

# CS 277 (W26): Control and Reinforcement Learning

## Quiz 8: Inverse RL and Bounded RL

Due date: Monday, March 9, 2026 (Pacific Time)

Roy Fox

<https://royf.org/crs/CS277/W26>

**Instructions:** please solve the quiz in the marked spaces and submit this PDF to Gradescope.

**Question 1** The Inverse RL (IRL) algorithms we saw also find a good policy, in which sense they are like Imitation Learning (IL). Comparing IRL to IL (check all that hold):

- Learning both a reward function and a policy can be an easier problem than only learning a policy.
- IL methods that also learn a reward function are typically more robust to suboptimal demonstrations than those that don't.
- IL methods that also learn a reward function are typically more robust to conflicting or multi-modal demonstrations than those that don't.
- Pre-training with IRL in one environment typically provides a good starting point for IL in another environment with similar but different dynamics, such as in simulation-to-reality transfer.
- Pre-training with IRL in one task typically provides a good starting point for IL in a completely different task with the same environment dynamics.
- None of the above.

**Question 2** In Soft Q-Learning (SQL) (check all that hold):

- As  $\beta \rightarrow 0$ , the algorithm learns a value function  $Q^{\pi_0}$  that evaluates  $\pi_0$ .
- In large action spaces, we can obtain an unbiased estimate of the target value  $r + \frac{\gamma}{\beta} \log \mathbb{E}_{(a'|s') \sim \pi_0} [\exp \beta Q_\beta(s', a')]$  by replacing the expectation with a sample  $(a'|s') \sim \pi_0$ .
- The soft-optimal policy can also be beneficial for exploration.
- When  $\pi_0$  is uniform and  $\beta$  is finite,  $Q_\beta(s, a)$  penalizes actions that lead to future states in which some actions are much better than others.
- None of the above.

**Question 3** In Soft Actor–Critic (SAC) (check all that hold):

- Both SAC and TRPO can be effective for Offline RL because they can both constrain the KL-divergence from the data distribution.
- As  $\beta \rightarrow \infty$  and  $\pi_\theta \rightarrow$  deterministic, the SAC loss approaches the DDPG loss.
- SAC can be applied in large and continuous action spaces because the actor maintains an explicit policy, unlike SQL which is value-based.
- The SAC actor imitates the soft-greedy policy suggested by the critic,  $\frac{\pi_0(a|s) \exp \beta Q_\beta(s,a)}{\exp \beta V(s)}$ . Since the normalizer needs to integrate the numerator, which is infeasible in large or continuous action spaces, SAC needs to also maintain a  $V$  network.
- None of the above.