

MDPs + RL

1 MDPs: Micro-Blackjack

In micro-blackjack, you repeatedly draw a card (with replacement) that is equally likely to be a 2, 3, or 4. You can either Draw or Stop if the total score of the cards you have drawn is less than 6. If your total score is 6 or higher, the game ends, and you receive a utility of 0. When you Stop, your utility is equal to your total score (up to 5), and the game ends. When you Draw, you receive no utility. There is no discount ($\gamma = 1$). Let's formulate this problem as an MDP with the following states: 0, 2, 3, 4, 5 and a *Done* state, for when the game ends.

1. What is the transition function and the reward function for this MDP?

The transition function is

$$\begin{aligned}T(s, \text{Stop}, \text{Done}) &= 1 \\T(0, \text{Draw}, s') &= 1/3 \text{ for } s' \in \{2, 3, 4\} \\T(2, \text{Draw}, s') &= 1/3 \text{ for } s' \in \{4, 5, \text{Done}\} \\T(3, \text{Draw}, s') &= \begin{cases} 1/3 \text{ if } s' = 5 \\ 2/3 \text{ if } s' = \text{Done} \end{cases} \\T(4, \text{Draw}, \text{Done}) &= 1 \\T(5, \text{Draw}, \text{Done}) &= 1 \\T(s, a, s') &= 0 \text{ otherwise}\end{aligned}$$

The reward function is

$$\begin{aligned}R(s, \text{Stop}, \text{Done}) &= s, s \leq 5 \\R(s, a, s') &= 0 \text{ otherwise}\end{aligned}$$

2. Fill in the following table of value iteration values for the first 4 iterations.

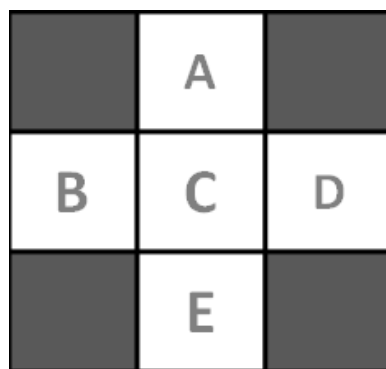
States	0	2	3	4	5
V_0	0	0	0	0	0
V_1	0	2	3	4	5
V_2	3	3	3	4	5
V_3	10/3	3	3	4	5
V_4	10/3	3	3	4	5

3. You should have noticed that value iteration converged above. What is the optimal policy for the MDP?

States	0	2	3	4	5
π^*	Draw	Draw	Stop	Stop	Stop

2 Learning in Gridworld

Consider the example gridworld that we looked at in lecture. We would like to use TD learning and q-learning to find the values of these states.



Passive Fixed Policy **Active**

1. (B, East, C, 2)

$$V(B) = 1 = (1-0.5)*0 + 0.5*(2+1*0)$$

2. (C, South, E, 4)

$$V(C) = 2 = (1-0.5)*0 + 0.5*(4+1*0)$$

3. (C, East, A, 6)

$$V(C) = 4 = (1-0.5)*2 + 0.5*(6+1*0)$$

3. (B, East, C, 2)

$$V(B) = 3.5 = (1-0.5)*1 + 0.5*(2+1*4)$$

Suppose that we have the following observed transitions:

(B, East, C, 2), (C, South, E, 4), (C, East, A, 6), (B, East, C, 2)

The initial value of each state is 0. Assume that $\gamma = 1$ and $\alpha = 0.5$.

1. What are the learned values from TD learning after all four observations?

$$V(B) = 3.5$$

$$V(C) = 4$$

All other states have a value of 0.

$$V^{\pi}(s) \leftarrow (1-\alpha)V^{\pi}(s) + \alpha(R(s') + \gamma V^{\pi}(s'))$$

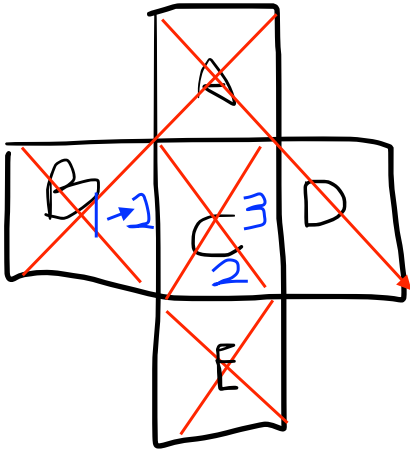
2. What are the learned Q-values from Q-learning after all four observations?

$$Q(B, East) = 3$$

$$Q(C, South) = 2$$

$$Q(C, East) = 3$$

All other q-states have a value of 0.



$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha(r + \gamma \max_{a'} Q(s', a'))$$

1. (B, East, C, 2)

$$Q(B, East) = 1 = (1-0.5)*0 + 0.5*(2+1*0)$$

2. (C, South, E, 4)

$$Q(C, South) = 2 = (1-0.5)*0 + 0.5*(4+1*0)$$

3. (C, East, A, 6)

$$Q(C, East) = 3 = (1-0.5)*0 + 0.5*(6+1*0)$$

3. (B, East, C, 2)

$$Q(B, East) = 3 = (1-0.5)*1 + 0.5*(2+1*3)$$