.5]
     = 0.5 [-1.5] + 0.5 [1.5]
     = -0.75 + 0.75
V3 (c) = 0.125 [4+0.5 V2(A)] + 0.5[8+0.5 V2(B)]+
     0.375 [4+0.5 Vice)]
     = 0.125 [4+0.5(-7)] + 0.5 [8+0.5(1)] +
       0.375 [4+0.5(7)]
     = 0.125 [4+(-3.5)] +0.5[8+0.5] +
          0-375 [4+3.5]
     = 0.125 [0.5] + 0.5 [8.5] + 0.375 [7.5]
  = 0.0625 + 4.25 + 2.8125
     = 7.1250 V (2.0++ 250.0
         0.125[5] 4.0.5[9] 4.0.275
```



Anjuman - I -Islam's

M.H SABOO SIDDIK COLLEGE OF ENGINEERING 8, Saboo Siddik Polytechnic Rd., Byculla, Mumbai – 400 008.

Department of Computer Science & Engineering (AI & ML) (2024-2025)

For Assignment [5 Marks]

Class : CSECAIME) Subject Name : RL Subject Code:

Assignment No:	2	
Title:	Module 4,5,6	
Date of Performance:	15/4/25	
Date of Submission:	FEEL CONTRACTOR OF THE PARTY OF	
Roll No:	211714	
Name of the Student:	Gram Mohd Ahmed Shaikh	

Evaluation:

Rubric	Marks
On time Submission & Completion (2)	a
Technical Content (2)	or
Presentation & Organization (1)	7
Total (5)	B
	On time Submission & Completion (2) Technical Content (2) Presentation & Organization (1)

Signature of the Teacher:

Date:

Assign-2 Explain key concept involved in Monte carlo Prediction ... Generate Episodes . . The agent follows a given policy (deterministic. or stochastic) to interact with the environment and generate complete episodes (from start to terminal. state), recording states, actions and rewards. 2 Sample Episodes Returns. For each state encountered during the episode calculate the return Gt, which is the total-discounted sum of rewards from that that state to the end of the episode. Gt = Rt+1 + 8 Rt+2 + 12 Rt+3 + ... + 8 T-t-1 RT 3 Estimate Value Function update the value function V(s) by averaging the returns observed for that State across multiple episodes. This can be done using: First-visit Mc: Use only the first occurrence of a state in each episode. Every visit Mc: Use all occurrence of the state. Repeat Over Many Episodes Continue generating episodes & updating estimates until the value function converges to a Stable approximation of the true value Function.

So, basically Monte Carlo prediction is a model free method that relies on complete episodes to estimate value functions based on actual returns, making it effective in environment with unknown dynamics. It is particularly useful when learning from experience with needing to know transition probabilities.

Assign 2 RL (2) Describe the update rule used an TD prediction. Temporal Difference prediction involves adjusting the value estimate of a state based on the observed reward and the estimated value of the next state. This allows learning to occur directly from raw experience, without needing a model of the environment. Temporal Difference Updated Rule -Temporal Difference Error: At each time Step, the TD prediction algorithm calculates a TD error, which represents the gap between the current prediction and the improved estimate. It is computed as: 6 = R+++ + 8. V(S++1) - V(5+) Rt+1 - Reward received after transitioning from state St to St+1. Y - discount factor (0 < Y < 1) which determines the importance of future reward. V(St) - current value estimate for the current State . V(St+1) - Value estimate for the next state 2 Update Equation : Once the TO error is computed, the value of the current state is updated using:

V(St) + V(St) + a. 6

	a: Learning rate - controls how much new information
	overndes the old.
	Oreman and the second s
3	Effect of Learning Rate (a)
	A high learning rate (x close to 1) causes fast learning
	but may lead to instability.
	A low learning rate provides more stable updates
	h low learning race provides those state
I The same of	but require more data to converge.
	Convergence:
4	Control
	Under appropriate condition (e.g the environment follows the
	Markov property), TD(0) prediction converges to the
	true value function over time.
N ill	more more
	That means the estimated values will become accurate
	as more transitions are observed and processed
	(3) V - (1:15) V - (-1:19 - 6)
	east - Reward second after hanstinging from
Milk alk his	The state St to Stri
	1 - Herewood farter (0 - 1 & 1) which debrances
	the manitoner of filers are and
	topolis the solution of the the course
A CONTRACTOR OF THE PARTY OF TH	State - State
	Start two others storales soler a factorite
	the Santage States
	core the 10 crear is competed the color of the
	private balabamies atore their
No.	Barrie Con Variable
TO SHARE THE PARTY OF THE PARTY	

Discuss a case study where dynamic channel allocation technique have been employed to improve the performance of wireless network

case Study: Dynamic Channel Allocation in IEEE 802.11 Wireless LANS

Background: Wireless Local Area Network (WLAMs),
especially those based on the IEEE 802.11 standard
(Wi-Fi), often Suffer from performance issues 30
high-density environment (eg. campuses, offices, malls) due to
co-channel interference and limited spectrum availability.
Traditional Static Channel assignment leads to Suboptimal
performance because it cannot adapt to varying traffic
loads and interference levels.

Problem: In a university campus Scenerio, multiple
Access Points (APs) were deployed across different
department. However, due to overlapping channels and
fixed allocations, Student & Staff experienced.

High packet loss
Increased latency
Poor throughput during peak usage.

Solution: Implementation of Dynamic Channel Allocation (DCA) -

The network administration implemented a DCA algorithm that:

Continuously monitored network conditions (eg. signal to-noise ratio, interference). Automatically reassigned channels to API based on current network load and interference patterns: Prioritized critical application (eg. online exam, video lecture) by dynamically reducing congestion on their channels Technique Used: A centralized DCA algorithm was employed which used real-time measurement from all APs 8: · Detected congested or noisy channels. · Reallocated APs to less congested frequencies. Balanced user load across APs and Channels. Results : After deploying DCA: Throughout improved by 40%. Packet loss reduced by 25%. · User experience became more consistent, especially during peak times The network adapted better to dynamic usage pattern (eg. class changeovers, events). Conclusion: This case study demonstrate how dynamic channel allocations technique can significantly enhance wireless network performance in real-world environment - By adapting to real-time condition. DCA provides better spectrum utilization, lower interference & improved quality of Service (Qos) for users.

Aspect Elevator Dispatching Dynamic Channel Job-Shop Allocation Scheduling. Key Managing unpredict- Handling Signal Allocation jobs to Challenges able, human traffic. Interference & machine with dynamic usage. Environ- Real-time, physical Wireless Communication Industrial, discrete ment world with user system with fluctur- System with job-behavior ting, demand. Specific roles. Objectives Minimize wait! Maximize bandwidth, Minimize completion travel time, save reduce interference. time, idle time & energy	Compare in elevi	ator dispatching	1 comomic at	
Aspect Elevator Dispatching Dynamic Channel Job-Shop Allocation Scheduling. Key Managing unpredict—Handling Signal Allocation jobs to Challenges able, human traffic. Interference 2 machine with dynamic usage. Environ—Real-time, physical Wireless Communication Industrial, discrete ment world with user system with fluctura—System with job-behavior ting, demand. Specific roles. Objectives Minimize wait! Maximize bandwidth, Minimize completion travel time, save reduce interference. time, idle time 2 delays. Constraints: Limited number of Limited frequency Machine capacit elevators, floor bands, legal job sequence requests. Common Multi-objective Need for real—combinational time adaptability complexity.		Jon-Shop	scheduling applie	cations.
Challenges able, human traffic. Interference 2 machine with dynamic usage. Multiple Constraints. Environ- Real-time, Physical Wireless Communication Industrial, discrete ment world with user system with fluctua- system with jobbehavior ting, demand. Specific roles. Objectives Minimize would maximize bandwidth, Minimize completion travel time, save reduce interference time, idle time 2 delays. Constraints: Limited number of Limited frequency Machine capacit elevators, floor bands, legal job sequence requests. Common Multi-objective Need for real- combinational time adaptability.	Aspect	Elevator Dispatching	Dynamic Channel	Joh-snop
ment world with user system with floctua- System with job- behavior ting, demand. Specific roles. objectives Minimize wait! Maximize bandwidth, Minimize completion travel time, save reduce interference. time, idle time & energy delays. Constraints: Limited number of Limited frequency Machine capacit elevators, floor bands, legal job sequence requests: regulations. Common Multi-objective Need for real- combinational time adaptability Complexity.			Interference 2	Allocation jobs to machine with
travel time, save reduce interference. time, idle time & energy eletays. Constraints: Limited number of Limited frequency Machine capacit elevators, floor bands, legal job sequence requests: Requests: regulations: Common Multi-objective Need for real- combinational time adaptability Complexity:		world with user	system with fluctua-	System with Job -
Constraints: Limited number of Limited frequency Machine capacit elevators, floor bands, legal job sequence requests: regulations. Common Multi-objective Need for real- time adaptability complexity.	o bjectives	travel time, save		time, idle time &
common Multi-objective Need for real- combinational time adaptability complexity.	constraints.	Limited number of elevators, floor	bands, legal	Machine capacita
Enamerges		Multi-objective	Need for real-	combinational complexity.
	Enallinges			