# Experiment 10: Solving a Markov Decision Process (MDP) in GridWorld

Anushka, 23/CS/071

November 2025

## 1 Introduction

In this experiment, we solve a stochastic GridWorld Markov Decision Process (MDP) using the Value Iteration algorithm. The GridWorld is a 3x4 grid where the agent must navigate to reach a goal state while avoiding a pit. Actions are stochastic: the agent may slip to the left or right of the intended direction. We introduce a living penalty to encourage shorter paths and analyze how this penalty affects the agent's optimal policy. The experiment aims to compute the optimal value function and policy, visualize them, and observe the effect of different living penalties on agent behavior.

## 2 Task 1: Defining the GridWorld MDP

### 2.1 States, Actions, and Rewards

- **States (S):** All cells except the wall at (1,1). Terminal states: Goal (0,3), Pit (1,3)

- **Actions (A):** ['up', 'down', 'left', 'right']

- **Rewards (R):**

    - Goal: +1
    - Pit: -1
    - All other states: living penalty (default -0.04)

- **Discount Factor:** $\gamma = 0.99$

### 2.2 Python Implementation — GridWorld Setup

Listing 1: GridWorld setup

```
import numpy as np

ROWS, COLS = 3, 4
WALL = (1, 1)
GOAL = (0, 3)
PIT  = (1, 3)
TERMINALS = {GOAL, PIT}
```

```
8
9  ACTIONS = ['up', 'down', 'left', 'right']
10 ACT_TO_DELTA = {
11     'up': (-1, 0),
12     'down': (1, 0),
13     'left': (0, -1),
14     'right': (0, 1),
15 }
16
17 SLIP_LEFT = {'up':'left','down':'right','left':'down','right':'up'}
18 SLIP_RIGHT= {'up':'right','down':'left','left':'up','right':'down'}
19
20 def in_bounds(s):
21     r, c = s
22     return 0 <= r < ROWS and 0 <= c < COLS
23
24 def step(s, a):
25     if s in TERMINALS:
26         return s
27     dr, dc = ACT_TO_DELTA[a]
28     ns = (s[0] + dr, s[1] + dc)
29     if not in_bounds(ns) or ns == WALL:
30         return s
31     return ns
32
33 def get_all_states():
34     return [(r,c) for r in range(ROWS) for c in range(COLS) if (r,c) !=
           WALL]
35
36 def make_rewards(living_penalty=-0.04):
37     R = {}
38     for s in get_all_states():
39         if s == GOAL:
40             R[s] = 1.0
41         elif s == PIT:
42             R[s] = -1.0
43         else:
44             R[s] = living_penalty
45     return R
46
47 def get_next_states(s, a):
48     if s in TERMINALS:
49         return [(1.0, s)]
50     intended = step(s, a)
51     left = step(s, SLIP_LEFT[a])
52     right = step(s, SLIP_RIGHT[a])
53     return [(0.8, intended), (0.1, left), (0.1, right)]
54
55 get_all_states()
```

# 3  Task 2: Value Iteration Algorithm

Listing 2: Value Iteration Algorithm

```
1 def value_iteration(gamma=0.99, theta=1e-4, living_penalty=-0.04):
2     R = make_rewards(living_penalty)
```

```
3    states = get_all_states()
4    V = {s: 0.0 for s in states}
5    V[GOAL] = 0.0
6    V[PIT]  = 0.0
7    def q_value(s, a):
8        return sum(p * (R[ns] + gamma * V[ns]) for p, ns in
             get_next_states(s, a))
9    while True:
10       delta = 0.0
11       for s in states:
12           if s in TERMINALS:
13               continue
14           v = V[s]
15           V[s] = max(q_value(s, a) for a in ACTIONS)
16           delta = max(delta, abs(v - V[s]))
17       if delta < theta:
18           break
19   return V
```

# 4 Task 3: Policy Extraction

Listing 3: Extract Optimal Policy

```
1  ARROWS = {'up':'^', 'down':'v', 'left':'<', 'right':'>'}
2
3  def extract_policy(V, gamma=0.99, living_penalty=-0.04):
4      R = make_rewards(living_penalty)
5      policy = {}
6      def q_value(s, a):
7          return sum(p * (R[ns] + gamma * V[ns]) for p, ns in
             get_next_states(s, a))
8      for s in get_all_states():
9          if s in TERMINALS:
10             policy[s] = None
11         else:
12             qs = [(a, q_value(s, a)) for a in ACTIONS]
13             policy[s] = max(qs, key=lambda x: x[1])[0]
14     return policy
```

# 5 Task 4: Visualization

Listing 4: Visualize Value Function and Policy

```
1
2  import matplotlib.pyplot as plt
3  import numpy as np
4
5  def values_to_grid(V):
6      grid = np.full((ROWS, COLS), np.nan)
7      for r in range(ROWS):
8          for c in range(COLS):
9              if (r,c) == WALL:
10                 continue
11             grid[r,c] = V.get((r,c), np.nan)
```

```
12      return grid
13
14  def plot_value_heatmap(V, title):
15      grid = values_to_grid(V)
16      plt.figure(figsize=(5,3))
17      plt.imshow(grid, interpolation='nearest')
18      plt.title(title)
19      plt.colorbar()
20      for (i, j), val in np.ndenumerate(grid):
21          if not np.isnan(val):
22              plt.text(j, i, f"{val:.2f}", ha='center', va='center')
23          else:
24              plt.text(j, i, '   ', ha='center', va='center')
25      plt.xticks(range(COLS))
26      plt.yticks(range(ROWS))
27      plt.show()
28
29  def policy_to_grid(policy):
30      grid = np.full((ROWS, COLS), '', dtype=object)
31      for r in range(ROWS):
32          for c in range(COLS):
33              s = (r,c)
34              if s == WALL: grid[r,c] = '   '
35              elif s == GOAL: grid[r,c] = 'G'
36              elif s == PIT:  grid[r,c] = 'P'
37              else:
38                  a = policy.get(s, None)
39                  grid[r,c] = ARROWS[a] if a else ''
40      return grid
41
42  def plot_policy_grid(policy, title):
43      grid = policy_to_grid(policy)
44      plt.figure(figsize=(5,3))
45      plt.imshow(np.zeros((ROWS, COLS)), vmin=0, vmax=1)
46      for (i, j), val in np.ndenumerate(grid):
47          plt.text(j, i, val, ha='center', va='center', fontsize=14)
48      plt.title(title)
49      plt.xticks(range(COLS))
50      plt.yticks(range(ROWS))
51      plt.show()
52
53  # Default run (R=-0.04)
54  V_default = value_iteration(gamma=0.99, theta=1e-4, living_penalty
        =-0.04)
55  pi_default = extract_policy(V_default, gamma=0.99, living_penalty
        =-0.04)
56  plot_value_heatmap(V_default, "Value Heatmap    Default (R=-0.04)")
57  plot_policy_grid(pi_default, "Policy    Default (R=-0.04)")
58
59  # Q1/Q2/Q3: change penalties
60  V_zero = value_iteration(gamma=0.99, theta=1e-4, living_penalty=0.0)
61  pi_zero = extract_policy(V_zero, gamma=0.99, living_penalty=0.0)
62  plot_value_heatmap(V_zero, "Value Heatmap    No Penalty (R=0.0)")
63  plot_policy_grid(pi_zero, "Policy    No Penalty (R=0.0)")
64
65  V_high = value_iteration(gamma=0.99, theta=1e-4, living_penalty=-0.5)
66  pi_high = extract_policy(V_high, gamma=0.99, living_penalty=-0.5)
67  plot_value_heatmap(V_high, "Value Heatmap    High Penalty (R=-0.5)")
```

```
68  plot_policy_grid(pi_high, "Policy      High Penalty (R=-0.5)")
```

# 6 Plots

## 6.1 Default Living Penalty (R=-0.04)


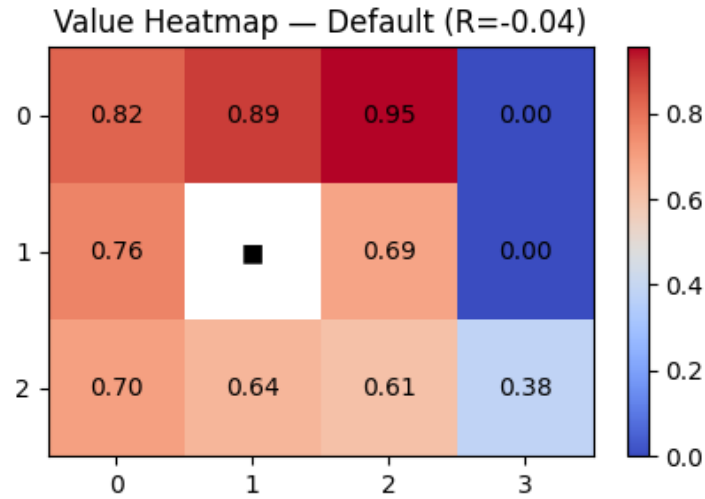
Figure 1: Value Heatmap — Default Living Penalty (R=-0.04)
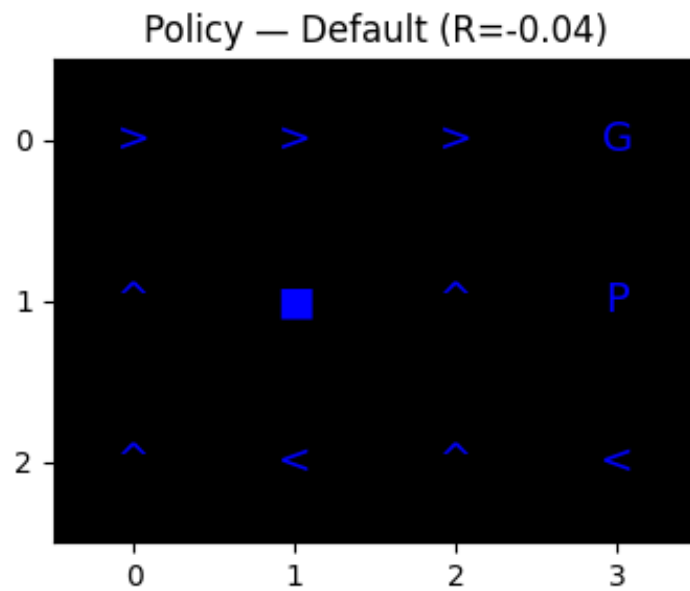


Figure 2: Optimal Policy — Default Living Penalty (R=-0.04)
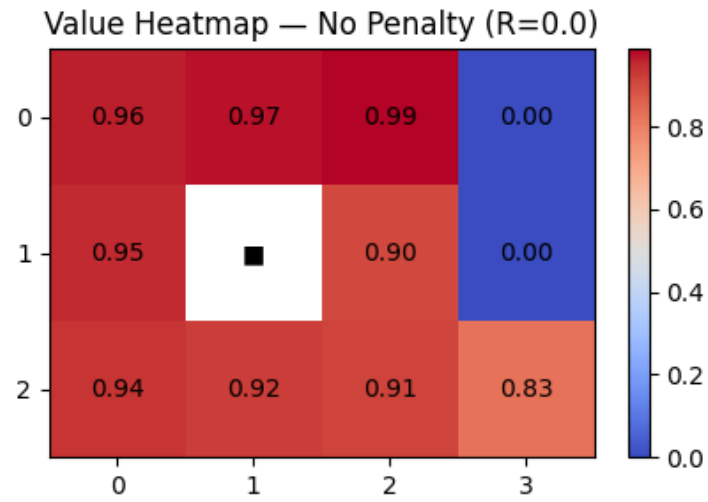
## 6.2 No Living Penalty (R=0.0)
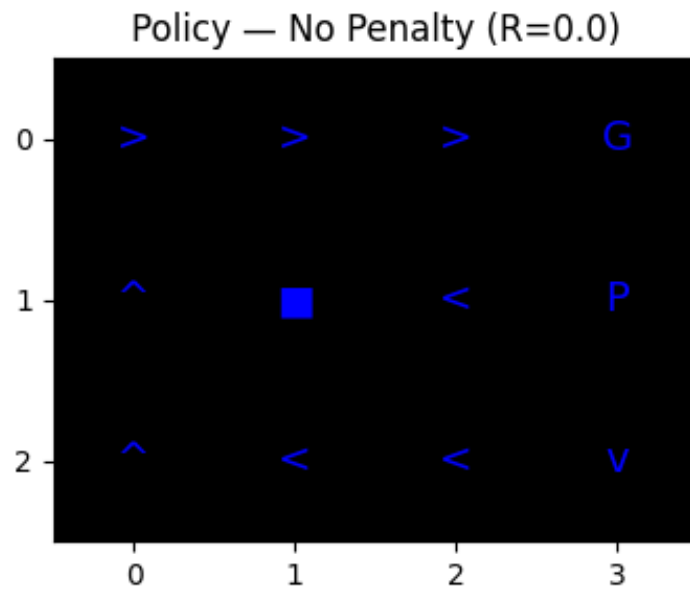


Figure 3: Value Heatmap — No Living Penalty (R=0.0)



Figure 4: Optimal Policy — No Living Penalty (R=0.0)
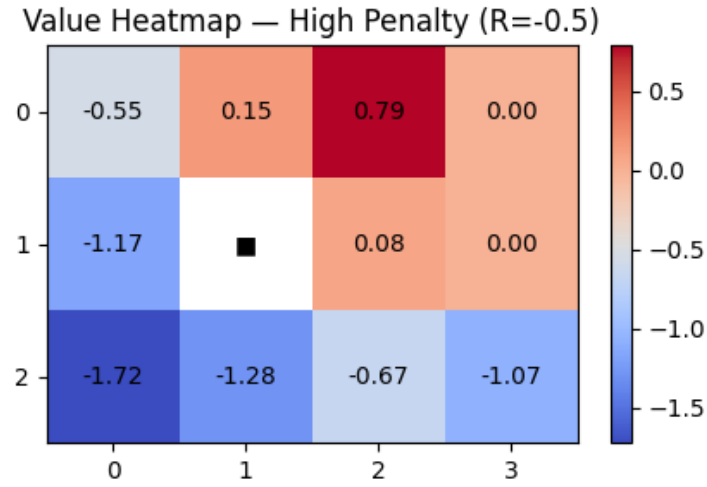
## 6.3 High Living Penalty (R=-0.5)



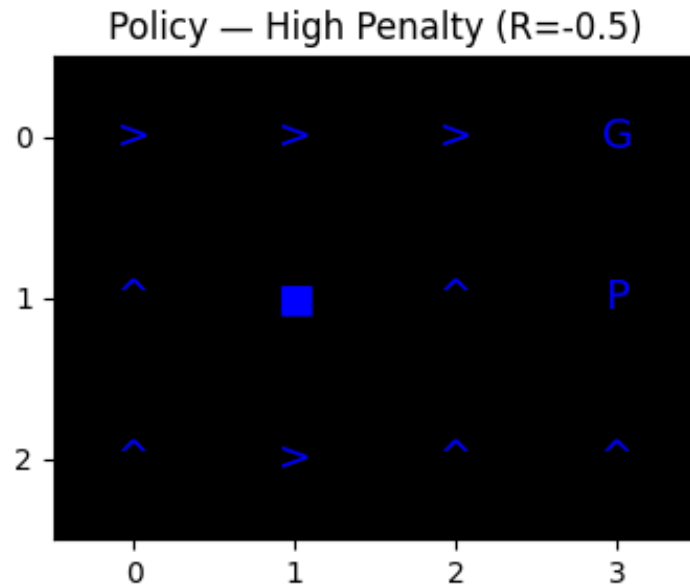Figure 5: Value Heatmap — High Living Penalty (R=-0.5)



Figure 6: Optimal Policy — High Living Penalty (R=-0.5)

# 7 Observations and Analysis

- **Default living penalty (-0.04):** Agent avoids the pit and finds the shortest path to the goal.

- **No living penalty (0.0):** Agent still reaches the goal but is less incentivized to take the shortest path; may wander more.

- **High living penalty (-0.5):** Agent aggressively chooses the shortest path to minimize cumulative penalty, even risking nearby pit.

- **Effect of stochasticity:** Slipping probabilities cause the agent to choose safer paths to avoid falling into the pit.

- **Value function interpretation:** Higher values correspond to states closer to the goal with lower risk of entering the pit.