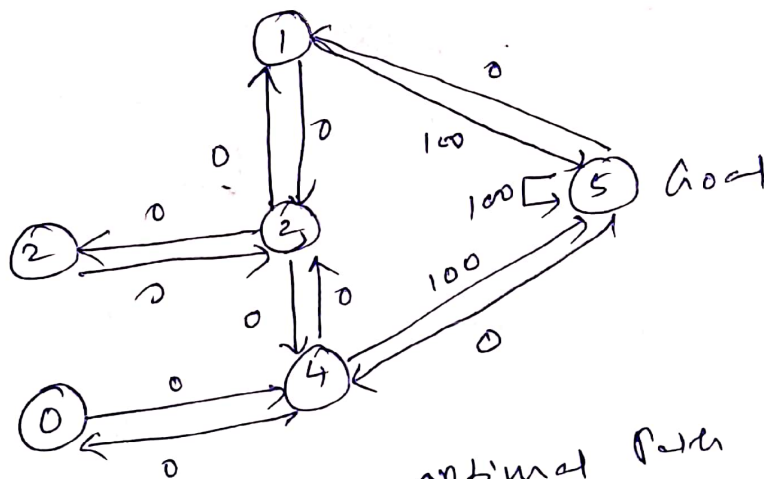


What is Q-Learning

Q-Learn is a reinforcement learning that will find the next best action, it is a model free MA reinforcement learning. The objective of this model is to find the best course of action given its current state.

Example:- Advertisement recommendation system.

In a normal ad recommendation system, the ads you get are based on your previous purchases, using Q-learning we can optimize the ad recommendation system to recommend products that are frequently brought together.



he need to find an optimal path from each state to goal state 5 where 5 is considered as goal state.

The actions which are leading to goal state are rewarded with 100, all the other actions are rewarded with 0.

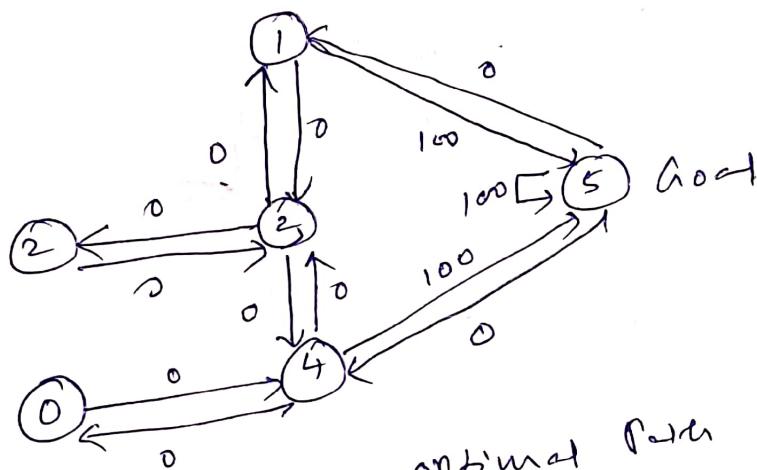
What is Q-Learning

①

Q-Learn is a reinforcement learning that will find the next best action, it is a model free MDP reinforcement learning. The objective of this model is to find the best course of action given its current state.

Example:- Advertisement recommendation system.

In a normal ad recommendation system, the ads you get are based on your previous purchases, using Q-learning we can optimize the ad recommendation system to recommend products that are frequently brought together.



he need to find an optimal path from each state to goal state 5 where 5 is considered as goal state.

The actions which are leading to goal state are rewarded with 100, all the other actions are rewarded with 0.

Q-learning algorithm is applied to get optimal Q

Table.

State

0 1 2 3 4 5

Actions

Reward matrix:

$$R = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & 1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix} \end{matrix}$$

when you are in state 1 you can perform
actions 3 & 5
(0) (100). \rightarrow indicates there is no
direct edge.

Learning Rate = 0.8, Initial State = 1

Now initialize Q Matrix

$$Q = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

Initial state is 1 there are two actions for
this state 3 & 5 with 0 & 100 rewards.
So will select next state as 5, immediate reward
5. from next state 5 you can perform an
action 1, 4, & 5 with rewards 0, 0, & 100.

Q-learning algorithm is applied to get optimal policy.

Reward matrix:

$$R = \begin{matrix} & \text{Action} \\ \text{State} & 0 & 1 & 2 & 3 & 4 & 5 \end{matrix}$$

0	-1	-1	-1	-1	0	-1
1	-1	-1	-1	0	-1	100
2	-1	-1	-1	0	-1	-1
3	-1	0	0	-1	0	-1
4	0	-1	-1	0	-1	100
5	-1	0	-1	-1	0	100

when you are in state 1 you can perform actions 3 & 5 (0) (100). \rightarrow indicates there is no direct edge.

Learning rate = 0.8, Initial State = 1

Now initialize Q matrix

$$Q = \begin{matrix} & \text{Action} \\ \text{State} & 0 & 1 & 2 & 3 & 4 & 5 \end{matrix}$$

0	0	0	0	0	0	0
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
5	0	0	0	0	0	0

Initial State is 1 there are two actions for this state 3 & 5 with 0 & 100 rewards. So I will select my next state as 5, immediate reward 5. from next state 5 you can perform an action 1, 4, & 5 with rewards 0, 0, & 100.

It has 5 possible actions 1, 4, 8, 5

(3)

Q-learning equation

$$Q(\text{state}, \text{action}) = R(\text{state}, \text{action}) +$$

$$\gamma \max_{\text{all actions}} (Q(\text{next state}, \text{action}))$$

$$Q(1, 5) = R(1, 5) + 0.8 \max(Q(5, (1, 4, 5)))$$

$$= 100 + 0.8 * 0 = \underline{\underline{100}}$$

For the next episode I will consider 3 as initial state.

It can go to 1, 2, 4.

So will select the next state as 1

$$Q(3, 1) = R(3, 1) + 0.8 \max(Q(1, (3, 5)))$$

$$= 0 + 0.8 \times 100 = \underline{\underline{80}} \checkmark$$

Next I will select 0 as initial state
0, 4.

$$Q(0, 4) = R(0, 4) + 0.8 \max(Q(4, (0, 3, 5)))$$

$$= 0 + 0.8 \times 100 = \underline{\underline{80}}$$

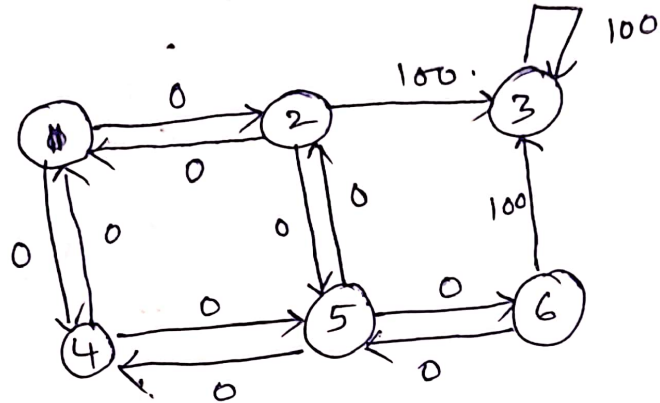
Next I will select 1 as my initial state

$$Q(1, 3) = R(1, 3) + 0.8 \max(Q(3, (1, 2, 4)))$$

	0	1	2	3	4	5
0	0	0	0	0	80	0
1	0	0	0	64	0	0
2	0	0	0	64	0	0
3	0	80	51	0	80	0
4	64	0	0	64	0	100
5	0	80	0	0	80	100

Q-Learning Example . Example - 2

5



R

	1	2	3	4	5	6
1	-1	0	-1	0	-1	-1
2	0	-1	100	-1	0	-1
3	-1	-1	100	-1	0	-1
4	0	-1	-1	0	-1	0
5	-1	0	-1	0	-1	-1
6	-1	-1	100	-1	0	-1

Q

	1	2	3	4	5	6
1	0	90	0	81	0	0
2	81	0	100	0	0	0
3	0	0	100	0	0	0
4	81	0	0	0	72	0
5	0	90	100	0	0	90
6	0	0	0	0	0	90

$$Q(\text{State}, \text{action}) = R(\text{State}, \text{action}) + \gamma * \max_{\text{Next State, action}} (Q(\text{State}, \text{action}))$$

	1	2	3	4	5	6
1		90		81		
2	81		100			
3						
4		81				
5			100			
6						

(1,2)

$$\begin{aligned}
 Q(1,2) &= R(1,2) + 0.9 * (\max[(2,1)(2,3)(2,5)]) \\
 &= 0 + 0.9 * (\max[0, 100, 0]) \\
 &= 0.9 * 100 = 90.
 \end{aligned}$$

(2,1)

$$\begin{aligned}
 Q(2,1) &= R(2,1) + 0.9 (\max[R(1,2)(1,4)]) \\
 &= 0 + 0.9 (\max(90, 0)) \\
 &= 81.
 \end{aligned}$$

(4,1)

$$\begin{aligned}
 Q(4,1) &= R(4,1) + 0.9 (\max(R(1,2)(1,4))) \\
 &= 0 + 0.9 * (\max(81, 90)) \\
 &= \underline{\underline{81}} \checkmark.
 \end{aligned}$$