



DIGITAL ASSIGNMENT: II

Q1. Learn an AI technique and understand how it works.

Q-Learning:

- Q-learning is typically used in Markov Decision Processes (MDPs), where an agent interacts with an environment in discrete time steps.
- The agent observes the current state (s) of the environment and takes an action (a) .
- After taking an action, the agent receives a reward (r) and transitions to a new state (s') .

- The goal of the agent is to learn a policy π , which maps states to actions, in order to maximize the cumulative reward over time.
- The value of a state-action pair $Q(s, a)$ represents the expected cumulative reward if the agent starts in state s , takes action a , and then follows the optimal policy thereafter.

- Initialize the Q-table: Start with an arbitrary initial Q-value for each state-action pair.

- Exploration vs. Exploitation: At each time step, the agent decides whether to explore new actions or exploit the current best action based on an exploration-exploitation trade-off strategy (e.g., epsilon-greedy).

- Update Q-values: When the agent takes an action and observes the reward and next state, it updates the Q-value of the current state-action pair using the Q-learning update rule:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

where:

- α is the learning rate, determining the extent to which new information overrides old information.
- γ is the discount factor, representing the importance of future rewards relative to immediate rewards.

4. Convergence:

- Q-learning is guaranteed to converge to the optimal Q-values under certain conditions, such as infinite exploration and appropriate decay schedules for the learning rate and exploration rate.

5. Policy Extraction:

- Once the Q-values converge, the agent can extract an optimal policy by selecting the action with the highest Q-value for each state.

6. Application:

- Q-learning can be applied to various problems, such as grid world navigation, robotic control, and game playing (e.g., Tic-Tac-Toe, Atari games).

Overall, Q-learning is a powerful technique for learning optimal policies in environments with discrete states and actions, enabling agents to make sequential decisions to maximize cumulative rewards over time.

Q2. Code it for an application in your programming language of preference

and upload your codes on your GITHUB repository.

Let's apply Q-learning to a simple grid world navigation problem. In this scenario, an agent needs to navigate through a grid world to reach a goal while avoiding obstacles. The agent receives a positive reward when it reaches the goal and a negative reward if it hits an obstacle. The objective is to learn an optimal policy that maximizes the cumulative reward.

Benefits:

Adaptability: Q-learning can adapt to changes in the environment or goals, making it suitable for dynamic scenarios.

Efficiency: Once trained, the agent can navigate the grid world efficiently, finding the optimal path to the goal.

Scalability: Q-learning can scale to larger grid worlds or more complex environments with minimal modifications.

```
main.py
1 import numpy as np
2
3 class QLearningAgent:
4     def __init__(self, num_states, num_actions, learning_rate=0.1, discount_factor=0.9, epsilon=0.1):
5         self.num_states = num_states
6         self.num_actions = num_actions
7         self.learning_rate = learning_rate
8         self.discount_factor = discount_factor
9         self.epsilon = epsilon
10        self.q_table = np.zeros((num_states, num_actions))
11
12    def choose_action(self, state):
13        if np.random.uniform(0, 1) < self.epsilon:
14            return np.random.choice(self.num_actions) # Exploration
15        else:
16            return np.argmax(self.q_table[state, :]) # Exploitation
17
18    def update_q_table(self, state, action, reward, next_state):
19        best_next_action = np.argmax(self.q_table[next_state, :])
20        td_target = reward + self.discount_factor * self.q_table[next_state, best_next_action]
21        td_error = td_target - self.q_table[state, action]
22        self.q_table[state, action] += self.learning_rate * td_error
23
24    def simulate_environment(agent, num_episodes, max_steps):
25        for episode in range(num_episodes):
26            state = np.random.randint(0, agent.num_states) # Start from a random state
27            for _ in range(max_steps):
28                action = agent.choose_action(state)
29                # Simulate taking action in the environment
30                next_state, reward = simulate_step(state, action)
31                agent.update_q_table(state, action, reward, next_state)
32                state = next_state
33
34    def simulate_step(state, action):
35        # Define transition dynamics and rewards for the grid world
36        # For simplicity, assume deterministic transitions and fixed rewards
37        if action == 0: # Move up
38            next_state = max(state - 3, 0)
39        elif action == 1: # Move down
40            next_state = min(state + 3, 8)
41        elif action == 2: # Move Left
42            next_state = max(state - 1, 0)
43        else: # Move right
44            next_state = min(state + 1, 8)
45
46        if next_state == 8: # Goal state
47            reward = 1
48        elif next_state in [3, 5, 7]: # Obstacle states
49            reward = -1
50        else:
51            reward = 0
```

```

46 ▾     if next_state == 8: # Goal state
47         reward = 1
48 ▾     elif next_state in [3, 5, 7]: # Obstacle states
49         reward = -1
50 ▾     else:
51         reward = 0
52
53     return next_state, reward
54
55 # Main
56 num_states = 9 # 3x3 grid world
57 num_actions = 4 # Up, down, left, right
58 agent = QLearningAgent(num_states, num_actions)
59 simulate_environment(agent, num_episodes=1000, max_steps=100)
60
61 # Extract Learned policy
62 optimal_policy = np.argmax(agent.q_table, axis=1)
63 print("Learned Optimal Policy:")
64 print(optimal_policy.reshape((3, 3)))

```

```

16         return np.argmax(self.q_table[state, :]) # Exploitation
17
18 ▾     def update_q_table(self, state, action, reward, next_state):
19         best_next_action = np.argmax(self.q_table[next_state, :])
20         td_target = reward + self.discount_factor * self.q_table[next_state, best_next_action]
21         td_error = td_target - self.q_table[state, action]
22         self.q_table[state, action] += self.learning_rate * td_error
23
24 ▾     def simulate_environment(agent, num_episodes, max_steps):
25         for episode in range(num_episodes):
26             state = np.random.randint(0, agent.num_states) # Start from a random state
27             for _ in range(max_steps):
28                 action = agent.choose_action(state)
29                 # Simulate taking action in the environment
30                 next_state, reward = simulate_step(state, action)
31                 agent.update_q_table(state, action, reward, next_state)
32                 state = next_state
33
34 ▾     def simulate_step(state, action):
35         # Define transition dynamics and rewards for the grid world
36         # For simplicity, assume deterministic transitions and fixed rewards
37 ▾         if action == 0: # Move up
38             next_state = max(state - 3, 0)
39 ▾         elif action == 1: # Move down
40             next_state = min(state + 3, 8)

```

Learned Optimal Policy:

```

[[1 0 0]
 [1 3 1]
 [1 1 1]]

```

...Program finished with exit code 0
Press ENTER to exit console. □

[GITHUB LINK](#) - I have hyperlinked the text on MS Word

Q3. Describe the working of your code and the result.

The provided code implements Q-learning, a reinforcement learning technique, for solving a simple grid world navigation problem. Here's a detailed description of how the code works and the expected results:

****Working of the Code**:**

1. ****Q-Learning Agent Initialization**:**

- The `QLearningAgent`` class is initialized with parameters such as the number of states (``num_states``), the number of actions (``num_actions``), learning rate (``learning_rate``), discount factor (``discount_factor``), and exploration rate (``epsilon``).
- The Q-table is initialized with zeros, representing the expected cumulative reward for each state-action pair.

2. ****Action Selection**:**

- The agent chooses actions based on an exploration-exploitation strategy.
- With probability ``epsilon``, the agent explores by selecting a random action.
- Otherwise, it exploits by selecting the action with the highest Q-value for the current state.

3. ****Q-Value Update**:**

- After taking an action and observing the reward and next state, the agent updates the Q-value of the current state-action pair using the Q-learning update rule.
- The update is based on the observed reward, the Q-value of the next state, and parameters such as the learning rate and discount factor.

4. ****Environment Simulation**:**

- The ``simulate_environment`` function simulates the agent's interaction with the environment over multiple episodes.
- In each episode, the agent starts from a random state and takes actions until a terminal state is reached or a maximum number of steps (``max_steps``) is reached.
- At each step, the agent chooses an action, receives a reward, and updates its Q-values accordingly.

5. ****Policy Extraction**:**

- After training, the learned optimal policy is extracted from the Q-table.
- The optimal policy specifies the best action to take in each state based on the learned Q-values.

****Result**:**

- After running the code, the learned optimal policy is printed, showing the best action to take in each state of the grid world.
- The optimal policy should indicate the path that maximizes the cumulative reward, navigating around obstacles to reach the goal state.
- By observing the learned optimal policy, we can evaluate the effectiveness of the Q-learning algorithm in solving the grid world navigation problem.

Overall, the code demonstrates the application of Q-learning in solving a simple navigation problem, highlighting the agent's ability to learn an optimal policy through interaction with the environment and updating Q-values based on observed rewards.

Q4. Develop an idea of how this can be implemented in hardware.

Implementing Q-learning in hardware presents some unique challenges and opportunities. Here's an idea of how it could be implemented:

****Hardware Architecture**:**

1. ****Custom Accelerator**:**

- Design a custom hardware accelerator optimized for Q-learning computations.
- The accelerator should be capable of performing matrix operations efficiently, as Q-learning involves updating Q-values stored in a large Q-table.
- Implement specialized hardware units for tasks such as matrix multiplication, element-wise operations, and memory access.

2. ****Parallelism**:**

- Exploit parallelism in hardware to speed up computations.
- Use parallel processing units to perform multiple Q-value updates simultaneously, leveraging the inherent parallel nature of Q-learning algorithms.

3. ****Memory Hierarchy**:**

- Design a memory hierarchy optimized for Q-learning workloads.
- Utilize on-chip memory for storing Q-values and intermediate computation results, reducing latency and energy consumption associated with off-chip memory access.
- Implement efficient memory access patterns to minimize data movement and maximize memory bandwidth utilization.

4. ****Scalability**:**

- Design the hardware accelerator to be scalable, allowing for parallel execution of multiple Q-learning agents or larger Q-tables.
- Support configurable parameters such as the number of states, actions, and Q-table size to accommodate different problem sizes and complexities.

5. ****Interface**:**

- Include interfaces for communication with external components, such as sensors and actuators, to interact with the environment.
- Provide interfaces for configuring the Q-learning algorithm parameters and monitoring performance metrics.

6. ****Power Efficiency**:**

- Optimize the hardware design for power efficiency to enable deployment in resource-constrained environments.
- Utilize techniques such as voltage scaling, clock gating, and power gating to minimize energy consumption during idle periods and low-utilization phases.

****Integration with Embedded Systems**:**

1. ****Embedded Platform**:**

- Integrate the hardware accelerator with embedded platforms such as FPGAs or ASICs to enable real-time Q-learning inference in embedded systems.
- Provide software APIs and drivers for seamless integration with higher-level application software running on the embedded platform.

2. ****Deployment in Edge Devices**:**

- Deploy the hardware-accelerated Q-learning solution in edge devices such as autonomous robots, drones, or smart IoT devices.
- Enable on-device learning and decision-making capabilities, reducing the need for frequent communication with centralized servers and enhancing privacy and security.

3. ****Application Examples**:**

- Use cases include autonomous navigation for robots or drones, adaptive control systems for industrial automation, and intelligent IoT devices for smart home applications.
- The hardware-accelerated Q-learning solution can enable efficient and autonomous operation in real-world environments with limited computational resources and power constraints.

In summary, implementing Q-learning in hardware requires designing a custom accelerator with optimized computational units, memory hierarchy, and interfaces for efficient and scalable deployment in embedded systems. Integration with embedded platforms enables on-device learning and decision-making capabilities for a wide range of applications in robotics, automation, and IoT.