

# 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	3
4	Solving a Discounted Grid-World	3
5	Off-Policy Evaluation in Episode-Based River-Swim	3

# 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:** “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  $\hat{V}_{\text{wIS}}$  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
- 3 Policy Evaluation in RiverSwim
- 4 Solving a Discounted Grid-World
- 5 Off-Policy Evaluation in Episode-Based River-Swim