

---

## Optimality of solutions via LMDPs

Do these two paths lead to the same place?

One of the main questions we have not addressed yet is; if we solve the MDP directly, or solve it via our linear abstraction (linearise, solve and project), do we end up in the same place? This is a question about the completeness of our abstraction. Can our abstraction represent (and find) the same solutions that the original can?

Why does this matter? When we apply our linear abstraction, we want to know: can I trust the answer it has given? If I follow the actions specified by optimal policy, am I going to get rewards?

$$\| V_{\pi^*} - V_{\pi_{u^*}} \|_{\infty} = \epsilon \quad (1)$$

$$= \| (I - \gamma P_{\pi^*})^{-1} r_{\pi^*} - (I - \gamma P_{\pi_{u^*}})^{-1} r_{\pi_{u^*}} \|_{\infty} \quad (2)$$

$$\leq \| (I - \gamma P_{\pi^*})^{-1} r_{\max} - (I - \gamma P_{\pi_{u^*}})^{-1} r_{\min} \|_{\infty} \quad (3)$$

$$= \| \left( (I - \gamma P_{\pi^*})^{-1} - (I - \gamma P_{\pi_{u^*}})^{-1} \right) \Delta r \|_{\infty} \quad (4)$$

$$\leq \Delta r_{\max} \| (I - \gamma P_{\pi^*})^{-1} - (I - \gamma P_{\pi_{u^*}})^{-1} \|_{\infty} \quad (5)$$

$$= \Delta r_{\max} \left\| \sum_{t=0}^{\infty} \gamma^t P_{\pi^*} - \sum_{t=0}^{\infty} \gamma^t P_{\pi_{u^*}} \right\|_{\infty} \quad (6)$$

$$= \Delta r_{\max} \left\| \sum_{t=0}^{\infty} \gamma^t (P_{\pi^*} - P_{\pi_{u^*}}) \right\|_{\infty} \quad (7)$$

$$= \frac{\Delta r_{\max}}{1 - \gamma} \| P_{\pi^*} - P_{\pi_{u^*}} \|_{\infty} \quad (8)$$

(1)

- (1) We want to compare the value achieved by the optimal policy and the value achieved by the optimal linearised solution.
- (2) Assume that there exists a policy that can generate the optimal control dynamics (as given by the LMDP). In that case we can set  $P_{\pi_{u^*}} = U^*$ .
- (3)  $r_{u^*}$  doesn't really make sense as the reward is action dependent. We could calculate it as  $r_{\pi_{u^*}}$ , but we don't explicitly know  $\pi_{u^*}$ .  $(I - \gamma P_{\pi^*})^{-1} r$  represents the action-values, or  $Q$  values. By doing this exchange, we might over estimate the difference under the infinity norm as two non-optimal actions may have larger difference. Also, use the element wise infinity norm.
- (4) Let's assume the optimal policy picks  $\max_a r(s, a)$  at every step. Then the worst case is that we pick  $\min_a r(s, a)$ . Write  $\max_a r(s, a) - \min_a r(s, a) = \Delta r$
- (5) No. Cannot do that... Assumes  $(I - \gamma P_{\pi^*})^{-1} = (I - \gamma P_{\pi_{u^*}})$ . Or that the eqn can be factored.

---

## Notes

- why are we using the infinity norm?!!
- What does  $\delta \geq \text{KL}(U^* \parallel P_{\pi_{u^*}})$  imply about  $\|P_{\pi^*} - P_{\pi_{u^*}}\|_\infty$ ?
- Is it possible to relate  $U^*$  to  $\pi^*$ ?  $\text{KL}(U^* \parallel P_{\pi^*})$
- Use  $r \in [0, 1]$ . Like in other papers. This would simplify and allow us to do (4)?? But need to check it generalises.

---

Alternative.

$$\begin{aligned} & \|V_{\pi^*} - V_{\pi_{u^*}}\|_\infty = \epsilon \quad (1) \\ & = \left\| \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] \right] - \mathbb{E}_{a \sim \pi_{u^*}(\cdot|s)} r(s, a) - \gamma \mathbb{E}_{s' \sim \sum_a P(\cdot|s, a) \pi_{u^*}(a|s)} [V(s')] \right\|_\infty \quad (2) \\ & = \left\| \max_a \left[ r(s, a) - \mathbb{E}_{\tilde{a} \sim \pi_{u^*}(\cdot|s)} r(s, \tilde{a}) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] - \gamma \mathbb{E}_{s' \sim \sum_{\tilde{a}} P(\cdot|s, \tilde{a}) \pi_{u^*}(\tilde{a}|s)} [V(s')] \right] \right\|_\infty \quad (2) \\ & \quad (2) \end{aligned}$$


---

Relies a lot on the result of the projection from optimal state distributions to state-action policies.

$$\begin{aligned} P_\pi(\cdot|s) &= \sum_a P(\cdot|s, a) \pi(a|s) \quad (3) \\ \pi &= \underset{\pi}{\text{argmin}} \text{KL}(u(\cdot|s) \parallel P_\pi(\cdot|s)) \quad (4) \\ \delta &= \text{KL}(u^*(\cdot|s) \parallel P_\pi(\cdot|s)) \quad (5) \\ \delta &= - \sum_{s'} u^*(s'|s) \log \frac{P_\pi(s'|s)}{u^*(s'|s)} \quad (6) \end{aligned}$$

Small delta implies ????. Everywhere  $u^*$  has non-zero probability,  $P_\pi(s'|s)$ , is  $\approx u(s'|s)$

KL

the expected number of extra bits required to code samples from  $P$  using a code optimized for  $Q$  rather than the code optimized for  $P$

Always gives more weight to the high probability (under  $u$ ) states. And because the two distributions

---

are normalised. For all mass that  $P_\pi > u$ , there will be equal mass where  $P_\pi < u$ . But the latter is weighted more. (kinda)

---

want to know the difference between the two policies.

$$\delta = \text{KL}(\pi^* \parallel \pi_{u^*}) \quad (7)$$

(8)

Could then use this to relate  $r_{\pi^*}$  and  $r_{\pi_{u^*}}$ ? But how!?!