

Answer any THREE questions. Each question is worth 20 marks. Use separate answer books for PART A and PART B. **Gatsby PhD students only:** answer *either* TWO questions from PART A and ONE question from PART B; *or* ONE question from PART A and TWO questions from PART B.

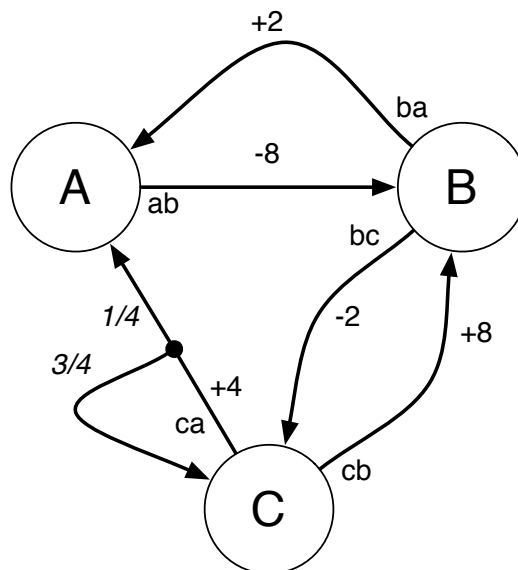
Marks for each part of each question are indicated in square brackets

Calculators are NOT permitted

Part A: Kernel Methods

Part B: Reinforcement Learning

- Consider the following Markov Decision Process (MDP) with discount factor $\gamma = 0.5$. Upper case letters A, B, C represent states; arcs represent state transitions; lower case letters ab, ba, bc, ca, cb represent actions; signed integers represent rewards; and fractions represent transition probabilities.



✓

- Define the *state-value function* $V^\pi(s)$ for a discounted MDP

$$V^\pi(s) = \mathbb{E}_\pi [R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s] \quad [1 \text{ marks}]$$

✓

- Write down the *Bellman expectation equation* for state-value functions

$$V^\pi(s) = \mathbb{E} [R_{t+1} + \gamma V^\pi(S_{t+1}) | S_t = s] \quad [2 \text{ marks}]$$

- Consider the uniform random policy $\pi_1(s,a)$ that takes all actions from state s with equal probability. Starting with an initial value function of $V_1(A) = V_1(B) = V_1(C) = 2$, apply one synchronous iteration of iterative policy evaluation (i.e. one backup for each state) to compute a new value function $V_2(s)$

$$V_2^{\pi_1}(A) = -8 + 2 = -6; V_2^{\pi_1}(B) = \frac{1}{2}(+2+2) + \frac{1}{2}(-2+2) = 2$$

- Apply one iteration of greedy policy improvement to compute a new, deterministic

$$\pi_2(A) = (0, 1, 0); \pi_2(B) = (0, 0, 1); \pi_2(C) = (0, 1, 0)$$

- Consider a deterministic policy $\pi(s)$. Prove that if a new policy π' is greedy with respect to V^π then it must be better than or equal to π , i.e. $V^{\pi'}(s) \geq V^\pi(s)$ for all s ; and that if $V^{\pi'}(s) = V^\pi(s)$ for all s then π' must be an optimal policy.

$$\exists \pi' \text{ is greedy} \rightarrow \forall s \in S: V_{\pi'}(s) \geq \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s')$$

- Define the optimal state-value function $V^*(s)$ for an MDP

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

- Write down the Bellman optimality equation for state-value functions

$$V_*(s) = \max_a (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s'))$$

- Starting with an initial value function of $V_1(A) = V_1(B) = V_1(C) = 2$, apply one synchronous iteration of value iteration (i.e. one backup for each state) to compute a new value function $V_2(s)$.

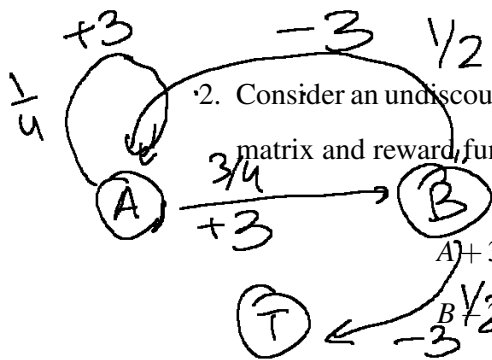
$$V(A) = -8 + 2 = -6; V(B) = \max(2+2, -2+2) = 4; V(C) = \max(4+2, -8+2) = 10$$

- Is your new value function $V_2(s)$ optimal? Justify your answer.

$$\begin{aligned} V_3(A) &= -8 + 4 = -4 \\ V_3(B) &= \max(+2-6, -2+10) = 8 \\ V_3(C) &= \max(+8+4, 4 + \frac{1}{4} \cdot (-6) + \frac{3}{4} \cdot 10) = 12 \end{aligned}$$

So it's not optimal

14/20



2. Consider an undiscounted Markov Reward Process with two states A and B . The transition matrix and reward function are unknown, but you have observed two sample episodes:

$A + 3 \rightarrow A + 2 \rightarrow B - 4 \rightarrow A + 4 \rightarrow B - 3 \rightarrow \text{terminate}$

$B + 2 \rightarrow A + 3 \rightarrow B - 3 \rightarrow \text{terminate}$

In the above episodes, sample state transitions and sample rewards are shown at each step, e.g. $A + 3 \rightarrow A$ indicates a transition from state A to state A , with a reward of $+3$.

Using first-visit Monte-Carlo evaluation, estimate the state-value function $V(A), V(B)$

$$V(A) = \frac{2 + (-2)}{2} = 0 \quad V(B) = 0 \quad X \quad [2 \text{ marks}]$$

Using every-visit Monte-Carlo evaluation, estimate the state-value function $V(A), V(B)$

$$V(A) = \frac{2(3) + (-2)(1)}{4} = 1 \quad ; \quad V(B) = \frac{2(2) + (-2)(2)}{4} = 0 \quad X \quad [2 \text{ marks}]$$

Draw a diagram of the Markov Reward Process that best explains these two episodes (i.e. the model that maximises the likelihood of the data - although it is not necessary to prove this fact). Show rewards and transition probabilities on your diagram.

[4 marks] ✓

Define the Bellman equation for a Markov reward process

$$V(s) = \mathbb{E}[R_{t+1} + \gamma V(s_{t+1}) | s_t = s] \quad V \quad [2 \text{ marks}]$$

Solve the Bellman equation to give the true state-value function $V(A), V(B)$. Hint: solve the Bellman equations directly, rather than iteratively. $V(A) = 2; V(B) = -2$ ✓

What value function would batch TD(0) find, i.e. if TD(0) was applied repeatedly to these two episodes?

$$V(A) = 2; V(B) = -2 \quad V \quad [2 \text{ marks}]$$

What value function would batch TD(1) find, using accumulating eligibility traces?

$$V(A) = 0; V(B) = 0 \quad \text{as it is MC} \quad X \quad [2 \text{ marks}]$$

What value function would LSTD(0) find?

$$V(A) = 2; V(B) = -2 \quad V \quad [2 \text{ marks}]$$

[Total 20 marks]

GI13

3

14/20

TURN OVER

Total: 28/40

This one is done for myself, not for eval.

3. A rat is involved in an experiment. It experiences one episode. At the first step it hears a bell. At the second step it sees a light. At the third step it both hears a bell and sees a light. It then receives some food, worth +1 reward, and the episode terminates on the fourth step. All other rewards were zero. The experiment is undiscounted.

- Represent the rat's state s by a vector of two binary features, $bell(s) \in \{0, 1\}$ and $light(s) \in \{0, 1\}$. Write down the sequence of feature vectors corresponding to this episode.

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix}; \begin{pmatrix} 0 \\ 1 \end{pmatrix}; \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \checkmark$$

[3 marks]

- Approximate the state-value function by a linear combination of these features with two parameters: $b \cdot bell(s) + l \cdot light(s)$. If $b = 2$ and $l = -2$ then write down the sequence of approximate values corresponding to this episode.

$$2; -2; 0 \quad \checkmark$$

[3 marks]

- Define the λ -return v_t^λ

$$v_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} v_t^{(n)} \quad \checkmark$$

[1 marks]

- Write down the sequence of λ -returns v_t^λ corresponding to this episode, for $\lambda = 0.5$ and $b = 2, l = -2$

$$1/2(2) = \frac{1}{2}; \quad 1/2(2) + (1/2)(-2) = 0; \quad 0 + (1/2)^2(0) = 0 \quad \text{[3 marks]} \quad \times$$

- Using the forward-view TD(λ) algorithm and your linear function approximator, what are the sequence of updates to weight b ? What is the total update to weight b ? Use $\lambda = 0.5, \gamma = 1, \alpha = 0.5$ and start with $b = 2, l = -2$

[3 marks]

- Define the TD(λ) *accumulating eligibility trace* e_t when using linear value function approximation

$$e_t = \gamma \lambda e_{t-1} + \phi(s) \quad \checkmark$$

[1 marks]

- Write down the sequence of eligibility traces e_t corresponding to the bell, using $\lambda = 0.5, \gamma = 1$

[3 marks]

- Using the backward-view TD(λ) algorithm and your linear function approximator, what are the sequence of updates to weight b ? (Use offline updates, i.e. do not actually change your weights, just accumulate your updates). What is the total update to weight b ? Use $\lambda = 0.5, \gamma = 1, \alpha = 0.5$ and start with $b = 2, l = -2$

GI13

Answer is on the next page? CONTINUED

[3 marks]

$$\Delta b_1 = \alpha \delta_1 e_1 = 0.5(0 + -2 - 2)1 = -2$$

$$\Delta b_2 = \alpha \delta_2 e_2 = 0.5(0 + 0 - -2)1/2 = 1/2$$

$$\Delta b_3 = \alpha \delta_3 e_3 = 0.5(1 + 0 - 0)5/4 = 5/8$$

$$\sum \Delta b = (-2 + 1/2 + 5/8) = -7/8$$

[Total 20 marks]

END OF PAPER