

21' APRIL

03

DAY 093-272 | Wk 14

SATURDAY

M	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

MARCH '21

M	T	W	T	F	S	S
			1	2	3	4
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30		

APRIL '21

Reinforcement learning

9

10

11

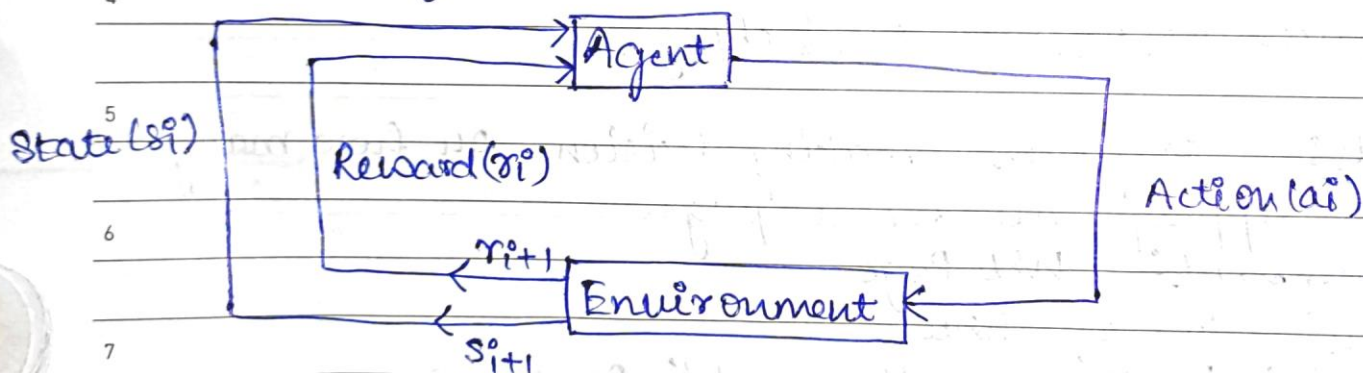
A subfield of ML that teaches agent how to choose an action from its action space within particular environment in order to maximize reward over time.

4 Elements:

12

- 1) Agent: The program you train, with the aim of doing job.
- 2) Environment: The world real or virtual in which the agent perform actions.
- 3) Action: A move made by agent which causes a status change in environment.
- 4) Rewards: The evaluation of action, positive or negative.

4



04 SUNDAY

Characteristics of RL

- 1) No supervisor (Reward signal)
- 2) Sequential Decision Making
- 3) Time plays crucial Role
- 4) Feed back is delayed
- 5) Agent's actions determine the subsequent data it receives.

<u>Basic</u>	<u>Positive Reinforcement</u>	<u>Negative Reinforcement</u>
1) <u>Def.</u>	Process of introducing a stimulus, to increase the <u>probability of occurrence of pattern or behaviour.</u>	Process of <u>removal of unfavourable stimulus</u> for the purpose of encouraging good behaviours.
2) <u>Stimuli</u>	Added	Removed.
3) <u>Consequences</u>	Pleasant	Unpleasant
4) <u>Reinforcer acts as</u>	Reward	Penalty.
5) <u>Result in</u>	<u>Strengthening and Maintaining Responses.</u>	<u>Avoiding and escaping responses.</u>
6) <u>Drawbacks</u> <u>DSRB</u>	<u>Maximize performance</u>	<u>Maintain minimum performance.</u>

S	M	T	W	T	F	S
				1	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	

MAY '21

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

MARCH '21

Wk 13 | DAY 086-279

SATURDAY

27

Applications of Reinforcement Learning

- 1) Robotics for Industrial Automation
- 2) Business strategy planning
- 3) Machine learning and Data processing
- 4) Aircraft control and Robot Motion control
- 5) Manufacturing Optimization problems
- 6) Sequential Scheduling Problems

Approaches:

1) Value Based: In this, we try to maximize a value function $V(s)$. In this method agent is expecting a long term return of current states under policy π .

2) Policy Based: In this we try to come up with such policy that the action performed in every state helps you to gain Max reward in future.

(i) Deterministic:

For every state same action is produced by policy π .

(ii) Stochastic:

Every action has a certain probability

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

SUNDAY 28

3) Model Based: In this we need to create a virtual model for each environment. The agent learns to perform in that specific environment.

Terminologies:

- 1) Agent: It is an assumed entity which performs actions in an envt. to gain some reward.
- 2) Environment (e): A scenario that an agent has to face.
- 3) Reward (R): An immediate return given to an agent when they perform specific task.
- 4) State (s): Current situation returned by environment.
- 5) Policy (π): strategy applied by agent to decide the next action based on current state. $\pi(s) \rightarrow a$
- 6) Value (V): It is expected long-term return with discount as compared to short term reward.
- 7) Value function: It specifies the value of a state that is a total amount of reward. It is an agent which should be expected beginning from the state.
- 8) Q-Value / Action Value: Q is quite similar to value. It takes additional parameters as a current action.

APRIL '21	S	M	T	W	T	F	S
					1	2	3
4	5	6	7	8	9	10	
11	12	13	14	15	16	17	
18	19	20	21	22	23	24	
25	26	27	28	29	30		

MAY '21	S	M	T	W	T	F	S
30	31						
2	3	4	5	6	7	8	
9	10	11	12	13	14	15	
16	17	18	19	20	21	22	
23	24	25	26	27	28	29	

learning → MDP
Model → 0

MARCH '21

Wk 14 | DAY 089-276

TUESDAY

30

Markov Decision Process

Markov property: The future is independent of the past given the present.

A state s_t is Markov iff

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

MDP consist of

1) Set of state S ,

2) Set of Action, A

3) Transition function $T(s, a, s')$

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

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

5) Policy: $\pi(s) \rightarrow a$
 π^*

↳ Transition Model gives an action's effect in a state.

↳ Policy is a solution to MDP. A policy is the mapping from s to a . It indicates action a to be taken while in state ' s '.

MDP is Mathematical framework to describe an environment in Reinforcement Learning.

Q-learning

A value based method of supplying information to inform which action an agent should take
[value based algorithm updates value for using Bellman eqn]

Q-function: $Q(s, a)$

→ Its value is the maximum discounted cumulative reward achieved by starting from state s, applying action a as first action.

→ Q is the reward received immediately upon executing action a from state s , plus the value of following optimal policy thereafter.

$$Q(s, a) \equiv r(s, a) + \gamma V^*(s, a)$$

$$\pi^*(s) = \arg \max_a Q(s, a)$$

$\pi^*(s)$, optimal action in state s .

$r(s, a)$, reward

V^* value of immediate successor state

γ discount.

Q-learning Algorithm:

1) for each s, a initialize table entry $\hat{Q}(s, a)$ to zero

2) Observe current state

3) Do forever:

↳ Select an action and execute it

↳ Receive Immediate reward r

↳ Observe new state s'

↳ Update table entry for $\hat{Q}(s, a)$ as:

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_a \hat{Q}(s', a')$$

↳ $s \leftarrow s'$

1) for each (s, a) initialize table entry $\hat{Q}(s, a)$ to zero

2) Observe current state

3) Do forever:

(i) Select an action and execute it

(ii) Receive Immediate reward r

(iii) Observe new state s'

(iv) Update table entry for $\hat{Q}(s, a)$ as:

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_a \hat{Q}(s', a')$$

(v) $s \leftarrow s'$

Examples

1) controlling a walking Robot:

✓ Agent: Program controlling a Walking Robot.

✓ Environment: Real World

✓ Action: One out of 4 moves.

(1) Forward

(2) Backward

(3) left

(4) Right

✓ Reward: Positive when it approaches the target destination direction; negative when it wastes time and goes in wrong direction.

2) Placement of Ads on Webpage:

✓ Agents: Program making decision on how many ads are appropriate for a page.

✓ Environment: Web Page

✓ Action: One of three

(i) Putting another add on the pg.

(ii) Drop add from the pg.

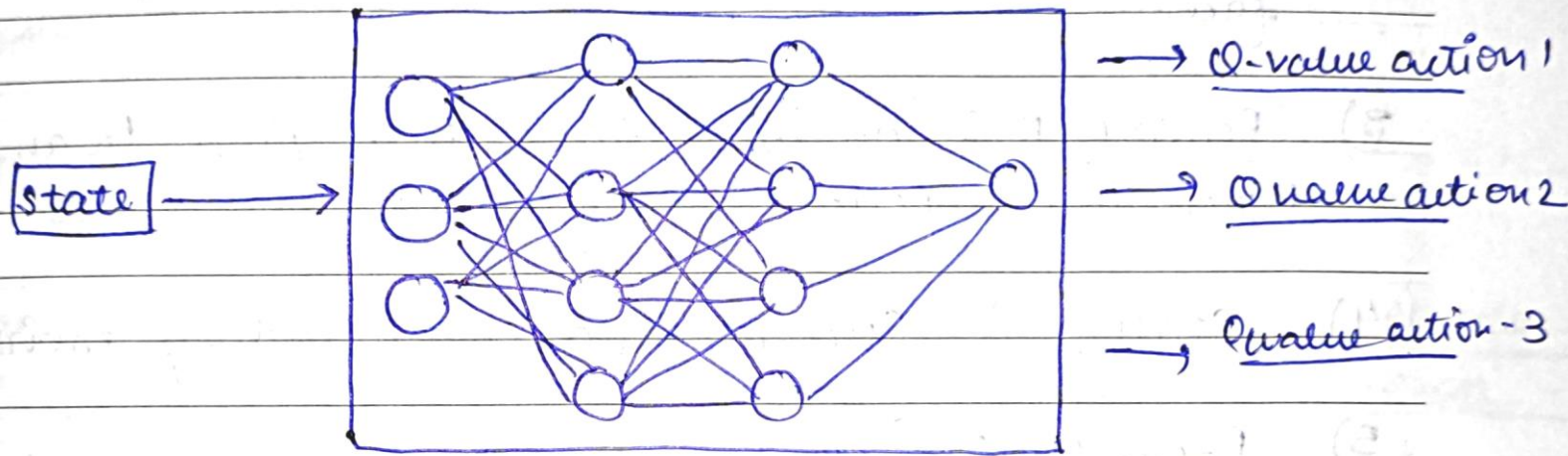
(iii) Neither add nor remove

✓ Reward: Positive when revenue increases
Negative when Revenue drops.

Deep Q-learning

In deep learning we use Neural Network to approximate the Q-value function.

The state is given as I/P and Q values of all possible actions is generated as the O/P.



Step:1 Store all the past experience in memory.

Step:2 Next action is determined by Max O/P of Q-network.

Step:3 Loss function = Mean sq. error of predicted Q-value and target Q-value. - Q^*

Bellman Equation:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$