

Assignment 1 – Q3

Majid Ghasemi

September 2024

1 Q3

The general form of Bellman Equation for the state-value function is:

$$v_{\pi}(s) = \mathbb{E}_{\pi} \{G_t \mid S_t = s\} \quad (1)$$

Which can be written as:

$$\mathbb{E}_{\pi} \{R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s\} \quad (2)$$

or without the expectation operator as:

$$v_{\pi}(s) = \sum_{a \in A} \pi(a \mid s) \left[\sum_{s' \in S} \left(\sum_{r \in R} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')] \right) \right] \quad (3)$$

In this question, I will use Eq. 4 as another form of Bellman Equation to answer the question.

$$v_{\pi}(s) = \sum_{a \in A} \pi(a \mid s) \left[\sum_{s' \in S} p(s' \mid s, a) [r + \gamma v_{\pi}(s')] \right] \quad (4)$$

1.1 Question (a)

As states in the question, we have the equiprobable policy, and the environment is deterministic, in which choosing actions are taking place deterministically. γ is 0.9. Here is a list of notations used for this question:

- N : North
- S : South
- W : West
- E : East

We can define the reward function as follows as well:

$$r(s, a, s') \quad (5)$$

Where s is the current state, a is the taken action, and s' is the next state that the agent ends up there after taking action a .

Now, we can rewrite the Bellman Equation without any summations.

$$\begin{aligned} v_{\pi}(s_5) = & \frac{1}{4} [r(s_5, N, s_0) + \gamma v_{\pi}(s_0)] + \frac{1}{4} [r(s_5, S, s_{10}) + \gamma v_{\pi}(s_{10})] \\ & + \frac{1}{4} [r(s_5, W, s_5) + \gamma v_{\pi}(s_5)] + \frac{1}{4} [r(s_5, E, s_6) + \gamma v_{\pi}(s_6)] \end{aligned} \quad (6)$$

Only the reward of hitting the wall and staying in s_5 (action W) is not zero (it is -1). Hence, we can rewrite it as follows:

$$v_\pi(s_5) = \frac{1}{4} [0 + \gamma v_\pi(s_0)] + \frac{1}{4} [0 + \gamma v_\pi(s_{10})] + \frac{1}{4} [-1 + \gamma v_\pi(s_5)] + \frac{1}{4} [0 + \gamma v_\pi(s_6)] \quad (7)$$

Substituting $\gamma = 0.9$, we have:

$$v_\pi(s_5) = \frac{1}{4} [0 + 0.9v_\pi(s_0)] + \frac{1}{4} [0 + 0.9v_\pi(s_{10})] + \frac{1}{4} [-1 + 0.9v_\pi(s_5)] + \frac{1}{4} [0 + 0.9v_\pi(s_6)] \quad (8)$$

We can simplify it further:

$$v_\pi(s_5) = \frac{1}{4} [[0 + 0.9v_\pi(s_0)] + [0 + 0.9v_\pi(s_{10})] + [-1 + 0.9v_\pi(s_5)] + [0 + 0.9v_\pi(s_6)]] \quad (9)$$

We can even simplify it more to be more clean and neat looking equation:

$$v_\pi(s_5) = \frac{1}{4} [0.9 [(v_\pi(s_0) + v_\pi(s_{10}) + v_\pi(s_6)) + (-1 + v_\pi(s_5))]] \quad (10)$$

I factored the gamma and put the states without reward in one place, and s_5 in another parenthesis.

Similarly, for s_{12} , we can have this series of equation. However, to save time for the readers, I will just give the final equation. All the rewards are zero, and by taking action N, agent goes to s_6 , by action S goes to s_{17} , action W takes it to s_{11} , and action E takes it to be in s_{13} :

$$v_\pi(s_{12}) = \frac{1}{4} [0.9 [(v_\pi(s_6) + v_\pi(s_{17}) + v_\pi(s_{11}) + v_\pi(s_{13}))]] \quad (11)$$

1.2 Question (b)

Using Eq. 6, we can answer this question. However, it must be noted that the value of $\pi(a|s)$ is different than the previous policy, and is given in the question.

Rewriting the Eq. 6, we have:

$$v_\pi(s_5) = 0.2 [r(s_5, N, s_0) + \gamma v_\pi(s_0)] + 0.2 [r(s_5, S, s_{10}) + \gamma v_\pi(s_{10})] + 0.1 [r(s_5, W, s_5) + \gamma v_\pi(s_5)] + 0.5 [r(s_5, E, s_6) + \gamma v_\pi(s_6)] \quad (12)$$

Again, only the reward of hitting the wall is not zero, and we can rewrite the equation above as:

$$v_\pi(s_5) = 0.2 [0 + \gamma v_\pi(s_0)] + 0.2 [0 + \gamma v_\pi(s_{10})] + 0.1 [-1 + \gamma v_\pi(s_5)] + 0.5 [0 + \gamma v_\pi(s_6)] \quad (13)$$

By substituting the numbers from the given table as v values, and gamma, we have:

$$v_\pi(s_5) = 0.2 [0 + 0.9 \times 3.3] + 0.2 [0 + 0.9 \times 0.1] + 0.1 [-1 + 0.9 \times 1.5] + 0.5 [0 + 0.9 \times 3] \quad (14)$$

Is equal to: 1.997.

1.3 Question (c)

Similar to question (b), for $v_\pi(s_{12})$, after simplifications, we have:

$$v_\pi(s_{12}) = 0.6 [0 + 0.9 \times 2.3] + 0.02 [0 + 0.9 \times -0.4] + 0.08 [0 + 0.9 \times 0.4] + 0.3 [0 + 0.9 \times 0.7] \quad (15)$$

Where it is equal to: 1.4526.

1.4 Question (d)

Based on the achieved results in the previous sections, we can see that the results under different policies, outperformed our initial equiprobable policy.

SIGNED AS: MAJID GHASEMI , September 22