

Homework 8

April 1, 2019 10:06 PM

Spend or Save

You are facing a small crisis where you have to decide whether you should spend your money or save it in your bank account. If the balance of your bank account is low, you can only buy something that gives you a little bit of joy. If you save your money, your account's balance increases, but you actually feel a bit sad. Later in the future, however, you can buy something more expensive that gives you more joy. Your dilemma is what the optimal decision should be depending on your current financial state (low or high).

This problem can be formulated as a discounted MDP where you have two states s_1 , which corresponds to low amount of money in your bank account, and s_2 , which corresponds to having a lot of money in your bank account. At each state you have two actions:

- a_1 : Save money
- a_2 : Spend money

Depending on the current state and the selected action, your financial state in the next time step might change. You also receive some amount of reward, which is an indicator of your level of joy.

To formulate this problem as a discounted MDP, we have to define transition probabilities, reward function, and the discount factor. The discount factor $0 \leq \gamma < 1$ determines how myopic/farsighted you are.¹ Figure 1 describes the dynamics and the reward. In this figure, all the transitions are deterministic, e.g., if you are at state s_1 and choose action a_1 , you will definitely move to state s_2 . In mathematical term, $\mathcal{P}(s_2|s_1, a_1) = 1$.

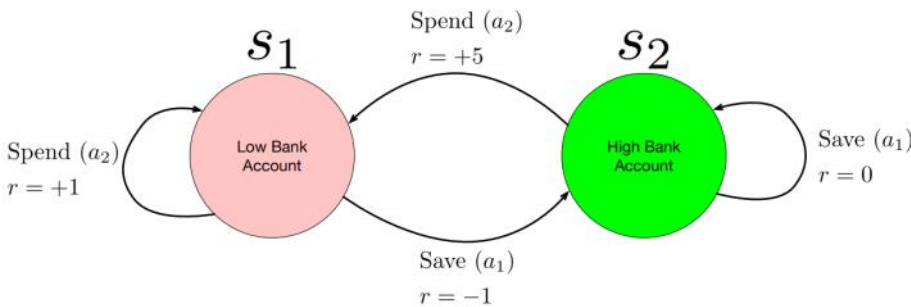


Figure 1: Spend or Save

The purpose of this question is to make you feel comfortable solving simple MDPs. Answer the following questions:

1. Write down $\mathcal{P}(s'|s, a)$ for all $s, s' = \{s_1, s_2\}$ and $a \in \{a_1, a_2\}$.

$$\begin{aligned} P(s'=1 | s=1, a=1) &= 0 & P(s'=1 | s=1, a=2) &= 1 \\ P(s'=2 | s=1, a=1) &= 1 & P(s'=2 | s=1, a=2) &= 0 \\ P(s'=1 | s=2, a=1) &= 0 & P(s'=1 | s=2, a=2) &= 1 \\ P(s'=2 | s=2, a=1) &= 1 & P(s'=2 | s=2, a=2) &= 0 \end{aligned}$$

2. Consider the following two policies:

- $\pi_{\text{save}}(s) = a_1$
- $\pi_{\text{spend}}(s) = a_2$

Write down the Bellman equations for $Q^{\pi_{\text{save}}}(s, a)$ and $Q^{\pi_{\text{spend}}}(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

$$\begin{aligned} Q^{\pi_{\text{save}}}(s_1, a_1) &= -1 + \gamma Q^{\pi_{\text{save}}}(s_2, a_1) \\ Q^{\pi_{\text{save}}}(s_1, a_2) &= 1 + \gamma Q^{\pi_{\text{save}}}(s_1, a_1) \end{aligned}$$

$$Q^{\pi_{\text{save}}}(s_2, a_1) = 0 + \gamma Q^{\pi_{\text{save}}}(s_2, a_1) \quad (1)$$

$$Q^{\pi_{\text{save}}}(s_2, a_2) = 5 + \gamma Q^{\pi_{\text{save}}}(s_1, a_1)$$

from (1),

$$\begin{aligned} Q^{\pi_{\text{save}}}(s_2, a_1) &= 0 \\ Q^{\pi_{\text{save}}}(s_1, a_1) &= -1 \\ Q^{\pi_{\text{save}}}(s_2, a_2) &= 5 - \gamma \\ Q^{\pi_{\text{save}}}(s_1, a_2) &= 1 - \gamma \end{aligned}$$

$$Q^{\pi_{\text{spend}}}(s_1, a_1) = -1 + \gamma Q^{\pi_{\text{spend}}}(s_2, a_2)$$

$$Q^{\pi_{\text{spend}}}(s_1, a_2) = 1 + \gamma Q^{\pi_{\text{spend}}}(s_1, a_2) \quad (2)$$

$$Q^{\pi_{\text{spend}}}(s_2, a_1) = 0 + \gamma Q^{\pi_{\text{spend}}}(s_2, a_2)$$

$$Q^{\pi_{\text{spend}}}(s_2, a_2) = 5 + \gamma Q^{\pi_{\text{spend}}}(s_1, a_2)$$

from (2),

$$\begin{aligned} Q^{\pi_{\text{spend}}}(s_1, a_2) &= \frac{1}{1-\gamma} \\ Q^{\pi_{\text{spend}}}(s_2, a_2) &= 5 + \frac{1}{1-\gamma} \\ Q^{\pi_{\text{spend}}}(s_1, a_1) &= -1 + \gamma \left(5 + \frac{1}{1-\gamma} \right) \\ Q^{\pi_{\text{spend}}}(s_2, a_1) &= \gamma \left(5 + \frac{1}{1-\gamma} \right) \end{aligned}$$

4. Write a simple program that computes the *optimal* action-value function and *optimal* policy for a given γ .

See HW8.py.

5. How does the choice of γ affect your optimal policy?

The larger the γ , the more we prioritize on future rewards.

The smaller the γ , the more we prioritize on immediate rewards.