# Title:  Reinforcement Learning with Q-Learning: Solving a Grid World Problem

## Abstract

This project explores Q-learning in a grid world with obstacles and a goal. The agent learns to navigate from a starting point to the goal while avoiding obstacles. Key aspects include defining the environment, implementing Q-learning, and employing visualizations like the step-to-go curve and intermediate Q-tables to showcase the agent's learning progress. This work offers valuable insights into Q-learning's performance in complex grid world scenarios.

## Introduction:

Reinforcement Learning (RL) stands as a subfield within machine learning that centralizes decision-making within sequential contexts. Within RL, agents dynamically interact with their environments, refining their decision-making through experimentation. At its core, RL grapples with challenges posed by scenarios like grid world navigation—a quintessential problem. Here, agents must master the art of transitioning from a predefined starting point to a target destination while skilfully obstacles.

This report embarks on an exploration of Q-learning, a widely recognized RL algorithm, in the context of resolving grid world navigation dilemmas. Our endeavour involves not only defining the intricacies of this challenge but also delineating our objectives and goals. Our paramount goal revolves around the successful implementation and scrutiny of Q-learning's efficacy when confronted with grid worlds containing obstacles. Through this investigation, we aim to cast a revealing light on its capacity to acquire optimal strategies within intricate environments.

## Problem Description:

The grid world problem under consideration involves an agent's navigation through a gridded environment, aiming to reach a predefined goal while encountering obstacles.

Description--

Grid Layout: The environment is represented as a two-dimensional grid, with each cell signifying a unique state. For instance, a 5x5 grid world consists of 25 individual states.

Obstacles: Certain cells within the grid are designated as obstacles, indicated by a distinct symbol (e.g., '#'). These obstacles act as impassable barriers, necessitating the agent to navigate around them.

Rewards: Each state in the grid world is associated with a specific reward, reflecting its desirability. For instance, successfully reaching the goal state results in a high positive reward, while colliding with an obstacle incurs a negative reward.

Agent's Objective: The primary goal for the agent is to traverse from an initial starting state to a predefined goal state. Achieving this objective yields a positive reward, marking the accomplishment of an episode.

States: States encompass individual cells or locations within the grid, providing information about the agent's current position and its surroundings.

Actions: The agent can execute particular actions to transition between states, including moving up, down, left, or right within the grid. The choice of action influences the subsequent state and the cumulative reward.

Rewards: Rewards are numerical values linked to each state-action pair, signifying the immediate advantage or cost of taking a specific action in a given state. For instance, progressing toward the goal state results in a favourable reward, while encountering an obstacle leads to an adverse reward.

## Methodology:

In tackling the grid world navigation problem, our approach centres on employing the Q-learning algorithm.

Below, are the key elements--

Q-Learning Algorithm: Q-learning, a reinforcement learning technique, determines the optimal policy for an agent within a Markov decision process. In this context, it guides the agent to maximize cumulative rewards by selecting the best actions in various states.

Application of Q-learning: We iteratively applied the Q-learning algorithm across multiple episodes. Within each episode, the agent initiates from an initial state, executes actions, and updates its Q-values based on observed rewards and state transitions. The objective remains the acquisition of an optimal Q-value function that steers the agent from start to goal while circumventing obstacles.

Hyperparameters:

Learning Rate ($\alpha$): Governs the pace at which the agent adjusts Q-values, ensuring a stable learning process.

Discount Factor ($\gamma$): Influences the agent's regard for future rewards. A high value emphasizes long-term gains, while a low one prioritizes immediate rewards.

Exploration Rate ($\varepsilon$): Balances exploration and exploitation. Higher values favor exploration, while lower ones lean toward exploiting learned Q-values.

Number of Episodes: Specifies the frequency of the agent's learning iterations and Q-value updates.

Q-Table: The Q-table, a pivotal component, holds Q-values for state-action pairs.

In our case, it adopts a 5x5x4 structure:

The first dimension (5x5) maps to grid rows and columns, signifying states.

The second dimension (4) encapsulates the four possible actions: up, down, left, and right.

For instance, q_table [0, 0] houses Q-values for the initial state (0, 0), and q_table [0, 0, 0] represents the Q-value for taking an "up" action from that state.

## Implementation:

Here is an overview of the implementation:

Q-Learning Code Snippets: We've divided the Q-learning algorithm into distinct steps, providing code segments with explanations. These steps encompass state transitions, Q-value updates, and checks for goal achievement.

Q-Table Initialization: To establish the groundwork for learning optimal policies, the Q-table is initialized as a multi-dimensional array. It faithfully represents the state-action space of the grid world.

Optimal Path Visualization Code: Included is code for visualizing the agent's optimal path from the starting point to the goal. This visual representation illuminates the agent's acquired navigation prowess within the grid world.

Here are some important code snippets—

1. Q-Value Update:

```python
# Update the Q-value using the Q-learning formula
q_table[state[0], state[1], action] += learning_rate * (
    reward + discount_factor * np.max(q_table[next_state[0], next_state[1]]) - q_table[state[0], state[1], action])
```

2. State Transition:

```python
# Calculate the next state
next_state = (state[0] + actions[action][0], state[1] + actions[action][1])
```

3. Q-Table Initialization:

```python
# Initialize the Q-table with zeros

num_rows, num_cols = grid_world.shape
num_actions = len(actions)
q_table = np.zeros((num_rows, num_cols, num_actions))
```

4. Update Q-value using Formula:

```python
# Update the Q-value using the Q-learning formula
q_table[state[0], state[1], action] += learning_rate * (
    reward + discount_factor * np.max(q_table[next_state[0], next_state[1]]) - q_table[state[0], state[1], action])
```

**Experimental Results:**

This section presents the findings of our experiments and encompasses the following key components--

Learned Q-Table: We showcase the Q-table, which encapsulates the agent's knowledge. This table reveals the Q-values associated with various state-action pairs, illustrating the agent's policy.

Optimal Path Visualization: Our representation of the agent's optimal path within the grid world explain the path clearly. This visualization effectively demonstrates the agent's successful navigation from the initial position to the goal while circumventing obstacles.

Performance Metrics: We provide essential performance metrics such as the number of steps taken, total rewards accumulated, and the success rate. These metrics serve to quantify the agent's proficiency in achieving its objectives.

Learning Progress Charts: We incorporate graphical representations that chart the agent's learning progress throughout the training episodes. These visuals offer insights into the agent's evolving performance and facilitate a comprehensive analysis of its learning trajectory.
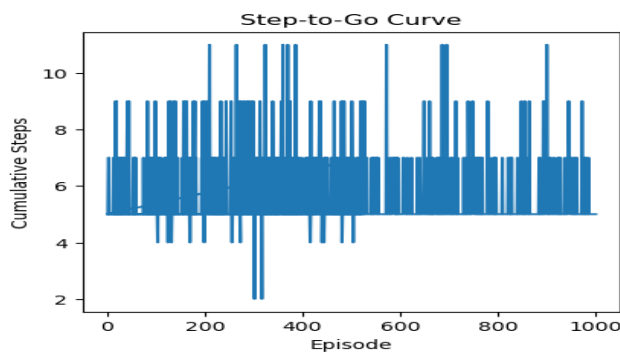
**Analysis**:

Within this section, we delve into the following critical aspects:

Q-Table Interpretation: We decipher the Q-table, elucidating the significance of Q-values associated with various state-action pairs. This interpretation unveils the rationale behind the agent's decisions and policies.
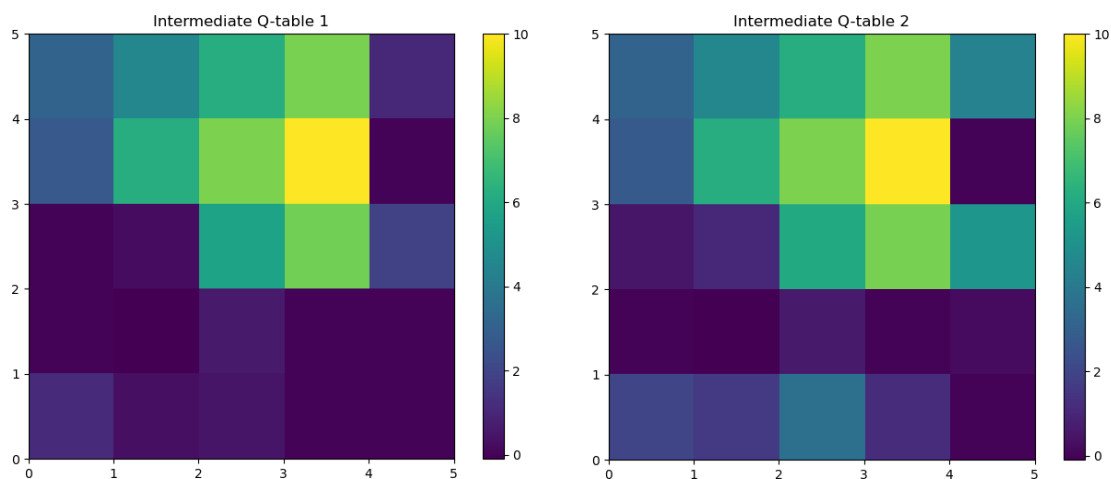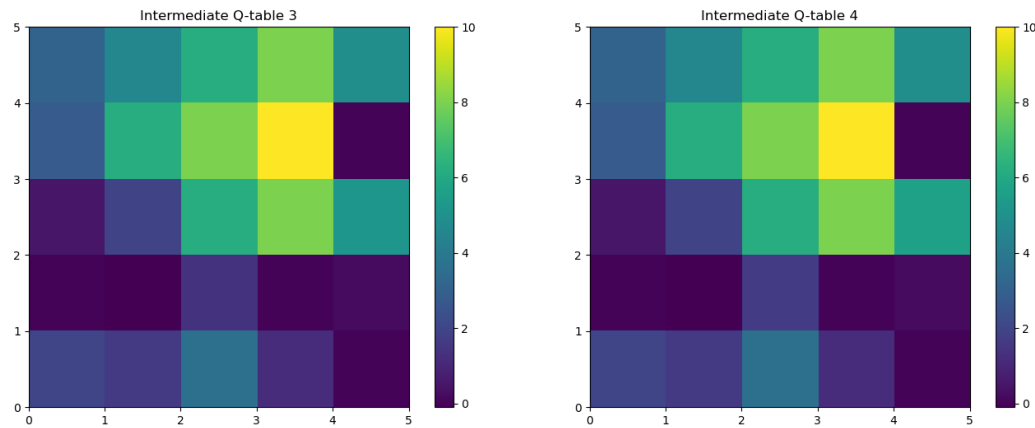
Code snippet:

```python
# Visualize intermediate Q-tables (if desired)
for i, q_table in enumerate(intermediate_q_tables):
    plt.figure(figsize=(8, 6))
    plt.imshow(np.max(q_table, axis=2), cmap='viridis', origin='upper', extent=(0, num_cols, 0, num_rows))
    plt.colorbar()
    plt.title(f'Intermediate Q-table {i + 1}')
    plt.show()
```

Performance Assessment: Insights and observed patterns, successes, and challenges encountered during the learning process.



Below we will see the Intermediate Q tables—

Intermediate Q-table 3        Intermediate Q-table 4

Hyperparameter Impact: Assessed the influence of hyperparameters on the learning journey. This evaluation provides valuable insights into how factors like learning rate, discount factor, and exploration rate affect the agent's behaviour.

```
# Hyperparameters

learning_rate = 0.1
discount_factor = 0.9
epsilon = 0.1
num_episodes = 1000
```

Comparative Analysis: Undertake a comparative exploration, juxtaposing the performance in different grid layouts or environments. This comparison sheds light on the algorithm's adaptability and effectiveness across diverse scenarios.

## Conclusion:

In summary, this study has produced several noteworthy findings while investigating the application of the Q-learning algorithm in the context of grid world problem-solving.

### Key Findings:

Effective Navigation: The Q-learning algorithm has demonstrated its proficiency in enabling an agent to successfully navigate a grid world. It has acquired optimal policies for guiding the agent from the initial position to the goal while skilfully avoiding obstacles.

Progressive Learning: Over the course of multiple episodes, the agent exhibits clear learning progress. It continually refines its policy, resulting in increasingly efficient navigation and a reduced number of steps needed to reach the goal.

Hyperparameter Influence: The selection of hyperparameters, encompassing learning rate, discount factor, and exploration rate, exerts significant impact on the learning process. Precisely tuning these parameters proves essential for achieving optimal outcomes.

Effectiveness of Q-Learning:

The Q-learning algorithm stands as a sturdy and adaptable solution for grid world navigation. Its capacity to learn and iteratively enhance its policy renders it a valuable asset for training agents in intricate, evolving environments. The agent's learned behaviour, encapsulated within the Q-table, attests to the algorithm's prowess in discovering optimal strategies.

**Limitations and Future Prospects:**

While Q-learning has demonstrated its capabilities, it is essential to acknowledge its constraints and avenues for prospective refinement:

Exploration-Exploitation Equilibrium: Striking the right equilibrium between exploration and exploitation, often regulated by the exploration rate (epsilon), remains a challenge. Research into dynamic exploration strategies holds the potential for more effective learning.

Complex Environments: The algorithm's performance in exceedingly complex or dynamic grid layouts could benefit from further enhancement through advanced techniques, such as deep reinforcement learning.

In conclusion, our exploration of Q-learning in the grid world problem underscores its potential as a robust and versatile reinforcement learning approach. While the results are promising, ongoing research and innovation within the realm of reinforcement learning are poised to address its limitations and elevate its efficacy in addressing even more intricate scenarios.