# CMPSCI 687 Homework 3

Due October 29, 2019, 11:55pm Eastern Time

**Instructions:** Collaboration is not allowed on any part of this assignment. Submissions must be typed (hand written and scanned submissions will not be accepted). You must use LaTeX. The assignment should be submitted as two documents: a .pdf with your written answers and a single .cpp file as described in the programming portion.

## Part One: Written (50 Points Total)

1. (2 Points) One day while working in the engineering department of the Starship Enterprise, your friend Geordi comes to you with an idea. He points out that the warp core (engine) uses a reinforcement learning algorithm to regulate its temperature. He hypothesizes the the value function that it uses would be easier to represent and/or faster to approximate in two distinct parts: one that estimates the value of a state given that the next state is safe (within desirable thresholds), and another that estimates the value of a state given that the next state is not safe. Working with Geordi, who of course uses the notation from this class, you decide to define $\mathcal{X}$ to be the set of safe states, and $\mathcal{X}^{\complement}$ to be the set of unsafe states, i.e., $\mathcal{X}^{\complement} = \mathcal{S} \setminus \mathcal{X}$.[1] In order to continue, you and Geordi decide to establish some notation. Specifically, you want to define $v_{\mathcal{Y}}^{\pi}(s)$ to be the expected discounted return given that the agent begins in state $s$, follows policy $\pi$, and the next state (but not necessarily the states after the next state) happens to be in $\mathcal{Y}$. Give a mathematical definition for $v_{\mathcal{Y}}^{\pi}$ like our definition for $v^{\pi}$:

   Ans 1.

   $$v_{\mathcal{Y}}^{\pi}(s) := \mathbf{E}[G_t | S_t = s, S_{t+1} \in \mathcal{Y}, \pi] \tag{1}$$

2. (5 Points) Having defined $v_{\mathcal{Y}}^{\pi}$, you decide to relate your new value functions, $v_{\mathcal{X}}^{\pi}$ and $v_{\mathcal{X}^{\complement}}^{\pi}$, to the standard value function, $v^{\pi}$. Derive an expression for $v^{\pi}(s)$ that *only* uses the following terms: $\pi, P, \mathcal{A}, \mathcal{S}, \mathcal{X}, v_{\mathcal{X}}^{\pi}$ and $v_{\mathcal{X}^{\complement}}^{\pi}$. Note: You may introduce variables when summing over sets, e.g., $x$ in $\sum_{x \in \mathcal{X}}$. Your final answer should not include expectations or any random variables like $S_t$ or $R_t$. You should begin with the definition of $v^{\pi}(s)$ and end with an expression that only contains the allowed terms. Show your work (show the steps, don't just jump to your final answer). You may want to derive some properties before proceeding with the derivation for $v^{\pi}(s)$—that is allowed.

   Ans 2.

$$v^{\pi}(s) = \mathbf{E}[G_t | S_t = s, \pi] \tag{2}$$

$$= \sum_{s' \in \mathcal{S}} \Pr(S_{t+1} = s' | S_t = s, \pi) \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] \quad (from \, Law \, of \, Total \, Expectation) \tag{3}$$

$$= \sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s' | S_t = s, \pi) \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] + \sum_{s' \in \mathcal{X}^{\complement}} \Pr(S_{t+1} = s' | S_t = s, \pi) \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] \tag{4}$$

(as $\mathcal{X} \, and \, \mathcal{X}^{\complement} \, are \, disjoint \, and \, \mathcal{X} \cup \mathcal{X}^{\complement} = \mathcal{S}$)

$$= \sum_{s' \in \mathcal{X}} \sum_{a \in \mathcal{A}} \Pr(A_t = a | S_t = s, \pi) \Pr(S_{t+1} = s' | S_t = s, A_t = a, \pi) \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] \tag{5}$$

$$+ \sum_{s' \in \mathcal{X}^{\complement}} \sum_{a \in \mathcal{A}} \Pr(A_t = a | S_t = s, \pi) \Pr(S_{t+1} = s' | S_t = s, A_t = a, \pi) \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] \tag{6}$$

$$= \sum_{s' \in \mathcal{X}} \sum_{a \in \mathcal{A}} \pi(s, a) \Pr(S_{t+1} = s' | S_t = s, A_t = a, \pi) \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] \tag{7}$$

$$+ \sum_{s' \in \mathcal{X}^{\complement}} \sum_{a \in \mathcal{A}} \pi(s, a) \Pr(S_{t+1} = s' | S_t = s, A_t = a, \pi) \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] \tag{8}$$

$$= \sum_{s' \in \mathcal{X}} \sum_{a \in \mathcal{A}} \pi(s, a) P(s, a, s') \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] + \sum_{s' \in \mathcal{X}^{\complement}} \sum_{a \in \mathcal{A}} \pi(s, a) P(s, a, s') \mathbf{E}[G_t | S_t = s, S_{t+1} = s', \pi] \tag{9}$$

$$= \sum_{s' \in \mathcal{X}} \sum_{a \in \mathcal{A}} \pi(s, a) P(s, a, s') v_{\mathcal{X}}^{\pi}(s) + \sum_{s' \in \mathcal{X}^{\complement}} \sum_{a \in \mathcal{A}} \pi(s, a) P(s, a, s') v_{\mathcal{X}^{\complement}}^{\pi}(s) \tag{10}$$

---

[1]In latex, here we are using the symbols \complement and \setminus for $\complement$ and $\setminus$ respectively.

3. (13 Points) Having related your new value functions to the standard value function, you now talk to Geordi about what to do next to design a reinforcement learning algorithm using your new value functions. Another friend named Data loads the course notes from CMPSCI 687 in Fall 2019. He finds that the next step towards developing an algorithm with this value function may be to write out a new Bellman equation for $v_{\mathcal{X}}^\pi$. Derive a Bellman-like equation for this new value function. You should begin with the definition of $v_{\mathcal{X}}^\pi$ according to your answer to the first question, and should end with a recursive expression for $v_{\mathcal{X}}^\pi$ that is written only in terms of $\mathcal{S}, \mathcal{A}, \mathcal{P}, R, d_0, \gamma, \pi, \mathcal{X}$, and $\mathcal{X}^{\complement}$. For this problem, use an alternate definition of $R$: $R(s, a, s') := \mathbf{E}[R_t | S_t = s, A_t = a, S_{t+1} = s']$. (Hint: Using font size "tiny", our answer spans two lines—do not expect a short answer).

Ans 3.

$$v_{\mathcal{X}}^\pi(s) = \mathbf{E}[G_t | S_t = s, S_{t+1} \in \mathcal{X}, \pi] \tag{11}$$

$$= \mathbf{E}[\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s, S_{t+1} \in \mathcal{X}, \pi] \tag{12}$$

$$= \mathbf{E}[R_t + \sum_{k=1}^{\infty} \gamma^k R_{t+k} | S_t = s, S_{t+1} \in \mathcal{X}, \pi] \tag{13}$$

$$= \mathbf{E}[R_t | S_t = s, S_{t+1} \in \mathcal{X}, \pi] + \mathbf{E}[\sum_{k=1}^{\infty} \gamma^k R_{t+k} | S_t = s, S_{t+1} \in \mathcal{X}, \pi] \tag{14}$$

$$= \mathbf{E}[R_t | S_t = s, S_{t+1} \in \mathcal{X}, \pi] + \mathbf{E}[\sum_{k=0}^{\infty} \gamma^{k+1} R_{t+k+1} | S_t = s, S_{t+1} \in \mathcal{X}, \pi] \tag{15}$$

$$= \mathbf{E}[R_t | S_t = s, S_{t+1} \in \mathcal{X}, \pi] + \gamma \mathbf{E}[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, S_{t+1} \in \mathcal{X}, \pi] \tag{16}$$

$$= \mathbf{E}[R_t | S_t = s, S_{t+1} \in \mathcal{X}, \pi] + \gamma \mathbf{E}[G_{t+1} | S_t = s, S_{t+1} \in \mathcal{X}, \pi] \tag{17}$$

Say term $1 = \mathbf{E}[R_t | S_t = s, S_{t+1} \in \mathcal{X}, \pi]$ and term $2 = \mathbf{E}[G_{t+1} | S_t = s, S_{t+1} \in \mathcal{X}, \pi]$

Our result is $v_{\mathcal{X}}^\pi(s) = term1 + \gamma * term2$.

Calculate for term 1:

$$\mathbf{E}[R_t|S_t = s, S_{t+1} \in \mathcal{X}, \pi] = \sum_{a \in \mathcal{A}} \Pr(A_t = a|S_t = s, S_{t+1} \in \mathcal{X}, \pi)\mathbf{E}[R_t|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{18}$$

$$= \sum_{a \in \mathcal{A}} \frac{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)\Pr(A_t = a|S_t = s, \pi)}{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, \pi)}\mathbf{E}[R_t|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{19}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a, \pi)\Pr(A_t = a|S_t = s, \pi)}{\sum_{a' \in \mathcal{A}} \Pr(A_t = a'|S_t = s, \pi)\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)}\mathbf{E}[R_t|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{20}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a', \pi)}\mathbf{E}[R_t|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{21}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')}\mathbf{E}[R_t|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{22}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')}\sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi) \tag{23}$$

$$\times \mathbf{E}[R_t|S_t = s, A_t = a, S_{t+1} = s', \pi] \tag{24}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')}\sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi)R(s, a, s') \tag{25}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \tag{26}$$

$$\times \sum_{s' \in \mathcal{X}} \frac{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, S_{t+1} = s', \pi)\Pr(S_{t+1} = s'|S_t = s, A_t = a, \pi)}{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)}R(s, a, s') \tag{27}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')}\sum_{s' \in \mathcal{X}} \frac{\Pr(S_{t+1} = s'|S_t = s, A_t = a, \pi)}{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)}R(s, a, s') \tag{28}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')}\sum_{s' \in \mathcal{X}} \frac{P(s, a, s')}{\sum_{s'' \in \mathcal{X}} P(s, a, s'')}R(s, a, s') \tag{29}$$

Calculate for term 2:

$$\mathbf{E}[G_{t+1}|S_t = s, S_{t+1} \in \mathcal{X}, \pi] = \sum_{a \in \mathcal{A}} \Pr(A_t = a|S_t = s, S_{t+1} \in \mathcal{X}, \pi)\mathbf{E}[G_{t+1}|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{30}$$

$$= \sum_{a \in \mathcal{A}} \frac{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)\Pr(A_t = a|S_t = s, \pi)}{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, \pi)}\mathbf{E}[G_{t+1}|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{31}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a, \pi)\Pr(A_t = a|S_t = s, \pi)}{\sum_{a' \in \mathcal{A}} \Pr(A_t = a'|S_t = s, \pi)\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)} \tag{32}$$

$$\times \mathbf{E}[G_{t+1}|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{33}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a', \pi)}\mathbf{E}[G_{t+1}|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{34}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')}\mathbf{E}[G_{t+1}|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi] \tag{35}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi) \tag{36}$$

$$\times \mathbf{E}[G_{t+1}|S_t = s, A_t = a, S_{t+1} = s', \pi] \tag{37}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi) \tag{38}$$

$$\times \mathbf{E}[G_{t+1}|S_{t+1} = s', \pi] \tag{39}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \Pr(S_{t+1} = s'|S_t = s, A_t = a, S_{t+1} \in \mathcal{X}, \pi)v^\pi(s') \tag{40}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \tag{41}$$

$$\times \sum_{s' \in \mathcal{X}} \frac{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, S_{t+1} = s', \pi)\Pr(S_{t+1} = s'|S_t = s, A_t = a, \pi)}{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)}v^\pi(s') \tag{42}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \frac{\Pr(S_{t+1} = s'|S_t = s, A_t = a, \pi)}{\Pr(S_{t+1} \in \mathcal{X}|S_t = s, A_t = a, \pi)}v^\pi(s') \tag{43}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \frac{P(s, a, s')}{\sum_{s'' \in \mathcal{X}} P(s, a, s'')}v^\pi(s') \tag{44}$$

$$= \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \frac{P(s, a, s')}{\sum_{s'' \in \mathcal{X}} P(s, a, s'')} \tag{45}$$

$$\times \left( \sum_{s'' \in \mathcal{X}} \sum_{a' \in \mathcal{A}} \pi(s', a')P(s', a', s'')v^\pi_{\mathcal{X}}(s') + \sum_{s'' \in \mathcal{X}^{\complement}} \sum_{a' \in \mathcal{A}} \pi(s', a')P(s', a', s'')v^\pi_{\mathcal{X}^{\complement}}(s') \right) \tag{46}$$

So,

$$v^\pi_{\mathcal{X}}(s) = \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \frac{P(s, a, s')}{\sum_{s'' \in \mathcal{X}} P(s, a, s'')}R(s, a, s') \tag{47}$$

$$+ \gamma \sum_{a \in \mathcal{A}} \frac{\sum_{s' \in \mathcal{X}} P(s, a, s')\pi(s, a)}{\sum_{a' \in \mathcal{A}} \pi(s, a')\sum_{s' \in \mathcal{X}} P(s, a', s')} \sum_{s' \in \mathcal{X}} \frac{P(s, a, s')}{\sum_{s'' \in \mathcal{X}} P(s, a, s'')} \tag{48}$$

$$\times \left( \sum_{s'' \in \mathcal{X}} \sum_{a' \in \mathcal{A}} \pi(s', a')P(s', a', s'')v^\pi_{\mathcal{X}}(s') + \sum_{s'' \in \mathcal{X}^{\complement}} \sum_{a' \in \mathcal{A}} \pi(s', a')P(s', a', s'')v^\pi_{\mathcal{X}^{\complement}}(s') \right) \tag{49}$$

4

4. (5 Points) Consider the following definition of an optimal policy:

> For any finite MDP with $\gamma < 1$ and precisely two actions, $a_1$ and $a_2$, for any two policies $\pi$ and $\pi'$, $\pi \geq \pi'$ iff $\forall s \in \mathcal{S}$, $q^\pi(s, a_1) \geq q^{\pi'}(s, a_1)$. A policy $\pi$ is optimal iff $\pi \geq \pi'$ for all policies $\pi'$.

Is this definition equivalent to the definition from Section 4.5 in the course notes? Prove that your answer is correct.

Ans 4. Say the definition proposed here is Definition A and the definition from class is Definition B. Definition A is not equivalent to Definition B. I am proving this by providing an example of an MDP for which Definition A does not imply Definition B.

**Proof**:

Consider an MDP in which $\mathcal{S} = \{s_1, s_2\}$, $\mathcal{A} = \{a_1, a_2\}$, $P(s_1, a_1, s_1) = P(s_1, a_2, s_1) = P(s_2, a_1, s_1) = P(s_2, a_2, s_1) = 0$, $P(s_1, a_2, s_2) = P(s_1, a_1, s_2) = P(s_2, a_1, s_2) = P(s_2, a_2, s_2) = 1$, $d_0(s_1) = 1, d_0(s_2) = 0$, $R(s_1, a_1) = +10, R(s_1, a_2) = -10, R(s_2, a_1) = 0, R(s_2, a_2) = 0, \gamma = 0$

$$q^\pi(s_1, a_1) = R(s_1, a_1) + \gamma \sum_{s' \in \mathcal{S}} P(s_1, a_1, s') \sum_{a' \in \mathcal{A}} \pi(s', a') q^\pi(s', a') \tag{50}$$

As $\gamma = 0$,

$$q^\pi(s_1, a_1) = R(s_1, a_1) \tag{51}$$
$$q^\pi(s_1, a_1) = +10 \tag{52}$$

$$q^\pi(s_2, a_1) = R(s_2, a_1) + \gamma \sum_{s' \in \mathcal{S}} P(s_2, a_1, s') \sum_{a' \in \mathcal{A}} \pi(s', a') q^\pi(s', a') \tag{53}$$

As $\gamma = 0$,

$$q^\pi(s_2, a_1) = R(s_2, a_1) \tag{54}$$
$$q^\pi(s_2, a_1) = 0 \tag{55}$$

We can see that for all policies $\pi$, $q^\pi(s_1, a_1) = +10$ and $q^\pi(s_2, a_1) = 0$, so under Definition A, all policies are optimal. This is because $q^\pi(s_1, a_1)$ and $q^\pi(s_2, a_1)$ do not depend on $\pi$ in this example: $q^\pi(s_1, a_1) = q^{\pi'}(s_1, a_1) = +10$ and $q^\pi(s_2, a_1) = q^{\pi'}(s_2, a_1) = 0$ for all $\pi$ and $\pi'$.

Although, consider a particular $\pi$: $\pi(s_1, a_1) = 1, \pi(s_1, a_2) = 0, \pi(s_2, a_1) = 0, \pi(s_2, a_2) = 1$ and a particular $\pi'$: $\pi'(s_1, a_1) = 0, \pi'(s_1, a_2) = 1, \pi'(s_2, a_1) = 0, \pi'(s_2, a_2) = 1$

$$v^\pi(s_1) = \sum_{a' \in \mathcal{A}} \pi(s_1, a') q^\pi(s_1, a') \tag{56}$$
$$v^\pi(s_1) = \pi(s_1, a_1) q^\pi(s_1, a_1) + \pi(s_1, a_2) q^\pi(s_1, a_2) \tag{57}$$
$$v^\pi(s_1) = 1 \times q^\pi(s_1, a_1) + 0 \times q^\pi(s_1, a_2) \tag{58}$$
$$v^\pi(s_1) = q^\pi(s_1, a_1) \tag{59}$$
$$v^\pi(s_1) = +10 \tag{60}$$

$$q^{\pi'}(s_1, a_2) = R(s_1, a_2) + \gamma \sum_{s' \in \mathcal{S}} P(s_1, a_2, s') \sum_{a' \in \mathcal{A}} \pi(s', a') q^\pi(s', a') \tag{61}$$

As $\gamma = 0$,

$$q^{\pi'}(s_1, a_2) = R(s_1, a_2) \tag{62}$$
$$q^{\pi'}(s_1, a_2) = -10 \tag{63}$$

5

$$v^{\pi'}(s_1) = \sum_{a' \in \mathcal{A}} \pi'(s_1, a') q^{\pi'}(s_1, a') \tag{64}$$

$$v^{\pi'}(s_1) = \pi'(s_1, a_1) q^{\pi'}(s_1, a_1) + \pi'(s_1, a_2) q^{\pi'}(s_1, a_2) \tag{65}$$

$$v^{\pi'}(s_1) = 0 \times q^{\pi'}(s_1, a_1) + 1 \times q^{\pi'}(s_1, a_2) \tag{66}$$

$$v^{\pi'}(s_1) = q^{\pi'}(s_1, a_2) \tag{67}$$

$$v^{\pi'}(s_1) = -10 \tag{68}$$

So, $v^{\pi}(s_1) > v^{\pi'}(s_1)$. This means that $\pi'$ is not an optimal policy under Defiition B, but it is an optimal policy under definition A. Thus, the two definitions are not equivalent. Hence proved.

5. (5 Points) Consider a different definition of $\geq$ for policies: $\pi \geq \pi'$ iff $\sum_{s \in \mathcal{S}} d_0(s) v^{\pi}(s) \geq \sum_{s \in \mathcal{S}} d_0(s) v^{\pi'}(s)$. Using this modified version of $\geq$, we can still define an optimal policy to be any policy $\pi$ such that $\pi \geq \pi'$ for all $\pi'$. Prove that using this definition of an optimal policy is equivalent to using our first definition:

$$\pi^* \in \arg\max_{\pi \in \Pi} J(\pi). \tag{69}$$

Ans 5. Both definitions are equivalent. To prove that both definitions are equivalent, consider the definition proposed here to be Definition A and consider the definition from class to be Definition B. We need to prove Definition A $\implies$ Definition B and Definition B $\implies$ Definition A.

**Proving Definition A $\implies$ Definition B**: We are given that $\pi \geq \pi'$ if and only if $\sum_{s \in \mathcal{S}} d_0(s) v^{\pi}(s) \geq \sum_{s \in \mathcal{S}} d_0(s) v^{\pi'}(s)$ for all $\pi$ then $\pi'$ is optimal.

$$\sum_{s \in \mathcal{S}} d_0(s) v^{\pi}(s) \geq \sum_{s \in \mathcal{S}} d_0(s) v^{\pi'}(s) \tag{70}$$

$$\implies \sum_{s \in \mathcal{S}} Pr(S_0 = s) v^{\pi}(s) \geq \sum_{s \in \mathcal{S}} Pr(S_0 = s) v^{\pi'}(s) \tag{71}$$

$$\implies \sum_{s \in \mathcal{S}} Pr(S_0 = s) \mathbf{E}[G_t | S_0 = s, \pi] \geq \sum_{s \in \mathcal{S}} Pr(S_0 = s) \mathbf{E}[G_t | S_0 = s, \pi'] \tag{72}$$

$$\implies \mathbf{E}[G_t | \pi] \geq \mathbf{E}[G_t | \pi'] \quad (from Law of Total Expectation) \tag{73}$$

$$\implies J(\pi) \geq J(\pi') \tag{74}$$

So, this implies that $\pi \geq \pi'$ if and only if $J(\pi) \geq J(\pi')$ for all $\pi'$ then $\pi$ is an optimal policy. Renaming the variables, we get:
For an optimal policy $\pi^*$, $\pi^* \geq \pi$ if and only if $J(\pi^*) \geq J(\pi)$ for all $\pi$. Therefore,

$$\pi^* \in \arg\max_{\pi \in \Pi} J(\pi) \tag{75}$$

Therefore, we have proved that Definition A $\implies$ Definition B. For the complete proof, we also need to prove that Definition B $\implies$ Definition A.

**Proving Definition B $\implies$ Definition A**:

$$\pi^* \in \arg\max_{\pi \in \Pi} J(\pi) \tag{76}$$

We have that $J(\pi^*) \geq J(\pi)$ for all $\pi$ then $\pi^*$ is optimal. Renaming the variables, we get: $J(\pi) \geq J(\pi')$ for all $\pi'$ then $\pi$ is optimal. Therefore,

$$J(\pi) \geq J(\pi') \tag{77}$$

$$\implies \mathbf{E}[G_t|\pi] \geq \mathbf{E}[G_t|\pi'] \quad (from Law of Total Expectation) \tag{78}$$

$$\implies \sum_{s \in \mathcal{S}} Pr(S_0 = s)\mathbf{E}[G_t|S_0 = s, \pi] \geq \sum_{s \in \mathcal{S}} Pr(S_0 = s)\mathbf{E}[G_t|S_0 = s, \pi'] \tag{79}$$

$$\implies \sum_{s \in \mathcal{S}} Pr(S_0 = s)v^{\pi}(s) \geq \sum_{s \in \mathcal{S}} Pr(S_0 = s)v^{\pi'}(s) \tag{80}$$

$$\implies \sum_{s \in \mathcal{S}} d_0(s)v^{\pi}(s) \geq \sum_{s \in \mathcal{S}} d_0(s)v^{\pi'}(s) \tag{81}$$

So, this implies that $\pi \geq \pi'$ iff $\sum_{s \in \mathcal{S}} d_0(s)v^{\pi}(s) \geq \sum_{s \in \mathcal{S}} d_0(s)v^{\pi'}(s)$ for all $\pi'$ then $\pi$ is an optimal policy.

Therefore, I have proven that both definitions are equivalent.

6. (20 Points) In class we proved that the Bellman operator is a contraction, and used this to show that value iteration converges to a unique fixed point. In this problem you will prove that the dynamic programming policy evaluation operator is a contraction, and so the policy evaluation algorithm converges to a unique fixed-point. (From the Bellman equation, it should then be clear that this fixed point is $v^{\pi}$, establishing that our dynamic programming policy evaluation algorithm converges to $v^{\pi}$.) Let $f$ denote the dynamic programming policy evaluation operator (this is currently equation (200) in the course notes, viewed as an operator on value function approximations):

$$fv(s) = \sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')\big(R(s, a) + \gamma v(s')\big). \tag{82}$$

Notice that the definition of $f$ relies on a specific policy $\pi$—this is the policy being evaluated by the policy evaluation algorithm. Prove that $f$ is a contraction under the $L^{\infty}$ norm (the same max norm used in our proof that the Bellman operator is a contraction).

Ans 6. To prove that $f$ is a contraction under the $L^{\infty}$ norm, we need to prove $||fv(s) - fv'(s)||_{\infty} \leq \gamma||v(s) - v'(s)||_{\infty}$

$$||fv(s) - fv'(s)||_{\infty} = \max_{s \in \mathcal{S}} |fv(s) - fv'(s)| \tag{83}$$

$$= \max_{s \in \mathcal{S}} |\sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')\big(R(s, a) + \gamma v(s')\big) - \sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')\big(R(s, a) + \gamma v'(s')\big)|$$
$$\tag{84}$$

$$= \max_{s \in \mathcal{S}} |\sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')\gamma v(s') - \sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')\gamma v'(s')| \tag{85}$$

$$= \gamma \max_{s \in \mathcal{S}} |\sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')v(s') - \sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')v'(s')| \tag{86}$$

$$= \gamma \max_{s \in \mathcal{S}} |\sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')(v(s') - v'(s'))| \tag{87}$$

As with the modulo, the terms would be cancelling each other's magnitudes,

$$\leq \gamma \max_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \pi(s, a)|\sum_{s' \in \mathcal{S}} P(s, a, s')(v(s') - v'(s'))| \tag{88}$$

$$\leq \gamma \max_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \pi(s, a) \sum_{s' \in \mathcal{S}} P(s, a, s')|v(s') - v'(s')| \tag{89}$$

$$\leq \gamma \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \max_{s' \in \mathcal{S}} |v(s') - v'(s')| \tag{90}$$

$$= \gamma \max_{s' \in \mathcal{S}} |v(s') - v'(s')| \tag{91}$$

$$= \gamma||v(s') - v'(s')||_{\infty} \tag{92}$$

Thus, I have proved that $f$ is a contraction under the $L^{\infty}$ norm. Also, as $f$ is a contraction mapping on a non-empty complete normed vector space, $f$ has a unique fixed point to which $f$ converges (from Banach Fixed-Point Theorem). Thus, I have proved that the dynamic programming policy evaluation algorithm converges to a unique fixed point.

# Part Two: Programming (25 Points Total)

For this part of the assignment, you will implement value iteration (modified to terminate when the value function estimate has not changed significantly between two iterations). Your program will read an MDP from a file, run value iteration on the MDP, and output the final estimate of the optimal value function and the policies that are greedy with respect to this value function. As a soft introduction to C++, we are providing you with most of the code here: your job is to fill in the missing lines in the function `valueIteration`, marked with a comment saying "TODO". Do not change the code logic outside of the valueIteration function (you may add new functions if you like, but do not modify any of the other functions in your final submission or it may fail to run as expected in our auto-grader).

You are free to use any IDE or toolchain you would like to program in C++. If you are not familiar with C++, we have provided two different systems for opening and working with this C++ code. If you are using Windows, you should download Microsoft Visual Studio. The community version is perfectly sufficient, and is free online (in my opinion, this is the best C++ experience out there). Clicking on the .sln file in HW3/build/VisualStudio will open the project. On the left you should see main.cpp—open this file to see all of the code for this assignment. If you are using Mac or Linux, we have provided a CLion project. CLion is free for students. To open this project, select "Open" when launching CLion. Select the file HW3/build/CLion/CMakeLists.txt. When prompted, select "Open as Project". If main.cpp does not immediately open, on the left click on HW3/main.cpp.

This assignment is your chance to begin to familiarize yourself with C++. Please look over all of the provided code, and feel free to ask if you have questions about what some portion of the code is doing. Also, take this opportunity to familiarize yourself with the debugger in your IDE—developing simple programs in C++ is a breeze when you are familiar with how to use the different capabilities of your debugger.

We have provided you (within the provided code) with 687Gridworld.txt, a text file containing the MDP we have been using in class. We will evaluate your program on other MDPs that we are not providing to you. You are welcome to create your own test MDPs, but do not share these with others.

You must submit your main.cpp file. A correct implementation is worth 20 points. *Any* incorrect output (beyond numerical issues) will result in 0/20 points. In the .pdf that you submit, answer the following questions.

1. (2 Points) Did your final code compile on your machine? (Yes or no).

   Ans 1. Yes.

2. (3 Points) Comment on your experience with this problem. Did your first implementation work, or did you introduce a bug at first? Was there anything we could do to smooth your introduction to C++? Did you implement any additional test MDPs (you do not have to in order to get full credit). Did the number of iterations required by value iteration surprise you? Do you have any other comments on this problem?

   Ans 2. I had a nice experience with this problem, as it was easy to solve and insightful. I did not introduce a bug as I was able to do this question in a single run of the program. I would say that the introduction to C++ was really nice, and I do not think anything more needs to be done to smoothen the introduction to C++ further. I did not implement any additional test MDPs apart from the 687Gridworld and MoreWatery687Gridworld. The number of iterations did surprise me, as the program took really few number of iterations to reach the optimal value functions for both, the 687Gridworld and the MoreWatery687Gridworld MDPs. I do not have any other comments on this problem.