

Design Defense

1. Examine the Distinctions Between Human and Machine Methods of Problem Solving

Human Approach to Solving the Maze:

- **Observation:** A human would start by observing the maze to understand its layout, identifying obstacles and possible paths.
- **Planning:** The human would plan a route mentally, considering the shortest or safest paths to the treasure while avoiding obstacles.
- **Trial and Error:** They might use a trial-and-error approach, trying different paths and backtracking if they encounter dead ends.
- **Learning and Memory:** If attempting the maze multiple times, the human would learn from previous attempts, remembering successful paths and avoiding past mistakes.
- **Goal-Oriented Actions:** The human remains focused on reaching the treasure, adjusting their strategy based on new observations.

Intelligent Agent's Approach:

- **Initialization:** The agent initializes its Q-table or neural network weights and begins without any past knowledge of the maze.

- **Exploration:** Initially, the agent explores the environment randomly to gather information about the states and rewards associated with different actions. For example, the agent uses **epsilon-greedy** policy to explore the maze:

```
def choose_action(state, epsilon):  
    if np.random.rand() < epsilon:  
        return np.random.randint(0, 4)  
    else:  
        q_values = model.predict(state)  
        return np.argmax(q_values[0])
```

- **Q-Learning Algorithm:** Based on the incentives it receives, the agent updates its knowledge using the Q-learning algorithm. It uses the Bellman equation to update the Q-values continuously:

```

def q_learning_train(maze, agent, episodes, gamma, epsilon):
    for episode in range(episodes):
        state = maze.reset()
        done = False
        while not done:
            action = choose_action(state, epsilon)
            next_state, reward, done = maze.step(action)
            target = reward + gamma *
np.max(agent.predict(next_state))
            q_values = agent.predict(state)
            q_values[0][action] = target
            agent.fit(state, q_values, epochs=1, verbose=0)
            state = next_state

```

- **Exploitation and Exploration:** The agent strikes a balance between using well-traveled routes that provide large rewards (exploitation) and discovering new paths (exploration).
- **Training and Convergence:** The agent's policy eventually converges to an ideal course of action through a series of events, decreasing the distance to the treasure while dodging hazards.

Similarities and Differences:

- **Similarities:**
 - Both humans and the agent use trial and error to find the optimal path.
 - Both approaches involve learning from past experiences to improve future decisions.
 - Both aim to reach the goal (treasure) while avoiding obstacles.
- **Differences:**
 - Humans rely on cognitive abilities and memory, whereas the agent relies on algorithms and computational processes.

- Humans can plan and visualize potential routes beforehand, while the agent iteratively improves its strategy through reinforcement learning.
- The agent requires a defined reward structure to learn effectively, whereas humans can adapt to various goals and rules intuitively.

2. Assess the Purpose of the Intelligent Agent in Pathfinding

Exploitation vs. Exploration:

- **Exploitation:** Use the previously acquired knowledge to execute the most well-known action. Based on prior results, this focuses on optimizing the immediate payoff.
- **Exploration:** Trying out new actions to discover their potential rewards. This is crucial for learning and improving the policy.
- **Ideal Proportion:** The ideal proportion of exploitation and exploration is often achieved through a decaying epsilon strategy in epsilon-greedy algorithms. Initially, a higher exploration rate (e.g., epsilon = 1) is used to gather information, which gradually decreases to a lower value (e.g., epsilon = 0.1) to favor exploitation as the agent becomes more confident in its knowledge. This equilibrium guarantees that the agent sufficiently investigates to gain an understanding of the surroundings, all the while utilizing existing knowledge to maximize efficiency:

```
epsilon = max(epsilon_min, epsilon_decay * epsilon)
```

Reinforcement Learning for Pathfinding:

- **Policy Learning:** By linking states with the most effective actions, reinforcement learning assists the agent in learning the optimum course of action to traverse the maze.
- **Reward Signals:** When the agent overcomes hurdles and reaches the treasure, they are rewarded. The learning process is guided by these reward signals:

```
reward = -1

if state == treasure_state:
    reward = 100
```

- **Q-Learning Updates:** Based on the Bellman equation, which considers both the reward that was received and the maximum reward that is anticipated in the future, the agent adjusts its Q-values. This aids in the agent's steady learning of the worth of various states and activities:

```
target = reward + gamma * np.max(agent.predict(next_state))
```

3. Evaluate the Use of Algorithms to Solve Complex Problems

Implementation of Deep Q-Learning:

- **Neural Network Structure:** The Q-value function is approximated using a neural network. An input layer, hidden layers, and an output layer that corresponds to the action space make up the network:

```
model = Sequential()  
model.add(Dense(24, input_dim=state_size, activation='relu'))  
model.add(Dense(24, activation='relu'))  
model.add(Dense(action_size, activation='linear'))  
model.compile(loss='mse', optimizer=Adam(lr=learning_rate))
```

- **State Representation:** The neural network is fed with the current state of the maze, which includes the agent's position as well as the locations of the obstacles and treasure.
- **Action Selection:** For choosing an action and weighing exploration versus exploitation, an epsilon-greedy policy is employed:

```
if np.random.rand() < epsilon:  
    action = np.random.randint(0, action_size)  
else:  
    action = np.argmax(model.predict(state)[0])
```

- **Experience Replay:** In a replay buffer, the agent records its experiences (state, action, reward, and next state). The Q-network is updated during training using random samples from this buffer, which increases learning efficiency and stability:

```
replay_buffer.append((state, action, reward, next_state, done))  
if len(replay_buffer) > batch_size:  
    minibatch = random.sample(replay_buffer, batch_size)
```

```

for state, action, reward, next_state, done in minibatch:
    target = reward
    if not done:
        target = reward + gamma *
np.amax(model.predict(next_state)[0])
    target_f = model.predict(state)
    target_f[0][action] = target
    model.fit(state, target_f, epochs=1, verbose=0)

```

- **Target Network:** Divergence in the learning process is minimized by computing the target Q-values using a different target network.
- **Training Loop:** The agent iterates through episodes, collecting experiences, updating the Q-network, and gradually improving its policy to find the optimal path to the treasure.

References

- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction* (2nd ed.). MIT Press.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529-533.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.