



AU332 Quiz9

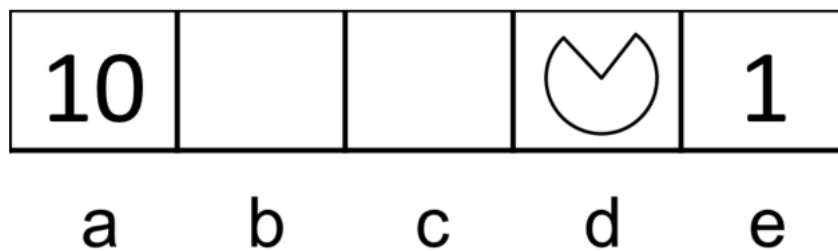
* 基本信息:

姓名:

学号:

*1. Solving MDPs

Consider the gridworld MDP for which Left and Right actions are 100% successful. Specifically, the available actions in each state are to move to the neighboring grid squares. From state a, there is also an exit action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state e, the reward for the exit action is 1. Exit actions are successful 100% of the time. There is no living reward.



Let the discount factor $\gamma = 1$. Fill in the following blanks.

$V_0(d) =$

$V_1(d) =$

$V_2(d) =$

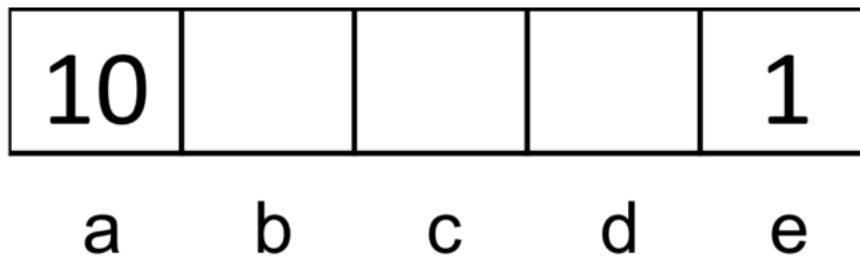
$V_3(d) =$

$V_4(d) =$

$V_5(d) =$

*2. Value Iteration Convergence Values

Consider the gridworld where Left and Right actions are successful 100% of the time. Specifically, the available actions in each state are to move to the neighboring grid squares. From state a, there is also an exit action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state e, the reward for the exit action is 1. Exit actions are successful 100% of the time. There is no living reward.



Let the discount factor $\gamma = 0.2$. Fill in the following blanks.

$$V^*(a) = V_\infty(a) = \underline{\hspace{2cm}}$$

$$V^*(b) = V_\infty(b) = \underline{\hspace{2cm}}$$

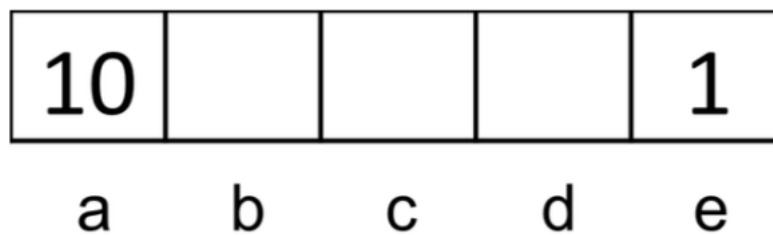
$$V^*(c) = V_\infty(c) = \underline{\hspace{2cm}}$$

$$V^*(d) = V_\infty(d) = \underline{\hspace{2cm}}$$

$$V^*(e) = V_\infty(e) = \underline{\hspace{2cm}}$$

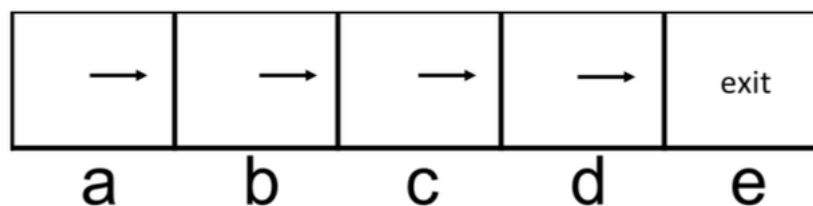
*3. Policy Evaluation

Consider the grid world where Left and Right actions are successful 100% of the time. Specifically, the available actions in each state are to move to the neighboring grid squares. From state a, there is also an exit action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state e, the reward for the exit action is 1. Exit actions are successful 100% of the time. The discount factor (γ) is 1.



Part 1

Consider the policy π_1 shown below, and evaluate the following quantities for this policy.



$$V^{\pi_1}(a) = \underline{\hspace{2cm}}$$

$$V^{\pi_1}(b) = \underline{\hspace{2cm}}$$

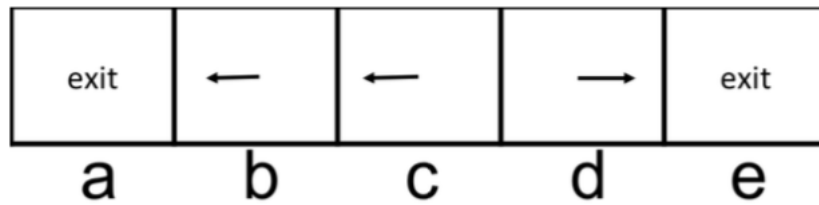
$$V^{\pi_1}(c) = \underline{\hspace{2cm}}$$

$$V^{\pi_1}(d) = \underline{\hspace{2cm}}$$

$$V^{\pi_1}(e) = \underline{\hspace{2cm}}$$

Part 2

Consider the policy π_2 shown below, and evaluate the following quantities for this policy.



$$V^{\pi^2}(a) = 10$$

$$V^{\pi^2}(b) = 10$$

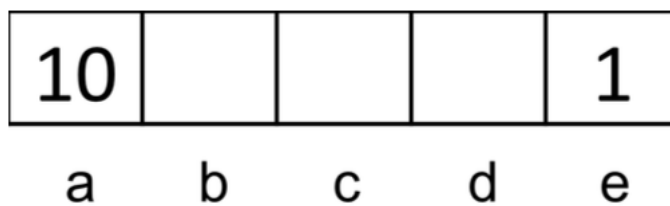
$$V^{\pi^2}(c) = 10$$

$$V^{\pi^2}(d) = 1$$

$$V^{\pi^2}(e) = 1$$

*4. Policy Iteration

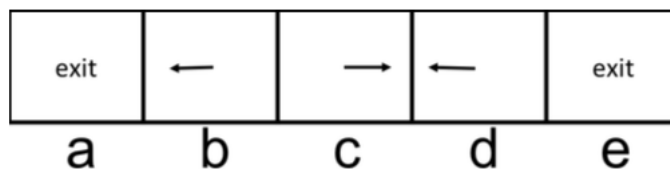
Consider the grid world where Left and Right actions are successful 100% of the time. Specifically, the available actions in each state are to move to the neighboring grid squares. From state a, there is also an exit action available, which results in going to the terminal state and collecting a reward of 10. Similarly, in state e, the reward for the exit action is 1. Exit actions are successful 100% of the time. The discount factor (γ) is 0.9.



We will execute one round of policy iteration.

Part 1 Policy Evaluation

Consider the policy π_i shown below, and evaluate the following quantities for this policy.



$$V^{\pi_i}(a) = 10$$

$$V^{\pi_i}(b) = 9$$

$$V^{\pi_i}(c) = 0$$

$$V^{\pi_i}(d) = 0$$

$$V^{\pi_i}(e) = 1$$

Part 2: Policy Improvement

Perform a policy improvement step. The current policy's values are the ones from Part 1.

*5. $\pi_{i+1}(a) =$

☒ Exit

☐ Right

*6. $\pi_{i+1}(\mathbf{b})=$

☒ Left

Right

*7. $\pi_{i+1}(\mathbf{c})=$

☒ Left

Right

*8. $\pi_{i+1}(\mathbf{d})=$

Left

☒ Right

*9. $\pi_{i+1}(\mathbf{e})=$

Left

☒ Exit

提交

问卷星 提供技术支持