

CS 542 Stats RL Homework 1

September 17, 2024

- **Submission deadline: Oct. 2 before class.** Submission website: Canvas.
- **Your submission must be typeset with L^AT_EX.** Handwritten solutions will not be accepted. The L^AT_EX source file of this homework will be provided and you can use it as a template.
- You can discuss with anyone and consult any material, but (1) you still need to write the homework on your own, and (2) if you get help from anyone other than the course instructors or any material other than the course notes, you will need to mention them in your homework.

1. Evaluation error to decision loss (5 pts)

(1) There are K items, $1, 2, \dots, K$. The i -th item has value $v_i \in \mathbb{R}$. Let $v^* := \max_{i \in [K]} v_i$, where $[K] := \{1, 2, \dots, K\}$, and $i^* = \arg \max_{i \in [K]} v_i$.

An agent chooses the j -th item, where $j = \arg \max_{i \in [K]} u_i$, and $\{u_i\}_{i=1}^K$ are K real numbers. Let

$$\epsilon := \max_{i \in [K]} |u_i - v_i|.$$

Upper-bound $v^* - v_j$ as a function of ϵ . Prove your result.

- (2) If we further have $\forall i, u_i \leq v_i$, can you improve the bound? What about $\forall i, u_i \geq v_i$?
- (3) State and prove an upper bound of $|u_j - v^*|$.

2. Loss of using a smaller γ (5 pts)

Suppose we are given an MDP $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$. In the lecture we have seen that heavy discounting leads to faster convergence of planning, so we may run some planning algorithm using $\gamma' < \gamma$. The question is, how lossy is the obtained policy when we evaluate it in the original MDP?

More concretely, let's ignore the details of the planning algorithm and say we can compute the optimal policy for $M' = (\mathcal{S}, \mathcal{A}, P, R, \gamma')$, that is, a new MDP that is the same as M in all parameters except that its discount factor is γ' instead of γ . Let $\pi_{\gamma'}^*$ denote the optimal policy of M' . Prove a bound on $\|V_M^* - V_M^{\pi_{\gamma'}^*}\|_\infty$. Here the subscript M is not necessary and only used to emphasize that both value functions are defined in the original M (i.e., w.r.t. γ , not γ'). Your bound should scale with $\gamma - \gamma'$, that is, when γ' is close to γ , the loss will be small. (If your bound is correct but loose by significant factor(s), you may lose up to 1 point.)

3. Perturbation bound for d^π

Let π_1, π_2 be two policies, and ϵ be their distance, defined as follows:

$$\epsilon := \max_{s \in \mathcal{S}} \|\pi_1(\cdot|s) - \pi_2(\cdot|s)\|_1.$$

That is, we consider the L_1 difference between their action distributions (which is proportional to the TV distance) in each state, and then take the worst-case error over all possible states. You are asked to upper-bound $d^{\pi_1} - d^{\pi_2}$ using ϵ , where $d^\pi = (1 - \gamma)(d_0^\top (I - \gamma P^\pi)^{-1})^\top$ is the discounted state occupancy of π when the initial state s_1 is drawn from some initial distribution $d_0 \in \Delta(\mathcal{S})$.

(1) (5 pts) Use ϵ to upper-bound $\|d^{\pi_1} - d^{\pi_2}\|_\infty$. Your bound should go to 0 when $\epsilon = 0$.¹

Hint: Recall that $d^\pi(s) = (1 - \gamma) \cdot \mathbb{E}[\sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{I}[s_t = s] | s_1 \sim d_0, \pi]$ can be viewed as the expected return under an indicator reward function $R(\tilde{s}, a) = \mathbb{I}[\tilde{s} = s]$ (up to the normalization factor $(1 - \gamma)$), and the performance difference lemma can be used to bound the difference between the returns of two policies.

(2) Optional (5 pts) What if we want a bound on $\|d^{\pi_1} - d^{\pi_2}\|_1$? While we can obtain a bound by $\|d^{\pi_1} - d^{\pi_2}\|_1 \leq |\mathcal{S}| \|d^{\pi_1} - d^{\pi_2}\|_\infty$ and plugging the bound from (1) into the RHS, this will incur a dependence on $|\mathcal{S}|$. Can you prove a bound that does not have such a dependence?

Hint 1: d^π satisfies the *Bellman flow equation* $d^\pi = (1 - \gamma)d_0 + \gamma(P^\pi)^\top d^\pi$ ($P^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ is the transition matrix under policy π , where $[P^\pi]_{s,s'} := \sum_a \pi(a|s)P(s'|s, a)$).

Hint 2: You may find the *data processing inequality* useful and can directly use it, whose general form can be found on Wikipedia. As a special case that suffices for our problem, consider two stochastic vectors $q, q' \in \mathbb{R}^n$ (i.e., they represent distributions over n items) and a column-stochastic matrix $T \in \mathbb{R}^{m \times n}$. Note that Tq is a stochastic vector in \mathbb{R}^m , and can be interpreted as $Y \sim Tq \Leftrightarrow X \sim q, Y \sim T(\cdot|X)$, where $T(\cdot|x)$ is the x -th column of T (i.e., T is interpreted as a conditional distribution of $Y|X$). The data processing inequality states that $\|Tq - Tq'\|_1 \leq \|q - q'\|_1$. The intuition is that if two different information sources ($X \sim q$ and $X \sim q'$) are passed through the same information processing channel ($T(Y|X)$), their difference can only be “erased” and will never increase.

¹Note: this result is strictly inferior to the one in (2). If you choose to work on (2) and are confident about your solution, you can skip this question, as a correct solution to (2) will be automatically counted also as a correct solution to (1).