
CS 234 Midterm - Winter 2018-19

****Do not turn this page until you are instructed to do so.**

Instructions

Please answer the following questions to the best of your ability. Read all the questions first before answering. You have **80** minutes to complete the exam. The exam is closed-book and closed-internet. Additionally use of all electronic items during the exam is prohibited. However, you may use one one-sided letter-sized page of notes as reference. All of the intended answers can be written well within the space provided. Good luck!

Stanford University Honor Code

The following is a statement of the Stanford University Honor Code:

1. The Honor Code is an undertaking of the students, individually and collectively: (1) that they will not give or receive aid in examinations; that they will not give or receive unpermitted aid in class work, in the preparation of reports, or in any other work that is to be used by the instructor as the basis of grading; (2) that they will do their share and take an active part in seeing to it that others as well as themselves uphold the spirit and letter of the Honor Code.
2. The faculty on its part manifests its confidence in the honor of its students by refraining from proctoring examinations and from taking unusual and unreasonable precautions to prevent the forms of dishonesty mentioned above. The faculty will also avoid, as far as practicable, academic procedures that create temptations to violate the Honor Code.
3. While the faculty alone has the right and obligation to set academic requirements, the students and faculty will work together to establish optimal conditions for honorable academic work.

By signing your name below, you acknowledge that you have abided by the Stanford Honor Code while taking this exam.

Signature:

Name:

SUNet ID:

Grading (For Midterm Graders Only)

Question #	1	2	3	4	Total
Maximum Points	25	10	20	25	80
Student Grade					

Question 1 – True/False and Short Answers

[25 pts]

(A) Circle True or False. Provide a one-sentence justification for each question. [2 pts each]

1. **TRUE** False An optimal policy of an infinite-horizon MDP is always stationary (i.e. it does not vary with the time step).

Solution: True. Infinite-horizon MDPs are guaranteed to have a stationary optimal policy.

2. True **FALSE** When using a linear value function approximation, both Monte Carlo and TD(0) policy evaluations achieve the minimum mean squared error possible at convergence.

Solution: False. They may converge to different value functions unless the environment is Markovian and both of them achieve zero mean squared error (which is not necessarily true due to error from function approximation).

3. **TRUE** False All discounted MDPs with finite states and action spaces with bounded reward (i.e. $|R(s, a)| \leq R_{\max}$ for some $R_{\max} \geq 0$) have bounded returns.

Solution: True. Since the reward is bounded and $\gamma < 1$, return is a bounded geometric series.

4. **TRUE** False Monte Carlo policy evaluation does not require the environment to be Markovian.

Solution: True. Since MC policy evaluation does not bootstrap, it does not require the environment to be Markovian.

5. True **FALSE** Double Q-learning helps training by directly estimating the advantage function.

Solution: False. Double Q-learning helps training by reducing maximization bias.

(B) For the next few questions answer in 2 - 3 sentences.

[3 pts each]

1. Explain why Q-learning is off-policy and SARSA is on-policy.

Solution: Q learning's target policy is the optimal policy (i.e. it always tries to approximate the optimal Q-values). In SARSA, the target policy is the same as behavior policy.

2. What is the challenge with using DQN when the action space is continuous?

Solution: While we can still represent the value function by requiring both state and action as arguments, it becomes impossible or very difficult to take the maximum over all actions for an arbitrary Q network.

3. If you need an algorithm that will definitely reach a local optimum, should you use DQN or REINFORCE, and why? Assume that the state space is very high-dimensional and cannot be stored in memory.

Solution: REINFORCE. For such large state space, we have to use function approximation. In such scenario, REINFORCE will converge to a local optimum because it directly optimized the policy using SGD. However, DQN does not and lacks convergence guarantee.

4. State two advantages of using value function approximation compared to a tabular representation.

Solution:

- Generalizability to unseen states
- More efficient memory usage (i.e. more compact value function representation)

5. Suppose we gathered infinitely many trajectories for an MDP with finite state space. If we run tabular Monte Carlo and TD(0) policy evaluations, would the two methods converge to the same value function? Justify your answer.

Solution: Since the environment is Markovian (MDP), both methods are guaranteed to converge to the same value function that achieves zero MSVE (mean squared value error).

Question 2 – Value Iteration

[10 pts]

Suppose we have an MDP $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$, where \mathcal{S} is a finite state space and \mathcal{A} is a finite action space. This MDP is defined such that $R(s, a) \geq 0$ for all state action pairs (s, a) . Furthermore, suppose that for every state $s \in \mathcal{S}$, there is some action $a_s \in \mathcal{A}$ such that $P(s' = s \mid s, a_s) \geq p$, where $0 \leq p \leq 1$ is some constant probability.

Consider performing value iteration on this MDP. Let $V_t(s)$ be the value of state s after t iterations. We initialize to $V_0(s) = 0$ for all states s .

(a) Prove that for all states s and $t \geq 0$, $V_{t+1}(s) \geq p\gamma V_t(s)$. (5 pts)

Solution: Consider an arbitrary state $s \in \mathcal{S}$ and an arbitrary iteration $t \geq 0$. From the Value Iteration algorithm, we have:

$$\begin{aligned} V_{t+1}(s) &= \max_{a \in \mathcal{A}} R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s, a) V_t(s') \\ &\geq R(s, a_s) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s, a_s) V_t(s') \\ &\geq R(s, a_s) + \gamma p V_t(s) \\ &\geq \gamma p V_t(s) \end{aligned}$$

Where the 3rd line follows by the fact that $\forall s : V_t(s) \geq 0$ because $R(s, a) \geq 0$ and initialization is 0. Last line uses $R(s, a) \geq 0$. This completes the proof.

(b) During a single iteration of the Value Iteration algorithm, we typically iterate over the states in \mathcal{S} in some order to update $V_t(s)$ to $V_{t+1}(s)$ for all states s . Is it possible to do this iterative process in parallel? Explain why or why not. (5 pts)

Solution: If we think of $V_t(\cdot)$ and $V_{t+1}(\cdot)$ as vectors of length $|\mathcal{S}|$, then

$$V_{t+1}(s) = \max_{a \in \mathcal{A}} R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s, a) V_t(s') \quad (1)$$

implies that $V_{t+1}(\cdot)$ depends only on values in $V_t(\cdot)$. Therefore, it is possible to parallelize the calculations in Equation 1.

Question 3 – Q-Learning and Function Approximation [20 pts]

Consider an episodic, deterministic chain MDP with $n = 7$ states assembled in a line.

The possible actions are $a \in \{-1, 1\}$, and the transition function is deterministic such that $s' = s + a$. Note that as an exception, taking $a = -1$ from $s = 1$ keeps us in $s = 1$, and taking $a = 1$ from $s = 7$ keeps us in $s = 7$.

We have a special goal state, $g = 4$, such that taking any action **from** g ends the episode with a reward of $r = 0$. From all other states, any action incurs a reward of $r = -1$. We let $\gamma = 1$.

The chain MDP is pictured in Figure 1, with the goal state s_4 shaded in.

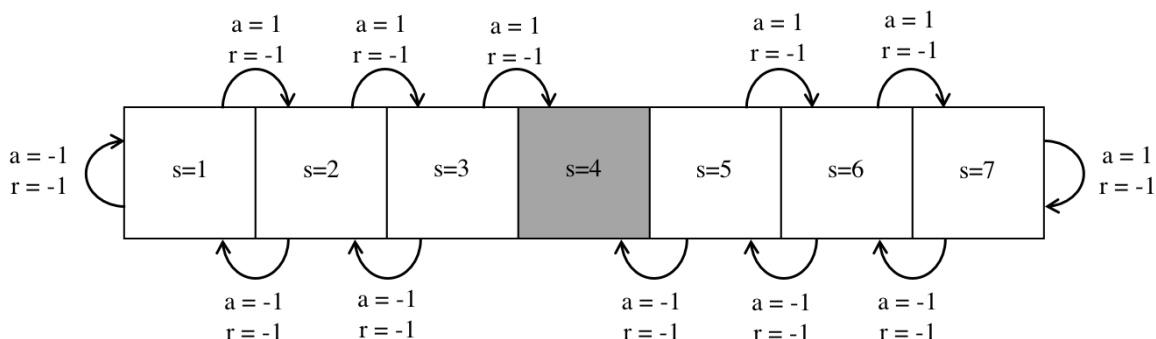


Figure 1: Chain MDP

By inspection, we see that $V^*(s) = -|s - 4|$.

(a) We would like to perform tabular Q-learning on this chain MDP. Suppose we observe the following 4 step trajectory (in the form $(state, action, reward)$):

$$(3, -1, -1), (2, 1, -1), (3, 1, -1), (4, 1, 0)$$

Suppose we initialize all Q values to 0. Use the tabular Q-learning update to give updated values for

$$Q(3, -1), Q(2, 1), Q(3, 1)$$

assuming we process the trajectory in the order given from left to right. Use the learning rate $\alpha = \frac{1}{2}$. (5 pts)

Solution: Given a tuple (s, a, r, s') , we use the update equation:

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

Using this equation with $\alpha = \frac{1}{2}, \gamma = 1$, we have:

$$Q(3, -1) \leftarrow 0 + \frac{1}{2}(-1 + \max_{a'} Q(2, a')) = \frac{1}{2}(-1 + 0) = -\frac{1}{2}$$

$$Q(2, 1) \leftarrow 0 + \frac{1}{2}(-1 + \max_{a'} Q(3, a')) = \frac{1}{2}(-1 + 0) = -\frac{1}{2}$$

$$Q(3, 1) \leftarrow 0 + \frac{1}{2}(-1 + \max_{a'} Q(4, a')) = \frac{1}{2}(-1 + 0) = -\frac{1}{2}$$

Now, we are interested in performing linear function approximation in conjunction with Q-learning. In particular, we have a weight vector

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \in \mathbb{R}^3$$

Given some state s and action $a \in \{-1, 1\}$, the featurization of this state, action pair is: $\begin{bmatrix} s \\ a \\ 1 \end{bmatrix}$

To approximate the Q-values, we compute

$$\hat{q}(s, a; w) = w^T \begin{bmatrix} s \\ a \\ 1 \end{bmatrix} = w_0 * s + w_1 * a + w_2$$

Given the parameters w and a single sample (s, a, r, s') , the loss function we will minimize is

$$J(w) = (r + \gamma \max_{a'} \hat{q}(s', a'; w^-) - \hat{q}(s, a; w))^2$$

where $\hat{q}(s', a'; w^-)$ is a target network parametrized by fixed weights w^- .

(b) Suppose we currently have a weight vectors

$$w = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}, \quad w^- = \begin{bmatrix} 1 \\ -1 \\ -2 \end{bmatrix}$$

and we observe a sample $(s = 2, a = -1, r = -1, s' = 1)$

Perform a single gradient update to the parameters w given this sample. Use the learning rate $\alpha = \frac{1}{4}$. Write out the gradient $\nabla_w J(w)$ as well as the new parameters w' . Show all work. (10 pts)

Solution:

$$\begin{aligned} \nabla_w J(w) &= -2(r + \gamma \max_{a'} \hat{q}(s', a'; w^-) - \hat{q}(s, a; w)) \nabla_w \hat{q}(s, a; w) \\ &= -2 \left(r + \max_{a'} (w^-)^T \begin{bmatrix} s' \\ a' \\ 1 \end{bmatrix} - w^T \begin{bmatrix} s \\ a \\ 1 \end{bmatrix} \right) \begin{bmatrix} s \\ a \\ 1 \end{bmatrix} \end{aligned}$$

Using this, the parameter update with a single sample (s, a, r, s') is:

$$\begin{aligned} w' &\leftarrow w - \alpha \nabla_w J(w) \\ &= w + \frac{1}{2} \left(r + \max_{a'} (w^-)^T \begin{bmatrix} s' \\ a' \\ 1 \end{bmatrix} - w^T \begin{bmatrix} s \\ a \\ 1 \end{bmatrix} \right) \begin{bmatrix} s \\ a \\ 1 \end{bmatrix} \end{aligned}$$

Using the sample $(2, -1, -1, 1)$ and the particular values of w and w^- yields:

$$\begin{aligned}
 w' &\leftarrow \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} + \frac{1}{2} \left(-1 + \max_{a'} \begin{bmatrix} 1 \\ -1 \\ -2 \end{bmatrix}^T \begin{bmatrix} 1 \\ a' \\ 1 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}^T \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} \right) \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} + \frac{1}{2} \left(-1 + \max_{a'} (1 - a' - 2) - (-2 - 1 + 1) \right) \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 0 \\ 1/2 \\ 3/2 \end{bmatrix}
 \end{aligned}$$

Note: Some students may have written the parameter update as:

$$w' \leftarrow w - \frac{1}{2} \alpha \nabla_w J(w)$$

This is fine, and the subsequent answer they should obtain is:

$$w' \leftarrow \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} + \frac{1}{4} \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1/2 \\ 3/4 \\ 5/4 \end{bmatrix}$$

Both answers are acceptable as long as the gradient update form used was written explicitly, and the subsequent w' matches that update.

(c) The optimal Q function $Q^*(s, a)$ is exactly representable by some neural network architecture \mathcal{N} . Suppose we perform Q-learning on this MDP using the architecture \mathcal{N} to represent the Q-values. Suppose we randomly initialize the weights of a neural net with architecture \mathcal{N} and collect infinitely many samples using an exploration strategy that is greedy in the limit of infinite exploration (GLIE). Are we guaranteed to converge to the optimal Q function $Q^*(s, a)$? Explain your answer. (5 pts)

Solution: Because this method of Q-learning involves function approximation, bootstrapping, and off-policy training, (the "deadly triad" according to Sutton & Barto), we are not guaranteed to converge to anything, which includes no guarantee of converging to the optimal Q function.

Question 4 – Certainty Equivalence Estimate

[25 pts]

Consider the discounted ($\gamma < 1$), infinite-horizon MDP $M = (S, A, P, R, \gamma)$, where we do not know the true reward function $R(s, a) \in \mathbb{R}$ and state transition probabilities $P(s'|s, a)$. With a slight abuse of notation, we will also write $P(s, a) \in \mathbb{R}^{|S|}$ to denote the vector of transition probabilities of size $|S|$, whose values sum up to 1. For a given policy π , V_M^π denotes the value function of π in M .

Fortunately, we are given estimates of R and P , namely, \hat{R} and \hat{P} , respectively, with the following properties:

$$\max_{s,a} |\hat{R}(s, a) - R(s, a)| < \epsilon_R$$

$$\max_{s,a} \|\hat{P}(s, a) - P(s, a)\|_1 < \epsilon_P$$

$$\text{where } \|\cdot\|_1 \text{ is the } L_1 \text{ norm for } \mathbf{x} \in \mathbb{R}^n: \|\mathbf{x}\| = \sum_{i=1}^n |x_i|$$

We can then define the *approximate MDP* $\hat{M} = (S, A, \hat{P}, \hat{R}, \gamma)$ and let $V_{\hat{M}}^\pi$ be the value function of policy π in \hat{M} . For simplicity, assume that both reward functions are bounded within $[0, R_{\max}]$.

(a) Show that for all policies π and states s , the value function is bounded as follows:

$$0 \leq V^\pi(s) \leq \frac{R_{\max}}{1 - \gamma} \quad (2)$$

(5 pts)

Solution: All state values are at least 0 because reward is nonnegative. Also reward is at most R_{\max} , so state values are bounded by the one that has maximal (R_{\max}) reward at every time step. Thus for any policy π and state s , we have

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \leq \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_{\max} \right] = R_{\max} \sum_{t=0}^{\infty} \gamma^t = \frac{R_{\max}}{1 - \gamma}$$

(b) Let V be the value function of an arbitrary policy. Show that for all state-action pairs (s, a) , the following holds:

$$\sum_{s' \in S} \left[V(s') \left(\hat{P}(s'|s, a) - P(s'|s, a) \right) \right] \leq \epsilon_P \frac{R_{\max}}{2(1 - \gamma)} \quad (3)$$

(7 pts)

Solution: Since $\sum_{s' \in S} \hat{P}(s'|s, a) = \sum_{s' \in S} P(s'|s, a) = 1$,

we can subtract $\frac{R_{\max}}{2(1-\gamma)} \sum_{s' \in S} \left(\hat{P}(s'|s, a) - P(s'|s, a) \right)$ without changing the value of the expression:

$$\begin{aligned}
& \sum_{s' \in S} \left[V(s') \left(\hat{P}(s'|s, a) - P(s'|s, a) \right) \right] \\
&= \sum_{s' \in S} \left[V(s') \left(\hat{P}(s'|s, a) - P(s'|s, a) \right) - \frac{R_{\max}}{2(1-\gamma)} \left(\hat{P}(s'|s, a) - P(s'|s, a) \right) \right] \\
&= \sum_{s' \in S} \left[\left(V(s') - \frac{R_{\max}}{2(1-\gamma)} \right) \left(\hat{P}(s'|s, a) - P(s'|s, a) \right) \right] \\
&\leq \sum_{s' \in S} \left| V(s') - \frac{R_{\max}}{2(1-\gamma)} \right| \left| \hat{P}(s'|s, a) - P(s'|s, a) \right| \\
&\leq \frac{R_{\max}}{2(1-\gamma)} \sum_{s' \in S} \left| \hat{P}(s'|s, a) - P(s'|s, a) \right| \quad \text{from part (a)} \\
&= \frac{R_{\max}}{2(1-\gamma)} \left\| \hat{P}(s, a) - P(s, a) \right\|_1 \leq \epsilon_P \frac{R_{\max}}{2(1-\gamma)}
\end{aligned}$$

(c) Using part (b), show that for any deterministic policy π , the error in state value is bounded as follows:

$$\|V_M^\pi - V_M^\pi\|_\infty \leq \frac{\epsilon_R}{1-\gamma} + \gamma \epsilon_P \frac{R_{\max}}{2(1-\gamma)^2} \quad (4)$$

where $\|\cdot\|_\infty$ is the L_∞ norm: $\|V_M^\pi - V_M^\pi\|_\infty = \max_s |V_M^\pi(s) - V_M^\pi(s)|$. (8 pts)

Solution: For any given state s let $a = \pi(s)$. From the definition of value function, we get:

$$\begin{aligned}
& V_M^\pi(s) - V_M^\pi(s) \\
&= \left(\hat{R}(s, a) + \gamma \sum_{s' \in S} \hat{P}(s'|s, a) V_M^\pi(s') \right) - \left(R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V_M^\pi(s') \right) \\
&= \left(\hat{R}(s, a) - R(s, a) \right) + \gamma \sum_{s' \in S} \left(\hat{P}(s'|s, a) V_M^\pi(s') - P(s'|s, a) V_M^\pi(s') \right) \\
&\leq \epsilon_R + \gamma \sum_{s' \in S} \left(\hat{P}(s'|s, a) V_M^\pi(s') - P(s'|s, a) V_M^\pi(s') + P(s'|s, a) V_M^\pi(s') - P(s'|s, a) V_M^\pi(s') \right) \\
&= \epsilon_R + \gamma \sum_{s' \in S} \left(\hat{P}(s'|s, a) - P(s'|s, a) \right) V_M^\pi(s') + \gamma \sum_{s' \in S} P(s'|s, a) \left(V_M^\pi(s') - V_M^\pi(s') \right) \\
&\leq \epsilon_R + \gamma \epsilon_P \frac{R_{\max}}{2(1-\gamma)} + \gamma \sum_{s' \in S} P(s'|s, a) \left\| V_M^\pi - V_M^\pi \right\|_\infty \quad \text{from part (b)} \\
&= \epsilon_R + \gamma \epsilon_P \frac{R_{\max}}{2(1-\gamma)} + \gamma \left\| V_M^\pi - V_M^\pi \right\|_\infty
\end{aligned}$$

Since this is true for all states, we can take the L_∞ norm of the left hand side:

$$\left\| V_M^\pi - V_M^\pi \right\|_\infty \leq \epsilon_R + \gamma \epsilon_P \frac{R_{\max}}{2(1-\gamma)} + \gamma \left\| V_M^\pi - V_M^\pi \right\|_\infty$$

Rearranging the expression and scaling by $\frac{1}{1-\gamma}$ gives the desired result.

(d) Let π^* and $\hat{\pi}^*$ be (deterministic) optimal policies in M and \widehat{M} , respectively. Using eq. (4), show that the following bound holds for all $s \in S$:

$$V_M^{\pi^*} - V_M^{\hat{\pi}^*} \leq 2 \left(\frac{\epsilon_R}{1-\gamma} + \gamma \epsilon_P \frac{R_{\max}}{2(1-\gamma)^2} \right) \quad (5)$$

(5 pts)

Note: What this means is that if we use a policy that is optimal in \widehat{M} , the amount we lose in value compared to the true optimal policy is bounded by the error in our approximation. In other words, the better the approximation, the closer we are to the true optimal policy!

Solution: From the definition of value function, we get:

$$\begin{aligned} & V_M^{\pi^*}(s) - V_M^{\hat{\pi}^*}(s) \\ &= V_M^{\pi^*}(s) - V_{\widehat{M}}^{\pi^*}(s) + V_{\widehat{M}}^{\pi^*}(s) - V_M^{\hat{\pi}^*}(s) \\ &\leq \left\| V_M^{\pi^*} - V_{\widehat{M}}^{\pi^*} \right\|_{\infty} + V_{\widehat{M}}^{\pi^*}(s) - V_M^{\hat{\pi}^*}(s) \\ &\leq \left\| V_M^{\pi^*} - V_{\widehat{M}}^{\pi^*} \right\|_{\infty} + V_{\widehat{M}}^{\hat{\pi}^*}(s) - V_M^{\hat{\pi}^*}(s) \quad \hat{\pi}^* \text{ is optimal in } \widehat{M} \\ &\leq \left\| V_M^{\pi^*} - V_{\widehat{M}}^{\pi^*} \right\|_{\infty} + \left\| V_{\widehat{M}}^{\hat{\pi}^*} - V_M^{\hat{\pi}^*} \right\|_{\infty} \\ &\leq 2 \left(\frac{\epsilon_R}{1-\gamma} + \gamma \epsilon_P \frac{R_{\max}}{2(1-\gamma)^2} \right) \quad \text{from part (c)} \end{aligned}$$