

Step 1: Grid Setup

1. **Create the Grid:** Generate a 100x100 grid where each cell can either be a free space, an obstacle, a start, or a goal point.
2. **Place Start and Goal Points:** Randomly assign two cells as start and goal points.
3. **Place Obstacles:** Randomly place obstacles in some cells, ensuring there's at least one path from the start to the goal.

```
import numpy as np
```

```
grid_size = 100
```

```
grid = np.zeros((grid_size, grid_size))
```

```
# Define start and goal points
```

```
start = (np.random.randint(0, grid_size), np.random.randint(0, grid_size))
```

```
goal = (np.random.randint(0, grid_size), np.random.randint(0, grid_size))
```

```
grid[start] = 1 # Start point
```

```
grid[goal] = 2 # Goal point
```

```
# Define obstacles
```

```
obstacle_count = int(0.2 * grid_size * grid_size) # 20% of cells as obstacles
```

```
for _ in range(obstacle_count):
```

```
    x, y = np.random.randint(0, grid_size), np.random.randint(0, grid_size)
```

```
    if (x, y) != start and (x, y) != goal:
```

```
        grid[x, y] = -1 # Obstacle
```

Step 2: Define the Markov Decision Process (MDP)

The MDP can be formulated with:

- **States:** Each cell in the grid.
- **Actions:** Moving up, down, left, or right.
- **Rewards:**

- +1 for reaching the goal.
- -1 for stepping into an obstacle or moving out of bounds.
- -0.1 for each regular step to encourage the agent to reach the goal faster.
- **Transition Probabilities:** With a deterministic environment, assume the agent successfully moves in the intended direction unless blocked.

Step 3: Implement an RL Agent

The agent's task is to learn an optimal policy that maximizes the expected cumulative reward. We can use Dynamic Programming (DP), Q-learning, and Deep Q-learning (DQN) for policy optimization.

Dynamic Programming Approach

1. Value Iteration:

- Initialize a value function for each state.
- Iteratively update the value function using Bellman equations until convergence.

2. Policy Extraction:

- After value iteration, derive a policy by selecting actions with the highest value.

import numpy as np

Initialize value function

values = np.zeros((grid_size, grid_size))

Define discount factor and threshold

gamma = 0.9

threshold = 0.01

Value iteration loop

while True:

 delta = 0

 for i in range(grid_size):

 for j in range(grid_size):

 if grid[i, j] == -1 or (i, j) == goal:

```
    continue
```

```
    old_value = values[i, j]
```

```
    new_value = max([reward + gamma * values[next_state]
```

```
                    for action, (next_state, reward) in transitions[(i, j)]])
```

```
    values[i, j] = new_value
```

```
    delta = max(delta, abs(old_value - new_value))
```

```
if delta < threshold:
```

```
    break
```

Q-learning Approach

Implement a Q-learning agent to explore and exploit the environment without needing a model of the environment's dynamics. This approach is more flexible than DP.

```
import random
```

```
# Initialize Q-table
```

```
Q = np.zeros((grid_size, grid_size, 4)) # 4 actions
```

```
# Define learning parameters
```

```
alpha = 0.1 # learning rate
```

```
epsilon = 0.1 # exploration rate
```

```
episodes = 1000
```

```
for episode in range(episodes):
```

```
    state = start
```

```
    while state != goal:
```

```
        if random.uniform(0, 1) < epsilon:
```

```
            action = random.choice(actions) # explore
```

```

else:

    action = np.argmax(Q[state]) # exploit

    next_state, reward = step(state, action)

    Q[state][action] = Q[state][action] + alpha * (reward + gamma * max(Q[next_state]) -
    Q[state][action])

    state = next_state

```

Deep Q-Network (DQN)

To handle large or complex environments, implement a DQN. This is a deep learning approach where a neural network approximates the Q-values for each state-action pair.

Step 4: Benchmark the Approaches

1. **Metrics:** Evaluate each method based on cumulative rewards, average path length to the goal, and convergence speed.
2. **Implementation and Testing:**
 - Use the Dynamic Programming approach as a baseline.
 - Compare with Q-learning and DQN on the same grid setup.

Sample Output:

Algorithm	Average Reward	Average Path Length	Convergence Speed
Dynamic Programming	0.85	50	Fast
Q-learning	0.80	55	Moderate
Deep Q-learning (DQN)	0.88	52	Slower