

CS 234 Winter 2020  
Assignment 1  
Due: January 22 at 11:59 pm

ksang  
kaitoy@qq.com

For submission instructions please refer to [website](#). For all problems, if you use an existing result from either the literature or a textbook to solve the exercise, you need to cite the source.

## 1 Gridworld [15 pts]

Consider the following grid environment. Starting from any unshaded square, you can move up, down, left, or right. Actions are deterministic and always succeed (e.g. going left from state 16 goes to state 15) unless they will cause the agent to run into a wall. The thicker edges indicate walls, and attempting to move in the direction of a wall results in staying in the same square (e.g. going in any direction other than left from state 16 stays in 16). Taking any action from the green target square (no. 12) earns a reward of  $r_g$  (so  $r(12, a) = r_g \forall a$ ) and ends the episode. Taking any action from the red square of death (no. 5) earns a reward of  $r_r$  (so  $r(5, a) = r_r \forall a$ ) and ends the episode. Otherwise, from every other square, taking any action is associated with a reward  $r_s \in \{-1, 0, +1\}$  (even if the action results in the agent staying in the same square). Assume the discount factor  $\gamma = 1$ ,  $r_g = +5$ , and  $r_r = -5$  unless otherwise specified.

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

- (a) (3pts) Define the value of  $r_s$  that would cause the optimal policy to return the shortest path to the green target square (no. 12). Using this  $r_s$ , find the optimal value for each square.

**Answer:**

$$r_s = -1$$

$$\begin{bmatrix} 0 & 1 & 2 & 3 \\ -5 & 2 & 3 & 4 \\ 2 & 3 & 4 & 5 \\ 1 & 0 & -1 & -2 \end{bmatrix}$$

- (b) (3pts) Lets refer to the value function derived in (a) as  $V_{old}^{\pi_g}$  and the policy as  $\pi_g$ . Suppose we are now in a new gridworld where all the rewards ( $r_s$ ,  $r_g$ , and  $r_r$ ) have +2 added to them. Consider still following  $\pi_g$  of the original gridworld, what will the new values  $V_{new}^{\pi_g}$  be in this second gridworld?

**Answer:**

$$\begin{bmatrix} 12 & 11 & 10 & 9 \\ -3 & 10 & 9 & 8 \\ 10 & 9 & 8 & 7 \\ 11 & 12 & 13 & 14 \end{bmatrix}$$

- (c) (3pts) Consider a general MDP with rewards, and transitions. Consider a discount factor of  $\gamma$ . For this case assume that the horizon is infinite (so there is no termination). A policy  $\pi$  in this MDP induces a value function  $V^\pi$  (lets refer to this as  $V_{old}^\pi$ ). Now suppose we have a new MDP where the only difference is that all rewards have a constant  $c$  added to them. Can you come up with an expression for the new value function  $V_{new}^\pi$  induced by  $\pi$  in this second MDP in terms of  $V_{old}^\pi$ ,  $c$ , and  $\gamma$ ?

**Answer:**

$$\begin{aligned} V_{old}^\pi(s) &= \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s \right] \\ V_{new}^\pi(s) &= \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t (r_t + c) | s_0 = s \right] \\ &= \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s \right] + c \sum_{t=0}^{\infty} \gamma^t \\ &= V_{old}^\pi(s) + \frac{c}{1 - \gamma} \end{aligned}$$

- (d) (2pts) Lets go back to our gridworld from (a) with the default values for  $r_g$ ,  $r_r$ ,  $\gamma$  and with the value you specified for  $r_s$ . Suppose we now derived a second gridworld by adding a constant  $c$  to all rewards ( $r_s$ ,  $r_g$ , and  $r_r$ ) such that  $r_s = +2$ . How does the optimal policy change (Just give a one or two sentence description)? What do the values of the unshaded squares become?

**Answer:**

It will become looping at states other than the terminate states, because each step will gain positive reward value.

Unshaded squares will become value  $+\infty$ .

- (e) (2pts) Now take the second gridworld from part (d) and change  $\gamma$  such that  $0 < \gamma < 1$ . Can the optimal policy change and does it depend on your choice of gamma? (A brief description is sufficient, no formal proof or mathematical analysis required).

**Answer:**

If  $\gamma$  is close to 1, the result is same, all unshaded squares will have value  $+\infty$  and optimal policy is keep looping.

if  $\gamma$  is equal to some value closing to 0, at some point optimal policy will be able to find a shortest path to green square.

- (f) (2pts) Lets go back to our gridworld from (a) with the default values for  $r_g$ ,  $r_r$ ,  $\gamma$  and with the value you specified for  $r_s$ . In this gridworld, our optimal policy from any unshaded square never terminates in the red square. Now suppose  $r_s$  can take on any real, non-infinite value and is not restricted to  $\{+1, 0, -1\}$  anymore. Give a value of  $r_s$  such that there are unshaded squares starting from which following the optimal policy results in termination in the red square.

**Answer:**

$$r_s \leq -5$$

## 2 Value of Different Policies [35 pts]

In many situations such as healthcare or education, we cannot run any arbitrary policy and collect data from running those policies for evaluation. In these cases, we may need to take data collected from following one policy and use it to evaluate the value of a different policy. The equality proved in the following exercise can be an important tool for achieving this. The purpose of this exercise is to get familiar on how to compare the value of different policies,  $\pi_1$  and  $\pi_2$ , on a fixed horizon MDP. A fixed horizon MDP is an MDP where the agent's state is reset after  $H$  timesteps;  $H$  is called the *horizon* of the MDP. There is no discount (i.e.,  $\gamma = 1$ ) and policies are allowed to be non-stationary, i.e., the action identified by a policy depends on the timestep in addition to the state. Let  $x_t \sim \pi$  denote the distribution over states at timestep  $t$  (for  $1 \leq t \leq H$ ) upon following policy  $\pi$  and  $V_t^\pi(x_t)$  denote the value function of policy  $\pi$  in state  $x_t$  and timestep  $t$ , and  $Q_t^\pi(x_t, a)$  denote the corresponding  $Q$  value associated to action  $a$ . As a clarifying example, we denote  $\mathbb{E}_{x_t \sim \pi_1} V(x_t)$  to represent the average value of the value function  $V(\cdot)$  over the states at timestep  $t$  encountered upon following policy  $\pi_1$ . Please show the following:

$$V_1^{\pi_1}(x_1) - V_1^{\pi_2}(x_1) = \sum_{t=1}^H \mathbb{E}_{x_t \sim \pi_2} \left( Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_1}(x_t, \pi_2(x_t, t)) \right)$$

**Intuition:** The above expression can be interpreted in the following way. For concreteness, assume that  $\pi_1$  is the better policy, i.e., achieving  $V_1^{\pi_1}(x_1) \geq V_1^{\pi_2}(x_1)$ . Suppose you're following policy  $\pi_2$  and you are at timestep  $t$  in state  $x_t$ . You have the option to follow  $\pi_1$  (the better policy) until the end of the episode, totalling  $Q_t^{\pi_1}(x_t, \pi_1(x_t, t))$  return from the current state-timestep; or you have the option to follow  $\pi_2$  for one timestep and then follow  $\pi_1$  instead until the end of the episode (you can follow many other policies of course). This would give you a "loss" of  $Q_t^{\pi_1}(x_t, \pi_1(x_t, t)) - Q_t^{\pi_1}(x_t, \pi_2(x_t, t))$  that originates from following the worse policy  $\pi_2$  instead of  $\pi_1$  in that timestep. Then the equation above means that the value difference of the two policies is the sum of all the losses induced by following the suboptimal policy for every timestep, weighted by the expected trajectory of the policy you're following.

### 3 Fixed Point [25 pts]

In this exercise we will use [Cauchy sequences](#) to prove that value iteration will converge to a unique fixed point (in this case, a value function  $V$ ) regardless of the starting point. An element  $V$  is a fixed point for an operator  $B$  (in this case the Bellman operator) if performance of  $B$  on  $V$  returns  $V$ , i.e.,  $BV = V$ . Recall that the Bellman backup operator  $B$  is defined as (in lecture 2):

$$V_{k+1} \stackrel{def}{=} BV_k = \max_a [R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_k^\pi(s')].$$

Additionally, in lecture 2, we proved that this Bellman backup is a contraction for  $\gamma < 1$  on the infinity norm

$$\|BV' - BV''\|_\infty \leq \gamma \|V' - V''\|_\infty$$

for any two value functions  $V'$  and  $V''$ , meaning if we apply it to two different value functions, the distance between value functions (in the  $\infty$  norm) shrinks after application of the operator to each element.

- (a) (5pts) Prove by induction that  $\|V_{n+1} - V_n\|_\infty \leq \gamma^n \|V_1 - V_0\|_\infty$
- (b) (10pts) Prove that for any  $c > 0$ ,  $\|V_{n+c} - V_n\|_\infty \leq \frac{\gamma^n}{1-\gamma} \|V_1 - V_0\|_\infty$

A *Cauchy sequence* is a sequence whose elements become arbitrarily close to each other as the sequence progresses. Formally a sequence  $\{a_n\}$  in metric space  $X$  with distance metric  $d$  is a Cauchy sequence if given an  $\varepsilon > 0$  there exists  $k$  such that if  $m, n > k$  then  $d(a_m, a_n) < \varepsilon$ . Real Cauchy sequences are convergent.

- (c) (2pts) Using this information about Cauchy sequences, argue that the sequence  $V_0, V_1, \dots$  is a Cauchy sequence and is therefore convergent and must converge to some element  $V$  and this  $V$  is a fixed point
- (d) (8pts) Show that this fixed point is unique.

## 4 Frozen Lake MDP [25 pts]

Now you will implement value iteration and policy iteration for the Frozen Lake environment from [OpenAI Gym](#). We have provided custom versions of this environment in the starter code.

- (a) **(coding)** (10 pts) Read through `vi_and_pi.py` and implement `policy_evaluation`, `policy_improvement` and `policy_iteration`. The stopping tolerance (defined as  $\max_s |V_{old}(s) - V_{new}(s)|$ ) is  $\text{tol} = 10^{-3}$ . Use  $\gamma = 0.9$ . Return the optimal value function and the optimal policy.
- (b) **(coding)** (10 pts) Implement `value_iteration` in `vi_and_pi.py`. The stopping tolerance is  $\text{tol} = 10^{-3}$ . Use  $\gamma = 0.9$ . Return the optimal value function and the optimal policy.
- (c) **(written)** (5 pts) Run both methods on the Deterministic-4x4-FrozenLake-v0 and Stochastic-4x4-FrozenLake-v0 environments. In the second environment, the dynamics of the world are stochastic. How does stochasticity affect the number of iterations required, and the resulting policy?