

Online and Reinforcement Learning (2025)

Home Assignment 2

Davide Marchi 777881

Contents

1	Short Questions	2
2	MDPs with Similar Parameters Have Similar Values	3
3	Policy Evaluation in RiverSwim	5
4	Solving a Discounted Grid-World	5
5	Off-Policy Evaluation in Episode-Based River-Swim	5

1 Short Questions

Determine whether each statement below is True or False and provide a very brief justification.

1. **Statement:** “In a finite discounted MDP, every possible policy induces a Markov Reward Process.”

Answer: False. This statement assumes that the policy depends only on the current state. If we allow policies to depend on the *entire* past history (*history-dependent* policies), then the resulting transitions in the state space may no longer satisfy the Markov property, since the chosen action at each step might be a function of all previous states and actions. Hence not *every* (fully history-dependent) policy necessarily induces a Markov Reward Process in the *original* state space.

2. **Statement:** “Consider a finite discounted MDP, and assume that π is an optimal policy. Then, the action(s) output by π does not depend on history other than the current state (i.e., π is necessarily stationary).”

Answer: False. While it is true that there *exists* an optimal policy which is stationary deterministic, it does not follow that *all* optimal policies must be so. In fact, multiple distinct policies (some stationary, others possibly history-dependent or randomized) can achieve exactly the same optimal value. Hence it is incorrect to say that *any* optimal policy π must be purely state-dependent (stationary).

3. **Statement:** “In a finite discounted MDP, a greedy policy with respect to optimal action-value function, Q^* , corresponds to an optimal policy.”

Answer: True. From the Bellman optimality equations for Q^* , a policy that selects

$$\arg \max_a Q^*(s, a)$$

at each state s is indeed an optimal policy. This policy attains the same value as Q^* itself, thus achieving the optimal value.

4. **Statement:** “Under the coverage assumption, the Weighted Importance Sampling Estimator \hat{V}_{wIS} converges to V^π with probability 1.”

Answer: True. The coverage assumption ensures that the target policy’s state-action probabilities are absolutely continuous w.r.t. the behavior policy. Under this assumption, Weighted Importance Sampling (though slightly biased) is a *consistent* estimator of V^π , meaning it converges almost surely to V^π as the sample size grows unbounded.

2 MDPs with Similar Parameters Have Similar Values

We recall the setting: two finite discounted MDPs

$$M_1 = (S, A, P_1, R_1, \gamma) \quad \text{and} \quad M_2 = (S, A, P_2, R_2, \gamma),$$

with the same discount factor $0 < \gamma < 1$ and finite state-action space. For each (s, a) we have:

$$|R_1(s, a) - R_2(s, a)| \leq \alpha, \quad \|P_1(\cdot | s, a) - P_2(\cdot | s, a)\|_1 \leq \beta, \quad R_1(s, a), R_2(s, a) \in [0, R_{\max}].$$

Let π be any fixed stationary policy (deterministic or randomized), and let V_1^π, V_2^π denote its value functions in M_1 and M_2 , respectively. We wish to show that, for every $s \in S$,

$$|V_1^\pi(s) - V_2^\pi(s)| \leq \frac{\alpha + \gamma R_{\max} \beta}{(1 - \gamma)^2}.$$

Bellman Operators. Define the Bellman operator T_m^π for each M_m ($m = 1, 2$) by

$$(T_m^\pi V)(s) = \sum_{a \in A} \pi(a | s) \left[R_m(s, a) + \gamma \sum_{s'} P_m(s' | s, a) V(s') \right].$$

Then V_m^π is the unique fixed point: $V_m^\pi = T_m^\pi V_m^\pi$, i.e.

$$V_m^\pi(s) = (T_m^\pi V_m^\pi)(s) = \sum_a \pi(a | s) \left[R_m(s, a) + \gamma \sum_{s'} P_m(s' | s, a) V_m^\pi(s') \right].$$

Step 1: Decompose the difference. For each $s \in S$, we have

$$V_1^\pi(s) - V_2^\pi(s) = (T_1^\pi V_1^\pi)(s) - (T_2^\pi V_2^\pi)(s).$$

Add and subtract $(T_1^\pi V_2^\pi)(s)$ inside the absolute value:

$$\begin{aligned} |V_1^\pi(s) - V_2^\pi(s)| &= |T_1^\pi V_1^\pi(s) - T_2^\pi V_2^\pi(s)| \\ &\leq \underbrace{|T_1^\pi V_1^\pi(s) - T_1^\pi V_2^\pi(s)|}_{(1) \text{ same operator, diff in values}} + \underbrace{|T_1^\pi V_2^\pi(s) - T_2^\pi V_2^\pi(s)|}_{(2) \text{ same value, diff in operators}}. \end{aligned}$$

(1) Same operator, difference in the value functions.

Since T_1^π has discount factor γ ,

$$|T_1^\pi V_1^\pi(s) - T_1^\pi V_2^\pi(s)| = \gamma \sum_a \pi(a | s) \left| \sum_{s'} P_1(s' | s, a) V_1^\pi(s') - \sum_{s'} P_1(s' | s, a) V_2^\pi(s') \right| \leq \gamma \max_x |V_1^\pi(x) - V_2^\pi(x)|$$

Define

$$\Delta = \sup_{s \in S} |V_1^\pi(s) - V_2^\pi(s)|.$$

Thus the first portion is at most $\gamma \Delta$.

(2) Same value function, difference in operators.

Next,

$$\left| T_1^\pi V_2^\pi(s) - T_2^\pi V_2^\pi(s) \right| = \left| \sum_a \pi(a | s) \left[R_1(s, a) + \gamma \sum_{s'} P_1(s' | s, a) V_2^\pi(s') \right] - \sum_a \pi(a | s) \left[R_2(s, a) + \gamma \sum_{s'} P_2(s' | s, a) V_2^\pi(s') \right] \right|$$

We can group the terms:

$$\leq \sum_a \pi(a | s) \left| \underbrace{R_1(s, a) - R_2(s, a)}_{\leq \alpha} + \gamma \sum_{s'} [P_1(s' | s, a) - P_2(s' | s, a)] V_2^\pi(s') \right|.$$

Hence

$$\left| T_1^\pi V_2^\pi(s) - T_2^\pi V_2^\pi(s) \right| \leq \sum_a \pi(a | s) \left[\alpha + \gamma \left| \sum_{s'} [P_1(s' | s, a) - P_2(s' | s, a)] V_2^\pi(s') \right| \right].$$

Since

$$\sum_{s'} |P_1(s' | s, a) - P_2(s' | s, a)| \leq \beta, \quad \text{and} \quad |V_2^\pi(s')| \leq \frac{R_{\max}}{1 - \gamma},$$

it follows that

$$\left| \sum_{s'} [P_1(s' | s, a) - P_2(s' | s, a)] V_2^\pi(s') \right| \leq \beta \frac{R_{\max}}{1 - \gamma}.$$

Hence

$$\left| T_1^\pi V_2^\pi(s) - T_2^\pi V_2^\pi(s) \right| \leq \alpha + \gamma \beta \frac{R_{\max}}{1 - \gamma}.$$

Combine (1) & (2).

Putting the two pieces together:

$$\left| V_1^\pi(s) - V_2^\pi(s) \right| \leq \underbrace{\gamma \Delta}_{\text{from (1)}} + \underbrace{\alpha + \gamma \beta \frac{R_{\max}}{1 - \gamma}}_{\text{from (2)}}.$$

Therefore, taking supremum over $s \in S$:

$$\Delta = \sup_s |V_1^\pi(s) - V_2^\pi(s)| \leq \gamma \Delta + \alpha + \gamma \beta \frac{R_{\max}}{1 - \gamma}.$$

Rearranging gives

$$(1 - \gamma) \Delta \leq \alpha + \gamma \beta \frac{R_{\max}}{1 - \gamma}, \implies \Delta \leq \frac{\alpha}{1 - \gamma} + \frac{\gamma \beta R_{\max}}{(1 - \gamma)^2}.$$

As in the usual derivations, since $\alpha/(1 - \gamma) \leq \alpha/(1 - \gamma)^2$ whenever $0 < \gamma < 1$, one can in fact write

$$\Delta \leq \frac{\alpha + \gamma \beta R_{\max}}{(1 - \gamma)^2}.$$

Thus, for every $s \in S$,

$$\left| V_1^\pi(s) - V_2^\pi(s) \right| \leq \Delta \leq \frac{\alpha + \gamma \beta R_{\max}}{(1 - \gamma)^2}.$$

This completes the alternate proof via Bellman equations.

3 Policy Evaluation in RiverSwim

4 Solving a Discounted Grid-World

5 Off-Policy Evaluation in Episode-Based River-Swim