

## Importance Sampling – 8 points

Recall the importance sampling ratio

$$\rho_{k,T} = \frac{\prod_{i=k}^T \pi(a_i | s_i)}{\prod_{i=k}^T b(a_i | s_i)}$$

for the target policy  $\pi$  and the behaviour policy  $b$ .

- a. Prove that unweighted importance sampling is unbiased, i.e.

$$\mathbb{E}_b \left[ \frac{\sum_{k \in \mathcal{T}(s_k)} \rho_{k:T(k)} g_k}{|\mathcal{T}(s_k)|} \right] = v_\pi(s_k)$$

where  $\mathcal{T}(s_k)$  are all episodes starting with state  $s_k$ .

- b. Prove that the weighted importance sampling converges, i.e.

$$\mathbb{E}_b \left[ \frac{\sum_{k \in \mathcal{T}(s_k)} \rho_{k:T(k)} g_k}{\sum_{k \in \mathcal{T}(s_k)} \rho_{k:T(k)}} \right]$$

as  $|\mathcal{T}(s_k)| \rightarrow \infty$  Hint: you can use the strong law of large numbers. It states that for  $X_1, X_2, \dots$ , iid samples then

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mu$$

almost surely, that is for almost all realizations except for a zero probability set. Consider suitable augmented numerator and denominator and use the law of large numbers on both of them.

- c. Why is the weighted importance sampling biased?

## Preliminaries for Implementing Learning Agents using the gym Environment

We will use two environments to try out this exercises control problems: (i) the easier `FrozenLake-v1` and (ii) the harder `Blackjack-v0` environment.

You can install both environments using `pip install gym` and load them in python using `env = gym.make('FrozenLake-v1')` and `env = gym.make('Blackjack-v0')` respectively. In order to make diagnostics easier, we will use `tqdm` for writing progress bars. Install with `pip` if package is not present in your system.

The following two exercises shall be implemented by completing the given code skeleton.

## MC Control – 8 points

Implement MC control with importance sampling. The behaviour policy should be the  $\epsilon$ -greedy version of the target policy. Evaluate the algorithm on both environments Save the convergence curve and final greedy performance for both environments and send them in.

## Expected SARSA – 8 points

Implement expected SARSA. Evaluate the algorithm on both environments Save the convergence curve and final greedy performance for both environments and send them in.