# Online and Reinforcement Learning (2025)
# Home Assignment 2

Davide Marchi 777881

## Contents

# 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 $\pi$ is an optimal policy. Then, the action(s) output by $\pi$ does not depend on history other than the current state (i.e., $\pi$ 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 $\pi$ must be purely state-dependent (stationary).

3. **Statement:** *"n 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 $\widehat{V}_{\mathrm{wIS}}$ converges to $V^\pi$ 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 $V^\pi$, meaning it converges almost surely to $V^\pi$ as the sample size grows unbounded.

# 2 MDPs with Similar Parameters Have Similar Values

We have 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 $\gamma \in (0, 1)$ and the same finite state–action space $S \times A$. The reward functions satisfy $R_m(s, a) \in [0, R_{\max}]$, and for all $(s, a)$:

$$\left| R_1(s, a) - R_2(s, a) \right| \leq \alpha, \qquad \left\| P_1(\cdot \mid s, a) - P_2(\cdot \mid s, a) \right\|_1 \leq \beta.$$

We let $\pi$ be any fixed *stationary deterministic* (or stationary randomized) policy, and write $V_m^\pi$ to denote its value function in $M_m$. We want to prove:

$$\left| V_1^\pi(s) - V_2^\pi(s) \right| \leq \frac{\alpha + \gamma R_{\max} \beta}{\left(1 - \gamma\right)^2} \quad \text{for every state } s \in S.$$

**Proof of (ii):** Fix $s \in S$. By definition of the value function under policy $\pi$, we have

$$V_m^\pi(s) = \sum_{a \in A} \pi(a \mid s) \left[ R_m(s, a) + \gamma \sum_{x \in S} P_m(x \mid s, a) V_m^\pi(x) \right] \quad \text{for } m = 1, 2.$$

Taking their difference:

$$\left| V_1^\pi(s) - V_2^\pi(s) \right| = \left| \sum_a \pi(a \mid s) \left[ R_1(s, a) + \gamma \sum_x P_1(x \mid s, a) V_1^\pi(x) - \left( R_2(s, a) + \gamma \sum_x P_2(x \mid s, a) V_2^\pi(x) \right) \right] \right|$$

Use the triangle inequality, plus linearity of the sum:

$$\leq \sum_a \pi(a \mid s) \left| \underbrace{R_1(s, a) - R_2(s, a)}_{\leq \alpha} + \gamma \sum_x P_1(x \mid s, a) V_1^\pi(x) - \gamma \sum_x P_2(x \mid s, a) V_2^\pi(x) \right|.$$

Hence

$$\left| V_1^\pi(s) - V_2^\pi(s) \right| \leq \sum_a \pi(a \mid s) \left[ \alpha + \gamma \left| \sum_x P_1(x \mid s, a) V_1^\pi(x) - \sum_x P_2(x \mid s, a) V_2^\pi(x) \right| \right].$$

We now split that big absolute difference into two parts:

$$\left| \sum_x P_1(x \mid s, a) V_1^\pi(x) - \sum_x P_2(x \mid s, a) V_2^\pi(x) \right|$$

$$\leq \left| \sum_x P_1(x \mid s, a) \left( V_1^\pi(x) - V_2^\pi(x) \right) \right| + \left| \sum_x \left( P_1(x \mid s, a) - P_2(x \mid s, a) \right) V_2^\pi(x) \right|$$

$$\leq \sum_x P_1(x \mid s, a) \left| V_1^\pi(x) - V_2^\pi(x) \right| + \sum_x \left| P_1(x \mid s, a) - P_2(x \mid s, a) \right| \left| V_2^\pi(x) \right|.$$

3

Since $\left|V_2^\pi(x)\right| \leq \frac{R_{\max}}{1-\gamma}$ for discounted MDPs, and $\sum_x \left|P_1(x \mid s,a) - P_2(x \mid s,a)\right| \leq \beta$, it follows that:

$$\left|\sum_x P_1(x \mid s,a)\, V_1^\pi(x) \;-\; \sum_x P_2(x \mid s,a)\, V_2^\pi(x)\right| \;\leq\; \sup_x \left|V_1^\pi(x) - V_2^\pi(x)\right| \;+\; \beta\, \frac{R_{\max}}{1-\gamma}.$$

Let

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

Then combining everything above,

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

Taking the supremum in $s$ on the left side gives

$$\Delta \;\leq\; \alpha \;+\; \gamma\,\Delta \;+\; \gamma\,\beta\,\frac{R_{\max}}{1-\gamma}.$$

Hence,

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

Finally, we can note that $\alpha \leq \alpha/(1-\gamma)$, or equivalently multiply out and observe

$$\frac{\alpha}{1-\gamma} \;=\; \frac{\alpha\,(1-\gamma)}{(1-\gamma)^2} \;\leq\; \frac{\alpha}{(1-\gamma)^2}\,.$$

So we get

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

Thus, for every state $s$,

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

This completes the proof of part (ii).

# 3 Policy Evaluation in RiverSwim

# 4 Solving a Discounted Grid-World

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