

## Q-Learning

- model free learning approach.

- S A R S'

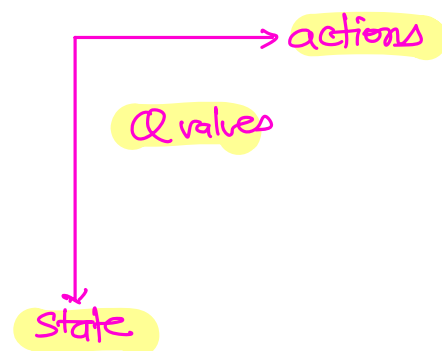
current state - S  
action - A  
reward - R  
new state - S'

- Goal of Q learning is to train Q Table

- There exists two HYPERPARAMETERS

0-1  $\gamma$  a) gamma (DISCOUNT FACTOR) ..... Applied on all Q learning methods  
0-1  $\alpha$  b) alpha (LEARNING RATE) ..... Applied on TD & SARSA

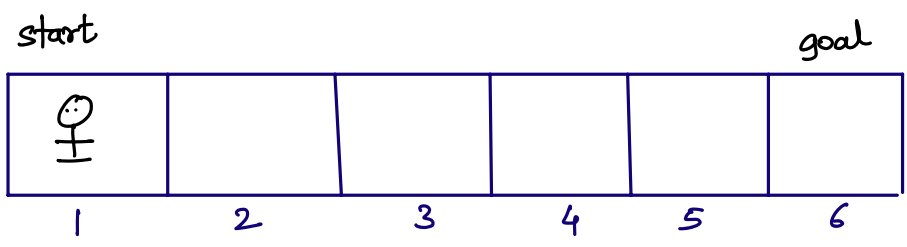
Q Table



Use-case

Given env

Goal: make agent reach the goal



Valid states → 1, 2, 3, 4, 5, 6  
Valid Action → Left, Right

Action \ state	Left	Right
1	0	0
2	0	0
3	0	0
4	0	0
5	0	0
6	0	0

Traditional Q learning\*

$$Q(s, A) = IR + (\text{gamma} * \max(s', a^*))$$

Annotations for Traditional Q learning:

- current state (points to s)
- current action (points to A)
- immediate reward (points to IR)
- discount factor (hyperparameters) (points to gamma)
- next state (points to s')
- all valid actions (points to a\*)

Q value of next state for all possible action and extract the max q value.

Modified Q learning\*

$$Q(s, A) = Q(s, A) + \alpha [R(s, A) + \gamma \max Q'(s', a^*) - Q(s, A)]$$

Annotations for Modified Q learning:

- current Q value from Q table (points to Q(s, A) on the left)
- learning rate (points to alpha)
- IR (points to R(s, A))
- discount factor (points to gamma)
- max q value of next state from all possible actions (points to max Q'(s', a\*))
- current Q value from Q table (points to Q(s, A) on the right)

## Algorithm to train the Q-Table

① initialize all Q-values in the Q-Table as 0.

② Observe the current state 's'.

③ Do forever

a. Select an action 'a' and execute it.

b. Receive an IR

c. Observe new state 's'.

d. Update the q-table entry for the given (S,a) using valid formula\*

until S is the terminal state | max episode count reaches

exploration or epsilon greedy.