
Policy Teaching via Environment Poisoning: Training-time Adversarial Attacks against Reinforcement Learning

Prakash K. Naikade
prna00001@stud.uni-saarland.de

1 Summary

The paper investigates an underlying reinforcement learning (RL) security threat in which an attacker poisons the environment, causing the agent to implement a target policy. This paper demonstrates that poisoning transition dynamics optimally is far more difficult than poisoning rewards, and that the attack may not always be possible. The paper lays the theoretical groundwork for environment poisoning against RL in a variety of new attack dimensions like, adversarial manipulation of transition dynamics; attack against RL agents maximizing average reward in an undiscounted infinite horizon; and analyzing different attack costs for offline planning and online learning settings. The paper also provides sufficient technical criteria for the assault to be successful, as well as lower and upper bounds on the attack cost. According to experimental data attacker can easily teach any target policy to the RL agents under mild conditions.

1.1 Overview of attacks

Paper presents attacks in two settings: an offline setting where the agent is doing planning in the poisoned environment and an online setting where the agent is learning a policy using a regret-minimization framework with poisoned feedback.

Goal of the attack against an offline planning agent is defined as,

$$\rho(\pi, \widehat{M}) \geq \rho(\pi, \widehat{M}) + \epsilon, \quad \forall \pi \neq \pi^\dagger \quad (1)$$

Where, ϵ is margin parameter and target policy is ϵ -robust optimal in the poisoned Markov Decision Process (MDP) \widehat{M} .

Cost of the attack against an offline planning agent is defined in terms of ℓ_p -norm of differences in reward functions (i.e., \widehat{R} and \overline{R}) and in transition dynamics (i.e., \widehat{P} and \overline{P}),

$$\|\widehat{R} - \overline{R}\|_p = \left(\sum_{s,a} \left(|\widehat{R}(s,a) - \overline{R}(s,a)| \right)^p \right)^{1/p}, \text{ for reward functions} \quad (2)$$

$$\|\widehat{P} - \overline{P}\|_p = \left(\sum_{s,a} \left(\sum_{s'} |\widehat{P}(s,a,s') - \overline{P}(s,a,s')| \right)^p \right)^{1/p}, \text{ for transition dynamics} \quad (3)$$

Goal of the attack against an online learning agent is defined as notion of *average mismatch* of learner's actions in time horizon T ,

$$(T) = \frac{1}{T} \cdot \left(\sum_{t=0}^{T-1} a_t \neq (s_t) \right) \quad (4)$$

The aim of the attacker is to ensure that (T) is $o(1)$, or, alternatively, the total number of time steps where there is a mismatch is $o(T)$.

Attacks in Offline Setting via Poisoning Rewards is formulated as an optimization problem,

$$\begin{aligned} \min_R \quad & \|R - \bar{R}\|_p \\ \text{s.t.} \quad & \hat{R} = R \text{ and } \hat{\mu} = \bar{\mu} \end{aligned} \quad (\text{P1})$$

The optimization problem (P1) is always feasible, and the cost of the optimum solution is bounded by,

$$\frac{\bar{\alpha}}{2} \cdot \|\bar{\chi}_\epsilon\|_\infty \leq \|\hat{R} - \bar{R}\|_p \leq \|\bar{\chi}_\epsilon\|_p \quad (5)$$

An optimization problem for **Attacks in Offline Setting via Poisoning Dynamics** is formulated as,

$$\min_{P, \mu, \mu^{sa}} \quad \|P - \bar{P}\|_p \quad (\text{P2})$$

If there exists a solution \hat{P} to the optimization problem (P2), its cost satisfies

$$\|\hat{P} - \bar{P}\|_p \cdot \|\bar{V}\|_\infty \geq \frac{\bar{\alpha}}{2} \cdot \|\bar{\chi}_0\|_\infty \quad (6)$$

If for every state s and action a , it holds that either $\bar{\beta}(s, a) \geq \epsilon \cdot (1 + \bar{D})$ OR $\bar{\chi}_\epsilon(s, a) = 0$, then the optimization problem (P2) has a solution \hat{P} whose cost is upper bounded by,

$$\|\hat{P} - \bar{P}\|_p \leq 2 \cdot \|\bar{\Lambda}\|_p \quad (7)$$

An attacker changes the environment feedback (reward r_t or state s_{t+1}) to attack an online learning agent.

2 Strengths

This paper presents environmental poisoning by manipulating both rewards or transition dynamics as compared to previous existing work which only manipulates rewards to perform poison attack on RL agents. This is important step in new direction as attacking specific applications could be efficiently done by manipulating transition dynamics instead of only rewards.

3 Weaknesses

The experimental results show that the examined attack models are effective, but they do not tell which types of learning algorithms are most sensitive to the attack tactics studied in the paper. Relaxing the assumptions about the attacker's understanding of the underlying MDP may result in more resilient attack tactics.

4 Possible Improvements

While the paper treats reward and transition attack models separately, attacking rewards and transitions at the same time may result in more robust and cost-effective solutions. This might have been solved by merging the corresponding constraints and specifying the cost function as a weighted sum of the cost functions used in (P1) and (P2).

5 Possible Extensions

Important possible extension could be expanding the attack models, such as attacking rewards and transitions simultaneously, and widening the range of attack goals, such as under partial specification of target policy.