

MOD 6

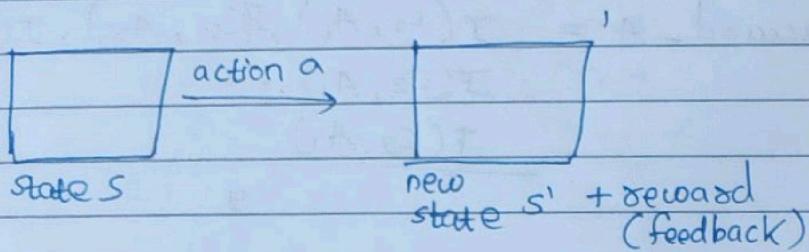
Reinforcement Learning

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

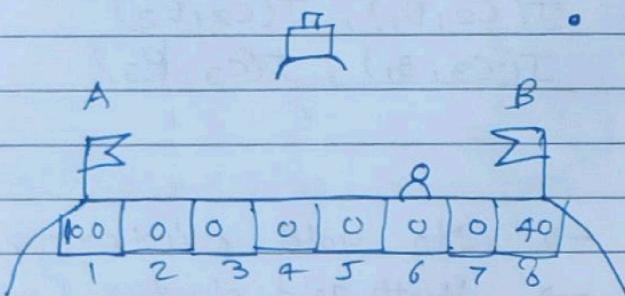
5/07/2021
Monday

- Reinforcement \rightarrow Strengthening something by adding some additional materials.

- No x, y



Mass Rover Example



- Discount factor γ (Gamma)

$0 \leftrightarrow 1$
↑
lesser cost of travelling.

- Return

$$\begin{aligned} \text{Return} &= 0 + \gamma \times 0 + \gamma^2 \times 0 + \gamma^3 \times 100 + \dots \\ &= (0.5)^3 \times 100 \\ &= 12.5 \end{aligned}$$

$$\begin{aligned} \text{Return} &= (0.99)^3 \times 100 \\ &= 97.02 \end{aligned}$$

$$\gamma = 0.99$$

γ return $\delta \gamma^n$

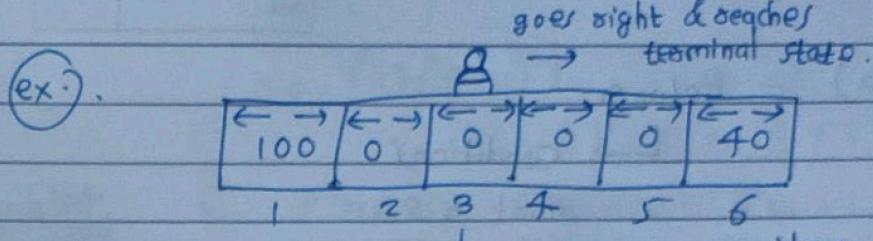
0.5 12.5 always left

0.99 97.02 always left.

0.5 20 always right

0.99 39.6 always right

$$\} \text{Return} = 40 \times \gamma$$



1, 6 \rightarrow terminal state. consider the state whose move goes as well.

$$r = 0.5$$

For state 3 & right.

$$\begin{aligned} \text{Return} &= r^0 \times 0 + r^1 \times 0 + r^2 \times 40 \\ &= 0.25 \times 40 \\ &= 2.5 \times 4 \\ &= 10 \end{aligned}$$

For state 2 & right.

$$\begin{aligned} \text{Return} &= r^3 \times 40 \\ &= 0.125 \times 40 \\ &= 1.25 \times 4 \\ &= 5 \end{aligned}$$

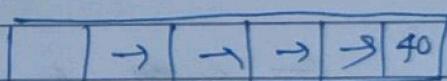
For state 3 &

0 $\leftarrow r \rightarrow 1$
 \uparrow
get closer reward tendency.
 \uparrow
get farther reward tendency.

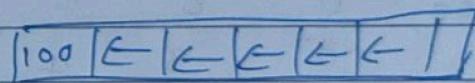
16/04/27
Tuesday

Policies (π)

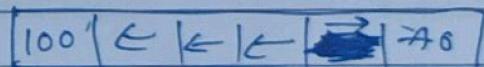
lower (right)



higher (left)

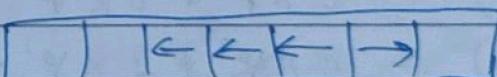


Neader



higher unless

(lower is 1 step away)



(Optimal)

$\pi(s) = a$
state action

$Q(s, a)$ = Return
 ↑
 State-action
 value function

↳ Conditions
 start in s
 takes action a
 behaves optimally

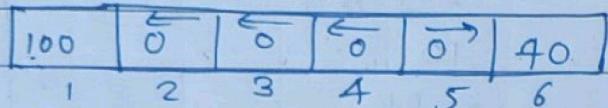
(eg).

$$Q(3, \rightarrow) = ?$$

↓
 $p = 0.5$

first move right

then follow
optimal path.



$$Q(3, \rightarrow) = \frac{0 + p \times 0 + p^2 \times 0 + p^3 \times 0 + p^4 \times 100}{(3 \rightarrow 4 \rightarrow 3 \rightarrow 2 \rightarrow 1)}$$

$$Q(5, \leftarrow) = 0 + p \times 0 + p^2 \times 0 + p^3 \times 0 + p^4 \times 100$$

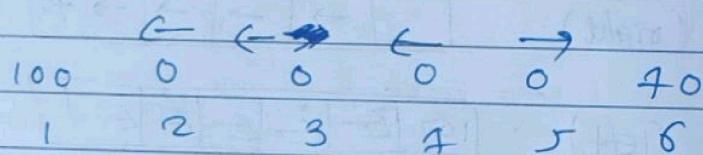
$$= (0.5)^4 \times 100$$

$$= 0.25 \times 25$$

$$= 6.25$$

$$(5 \rightarrow 4 \rightarrow 3 \rightarrow 2 \rightarrow 1)$$

(eg.2)



$$p = 0.5$$

$$Q(3, \leftarrow) = 0 \times p^0 + 0 \times p^1 + 100 \times p^2 = 25$$

$$Q(3, \rightarrow) = 0 \times p^0 + 0 \times p^1 + 0 \times p^2 + 40 \times p^3 =$$

~~$$= 40 \times 0.5 \times 25$$~~

~~$$= 20 \times 25$$~~

~~$$= 500$$~~

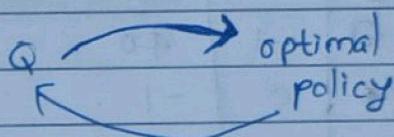
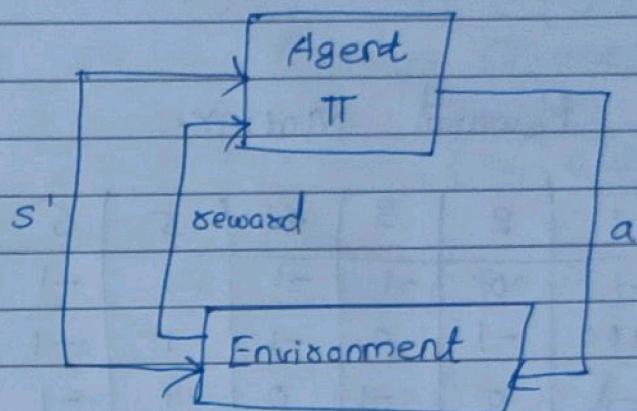
$$0 \times p^0 + 0 \times p^1 + 0 \times p^2 + 0 \times p^3 + 100 \times p^4$$

$$Q(3, \rightarrow) = 6.25$$

$$\max(Q(s, a))$$

18/04/24
Thursday

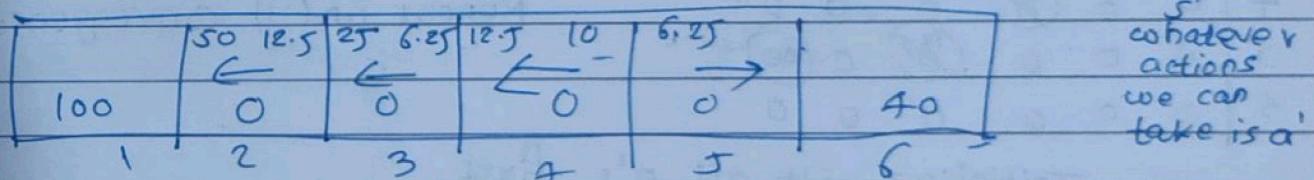
Markov Decision Process



Bellman Equation (guides all reinforcement algos).

$$Q(s, a) = R(s) + \gamma \max_{\text{all actions over } s'} Q(s', a')$$

Current state s
 action a
 after going to s'
 whatever actions we can take is a'



$$Q(4, \rightarrow) = 0 \times \gamma^0 + 0 \times \gamma^1 + 40 \times \gamma^2 = 40 \times 0.25 \quad \cancel{= -20 \times 0.25} \\ = 5 \cdot 10$$

Using Bellman equation.

$$Q(3, \rightarrow) = R(3) + (0.5) \max([12.5, 10])$$

$$Q(4, \rightarrow) = R(4) + (0.5) \max([6.25, 20])$$

$$Q(4, \leftarrow) = R(4) + (0.5) \max([25, 6.25]) = 0 + 12.5 = 12.5$$

$$\begin{aligned}
 Q(s, a) &= R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \\
 &= R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \dots]
 \end{aligned}$$

22/04/2024
Monday

(e)

Reward Matrix.

	1	2	3	4	5	6
1	-1	0	-1	-1	-1	-1
2	100	-1	0	-1	-1	-1
3	-1	0	-1	0	-1	-1
4	-1	-1	0	-1	0	-1
5	-1	-1	-1	0	-1	40
6	-1	-1	-1	-1	0	-1

Q-Matrix (Based on Bellman Equation)

	1	2	3	4	5	6
1	0	0	0	0	0	0
2	100	0	0	0	0	0
3	0	50	0	0	0	0
4	0	0	25	0	0	0
5	0	0	0	12.5	0	40
6	0	0	0	0	0	0

Choose action

Perform action

Calculate Q-matrix (Bellman eqn)

Repeat until convergence!

Initialize with all 0's first.

$$Q(2, 3) = R + \gamma \max_{\text{over all } a} Q(3, a)$$

$$= 0 + 0.5 \max([0, 0])$$

= 0

→ take from

Q-matrix.

$$Q(2,1) = 100 + 0.5 \times \max(\{10, ?\}) \\ = 100$$

(now Q matrix changes).

$$Q(3,2) = 0 + 0.5 \times \max(\{100, 0 ?\}) \\ = 50 \quad (\text{change Q-matrix}).$$

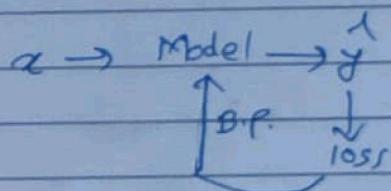
$$Q(4,3) = 0 + 0.5 \times \max(\{0, \dots, 0, 50 ?\}) \\ = 25$$

$$Q(5,4) = 0 + 0.5 \times \max(\{0, 0, \dots, 25 ?\}) \\ = 12.5$$

$$Q(5,6) = 40 + 0.5 \max(\{0, 0 ?\}) \\ = 40$$

Deep Reinforcement Learning

- Replay of data



- Exploration & Exploitation