

# Midterm Exam: Examples of Questions

## An Open-Notes Exam

CS4771 Reinforcement Learning (Spring 2026)

Other important notes:

- Write your answer in the **answer boxes**. Anything written outside the answer boxes or ~~crossed out~~ will not be graded.
  - Besides the answer box, there are several questions with a block of **Detailed calculation**. It is NOT necessary to provide detailed calculation. However, that could earn you partial credit if your answer is incorrect.
  - **Fractions** need NOT be simplified or converted to decimals in your answer.
  - The number of correct answers in **multiple choice questions** can range from 0 to 5. They will be graded as 5 independent true-or-false questions. If you think none of the choices are correct, put “None” in the answer box; leaving it blank will result in zero points.
1. In this problem, we use Markov decision process to model the operation of a machine.

The machine can be in one of two states: B (Bad condition), G (Good condition). In each state, there are two possible actions: N (Do nothing) and R (Repair). The machine operates indefinitely with a discount factor  $\gamma = 0.9$ . The rewards and transition probabilities are defined in the tables below.

$s$	$a$	$R(s, a)$
B	N	-5
B	R	-13
G	N	+10
G	R	+8

Rewards

$s$	$a$	$s'$	$P(s'   s, a)$
B	N	B	1.0
B	R	G	1.0
G	N	G	0.8
G	N	B	0.2
G	R	G	1.0

Transition Probabilities

Rationale behind the reward and transition functions are as follows:

- **Do nothing in Bad condition:** Incur a loss of  $-5$  due to lost productivity and remain in bad condition.
- **Repair in Bad condition:** Incur a total cost of  $-13$  (loss of productivity and repair cost) and transition to good condition.

- **Do nothing in Good condition:** Earn +10 from normal production, with a 20% chance of transitioning to bad condition.
- **Repair in Good condition:** Pay a maintenance cost of -2, resulting in a net reward of +8, and ensure the machine stays in good condition.

Suppose we start from  $V_0(s) = Q_0(s, a) = 0$  for all  $s, a$ . Let  $V_k(s)$  and  $Q_k(s, a)$  be the value functions after performing value iterations for  $k$  times. Let  $V^*(s) = \lim_{k \rightarrow \infty} V_k(s)$  and  $Q^*(s, a) = \lim_{k \rightarrow \infty} Q_k(s, a)$ .

- 1.1.** Write the Bellman equations (the relations that should be satisfied for  $V^*$  and  $Q^*$ ) for this MDP, using the specific numbers provided in the tables.

**Answer:**

- 1.2.** Starting with  $V_0(B) = 0$  and  $V_0(G) = 0$ , perform two iterations of value iteration. Compute  $V_1(B)$ ,  $V_1(G)$ ,  $V_2(B)$ , and  $V_2(G)$ .

$$V_1(B) = \boxed{\phantom{000}} \quad V_1(G) = \boxed{\phantom{000}} \quad V_2(B) = \boxed{\phantom{000}} \quad V_2(G) = \boxed{\phantom{000}}$$

- 1.3.** Determine the optimal policy  $\pi^*$  and the optimal discounted value function  $V^*$ .

$$\pi^*(B) = \boxed{\phantom{000}} \quad \pi^*(G) = \boxed{\phantom{000}} \quad V^*(B) = \boxed{\phantom{000}} \quad V^*(G) = \boxed{\phantom{000}}$$

2. Which of the following statements about contextual bandits are true?
- A. In contextual bandits, the action chosen by the learner may affect future contexts.
  - B. Unlike supervised learning, contextual bandits suffer from “bandit feedback,” meaning the learner only observes the reward for the chosen action and not the rewards of the unchosen actions.
  - C. In policy-based methods, we construct unbiased reward estimators to encourage exploration.
  - D. With  $\epsilon$ -greedy, we usually use decreasing  $\epsilon$  over time to balance exploitation and exploration.
  - E. If the learner has a reward estimation  $\hat{R}$  such that  $|\hat{R}(x, a) - R(x, a)| \leq \Delta$  for all  $x, a$  ( $R$  is the true underlying reward function), then the greedy-policy  $\pi(x) = \operatorname{argmax}_a \hat{R}(x, a)$  guarantees  $\max_a R(x, a) - R(x, \pi(x)) \leq 2\Delta$ .

**Answer:**

3. Recall the UCB algorithm for multi-armed bandits:

---

**Algorithm 1** UCB

---

**Parameters:** fixed constant  $c = 2$ .

**for**  $t = 1, 2, \dots, T$  **do**

Select arm

$$a_t = \operatorname{argmax}_a \left( \hat{R}_t(a) + c \sqrt{\frac{1}{N_t(a)}} \right) \quad (1)$$

where  $\hat{R}_t(a)$  is the empirical mean of arm  $a$  up to time  $t - 1$ , and  $N_t(a)$  is the number of draws to arm  $a$  up to time  $t - 1$ . In case of a tie, the algorithm selects the arm with the smallest index.

---

Suppose at time  $t = 30$ , the learner has collected the following statistics for three arms ( $A = 3$ ):

Arm ( $a$ )	Empirical Mean ( $\hat{R}_{30}(a)$ )	Number of Pulls ( $N_{30}(a)$ )
1	0.85	16
2	0.70	9
3	0.40	4

Which arm  $a_{30}$  will the algorithm select? Justify your answer.

**Answer:**

4. Consider the KL-regularized policy update algorithm for multi-armed bandits:

$$\pi_{t+1} = \operatorname{argmax}_{\pi} \left\{ \langle \pi, r_t \rangle - \beta \text{KL}(\pi, \pi_t) \right\}.$$

Here,  $r_t(a) = R(a) + w_t(a)$  where  $R(a)$  being the true expected reward of arm  $a$  and  $w_t(a)$  is a zero-mean noise. Explain the drawbacks of setting  $\beta$  too large or too small, respectively.

**Answer:**