

Question A

In chapter 4, we found that Eval_Trunc could be used in place of Eval. Briefly and clearly explain the following:

- a. Why is this important for Dynamic Programming?

Ans: Since Dynamic Programming includes policy evaluation in each iteration, but only the first couple iterations have large impact, therefore it's more efficient to truncate policy evaluation whenever it converges well enough.

- b. Why is this important for Reinforcement Learning? Explain the significance of this result.

This is important for RL because the idea can be applied to other algorithms, it can significantly save computation cost and its result won't lose convergence guarantees of policy iteration.

Question B:

(10 points) Clearly explain the importance of the Bellman equation in the context of formulating the Dynamic Programming algorithm for prediction? Use equations where needed for clarity.

The bellman equation is important because it allows us to formulate the value of a state and thus we can use dynamic programming to update the evaluation and find policy based on the state value. According to Bellman Equation

$$\begin{aligned} v_*(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma v_*(s') \right], \end{aligned}$$

we can see we can formulate the value of a state by the sum of next states' values and reward multiply by the possibility of reaching next state, or we can also use Bellman's state action equation

$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right], \end{aligned}$$

to evaluate the value of taking a specific action at certain state.

Question C:

Clearly answer the below.

- a. (10 points) Explain how exploration concretely manifests itself in Reinforcement Learning as presented for Dynamic Programming and Monte Carlo methods.

In Dynamic Programming since the environment is fully known the exploration is simply iterate over all states and evaluate them and update policy accordingly. But in Monte Carlo method we do not have complete knowledge over environment and thus we need a proactive exploration step, that is, e-greedy or e-soft, to take an exploration step and explore actions we have not tested yet, then we can improve our policy so that we can find optimal action.

- b. (5 points) Why is exploration needed?

Because if we always follow the policy we have we will not be able to try new actions, and some of the actions may lead to better result than current action and has potential to be the optimal action. Without the exploration we will only be updating existing states' or actions' values repeatedly.

Question D:

(10 points) In Figure 4.1 of the text, observe the $k=3$ value function. Clearly show the calculations that produced row 1 of this value function.

Denote row 1 at $k = 3$ as s_0, s_1, s_2, s_3

$S_0 = 0$ because it's the final state

$$\begin{aligned} S_1 &= \frac{1}{4} * (0 + -1) + \frac{1}{4} * (-1.7 + -1) + \frac{1}{4} * (-2 + -1) + \frac{1}{4} * (-2 + -1) \\ &= -1/4 + \frac{1}{4} * -2.7 + \frac{1}{4} * -3 + \frac{1}{4} * -3 \\ &= -2.425 = -2.4 \end{aligned}$$

$$\begin{aligned} S_2 &= \frac{1}{4} * (-1.7 + -1) + \frac{1}{4} * (-2 + -1) + \frac{1}{4} * (-2 + -1) + \frac{1}{4} * (-2 + -1) \\ &= -2.925 = -2.9 \end{aligned}$$

$$\begin{aligned} S_3 &= \frac{1}{4} * (-1.7 + -1) + \frac{1}{4} * (-2 + -1) + \frac{1}{4} * (-2 + -1) + \frac{1}{4} * (-2 + -1) \\ &= -2.925 = -2.9 \end{aligned}$$

Question E:

I have a policy with 10 possible actions $\{a_0, a_1, \dots, a_9\}$ in state s .

$$q(a_0|s) = q(a_1|s) > q(a_i | s) \text{ for all } a_i \in \{a_2, \dots, a_9\}$$

- a. (5 points) Show the optimization action selection policy for s .

$$P_i(s) = \operatorname{argmax}_{a_i} (q(a_i|s)) \text{ for all } a_i \text{ in } \{a_0, \dots, a_9\} = a_0 \text{ or } a_1$$

- b. (5 points) Show a fair epsilon-greedy action selection policy for s .

$$\text{Let } \epsilon = 0.5$$

$$P(q(a_i | s)) = \epsilon/10 \text{ for all } a_i \text{ in } \{a_2, \dots, a_9\} = 5\%$$

$$P(q(a_0|s)) = P(q(a_2|s)) = \frac{1}{2} * (1 - \epsilon + \epsilon/10 + \epsilon/10) = 30\%$$

- c. (5 points) Show a epsilon-soft action selection policy for s that is distinct from b

$$P(q(a_0|s)) = P(q(a_2|s)) = \dots = P(q(a_9|s)) = 10\%$$

- d. (5 points) What differentiates ϵ -soft and ϵ -greedy action selection policies?

ϵ -greedy action selection policy is an ϵ -soft policy, but a soft max is not required, and ϵ -greedy is also the closest to greedy policy in ϵ -soft policies.

Question F:

In MC, the off-policy methods generate an episode with b which is distinct from π .

- a. (5 points) How are we able to estimate Q_π using b as the episode-generating policy?

We need to make sure b is an ϵ -soft policy, that is all actions is possible to be explored at least occasionally, then we can say that every action taken under π is also taken at least occasionally under b , which is the assumption of converge.

We can calculate the importance-sampling ratio,

$$\rho_{t:T-1} = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k) p(S_{k+1}|S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k|S_k) p(S_{k+1}|S_k, A_k)} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}.$$

then we can estimate Q_π by the ratios and average the result with number of times visited.

$$V(s) = \frac{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1} G_t}{|\mathcal{T}(s)|}.$$

or we can use weighted average

$$V(s) = \frac{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1} G_t}{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1}},$$

- b. (5 points) In computing the returns, why does the iteration process episodes in reverse-time order?

Ans: Since with importance sampling we need next states information, in p110's function we will need next state's weight, value and cumulate sum to compute the return. Forward-time order will require recurrent function which is inefficient.

- c. (10 points) Suppose we wrote the iteration in forward-time order. Show the recurrent (sample update) expression for the computation of the return. That is, how do the first two lines of the inner loop on p 110 change.

Loop For each $t=0, \dots, T-1$, while $w \neq 0$

$$G = \sum_{t'=t+1, \dots, T-1} r^{t'} R_{t'+1}$$

$$C(s_t, a_t) = \pi \frac{\pi(A_{t+1} | S_{t+1})}{b(A_{t+1} | S_{t+1})} \text{ for } t' = t+1, \dots, T-1$$

- d. (5 points) Suppose we wrote the iteration in forward-time order. Show how the inner loop of p 111 would have to be modified to process the episodes.

Loop For each $t=0, \dots, T-1$ which $w \neq 0$

$$G' = \sum_{t'=t+1, \dots, T-1} r^{t'} R_{t'+1}$$

$$w' = \pi \frac{\pi(A_{t+1} | S_{t+1})}{b(A_{t+1} | S_{t+1})} \text{ for } t' = t+1, \dots, T-1$$

$$G = \gamma G' + R_{t+1}$$

$$C(S_t, A_t) = C(S_t, A_t) + w'$$

$$Q(S_t, A_t) = Q(S_t, A_t) + \frac{w}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \operatorname{argmax}_a (Q_{S_t}, a)$$

If $A_t \neq \pi(S_t)$ exit inner loop

Question G:

Suppose we are using on-policy MC on a simulated environment.

- a. (5 points) Explain the two approaches for GPI with on-policy MC.

On-policy MC uses the same policy to generate episode and update policy, whereas off-policy MC uses a behavior policy to generate episodes and aim to estimate the target policy.

- b. (5 points) Which of the two are most appropriate for our problem and why? Why would you choose this over the other?

In my case I would choose off policy, since I think off policy don't have to consider exploration step and thus it's more likely to explore all the states and return a better result

- c. (10 points) Suppose you used the less appropriate version, while your colleague used the more appropriate one. You both have waited a very long time for your algorithms to converge. After converging, you find that your policy is performing worse than your colleagues. Explain the reasons why this has likely occurred, and how would you adjust the algorithm you have chosen (note: you can't use your friend's algorithm) to improve performance.

If I'm using on-policy approach and it performs worse then it's probably because the ϵ -greedy method is not good enough, which means the approaches rarely explores and thus it did not find the optimal policy. I will adjust epsilon value so that the approach explores more often.