

# Deep Q-Learning



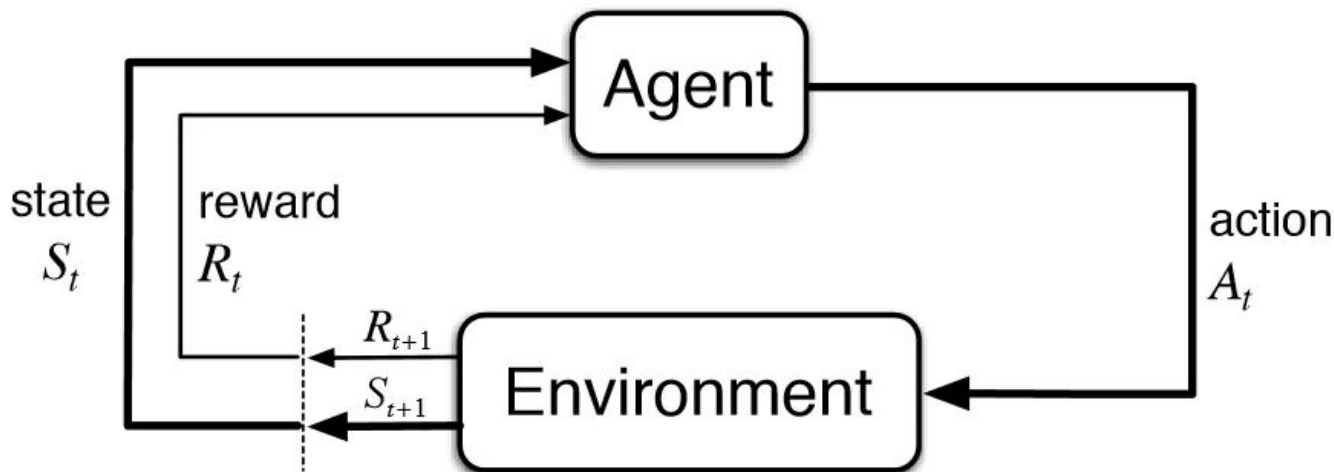
Presentation by Rahul Vishwakarma

# Content

- Reinforcement Learning
- Q-Learning
- Epsilon greedy policy
- Bellman's Equation
- Deep Q-Learning
- CartPole Game

# Reinforcement Learning

Training an agent to interact with an environment in order to maximize the cumulative reward over time.



# Q-Learning

The Q-learning algorithm uses a Q-table of State-Action Values. This Q-table has a row for each state and a column for each action. Each cell contains the estimated Q-value for the corresponding state-action pair.



	Left	Right	Up	Down
(1,1)	0	0	0	0
(1,2)	0	0	0	0
(1,3)	0	0	0	0
(2,1)	0	0	0	0
(2,2)	0	0	0	0
(2,3)	0	0	0	0
(3,1)	0	0	0	0
(3,2)	0	0	0	0
(3,3)	0	0	0	0

1

Initialise Q-Value (ie. State Action value) estimates with zero value, and pick the initial state.

	a1	a2	a3	a4
S1	0	0	0	0
S2	0	0	0	0
S3	0	0	0	0
S4	0	0	0	0
S5	0	0	0	0

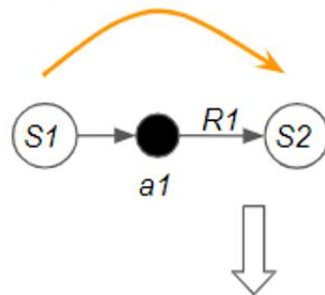


2

Agent picks an **action to execute** from the current state using  $\epsilon$ -greedy policy.

3

Agent obtains observation data from the environment ( $S1, a1, R1, S2$ )



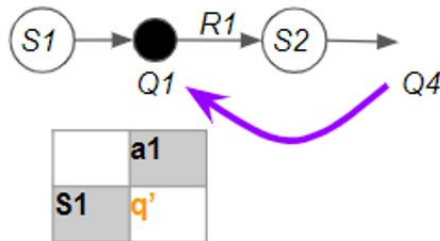
(S2)

Next state becomes the Current State

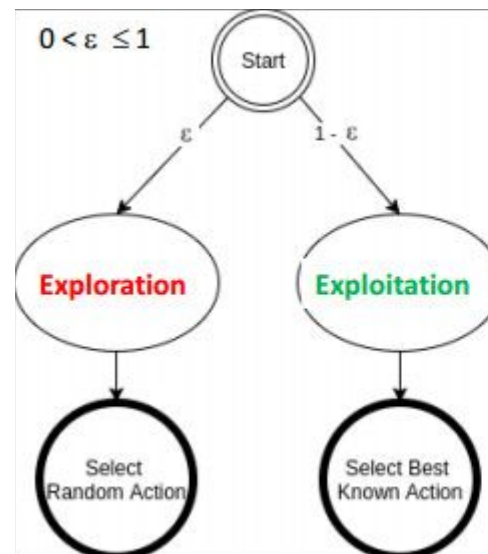
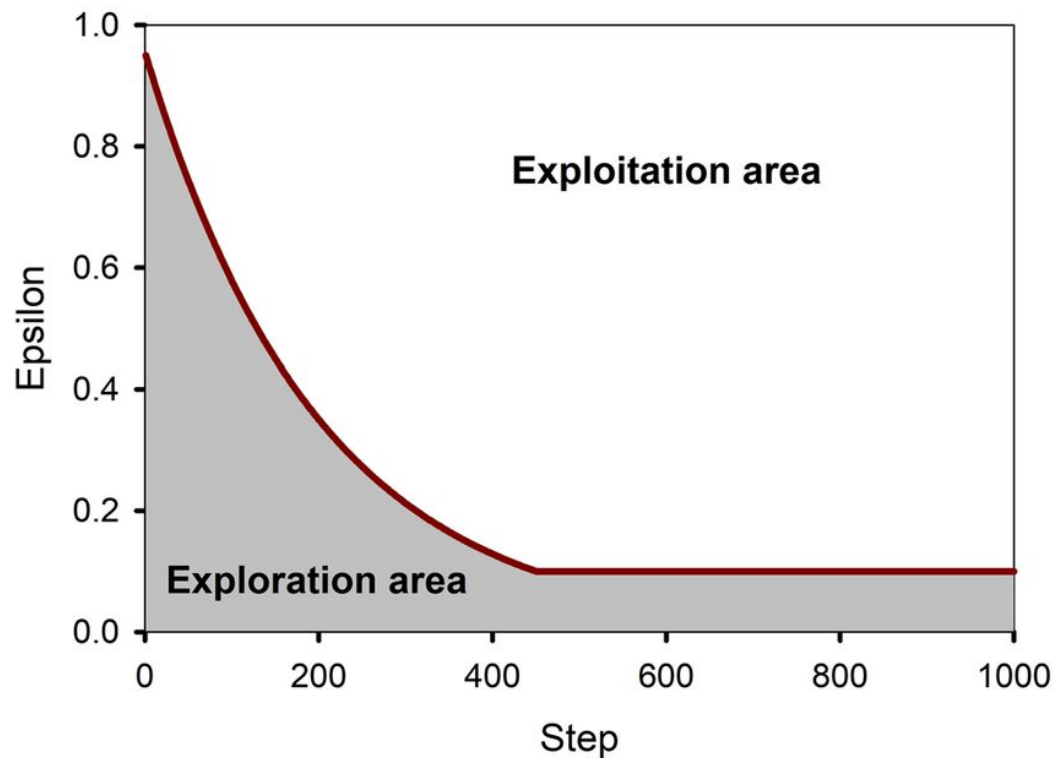
4

Update the **current Q-value** using the observed reward and the **target Q-value**. The target Q-value is the action with the max Q-value from the next state.

$$Q(S, A) = Q(S, A) + \alpha(R + \gamma \max_a Q(S', a) - Q(S, A))$$



# Epsilon Greedy Policy



# Bellman's Equation

$$\text{New } Q(s,a) = Q(s,a) + \alpha [R(s,a) + \gamma \max_{a'} Q'(s',a') - Q(s,a)]$$

Diagram illustrating the components of Bellman's Equation:

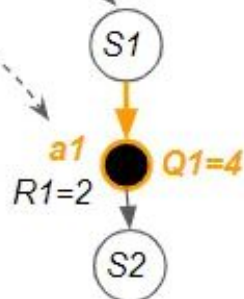
- New  $Q(s,a)$** : New Q value for the state and action (indicated by an arrow pointing to the first term on the left).
- $Q(s,a)$** : Current Q values (indicated by an arrow pointing to the second term).
- $\alpha$** : Learning Rate (indicated by an arrow pointing to the coefficient).
- $R(s,a)$** : Reward for taking an action in a state (indicated by an arrow pointing to the term inside the brackets).
- $\gamma$** : Discount Rate (indicated by an arrow pointing to the coefficient).
- $\max Q'(s',a')$** : Maximum expected future reward (indicated by an arrow pointing to the term inside the brackets).
- $Q(s,a)$** : Current Q values (indicated by an arrow pointing to the final subtraction term).

2

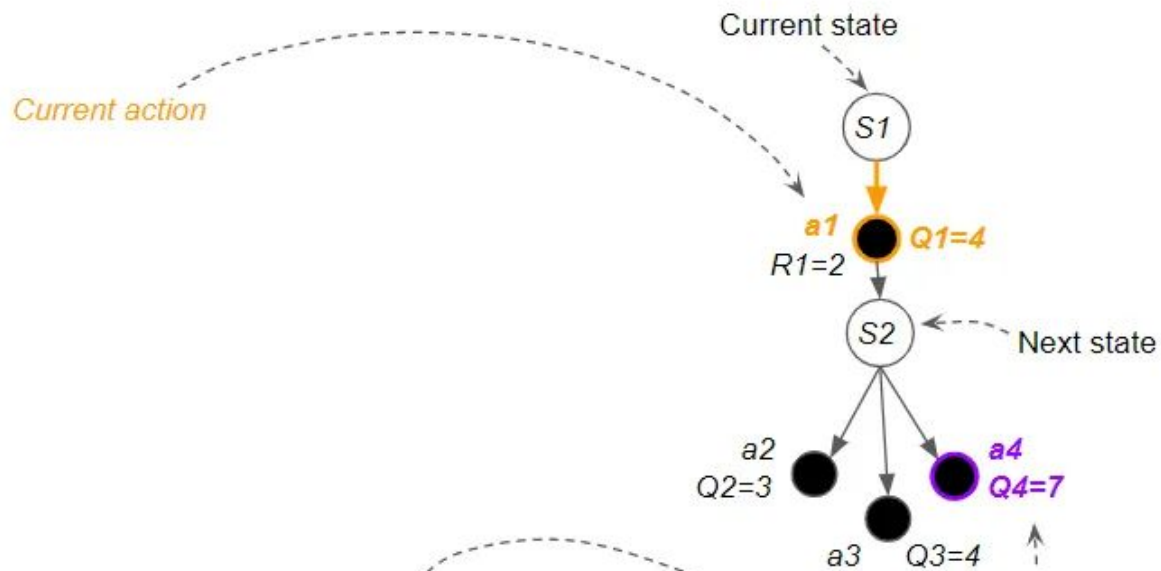
Current action picked using an  $\epsilon$ -greedy policy.

	a1	a2	a3	a4
S1	4	9	2	3
S2	0	3	4	7

Current state



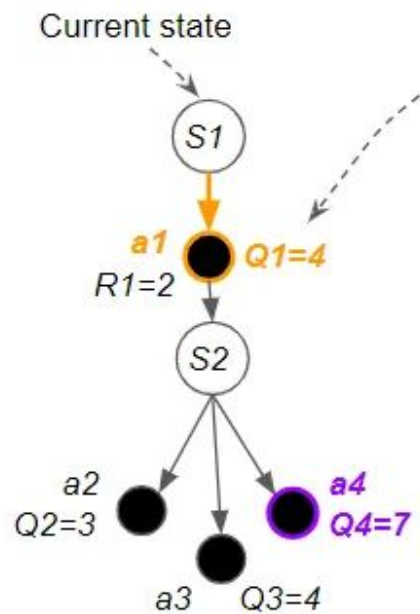




4a

The *target action* is the action from the next state with the highest Q-value

	a1	a2	a3	a4
S1	4	9	2	3
S2	0	3	4	7



**4b**

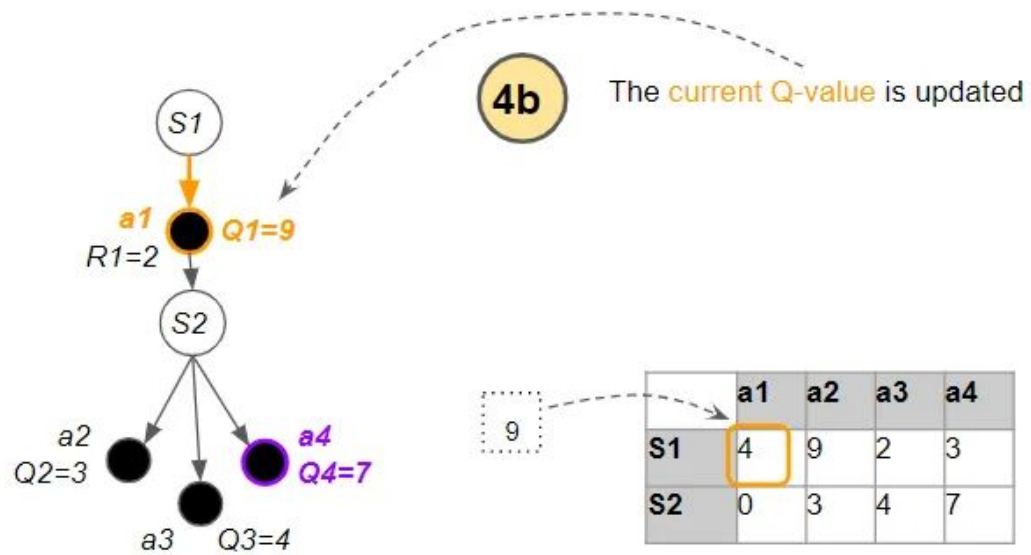
The **current Q-value** is updated

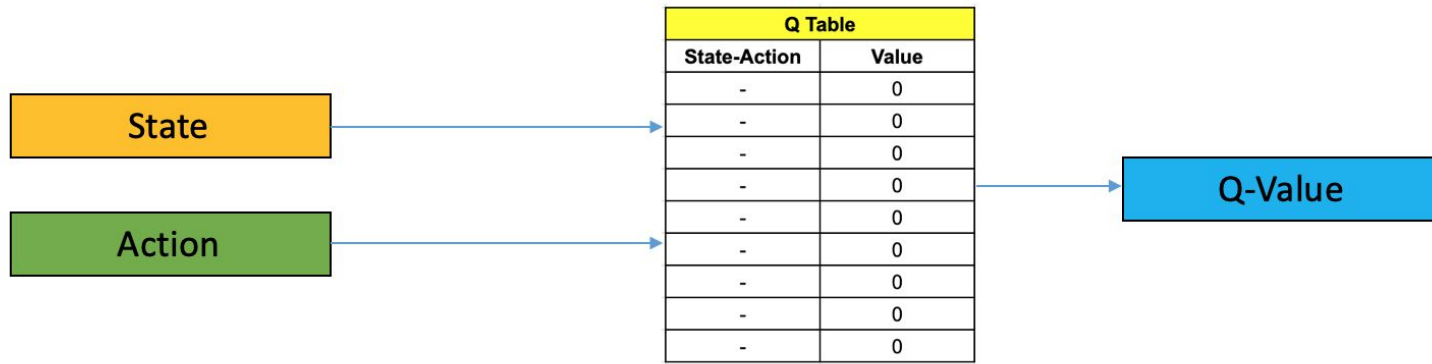
$$Q(S, A) = Q(S, A) + \alpha(R + \gamma \max Q(S', a) - Q(S, A))$$

$$Q1 = Q1 + \alpha(R1 + \gamma * Q4 - Q1)$$

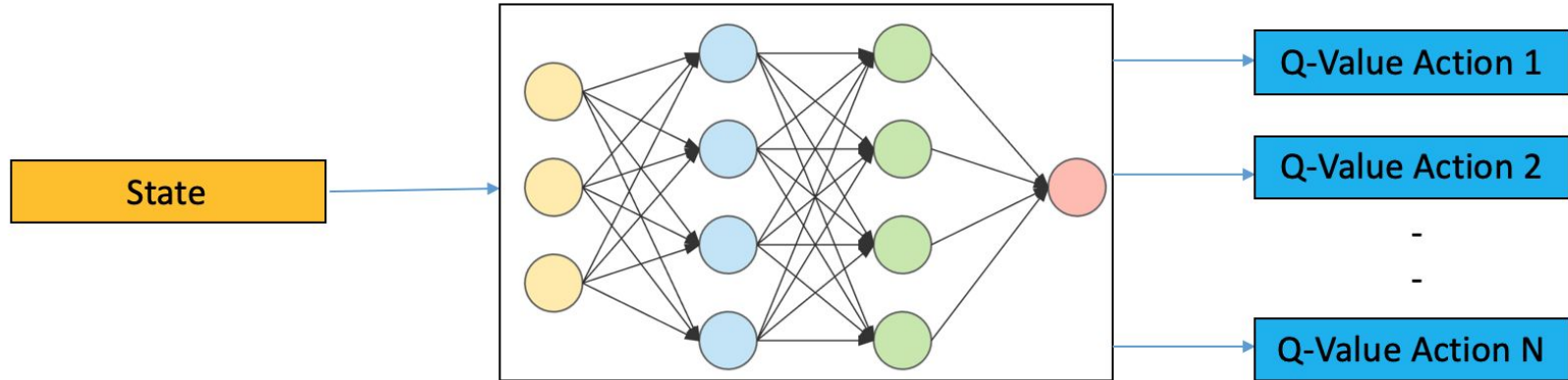
$$Q1 = 4 + (2 + 7 - 4) = 9$$

NB: taking  $\alpha = \gamma = 1$  for simplicity





## Q Learning



## Deep Q Learning

# Q Table vs Neural Network (DQN)

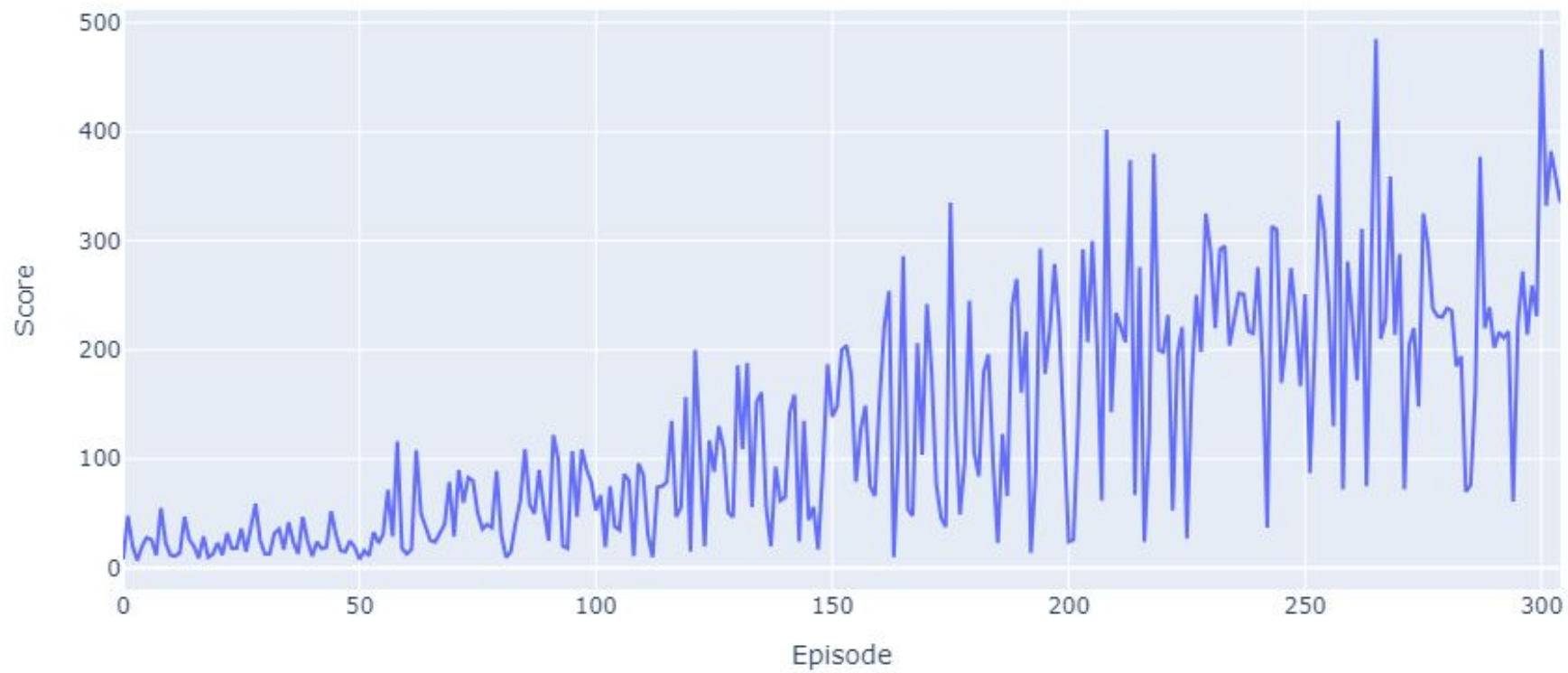
**Scalability:** Q-tables become impractical as the number of states and actions in the environment grows, making it infeasible to store and update all Q-values explicitly.

**Continuous action spaces:** In environments with continuous action spaces, it's challenging to create a Q-table that covers all possible actions finely.

**Memory requirements:** Neural networks typically require less memory than Q-tables when dealing with large state spaces because they store function approximations rather than explicit values for every state-action pair.

**Function approximation:** Q-tables can only represent linear relationships and are limited in their modeling capacity.

# **CartPole Game in Jupyter Notebook**



# References

- Towards Data Science: Reinforcement Learning Explained Visually (Part 4): Q Learning, step-by-step
- Analytics Vidhya: A Hands-On Introduction to Deep Q-Learning using OpenAI Gym in Python
- OpenAI Gymnasium Documentation:  
[https://gymnasium.farama.org/environments/classic\\_control/cart\\_pole/](https://gymnasium.farama.org/environments/classic_control/cart_pole/)



**Thank You**