

**Midterm Exam**

Name: Ha Manh Bui

**Issued.** Friday, October 27, 2022, 9:00 am EDT

**Due.** Friday, October 27, 2022, 11:59pm EDT

**Instructions (Read carefully!)**

1. This exam is open book and class notes.
2. No online sources besides the ones provided in blackboard are allowed.
3. Where proofs are required, give neat, step-by-step solutions that include all details and are based on first principles.
4. Where computations are required, give neat, step-by-step solutions that include all the details.
5. You **may NOT discuss the problems with anyone.**
6. Upload a unique pdf file that includes a signed copy of this cover page. **Good Luck!**

**Marks**

Question	Maximum	Marks
1	20	
2	20	
3	20	
4	20	
Total	80	

**Statement of Academic Honesty:** For this exam, I make the following truthful statements:

1. I have not received, I have not given, nor will I give or receive, any assistance to another student taking this exam.
2. I will not use any non-instructor approved sources to assist me on an exam.
3. I will not plagiarize someone else's work and turn it in as my own.
4. I understand that acts of academic dishonesty may be penalized to the full extent allowed by the Johns Hopkins University Student Conduct Code, including receiving a failing grade for the course. I recognize that I am responsible for understanding the provisions of the Johns Hopkins University Student Conduct Code as they relate to this academic exercise.

Place your initials here: 

### Problem 1.

Consider an infinite horizon discounted problem ( $0 < \gamma < 1$ ) in which the state space is finite, with  $n$  states, and there are only two possible decisions: stop or continue.

- If at time  $t$  you are at state  $s$ , and you decide to stop, you incur a stopping cost  $g(s)$  and move to a cost-free state, at which you stay forever.
- If at time  $t$  you are at state  $s$ , and you decide to continue, you incur a continuation cost  $c(s)$  and move to a next state  $s'$ , selected at random according to transition probabilities  $p(s, s') = \mathbb{P}(S_{t+1} = s' | S_t = s)$ .

Assume that all  $g(s)$  and  $c(s)$  are nonnegative. We use policy iteration to find the optimal policy that minimizes the discounted total cost  $\sum_{t=1}^{\infty} \gamma^t C_t$  for this MDP.

Note that a stationary policy  $\pi$  is completely specified in terms of the set  $S_\pi$  of states at which the policy decides to stop. Let  $\pi_0, \pi_1, \dots$  be the sequence of policies generated by the policy iteration algorithm, starting with a given arbitrary policy  $\pi_0$ . Suppose that any ties are broken in favor of the “continue” decision when a greedy policy is constructed in the policy improvement phase.

1. Given any value function  $v$ , what is its greedy policy?
2. Given a policy  $\pi$ , what is its value function  $v_\pi$ ? Find a closed form.  
Hint: consider (1) the states at which  $\pi$  stops and (2) the rest, write down their value function separately.
3. Is it always true that for any initial policy  $\pi_0$ , we have  $S_{\pi_1} \subseteq S_{\pi_0}$ ? Give the proof or a counterexample.
4. Suppose  $S_{\pi_0} = \mathcal{S}$ , i.e. the policy  $\pi_0$  stops at every state. It is obvious that  $S_{\pi_1} \subseteq S_{\pi_0}$ . Now show that

$$S_{\pi_2} \subseteq S_{\pi_1} \subseteq S_{\pi_0} .$$

You may use the fact that for any non-negative matrix  $M$  with  $\sum_j M_{ij} \leq 1$ ,  $(I - \gamma M)^{-1}$  is also non-negative.

**Problem 2.**

Given a MDP  $\mathcal{M}(\mathcal{S}, \mathcal{A}, p, r)$ , we assume that the reward function  $r(s, a)$  is known, deterministic and bounded  $|r(s, a)| \leq 1$ .

Suppose we have access to a generative model, which can provide us with a sample  $s' \sim p(\cdot|s, a)$  upon input of any state action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Suppose we call our simulator  $N$  times at each state action pair. Let  $\hat{p}$  be our empirical model, defined as follows:

$$\hat{p}(s'|s, a) = \frac{\text{count}(s', s, a)}{N}, \quad \forall s', (s, a)$$

where  $\text{count}(s', s, a)$  is the number of times the state-action pair  $(s, a)$  transitions to state  $s'$ . As the  $N$  is the number of calls for each state action pair, the total number of calls to our generative model is  $|\mathcal{S}||\mathcal{A}|N$ . We define  $\hat{\mathcal{M}}$  to be the empirical MDP that is identical to the original  $\mathcal{M}$ , except that it uses  $\hat{p}$  instead of  $p$  for the transition model. We let  $\hat{V}^\pi$  denote the value function of policy  $\pi$  under  $\hat{\mathcal{M}}$ .

1. Show that for any policy  $\pi$ , we have

$$\|V^\pi\|_\infty \leq \frac{1}{1-\gamma}.$$

2. Consider a fixed policy  $\pi$ ,

- (a) Define

$$\|P_\pi - \hat{P}_\pi\|_1 := \max_s \sum_{s'} \left| \sum_a p(s'|s, a)\pi(a|s) - \sum_a \hat{p}(s'|s, a)\pi(a|s) \right|,$$

show that

$$\|P_\pi - \hat{P}_\pi\|_1 \leq \max_s \max_a \sum_{s'} |p(s'|s, a) - \hat{p}(s'|s, a)|.$$

- (b) It can be shown that for any fixed  $(s, a)$ -pair, we have

$$\mathbb{P} \left( \sum_{s'} |p(s'|s, a) - \hat{p}(s', a)| \geq \epsilon \right) \leq \exp \left( -N \frac{\epsilon^2}{4|\mathcal{S}|} \right),$$

for any  $N \geq \frac{4|\mathcal{S}|}{\epsilon^2}$ .

Use this to show: Given  $\delta \in (0, 1)$  and  $\epsilon > 0$ . If

$$N \geq \frac{4|\mathcal{S}|}{\epsilon^2} \log \frac{|\mathcal{S}||\mathcal{A}|}{\delta},$$

then with probability  $1 - \delta$ , we have

$$\|P_\pi - \hat{P}_\pi\|_1 \leq \epsilon.$$

3. Show that for a fixed policy  $\pi$ ,

$$\|V^\pi - \hat{V}^\pi\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P_\pi - \hat{P}_\pi)V^\pi\|_\infty$$

4. Finally, show: Given  $\delta \in (0, 1)$  and  $\epsilon > 0$ . If

$$N \geq \frac{4\gamma^2|\mathcal{S}|}{(1-\gamma)^4\epsilon^2} \log \left( \frac{|\mathcal{S}||\mathcal{A}|}{\delta} \right),$$

then with at least probability  $1 - \delta$ , we have

$$\|V^\pi - \hat{V}^\pi\|_\infty \leq \epsilon.$$

**Problem 3.**

For the multi-arms bandit problem we discussed in the class. Suppose that we get return  $G_n$  at  $n$ -th time we do action  $a$ , and we follow the usual assumptions: 1)  $G_n$ 's are independent; 2)  $\mathbb{E}G_n = r$ ,  $n = 1, 2, \dots$ . Let  $Q_{n+1}$  be our estimates of  $r$  after we do action  $a$  the  $n$ -th time, and we have the following update rule

$$Q_{n+1} = Q_n + \alpha_n(G_n - Q_n), \quad Q_1 = 0. \quad (1)$$

We would like to show that, in this case, a weaker condition than  $\sum_n \alpha_n^2 < +\infty$  exists for convergence in expectation.

1. Show that

$$\prod_{k=2}^n \left(1 - \frac{1}{\sqrt{k}}\right) \leq \frac{1}{\sqrt{n}},$$

(Hint:  $\sqrt{k} - 1 \leq \sqrt{k-1}$  when  $k \geq 1$ )

2. Show that for any  $c > 0$ , there exists  $k_0 > 0$  that depends on  $c$ , such that

$$\sqrt{k} - c \leq \sqrt{k-1}, \quad \forall k > k_0.$$

3. Show that for any  $c > 0$ ,

$$\lim_{n \rightarrow \infty} \prod_{k=1}^n \left|1 - \frac{c}{\sqrt{k}}\right| = 0.$$

4. Suppose that we choose in (1)  $\alpha_n = \frac{c}{\sqrt{n}}$  for some  $c > 0$ , show that

$$\lim_{n \rightarrow \infty} |\mathbb{E}Q_n - r| = 0.$$

**Problem 4.** Suppose you have  $s_0$  items that you would like to sell ( $s_0 \in \mathbb{N}$ ). You can sell at most one item per day, and you have a total of  $N$  days in which you can sell them.

At the  $k$ -th day, given the remaining  $s_k$  items you have for sale, you can decide the unit price as  $a_k \in \mathbb{R}_{\geq 0}$ . After deciding the price, the probability of an item being sold is given by:

$$\lambda(a_k) = \beta e^{-a_k} \quad (2)$$

where  $0 < \beta \leq 1$ . The goal in this setting is to find the optimal way to set prices, with the goal of maximizing the expected total reward.

**1.** What is the MDP associated with this problem? Specify the horizon, state space, action space and state-reward transition kernel.

What is the objective function you are trying to maximize?

**2.** Let  $V_k(s_k)$  be the optimal value function for the  $k$ -th stage at state  $s_k$ . Write the Dynamic Programming equations for the optimal policy for all stages.

**3.** Assume the value functions  $V_k(s_k)$  are monotonically nondecreasing on  $s_k$ . Show that for all  $s_k > 0$ , the optimal price-setting policy satisfies:

$$\pi_k^*(s_k) = 1 + V_{k+1}(s_k) - V_{k+1}(s_k - 1) \quad (3)$$

Furthermore, show that:

$$V_k(s_k) = \beta e^{-\pi_k^*(s_k)} + V_{k+1}(s_k). \quad (4)$$

State clearly how you are using the given assumption.

**4.** Show that, for all  $k$ :

- (a)  $V_k(s_k)$  is indeed non-decreasing in  $s_k$ .
- (b) The optimal policy  $\pi_k^*(s_k)$  is non-increasing in  $s_k$ .
- (c) The closed-form solution for  $V_k$  is:

$$V_k(s_k) = \begin{cases} (N - k)\beta e^{-1} & \text{if } s_k \geq N - k, \\ \sum_{i=k}^{N-s_k} \beta e^{-\pi_i^*(s_k)} + s_k \beta e^{-1} & \text{if } 0 < s_k < N - k, \\ 0 & \text{if } s_k = 0 \end{cases}$$

*Hint: Use induction to simultaneously show (a) and (b).*

## Problem 1:

1, Under the infinite horizon discounted setting, by definition, we have two possible action "stop" or "continue", so the Q-function follows:

$$\begin{cases} q_{\pi}(s=s, A=\text{"stop"}) = g(s) \\ q_{\pi}(s=s, A=\text{"continue"}) = c(s) + \gamma \sum_{s'} p(s', s) v(s). \end{cases}$$

$\Rightarrow$  If  $v$ , the greedy policy  $\pi_\theta$ ,  $\pi_\theta(a|s) \geq 0$  is

$$\begin{cases} \text{"stop"} & \text{if } g(s) < c + \gamma \sum_{s'} p(s', s) v(s) \\ \text{"continue"} & \text{else} \end{cases}$$

2, If  $s \in S_\pi$ , then  $v_\pi(s) = g(s)$ . If  $s \notin S_\pi$ , we have

$$v_\pi(s) = c(s) + \gamma \sum_{s' \notin S_\pi} p(s'|s) v_\pi(s') + \gamma \sum_{s' \in S_\pi} p(s'|s) g(s') \quad (1)$$

Let  $P'(s'|s)$  be a  $m \times m$  transition matrix for  $s' \in S_\pi$ , where  $m = |S| - |S_\pi|$  &  $r_\pi(s) = c(s) + \gamma \sum_{s' \in S_\pi} p(s'|s) g(s')$ , then from (1), we get

$$v_\pi(s) = r_\pi(s) + \gamma \sum_{s' \in S_\pi} P'(s'|s) v_\pi(s')$$

$\Rightarrow$  the closed form of  $v_\pi(s)$  is

$$v_\pi = (\mathbb{I} - \gamma P')^{-1} r_\pi. \quad (2)$$

3, It is not always true that for any initial policy  $\pi_0$ , we have  $S_{\pi_1} \subseteq S_{\pi_0}$ .  
 One counterexample is:

Let  $S_{\pi_0} = \{s\}$ , the policy  $\pi_0$  decide to continue at any state.

Assume  $c(s) = 1$

$$\Rightarrow v_{\pi_0}(s) = q_{\pi_0}(s, "continue") = \sum_{t=1}^{\infty} V^t c_t = \frac{1}{1-V}$$

If  $g(s_1) < \frac{1}{1-V}$ , after the policy improvement step,  $\pi_1$  will stop at state  $s_1$ .

$$\Rightarrow S_{\pi_1} \supset S_{\pi_0} \text{ (contradict).}$$

4, consider at the beginning:  $v_{\pi_0}(s) = g(s)$ ,  $\forall s \in S$ .

Assume  $\pi_1(s_1) = "continue"$

$$\Rightarrow g(s_1) = q_{\pi_0}(s_1, "stop") \geq q_{\pi_0}(s_1, "continue").$$

Using the fact that  $P'$  is non-negative matrix  $P'$  with  $\leq P'_{i,j} \leq 1$ ,

$(I - VP')^{-1}$  is also non-negative. Apply for ②, we get

$$v_{\pi_0}(s_1) \geq v_{\pi_1}(s_1) \geq v_{\pi_2}(s_1).$$

$$\Rightarrow q_{\pi_2}(s_1, "continue") \geq q_{\pi_2}(s_1, "stop")$$

$$\Rightarrow g(s_1) \geq q_{\pi_2}(s_1, "continue") \geq q_{\pi_2}(s_1, "stop"),$$

i.e., when  $\pi_1$  change from "stop" to "continue",  $\pi_2$  can not change from "continue" to "stop". Combining with the fact that if  $\pi_2$  doesn't "continue", then  $S_{\pi_2} = S$  &  $\pi_2 = \pi_0 \Rightarrow S_{\pi_2} \subseteq S_{\pi_1}$ . As a result, we achieve

$$S_{\pi_2} \subseteq S_{\pi_1} \subseteq S_{\pi_0}. \quad \square$$

## Problem 2:

1, Under the infinite horizon discounted setting, by definition of the value function,  $\forall s \in S$ , we have:

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t R_t \mid s=s \right].$$

since  $|r(s,a)| \leq 1$ , we get!

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t R_t \mid s=s \right] \leq \sum_{t=0}^{\infty} \gamma^t = \lim_{t \rightarrow \infty} \frac{1-\gamma^t}{1-\gamma} = \frac{1}{1-\gamma}, \forall s \in S.$$

Combining with the fact that  $\|X\|_\infty = \max_i \|X_i\|$ , we obtain:

$$\|V^\pi\|_\infty \leq \frac{1}{1-\gamma} \quad \square \quad (3)$$

2a, Apply triangle inequality,  $\forall s \in S$ , we have:

$$\begin{aligned} & \sum_{s'} \left| \sum_a p(s'|s,a) \pi(a|s) - \sum_a \hat{p}(s'|s,a) \pi(a|s) \right| \\ & \leq \sum_{s'} \sum_a \left| p(s'|s,a) \pi(a|s) - \hat{p}(s'|s,a) \pi(a|s) \right| \\ & = \sum_a \pi(a|s) \sum_{s'} \left| p(s'|s,a) - \hat{p}(s'|s,a) \right| \\ & \leq \sum_a \pi(a|s) \left( \max_a \sum_{s'} \left| p(s'|s,a) - \hat{p}(s'|s,a) \right| \right) \\ & = \max_a \sum_{s'} \left| p(s'|s,a) - \hat{p}(s'|s,a) \right| \end{aligned}$$

$$\Rightarrow \max_a \sum_{s'} \left| \sum_a p(s'|s,a) \pi(a|s) - \sum_a \hat{p}(s'|s,a) \pi(a|s) \right| = \|P_\pi - \hat{P}_\pi\|_1$$

$$\leq \max_s \max_a \sum_{s'} \left| p(s'|s,a) - \hat{p}(s'|s,a) \right| \quad \square \quad (4)$$

2b, Given  $\epsilon > 0$ , by definition & from ⑤, we have

$$\|P_{\pi} - \hat{P}_{\pi}\|_1 \geq \epsilon \Rightarrow \max_s \max_a \sum_{s'} |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon.$$

Using the fact that if  $A \Rightarrow B$ , then  $P(A) \leq P(B)$

$$\Rightarrow P(\|P_{\pi} - \hat{P}_{\pi}\|_1 \geq \epsilon) \leq P\left(\max_s \max_a \sum_{s'} |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon\right) \quad (5)$$

On the other hand, fix a fixed  $(s,a)$ -pair s.t.  $\sum_{s'} |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon$ ,

then  $\max_s \max_a |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon$ , yielding

$$P\left(\exists s,a: \sum_{s'} |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon\right) = P\left(\max_s \max_a |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon\right).$$

$\Rightarrow$  If  $N \geq \frac{4|S|}{\epsilon^2}$ , applying union bound, we get:

$$P\left(\exists s,a: \sum_{s'} |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon\right) \leq \sum_{s,a} P\left(\sum_{s'} |p(s'|s,a) - \hat{p}(s'|s,a)| \geq \epsilon\right).$$

$$\leq \sum_{s,a} \exp(-N \frac{\epsilon^2}{2|S|}) = |S||A| \exp(-N \frac{\epsilon^2}{2|S|}) \leq \delta, \forall \delta \in (0,1).$$

$$\Rightarrow N \geq \frac{4|S|}{\epsilon^2} \log \frac{\delta}{|S||A|}.$$

Combining with ⑤  $\Rightarrow P(\|P_{\pi} - \hat{P}_{\pi}\|_1 \geq \epsilon) \leq \delta$ ,

i.e., if  $N \geq \frac{4|S|}{\epsilon^2} \log \frac{|S||A|}{\delta}$ , then with prob.  $1-\delta$ , we have

$$\|P_{\pi} - \hat{P}_{\pi}\|_1 \leq \epsilon. \square.$$

3, Using the Bellman Eq., we have

$$\begin{aligned}
 \|V^\pi - \hat{V}^\pi\|_\infty &= \|(R + \gamma P_\pi V^\pi) - (R + \gamma \hat{P}_\pi \hat{V}^\pi)\|_\infty \\
 &\leq \gamma \|P_\pi V^\pi - \hat{P}_\pi \hat{V}^\pi\|_\infty \\
 &= \gamma \|P_\pi V^\pi - \hat{P}_\pi V^\pi + \hat{P}_\pi V^\pi - \hat{P}_\pi \hat{V}^\pi\|_\infty \\
 &\leq \gamma \|(P_\pi - \hat{P}_\pi)V^\pi\|_\infty + \gamma \|\hat{P}_\pi(V^\pi - \hat{V}^\pi)\|_\infty. \quad \textcircled{6}
 \end{aligned}$$

Since  $V \in S'$ , we have.

$$\begin{aligned}
 \sum_s p(s'|s) \left( V^\pi(s') - \hat{V}^\pi(s') \right) &\leq \sum_s p(s'|s) \max_{s'} \left( V^\pi(s') - \hat{V}^\pi(s') \right) \\
 &= \max_{s'} \left( V^\pi(s') - \hat{V}^\pi(s') \right),
 \end{aligned}$$

i.e.,  $\|\hat{P}_\pi(V^\pi - \hat{V}^\pi)\|_\infty \leq \|V^\pi - \hat{V}^\pi\|_\infty$ . Applying this result to  $\textcircled{6}$ ,

We obtain:

$$\begin{aligned}
 \|V^\pi - \hat{V}^\pi\|_\infty &\leq \gamma \|(P_\pi - \hat{P}_\pi)V^\pi\|_\infty + \gamma \|\hat{P}_\pi(V^\pi - \hat{V}^\pi)\|_\infty \\
 &\leq \gamma \|(P_\pi - \hat{P}_\pi)V^\pi\|_\infty + \gamma \|V^\pi - \hat{V}^\pi\|_\infty \\
 &\leq \frac{\gamma}{1-\gamma} \|(P_\pi - \hat{P}_\pi)V^\pi\|_\infty \quad \text{D. } \textcircled{7}
 \end{aligned}$$

4, Using the result in ⑦, we have:

$$\begin{aligned}
 \|V^\pi - \hat{V}^\pi\|_\infty &\leq \frac{\gamma}{1-\gamma} \|(\mathbb{P}_\pi - \hat{\mathbb{P}}_\pi) V^\pi\|_\infty \\
 &\leq \frac{\gamma}{1-\gamma} \|\mathbb{P}_\pi - \hat{\mathbb{P}}_\pi\|_1 \|V^\pi\|_\infty \\
 &\leq \frac{\gamma}{(1-\gamma)^2} \|\mathbb{P}_\pi - \hat{\mathbb{P}}_\pi\|_1 \quad (\text{from ⑨})
 \end{aligned}$$

$\Rightarrow$  If  $\|V^\pi - \hat{V}^\pi\|_\infty \geq \varepsilon$ , then  $\|\mathbb{P}_\pi - \hat{\mathbb{P}}_\pi\|_1 \geq \frac{(1-\gamma)^2}{\gamma} \varepsilon$ , this yields

$$\mathbb{P}(\|V^\pi - \hat{V}^\pi\|_\infty \geq \varepsilon) \leq \mathbb{P}(\|\mathbb{P}_\pi - \hat{\mathbb{P}}_\pi\|_1 \geq \frac{(1-\gamma)^2}{\gamma} \varepsilon)$$

$$\leq \delta, \quad \forall \delta \in (0, 1),$$

i.e.,  $\mathbb{P}(\|V^\pi - \hat{V}^\pi\|_\infty \leq \varepsilon) \geq 1 - \delta$ . Applying the result in 2b,

this inequality holds, i.e., with at least prob.  $1 - \delta$ ,

we have  $\|V^\pi - \hat{V}^\pi\|_\infty \leq \varepsilon$  if  $N \geq \frac{\varepsilon \gamma^2 |S|}{(1-\gamma)^2 \delta^2} \log \left( \frac{|S| |\mathcal{A}|}{\delta} \right)$ .

D-

### Problem 3,

1, We have:

$$\prod_{h=2}^n \left(1 - \frac{1}{\sqrt{h}}\right) = \prod_{h=2}^n \left(\frac{\sqrt{h}-1}{\sqrt{h}}\right).$$

Using the fact that  $\sqrt{h}-1 \leq \sqrt{h-1}$  when  $h \geq 1$

$$\Rightarrow \prod_{h=2}^n \left(1 - \frac{1}{\sqrt{h}}\right) \leq \prod_{h=2}^n \left(\frac{\sqrt{h-1}}{\sqrt{h}}\right) = \frac{1}{\sqrt{n}} \quad \square.$$

2, If  $c > 0$ , we have:

$$\sqrt{h} - c \leq \sqrt{h-1}$$

$$\Leftrightarrow \sqrt{h} \leq \sqrt{h-1} + c$$

$$\Leftrightarrow h \leq h-1 + 2c\sqrt{h-1} + c^2$$

$$\Leftrightarrow 1 - c^2 \leq 2c\sqrt{h-1}.$$

If  $1 - c^2 < 0 \Rightarrow$  this inequality holds if  $h \geq 1$  (condition for  $\sqrt{h-1}$  is defined)  
otherwise, i.e.,  $1 - c^2 \geq 0$

$$\Rightarrow h \geq \frac{(1-c^2)^2}{4c^2} + 1$$

$\Rightarrow$  If  $c > 0$ , there is always exists  $k_0 > 0$  that depends on  $c$ ,

$$\text{s.t. } \sqrt{h} - c \leq \sqrt{h-1}, \text{ if } h \geq k_0. \quad \square.$$

3, we have

$$\begin{aligned}
 & \lim_{n \rightarrow \infty} \prod_{i=1}^n \left| 1 - \frac{c}{\sqrt{h}} \right| = \lim_{n \rightarrow \infty} \prod_{i=1}^n \left| \frac{\sqrt{h} - c}{\sqrt{h}} \right| = \lim_{n \rightarrow \infty} \prod_{k=1}^{h_0} \left| \frac{\sqrt{h} - c}{\sqrt{h}} \right| \cdot \prod_{k=h_0+1}^n \left| \frac{\sqrt{h} - c}{\sqrt{h}} \right| \\
 &= \lim_{n \rightarrow \infty} \prod_{k=1}^{h_0} \left| \frac{\sqrt{h} - c}{\sqrt{h}} \right| \left| \frac{\sqrt{h_0+1} - c}{\sqrt{n}} \right| \quad (\text{from 3.1, 3.2}) \\
 &= 0 \quad (\text{due to } \prod_{k=1}^{h_0} \left| \frac{\sqrt{h} - c}{\sqrt{h}} \right| < \infty \text{ when } \sqrt{h} - c \leq \sqrt{h})
 \end{aligned}$$

4, By definition & assumption (1) & (2), we have

$$\begin{aligned}
 \mathbb{E}\{Q_{n+1}\} &= \mathbb{E}\{Q_n + \alpha_n(G_n - Q_n)\} = \mathbb{E}\{Q_n\} + \alpha_n(\mathbb{E}\{G_n\} - \mathbb{E}\{Q_n\}) \\
 &= \mathbb{E}\{Q_n\} + \alpha_n(r - \mathbb{E}\{Q_n\}).
 \end{aligned}$$

Consider  $|\mathbb{E}\{Q_{n+1}\} - r|$ , we have

$$\begin{aligned}
 |\mathbb{E}\{Q_{n+1}\} - r| &= |\mathbb{E}\{Q_n\} + \alpha_n(r - \mathbb{E}\{Q_n\}) - r| = |(1 - \alpha_n)\mathbb{E}\{Q_n\} - (1 - \alpha_n)r| \\
 &= |1 - \alpha_n| |\mathbb{E}\{Q_n\} - r| = |1 - \alpha_n|^{n+2} |\mathbb{E}\{Q_n\} - r|.
 \end{aligned}$$

When  $n \rightarrow \infty$ , if  $\alpha_n = \frac{c}{\sqrt{n}}$ , then  $0 < \alpha_n < 2$  holds. This yields

$$\lim_{|1 - \alpha_n| < 1} |1 - \alpha_n|^{n+2} |\mathbb{E}\{Q_n\} - r| = 0 \Rightarrow \lim_{n \rightarrow \infty} |\mathbb{E}\{Q_n\} - r| = 0 \quad \square.$$

## Problem 5.

2, we can formulate by DP with Fixed-Horizon  $T = N$  stages corresponding to  $N$  days in which we can sell the items.

$\Rightarrow \forall h \in \{0, N-1\}$ , we have:

- †  $S_h \in S$  be the state r.v. at stage  $h \Rightarrow S = \{0, \dots, s_0; s_0 \in \mathbb{N}\}$
- †  $A_h \in A$  be the action r.v. at stage  $h$  represent for price  $a_h \Rightarrow A = \mathbb{R}_{\geq 0}$
- †  $R_{h+1} \in \mathbb{R}_{\geq 0}$  be the revenue reward of  $A_h$  at  $h$ .

Since we can sell at most 1 item per day with prob.  $\lambda(a_h)$ , so the state-reward transition kernel  $p(R_{h+1}, S_{h+1} | S_h, A_h)$  is

$$\begin{cases} p(R_{h+1} = a_h, S_{h+1} = s_{h+1} | S_h = s_h, A_h = a_h) = \lambda(a_h) \\ p(R_{h+1} = 0, S_{h+1} = s_h | S_h = s_h, A_h = a_h) = 1 - \lambda(a_h) \end{cases}$$

$\Rightarrow$  The objective function we want to maximize is the cumulative expected return!

$$V_{\pi}(s_0) = \mathbb{E}_{\pi} \left\{ \sum_{h=0}^{N-1} R_{h+1} | S = s_0 \right\}.$$

2, Using the backward recursion, we have the optimal policy  $\pi^*_{\pi}(a|s) \forall a$  for all stages are:

$$\begin{aligned} a &\in \arg\max_a \sum_{s' \in S} p(r, s' | s, a) \left[ r + V_{h+1}(s') \right] \\ &= \arg\max_{a_h \in \mathbb{R}_{\geq 0}} \left\{ \lambda(a_h) \left[ a_h + V_{h+1}(s_h - 1) \right] + (1 - \lambda(a_h)) V_{h+1}(s_h) \right\} \\ &= \arg\max_{a_h \in \mathbb{R}_{\geq 0}} \left\{ \beta e^{-\alpha a_h} \left[ a_h + V_{h+1}(s_h - 1) \right] + (1 - \beta e^{-\alpha a_h}) V_{h+1}(s_h) \right\} \end{aligned}$$

3, Consider  $S_h \geq 0$ , let  $f(d_h) = \beta e^{-d_h} (d_h + V_{h+1}(S_h - 1)) + (1 - \beta e^{-d_h}) V_{h+1}(S_h)$

since  $\Pi_h^*(S_h)$  is the optimal price-setting policy, then.

$$\Pi_h^*(S_h) = \underset{d_h \in \mathbb{R}_{\geq 0}}{\arg \max} f(d_h).$$

Take derivative  $\frac{\partial f(d_h)}{\partial d_h}$ , we have  $\frac{\partial f(d_h)}{\partial d_h} = 0$  when

$$e^{-d_h} \{1 + d_h - V_{h+1}(S_h - 1) + V_{h+1}(S_h)\} = 0$$

Since  $\ln(e) \neq 0$  &  $V_h(S_h)$  are monotonically non-decreasing on  $S_h$ , i.e.,  $V_{h+1}(S_h - 1) \leq V_{h+1}(S_h)$ , we get

$$d_h = 1 + V_{h+1}(S_h) - V_{h+1}(S_h - 1) \in \mathbb{R}_{\geq 0}.$$

$$\Rightarrow \Pi_h^*(S_h) = 1 + V_{h+1}(S_h) - V_{h+1}(S_h - 1). \quad (8)$$

In addition, we have

$$\begin{aligned} V_h(S_h) &= \max_{d_h \in \mathbb{R}_{\geq 0}} \left\{ \sum_{S' \in I} p(r, s' | S_h, d_h) (r + V_{h+1}(s')) \right\} \\ &= \max_{d_h \in \mathbb{R}_{\geq 0}} \{ \beta e^{-d_h} \{d_h + V_{h+1}(S_h - 1)\} + (1 - \beta e^{-d_h}) V_{h+1}(S_h) \} \\ &= \beta e^{-d_h} (-1 - V_{h+1}(S_h) + V_{h+1}(S_h - 1)) + V_{h+1}(S_h) \end{aligned}$$

$$(8) \Rightarrow \beta e^{-\Pi_h^*(S_h)} + V_{h+1}(S_h). \quad \square \quad (9)$$

4a, Since

$$V_h(S_h) = \max_{\alpha_h \in K_{\geq 0}} \{ \beta e^{-\alpha_h} [a_h + V_{h+1}(S_{h-1})] + (1 - \beta e^{-\alpha_h}) V_{h+2}(S_h) \}$$

and

$$V_h(S_h-1) = \max_{\alpha_h \in K_{\geq 0}} \{ \beta e^{-\alpha_h} [a_h + V_{h+1}(S_{h-2})] + (1 - \beta e^{-\alpha_h}) V_{h+2}(S_{h-1}) \}$$

$$\begin{aligned} \Rightarrow \forall \alpha_h \in K_{\geq 0}, \beta \in (0, 1), \text{ we obtain} \\ V_h(S_h) - V_h(S_h-1) &\geq \max \{ \beta e^{-\alpha_h} [a_h + V_{h+1}(S_{h-1})] + (1 - \beta e^{-\alpha_h}) V_{h+2}(S_h) \\ &\quad - \beta e^{-\alpha_h} [a_h + V_{h+1}(S_{h-2})] - (1 - \beta e^{-\alpha_h}) V_{h+2}(S_{h-1}) \} \\ &= \max_{\alpha_h \in K_{\geq 0}} \{ \beta e^{-\alpha_h} (2V_{h+1}(S_h-1) - V_{h+1}(S_h) - V_{h+2}(S_h) + V_{h+2}(S_h-1)) \} \\ &\geq 0, \forall h. \Rightarrow V_h(S_{h-1}) \leq V_h(S_h), \forall h, \end{aligned}$$

i.e.,  $V_h(S_h)$  is indeed non-decreasing on  $S_h$ .  $\square$ .

4b, Using the result from (8), we have.

$$\Pi_h^+(S_h) = 1 + V_{h+1}(S_h) - V_{h+2}(S_h-1)$$

and

$$\Pi_h^+(S_h-1) = 1 + V_{h+1}(S_h-1) - V_{h+2}(S_h-2).$$

since  $V_h(S_h)$  is non-decreasing on  $S_h$ , i.e.,

$$V_{h+2}(S_h-2) \leq V_{h+2}(S_h-1) \leq V_{h+2}(S_h),$$

we obtain

$$\Pi_h^+(S_h-1) \geq \Pi_h^+(S_h),$$

i.e.,  $\Pi_h^+(S_h)$  is non-increasing on  $S_h$ .  $\square$ .

9c, consider the term

$$V_h(S_h) = \max_{a_h \in R_{\geq 0}} \left\{ \beta e^{-\alpha_h} (a_h + V_{h+1}(S_h - 1)) + (1 - \beta e^{-\alpha_h}) V_{h+1}(S_h) \right\}$$

IF  $S_h = 0 \Rightarrow$  we don't have any item to sell, i.e.,  $\Pi_h(a_h | S_h=0) = 0$

$$\Rightarrow V_h(S_h=0) = 0.$$

IF  $S_h \geq N-h \Rightarrow$  we have  $\# \text{items} \geq \# \text{days}$  we can sell it.  
Due to each day, we can only sell at most 1 item

$$\Rightarrow V_h(S_h) = \mathbb{E}_{\Pi^*} \left[ \sum_{h=t}^N R_{h+1} | S=S_h \right] = (N-h) \beta e^{-1}. \quad (10)$$

IF  $0 < S_h < N-h$ , since (9), we have

$$V_h(S_h) = \beta e^{-\Pi^*(S_h)} + V_{h+1}(S_h),$$

We can set the optimal policy  $\Pi_i^*(S_h)$  until  $S_h \geq N-h$ , combining with the result from (10), we obtain:

$$V_h(S_h) = \sum_{i=h}^{N-h} \beta e^{-\Pi_i^*(S_h)} + S_h \beta e^{-1}.$$

As a result, the closed form solution for  $V_h$  is,

$$V_h(S_h) = \begin{cases} (N-h) \beta e^{-1} & \text{if } S_h \geq N-h \\ \sum_{i=h}^{N-h} \beta e^{-\Pi_i^*(S_h)} + S_h \beta e^{-1} & \text{if } 0 < S_h < N-h \\ 0 & \text{if } S_h = 0 \end{cases}$$

D-