

MEIC

PLANNING, LEARNING AND INTELLIGENT DECISION MAKING

Homework 2

Grupo 23:

André Ramos Loja Infante Barbosa (103008)

Artur Filipe de Menezes Crespo Ferreira (102482)

2024/2025 – 2º Semestre, P3

Conteúdo

1	Problem Definition	2
2	MDP Components	2
2.1	State Space	2
2.2	Action Space	2
2.3	Transition Matrices	2
3	Optimal Policy at State 2	4

1 Problem Definition

The problem is a Markov Decision Process (MDP) where a spider climbs a ladder with probabilistic transitions, influenced by weather conditions. A new action, **Stop**, allows the spider to build a web at state 2, preventing it from falling completely to the ground. The discount factor is given as $\gamma = 0.9$.

2 MDP Components

2.1 State Space

We introduce additional states to represent the web scenario:

$$X = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11\} \quad (1)$$

where states 6 – 11 correspond to the web-enabled ladder states.

2.2 Action Space

The spider has two available actions:

- **Play (P)**: Moves up the ladder based on die rolls.
- **Stop (S)**: Builds a web at state 2; stopping elsewhere results in falling.

2.3 Transition Matrices

The transition probability matrix for the **Play** action (extended to account for web states) is:

$$P_{\text{play}} = \begin{bmatrix} 0.2 & 0.4 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.2 & 0 & 0.4 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.2 & 0 & 0 & 0.4 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.2 & 0 & 0 & 0 & 0.4 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.2 & 0 & 0 & 0 & 0 & 0.8 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.2 & 0.4 & 0.4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.2 & 0.4 & 0.4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.2 & 0 & 0 & 0.4 & 0.4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.2 & 0 & 0 & 0 & 0.4 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.2 & 0 & 0 & 0 & 0 & 0.8 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

For the **Stop** action, the matrix of transitions is like follows:

$$P_{\text{stop}} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

Part (b): Cost-to-Go for Two Deterministic Policies in State 2

We compute the cost-to-go $J^\pi(s)$ for the two deterministic policies in state 2.

Policy 1: Always "Play" in State 2

For **Policy 1**, the spider always chooses to "Play" in state 2. The Bellman equation for state 2 is:

$$J^{\pi_1}(2) = \sum_{s'} P(2, \text{Play}, s') [c(2, \text{Play}, s') + \gamma J^{\pi_1}(s')]$$

Listing 1: Python Code for Perceptron Training

```

1 import numpy as np
2
3
4 # Discount factor
5 gamma = 0.9
6
7 # Number of states
8 n_states = 12
9
10 # Transition probability matrix for "Play" action (P_play)
11 P_play = np.array([
12     [0.2, 0.4, 0.4, 0, 0, 0, 0, 0, 0, 0, 0, 0],
13     [0.2, 0, 0.4, 0.4, 0, 0, 0, 0, 0, 0, 0, 0],
14     [0.2, 0, 0, 0.4, 0.4, 0, 0, 0, 0, 0, 0, 0],
15     [0.2, 0, 0, 0, 0.4, 0.4, 0, 0, 0, 0, 0, 0],
16     [0.2, 0, 0, 0, 0, 0.8, 0, 0, 0, 0, 0, 0],
17     [0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0],
18     [0, 0, 0, 0, 0, 0, 0.2, 0.4, 0.4, 0, 0, 0],
19     [0, 0, 0, 0, 0, 0, 0.2, 0.4, 0.4, 0, 0, 0],
20     [0, 0, 0, 0, 0, 0, 0.2, 0, 0, 0.4, 0.4, 0],
21     [0, 0, 0, 0, 0, 0, 0.2, 0, 0, 0, 0.4, 0.4],
22     [0, 0, 0, 0, 0, 0, 0.2, 0, 0, 0, 0, 0.8],
23     [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]

```

```

24 ]
25
26 # Transition probability matrix for "Stop" action (P_stop)
27 P_stop = np.array([
28     [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
29     [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
30     [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0],
31     [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
32     [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
33     [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
34     [0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0],
35     [0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0],
36     [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0],
37     [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0],
38     [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0],
39     [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0]
40 ])
41
42 # Cost vector (1 everywhere except final states 5 and 11 where cost is 0)
43 c = np.ones(n_states)
44 c[5] = 0
45 c[11] = 0
46
47 # Compute J_pi for Play
48 I = np.eye(n_states)
49 J_play = np.linalg.solve(I - gamma * P_play, c)
50
51 # Compute J_pi for Stop
52 J_stop = np.linalg.solve(I - gamma * P_stop, c)
53
54 # Compute Q-functions
55 Q_play = c + gamma * np.dot(P_play, J_play)
56 Q_stop = c + gamma * np.dot(P_stop, J_stop)
57
58 # Compute the optimal cost-to-go function J*(s) (minimum of Q-functions)
59 J_star = np.minimum(Q_play, Q_stop)
60
61 # Print results
62 print("Optimal Cost-to-Go Function J* considering both Play and Stop:")
63 print(J_star)

```

The output of the program is the following:

Optimal Cost – to – Go Function J* considering both Play and Stop :

```

[4.35671816  3.84578458  3.30007346  2.42652461  1.78420927  0.
 4.89747459  4.89747459  3.48010642  2.55890178  1.88154543  0.]

```

3 Optimal Policy at State 2

The cost-to-go for the two policies in state 2 are: - **Policy 1 (Always "Play")**: $J^{\pi_1}(2) = 3.30$ - **Policy 2 (Always "Stop")**: $J^{\pi_2}(2) = 3.48$

Thus, **Policy 1** is better in state 2 because it has a lower cost-to-go.