Adversarial Sequential Decision Making

Goran Radanović, Adish Singla, Wen Sun, Xiaojin Zhu

International Joint Conference on AI (IJCAI) 2022







Outline

- Preliminaries
- Test-time Attacks and Defenses in RL
- Training-time Attacks in RL
- Training-time Defenses in RL
- Adversarial Attacks in Multi-agent RL
- Concluding Remarks

Markov Games



Markov Games



Markov Games



Agents optimize their returns \Rightarrow find an equilibrium

Adversarial Multi-Agent Setting



Adversarial Multi-Agent Setting



Adversarial Policies



[Gleave et al., 2020] See also [Guo et al., 2021]

Adversarial Policies

Test-time attack



Adversarial Policies



Adversary trained with < 3% or time-steps used for training Victim.



Masking adversary's position helps: the victim's win rate increases.

[Gleave et al., 2020]

Backdoor Attacks



[Wang et al., 2021]

Backdoor Attacks



Backdoor Attacks

Victim's winning rate reduces by 17%-37% when the backdoor is triggered

Environment	Triggered			Not Triggered			Benign Baseline		
	Failing	Tie	Winning	Failing	Tie	Winning	Failing	Tie	Winning
Run To Goal (Ants)	73.8%	2.4%	23.8%	45.0%	5.1%	49.9%	46.0%	3.0%	51.0%
Run To Goal (Humans)	20.8%	69.3%	9.9%	52.2%	0.7%	47.1%	51.2%	1.4%	47.4%
You Shall Not Pass (Humans)	83.0%	0.0%	17.0%	50.1%	0.0%	49.9%	50.5%	0.0%	49.5%
Sumo (Humans)	34.4%	54.7%	10.9%	29.7%	42.2%	28.1%	30.1%	34.4%	35.5%

Table 2: The failing/tie/winning rate of the victim agent when the backdoor is triggered (or not). Benign baselines are measured on two normal agents.

Environment/Failing Rate	BACKDOORL	Fine-tuned
Run to Goal (Ants)	23.8%	39.0%
Run to Goal (Humans)	9.9%	5.0%
You Shall Not Pass	17.0%	23.8%
Sumo	10.9%	22.5%

Fine tuning defense not fully successful

Table 6: Winning rate before/after fine-tuning when facing the trigger. Bolded numbers are the higher winning rates.

Adversarial Multi-Agent Setting



Reward Poisoning Attacks





Force a joint target policy π_{\dagger} : Minimally change R^i s.t. π_{\dagger} is an ε -strict Markov perfect **dominant strategy equilibrium**

Reward Poisoning Attacks

- Offline finite-horizon setting:
 - Attacker: Modifies the rewards in a given dataset
 - Agents/Learners: Estimate the parameters of Markov game from the poisoned data
- Q-confidence bound backward induction minimize cost $C(r, r_0)$ while satisfying

$$\underline{Q}_{i,h}(s, (\pi^{i,h}_{\dagger}(s), a^{-i})) \ge \overline{Q}_{i,h}(s, (a^{i}, a^{-i})) + \varepsilon \quad \longleftarrow \quad \begin{array}{c} Equilibrium \\ condition \end{array}$$

- Exponential in #Agents \rightarrow Greedy backward-induction

Summary

- MARL as a framework for physical adversarial attacks
- Attacks and defense in MARL largely unexplored
- Byzantine Attacks in Distributed RL
 - Fault-tolerant Federated RL [Fan et al., 2021]
 - Byzantine-Robust Distributed RL [Chen et al., 2022]



Outline

- Preliminaries
- Test-time Attacks and Defenses in RL
- Training-time Attacks in RL
- Training-time Defenses in RL
- Adversarial Attacks in Multi-agent RL
- Concluding Remarks

Outline

- Preliminaries
- Test-time Attacks and Defenses in RL
- Training-time Attacks in RL
- Training-time Defenses in RL
- Adversarial Attacks in Multi-agent RL
- Concluding Remarks

Concluding Remarks



Concluding Remarks

Attacks/Defenses









تابه کدید Microsoft Created a Twitter Bot to Learn From Users. It Quickly Became a Racist Jerk.

🛱 Give this article 🔗 🗍



Concluding Remarks

Attacks/Defenses









قheكicutHork Eines Microsoft Created a Twitter Bot to Learn From Users. It Quickly Became a Racist Jerk.

🛱 Give this article





Steering/Teaching







Author: Mtnman79

References

- Gleave et al., Adversarial Policies: Attacking Deep Reinforcement Learning, 2020.
- Guo et al., Adversarial Policy Learning in Two-player Competitive Games, 2021.
- Wang et al., BACKDOORL: Backdoor Attack against Competitive Reinforcement Learning, 2021.
- Fan et al., Fault-Tolerant Federated Reinforcement Learning with Theoretical Guarantee, 2021.
- Wu et al., Reward Poisoning Attacks on Offline Multi-Agent Reinforcement Learning, 2022.
- Chen et al., Byzantine-Robust Online and Offline Distributed Reinforcement Learning, 2022.