

Attacker Setting

- Attacker setting:

- The attacker wants the victims to take a target (deterministic) policy $\pi^\dagger = (\pi_1^\dagger, \pi_2^\dagger, \dots, \pi_n^\dagger)$ as often as

possible, i.e. maximize $\sum_{t=1}^T \mathbb{1}_{(a_t = \pi^\dagger)}$.

- The attacker can modify the rewards that the victims see from $r^o(a)$ to $r^\dagger(a)$.
- The attacker wants sublinear design cost

$$\sum_{t=1}^T \left\| r^o(a_t) - r_t^\dagger(a_t) \right\|_p.$$

Interior Design Algorithm

- Given $r^o(a) \in [-1, 1]$, first consider the interior case when $r^o(\pi^\dagger) > -1$.
- Assumption: $r^o(\pi^\dagger) \geq -1 + \rho$, for some $\rho > 0$.

- Attack: $r_{i,t}^\dagger(a) = \begin{cases} r_i^o(\pi^\dagger) + \left(1 - \frac{d(a_t)}{n}\right) \rho & \text{if } a_{i,t} = \pi_i^\dagger \\ r_i^o(\pi^\dagger) - \frac{d(a_t)}{n} \rho & \text{if } a_{i,t} \neq \pi_i^\dagger \end{cases}$,

where $d(a_t) = \sum_{i=1}^n \mathbb{1}_{\{a_{i,t} = \pi_i^\dagger\}}$.

Interior Design Result

Theorem

Using the interior design, π^\dagger is used $T - O(nT^\alpha)$ times while incurring design cost $O(n^{1+1/p}T^\alpha)$.

- For example, EXP3. P with L_1 cost can achieve π^\dagger being used $T - O(n\sqrt{T})$ times with cost $O(n^2\sqrt{T})$.

Interior Design Proof Sketch

- Under this attack, we have,

$$r_{i,t}^{\dagger}(a) = \begin{cases} r_i^{\circ}(\pi^{\dagger}) + \left(1 - \frac{d(a_t)}{n}\right) \rho & \text{if } a_{i,t} = \pi_i^{\dagger} \\ r_i^{\circ}(\pi^{\dagger}) - \frac{d(a_t)}{n} \rho & \text{if } a_{i,t} \neq \pi_i^{\dagger} \end{cases}.$$

- π^{\dagger} is strictly dominant:

$$r_{i,t}^{\dagger}(\pi_{i,t}^{\dagger}, a_{-i,t}) = r_{i,t}^{\dagger}(a_{i,t}, a_{-i,t}) + \left(1 - \frac{1}{n}\right) \rho, \forall a_{i,t} \neq \pi_{i,t}^{\dagger}.$$

- π^{\dagger} rewards remain unchanged: $r_{i,t}^{\dagger}(\pi^{\dagger}) = r_i^{\circ}(\pi^{\dagger})$.

- No-regret learners will use the optimal action profile π^{\dagger} in all but $O(T^{\alpha})$ rounds while incurring $O(T^{\alpha})$ design cost.

Boundary Design Example

- When $r^o(\pi^\dagger) = -1$, it is impossible to decrease other entries below -1 : another design is needed.
- Suppose again $\pi^\dagger = (1, 1)$, then,

$$r^o = \begin{bmatrix} (-1, -1) & (-1, \boxed{1}) & (\boxed{1}, -1) \\ (\boxed{1}, -1) & (-1, -1) & (-1, \boxed{1}) \\ (-1, \boxed{1}) & (\boxed{1}, -1) & (-1, -1) \end{bmatrix},$$

$$r_1^\dagger \approx \begin{bmatrix} (\boxed{-0.8}, \boxed{-0.8}) & (\boxed{-0.7}, -0.9) & (\boxed{-0.7}, -0.9) \\ (-0.9, \boxed{-0.7}) & (-1, -1) & (-1, -1) \\ (-0.9, \boxed{-0.7}) & (-1, -1) & (-1, -1) \end{bmatrix},$$

Boundary Design Algorithm

- Assumption: $r^o(\pi^\dagger) = -1$.
- Attack: $r_{i,t}^\dagger(a) = w_t r_{i,\text{interior}}^\dagger(a) + (1 - w_t) r^o(\pi^\dagger)$, where $w_t = t^{\alpha+\varepsilon-1}$, for some $\varepsilon \in (0, 1 - \alpha]$.

Boundary Design Result

Theorem

Using the boundary design with $\varepsilon = \frac{1 - \alpha}{2}$, π^\dagger is used $T = O(nT^{(1+\alpha)/2})$ times while incurring design cost $O(n^{1/p}(1+n)T^{(1+\alpha)/2})$.

Boundary Design Proof Sketch

- Under this attack, we have,

$$r_{i,t}^\dagger(a) = w_t r_{i,\text{interior}}^\dagger(a) + (1 - w_t) r_i^o(\pi^\dagger), \text{ where } w_t = t^{\alpha+\varepsilon-1}.$$

- 1 π^\dagger is strictly dominant:

$$r_{i,t}^\dagger(\pi_{i,t}^\dagger, a_{-i,t}) = r_{i,t}^\dagger(a_{i,t}, a_{-i,t}) - \left(1 - \frac{1}{n}\right) \rho w_t, \forall a_{i,t} \neq \pi_{i,t}^\dagger.$$

- 2 π^\dagger rewards are almost unchanged:

$$\left\| r_{i,t}^\dagger(\pi^\dagger) - r_i^o(\pi^\dagger) \right\|_p \leq 2bn^{1/p} w_t.$$

- No-regret learners will use the optimal action profile π^\dagger in all but $O(T^{(1+\alpha)/2})$ rounds while incurring $O(T^{(1+\alpha)/2})$ design cost.

Nash Attack

- Victim setting:

- 1 The victims are uncertainty-aware offline learners that use additive bonus terms β when estimating the Q function, i.e. $Q = \hat{R} - \beta + \mathbb{E}_{\hat{P}}[V']$.

- 2 The victims learn a two-player zero-sum Markov game from a training set $\left\{ \left(\left(s_t^{(k)}, a_t^{(k)}, r_t^{(k)} \right)_{t=1}^T \right) \right\}_{k=1}^K$, with $r_t^{(k)} \in [0, 1]$.

Attacker Setting

- Attacker setting:

- 1 The attacker wants the victims to learn a target (deterministic) policy π^\dagger as the unique Markov perfect (Nash) equilibrium.

- 2 The attacker can modify the rewards in the training set from r^o to r^\dagger .

- 3 The attacker minimizes the reward modification cost

$$\|r^\dagger - r^o\|, \text{ e.g. } \sum_{k=1}^K \sum_{t=1}^T \|r_t^{\dagger, (k)} - r_t^{o, (k)}\|_1.$$

- 4 The attacker does not know \hat{R} and \hat{P} , but assumes $\|\hat{R} - R^{(\text{MLE})}\| < \rho^{(R)}$ and $\|\hat{P} - P^{(\text{MLE})}\|_1 < \rho^{(P)}$.

iNash Formulation

- The attack can be formulated as

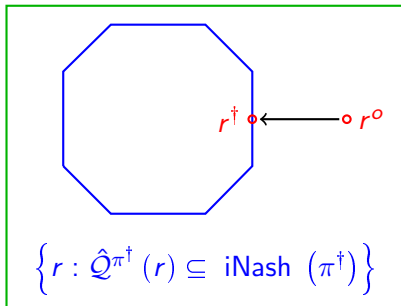
$$\begin{aligned} \min_{r^\dagger} & \|r^\dagger - r^o\| \\ \text{s.t.} & \hat{Q}^{\pi^\dagger}(r^\dagger; \rho^{(R)}, \rho^{(P)}) \subseteq \text{iNash}(\pi^\dagger), \end{aligned}$$

where,

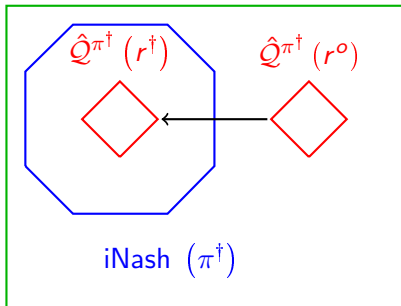
- 1 $\hat{Q}^\pi(r)$ is the set of plausible Q functions computed based on r evaluated on π ,
- 2 $\text{iNash}(\pi)$ is the inverse Nash polytope of Q functions such that π is the strict Markov perfect (Nash) equilibrium.

iNash Diagram

Space of Data Sets



Space of Q Functions



Feasibility

Theorem

The attack is feasible if $\rho_t^{(R)}(s, a) + |\beta_t(s, a)| < \frac{1}{4T}$, $\forall t, s$, and actions a such that $a_1 = \pi_{1,t}^\dagger(s)$ or $a_2 = \pi_{2,t}^\dagger(s)$.

- For example, if $\rho^{(R)} = 0$ and $\beta = \frac{c}{\sqrt{N_t(s, a)}}$, then the condition is a data coverage condition, $N_t(s, a) > 16cT^2$ for actions profiles in the same row or column as π^\dagger in the stage game matrices.

Feasible Example

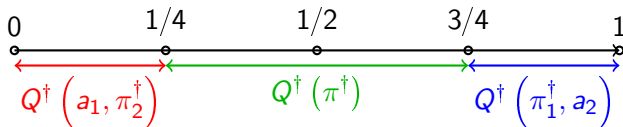
- Suppose $\pi^\dagger = (1, 1)$ in a stage game, then the following attack is feasible under the previous feasibility condition,

$a_1 \backslash a_2$	1	2	3	4
1	0.5	1	1	1
2	0	-	-	-
3	0	-	-	-
4	0	-	-	-

- Unspecified cells' corresponding rewards do not need to be poisoned.

Feasibility Proof Sketch

- The condition $\rho_t^{(R)}(s, a) + |\beta_t(s, a)| < 1/(4T)$ implies that the cumulated confidence interval width for R and P in the future periods is bounded by $1/4$.
- In period t , state s , for every $a_1 \neq \pi_1^\dagger$ and $a_2 \neq \pi_2^\dagger$, the Q values have the following relationship.



- Therefore, $\pi_t^\dagger(s)$ is the strict, thus unique, Nash equilibrium in every stage game (t, s) .

Linear Program Formulation

- The attacker's problem is given by,

$$\min_{r^\dagger} \sum_{k=1}^K \sum_{t=1}^T \left\| r_t^{\dagger, (k)} - r_t^{o, (k)} \right\|_1$$

s.t. for every t, s , and $Q_t^\dagger \in \hat{Q}^{\pi^\dagger}(r^\dagger)$,

$$Q_t^\dagger(s, \pi_t^\dagger(s)) > Q_t^\dagger(s, (a_1, \pi_{t,2}^\dagger(s))), \forall a_1 \neq \pi_{t,1}^\dagger(s),$$

$$Q_t^\dagger(s, \pi_t^\dagger(s)) < Q_t^\dagger(s, (\pi_{t,1}^\dagger(s), a_2)), \forall a_2 \neq \pi_{t,2}^\dagger(s).$$

- Since $\hat{Q}^\pi(r)$ are polytopes, this problem can be formulated as a linear program and solved efficiently.

Multi Attacker

- Incomplete, joint work ($\approx 75\%$ contribution) with Elliot Pickens, Jin-Yi Cai, and Jerry Zhu.
- Victim setting:
 - ① Single or multiple identical victims that estimates the mean $\hat{\mu}$ of a data set, based on a training provided by the attackers.

Attacker Setting, Direction, Continuous

- Attacker Setting 1.1:

- ① Each of K attackers has a target direction x_k^\dagger with the goal of minimizing $(x_k^\dagger)^T \hat{\mu}$.
 - ② Each attacker creates a training set with X_k from a convex and compact domain X , and the (disjoint) union of $\{X_k\}_{k=1}^K$ is given to the victim.
- The game has a (weakly) dominant strategy equilibrium, in which the attackers choose the most extreme points in X in the x^\dagger directions.

Attacker Setting, Direction, Discrete

- Attacker Setting 1.2:

- ① Each of K attackers has a target direction x_k^\dagger with the goal of minimizing $(x_k^\dagger)^T \hat{\mu}$.
 - ② Each attacker creates a training set with X_k from n existing points X , and the (disjoint) union of $\{X_k\}_{k=1}^K$ is given to the victim.
- The game has a (weakly) dominant strategy equilibrium, in which the attackers choose the most extreme points in X in the x^\dagger directions.

Attacker Setting, Point, Continuous

- Attacker Setting 2.1:
 - 1 Each of K attackers has a target point x_k^\dagger with the goal of minimizing $\|x_k^\dagger - \hat{\mu}\|^2$.
 - 2 Each attacker creates a training set with X_k from a convex and compact domain X , and the (disjoint) union of $\{X_k\}_{k=1}^K$ is given to the victim.
- The game has at least one pure strategy Nash equilibrium, and it can be found using
 - 1 Best response dynamics,
 - 2 Maximizing a (weakly) concave potential function on convex and compact X .

Attacker Setting, Point, Discrete

- Attacker Setting 2.2:
- ① Each of K attackers has a target point x_k^\dagger with the goal of minimizing $\|x_k^\dagger - \hat{\mu}\|^2$.
- ② Each attacker creates a training set with X_k from n existing points X , and the (disjoint) union of $\{X_k\}_{k=1}^K$ is given to the victim.
- The game has at least one pure strategy Nash equilibrium, and it can be found using
 - ① Best response dynamics,
 - ② Maximizing a potential function on finite X .

Potential Function Formulation

- The payoff to attacker k can be written as

$$- \left\| x_k^\dagger - \left(x_0 + \sum_{k=1}^K x_k \right) \right\|^2,$$

where x_k is the centroid of the points provided by attacker k .

- The potential function is

$$- \sum_{k=1}^K \|w_k x_k\|^2$$
$$- 2 \sum_{i \neq j} \left(x_i^\dagger - \left(x_0 + \sum_{k \neq i} w_k x_k \right) \right) \left(x_j^\dagger - \left(x_0 + \sum_{k \neq j} w_k x_k \right) \right).$$

References I

- [1] Pieter Abbeel and Andrew Y Ng.
Apprenticeship learning via inverse reinforcement learning.
In Proceedings of the twenty-first international conference on Machine learning, page 1, 2004.
- [2] Natalia Akchurina.
Multiagent reinforcement learning: algorithm converging to nash equilibrium in general-sum discounted stochastic games.
In Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems-Volume 2, pages 725–732, 2009.

References II

- [3] Ashton Anderson, Yoav Shoham, and Alon Altman.
Internal implementation.
In *Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems: volume 1-Volume 1*, pages 191–198. Citeseer, 2010.
- [4] Gautam Appa.
On the uniqueness of solutions to linear programs.
Journal of the Operational Research Society, 53:1127–1132, 2002.
- [5] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire.
The nonstochastic multiarmed bandit problem.
SIAM journal on computing, 32(1):48–77, 2002.

References III

- [6] Avraham Bab and Ronen I Brafman.
Multi-agent reinforcement learning in common interest and fixed sum stochastic games: An experimental study.
Journal of Machine Learning Research, 9(12), 2008.
- [7] Kiarash Banhashem, Adish Singla, and Goran Radanovic.
Defense against reward poisoning attacks in reinforcement learning.
arXiv preprint arXiv:2102.05776, 2021.
- [8] Kiarash Banhashem, Adish Singla, Jiarui Gan, and Goran Radanovic.
Admissible policy teaching through reward design.
arXiv preprint arXiv:2201.02185, 2022.

References V

- [11] Ilija Bogunovic, Arpan Losalka, Andreas Krause, and Jonathan Scarlett.
Stochastic linear bandits robust to adversarial attacks.
In *International Conference on Artificial Intelligence and Statistics*, pages 991–999. PMLR, 2021.
- [12] Michael Bowling.
Convergence problems of general-sum multiagent reinforcement learning.
In *ICML*, pages 89–94, 2000.
- [13] Michael Bowling and Manuela Veloso.
Rational and convergent learning in stochastic games.
In *International joint conference on artificial intelligence*, volume 17, pages 1021–1026. Lawrence Erlbaum Associates Ltd, 2001.

References VI

- [14] David Brandfonbrener, Will Whitney, Rajesh Ranganath, and Joan Bruna.
Offline rl without off-policy evaluation.
Advances in Neural Information Processing Systems, 34:4933–4946, 2021.
- [15] Noam Brown, Tuomas Sandholm, and Strategic Machine.
Libratus: The superhuman ai for no-limit poker.
In *IJCAI*, pages 5226–5228, 2017.
- [16] Noam Brown and Tuomas Sandholm.
Superhuman ai for multiplayer poker.
Science, 365(6456):885–890, 2019.

References VIII

- [21] William J. Clancey.
Transfer of Rule-Based Expertise through a Tutorial Dialogue.
Ph.D. diss., Dept. of Computer Science, Stanford Univ.,
Stanford, Calif., 1979.
- [22] William J. Clancey.
Communication, Simulation, and Intelligent Agents:
Implications of Personal Intelligent Machines for Medical
Education.
*In Proceedings of the Eighth International Joint Conference
on Artificial Intelligence (IJCAI-83)*, pages 556–560, Menlo
Park, Calif, 1983. IJCAI Organization.

References X

- [27] Gregory Dudek, Michael RM Jenkin, Evangelos Miliос, and David Wilkes.
A taxonomy for multi-agent robotics.
Autonomous Robots, 3(4):375–397, 1996.
- [28] Robert Engelmοre and Anthony Morgan, editors.
Blackboard Systems.
Addison-Wesley, Reading, Mass., 1986.
- [29] Wikipedia contributors.
Volunteer’s dilemma — Wikipedia, the free encyclopedia,
2021.
[Online; accessed 16-September-2021].

References XI

- [30] Wei Fu, Chao Yu, Zelai Xu, Jiaqi Yang, and Yi Wu.
Revisiting some common practices in cooperative multi-agent reinforcement learning.
arXiv preprint arXiv:2206.07505, 2022.

- [31] Evrard Garcelon, Baptiste Roziere, Laurent Meunier, Olivier Teytaud, Alessandro Lazaric, and Matteo Pirotta.
Adversarial attacks on linear contextual bandits.
arXiv preprint arXiv:2002.03839, 2020.

- [32] Adam Gleave, Michael Dennis, Cody Wild, Neel Kant, Sergey Levine, and Stuart Russell.
Adversarial policies: Attacking deep reinforcement learning.
arXiv preprint arXiv:1905.10615, 2019.

References XII

- [33] RA Good.
f-finger morra.
SIAM Review, 7(1):81–87, 1965.
- [34] Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine.
Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates.
In 2017 IEEE international conference on robotics and automation (ICRA), pages 3389–3396. IEEE, 2017.

References XV

- [41] Pablo Hernandez-Leal and Michael Kaisers.
Towards a fast detection of opponents in repeated stochastic games.
In International Conference on Autonomous Agents and Multiagent Systems, pages 239–257. Springer, 2017.
- [42] GA Heuer.
Uniqueness of equilibrium points in bimatrix games.
International Journal of Game Theory, 8:13–25, 1979.
- [43] Junling Hu and Michael P Wellman.
Nash q-learning for general-sum stochastic games.
Journal of machine learning research, 4(Nov):1039–1069, 2003.

References XVII

- [46] Lantao Yu, Jiaming Song, and Stefano Ermon. Multi-agent adversarial inverse reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7194–7201. PMLR, 09–15 Jun 2019.
- [47] Haoqi Zhang and David Parkes. Value-based policy teaching with active indirect elicitation. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 1*, AAI'08, page 208–214. AAAI Press, 2008.

References XIX

- [51] Ying Jin, Zhuoran Yang, and Zhaoran Wang.
Is pessimism provably efficient for offline rl?
In International Conference on Machine Learning, pages 5084–5096. PMLR, 2021.
- [52] Kwang-Sung Jun, Lihong Li, Yuzhe Ma, and Jerry Zhu.
Adversarial attacks on stochastic bandits.
Advances in Neural Information Processing Systems, 31:3640–3649, 2018.
- [53] Panagiota Kiourti, Kacper Wardega, Susmit Jha, and Wenchao Li.
Trojdr: evaluation of backdoor attacks on deep reinforcement learning.
In 2020 57th ACM/IEEE Design Automation Conference (DAC), pages 1–6. IEEE, 2020.

References XX

- [54] Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- [55] Jernej Kos and Dawn Song. Delving into adversarial attacks on deep policies. *arXiv preprint arXiv:1705.06452*, 2017.
- [56] Erich Kutschinski, Thomas Uthmann, and Daniel Polani. Learning competitive pricing strategies by multi-agent reinforcement learning. *Journal of Economic Dynamics and Control*, 27(11-12):2207–2218, 2003.

References XXI

- [57] Jae Won Lee and Jangmin O.
A multi-agent q-learning framework for optimizing stock trading systems.
In International Conference on Database and Expert Systems Applications, pages 153–162. Springer, 2002.
- [58] Jae Won Lee, Jonghun Park, O Jangmin, Jongwoo Lee, and Euyseok Hong.
A multiagent approach to q -learning for daily stock trading.
IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans, 37(6):864–877, 2007.

References XXII

- [59] Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel.
Multi-agent reinforcement learning in sequential social dilemmas.
arXiv preprint arXiv:1702.03037, 2017.
- [60] Xiaomin Lin, Peter A Beling, and Randy Cogill.
Multiagent inverse reinforcement learning for two-person zero-sum games.
IEEE Transactions on Games, 10(1):56–68, 2017.
- [61] Yen-Chen Lin, Zhang-Wei Hong, Yuan-Hong Liao, Meng-Li Shih, Ming-Yu Liu, and Min Sun.
Tactics of adversarial attack on deep reinforcement learning agents.
arXiv preprint arXiv:1703.06748, 2017.

References XXV

[67] Guanlin Liu and Lifeng Lai.
Provably efficient black-box action poisoning attacks against reinforcement learning.
Advances in Neural Information Processing Systems, 34, 2021.

[68] Yunlong Lu and Kai Yan.
Algorithms in multi-agent systems: a holistic perspective from reinforcement learning and game theory.
arXiv preprint arXiv:2001.06487, 2020.

[69] Shiyin Lu, Guanghui Wang, and Lijun Zhang.
Stochastic graphical bandits with adversarial corruptions.
In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 8749–8757, 2021.

References XXVI

- [70] Thodoris Lykouris, Max Simchowitz, Alex Slivkins, and Wen Sun.
Corruption-robust exploration in episodic reinforcement learning.
In Conference on Learning Theory, pages 3242–3245. PMLR, 2021.
- [71] Yuzhe Ma, Kwang-Sung Jun, Lihong Li, and Xiaojin Zhu.
Data poisoning attacks in contextual bandits.
In International Conference on Decision and Game Theory for Security, pages 186–204. Springer, 2018.
- [72] Yuzhe Ma, Xuezhou Zhang, Wen Sun, and Jerry Zhu.
Policy poisoning in batch reinforcement learning and control.
Advances in Neural Information Processing Systems, 32:14570–14580, 2019.

References XXVII

[73] Yuzhe Ma, Young Wu, and Xiaojin Zhu.
Game redesign in no-regret game playing.
arXiv preprint arXiv:2110.11763, 2021.

[74] Liam MacDermed, Charles Isbell, and Lora Weiss.
Markov games of incomplete information for multi-agent reinforcement learning.
In Workshops at the Twenty-Fifth AAAI Conference on Artificial Intelligence, 2011.

[75] Olvi Mangasarian.
Uniqueness of solution in linear programming.
Technical report, University of Wisconsin-Madison Department of Computer Sciences, 1978.

References XXVIII

- [76] Patrick Mannion, Karl Mason, Sam Devlin, Jim Duggan, and Enda Howley.
Dynamic economic emissions dispatch optimisation using multi-agent reinforcement learning.
In Proceedings of the Adaptive and Learning Agents workshop (at AAMAS 2016), 2016.
- [77] Eric Maskin and Jean Tirole.
Markov perfect equilibrium: I. observable actions.
Journal of Economic Theory, 100(2):191–219, 2001.

References XXIX

- [78] Linghui Meng, Muning Wen, Yaodong Yang, Chenyang Le, Xiyun Li, Weinan Zhang, Ying Wen, Haifeng Zhang, Jun Wang, and Bo Xu.
Offline pre-trained multi-agent decision transformer: One big sequence model conquers all starcraftii tasks.
arXiv preprint arXiv:2112.02845, 2021.
- [79] CB Millham.
Constructing bimatrix games with special properties.
Naval Research Logistics Quarterly, 19(4):709–714, 1972.
- [80] Junichi Minagawa.
On the uniqueness of nash equilibrium in strategic-form games.
Journal of Dynamics & Games, 7(2):97, 2020.

References XXX

- [81] Lin Yang, Mohammad Hajiesmaili, Mohammad Sadegh Talebi, John Lui, and Wing Shing Wong.
Adversarial bandits with corruptions: Regret lower bound and no-regret algorithm.
In Advances in Neural Information Processing Systems (NeurIPS), 2021.
- [82] Xiaomin Lin, Stephen C. Adams, and Peter A. Beling.
Multi-agent inverse reinforcement learning for certain general-sum stochastic games.
Journal of Artificial Intelligence Research, 66:473–502, oct 2019.

References XXXI

- [83] Sharada Mohanty, Erik Nygren, Florian Laurent, Manuel Schneider, Christian Scheller, Nilabha Bhattacharya, Jeremy Watson, Adrian Egli, Christian Eichenberger, Christian Baumberger, et al.
Flatland-rl: Multi-agent reinforcement learning on trains.
arXiv preprint arXiv:2012.05893, 2020.
- [84] Dov Monderer and Moshe Tennenholtz.
k-implementation.
In Proceedings of the 4th ACM conference on Electronic Commerce, pages 19–28, 2003.
- [85] Dov Monderer and Moshe Tennenholtz.
k-implementation.
Journal of Artificial Intelligence Research, 21:37–62, 2004.

References XXXII

- [86] John F Nash Jr.
Equilibrium points in n-person games.
Proceedings of the national academy of sciences,
36(1):48–49, 1950.
- [87] Norihiko Ono and Kenji Fukumoto.
A modular approach to multi-agent reinforcement learning.
In *Workshop on Learning in Distributed Artificial Intelligence
Systems, Workshop on Learning, Interaction, and
Organization in Multiagent Environments*, pages 25–39.
Springer, Berlin, Heidelberg, 1997.
- [88] Martin J Osborne.
An introduction to game theory, volume 3.
Oxford university press New York, 2004.

References XXXIII

- [89] Ling Pan, Longbo Huang, Tengyu Ma, and Huazhe Xu. Plan better amid conservatism: Offline multi-agent reinforcement learning with actor rectification. In *International Conference on Machine Learning*, pages 17221–17237. PMLR, 2022.
- [90] Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommaman, and Girish Chowdhary. Robust deep reinforcement learning with adversarial attacks. *arXiv preprint arXiv:1712.03632*, 2017.
- [91] HL Prasad and Shalabh Bhatnagar. A study of gradient descent schemes for general-sum stochastic games. *arXiv preprint arXiv:1507.00093*, 2015.

References XXXIV

- [92] HL Prasad, Prashanth LA, and Shalabh Bhatnagar.
Two-timescale algorithms for learning nash equilibria in general-sum stochastic games.
In Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems, pages 1371–1379, 2015.
- [93] Luis G Quintas.
Uniqueness of Nash equilibrium points in bimatrix games.
Center for Mathematical Studies in Economics and Management Science., 1988.
- [94] Arthur L. Robinson.
New ways to make microcircuits smaller.
Science, 208(4447):1019–1022, 1980.

References XXXV

- [95] Arthur L. Robinson.
New Ways to Make Microcircuits Smaller—Duplicate Entry.
Science, 208:1019–1026, 1980.
- [96] James Rice.
Poligon: A System for Parallel Problem Solving.
Technical Report KSL-86-19, Dept. of Computer Science,
Stanford Univ., 1986.
- [97] TES Raghavan.
Zero-sum two-person games.
Handbook of game theory with economic applications,
2:735–768, 1994.

References XXXVI

- [98] Amin Rakhsha, Goran Radanovic, Rati Devidze, Xiaojin Zhu, and Adish Singla.
Policy teaching via environment poisoning: Training-time adversarial attacks against reinforcement learning.
In International Conference on Machine Learning, pages 7974–7984. PMLR, 2020.
- [99] Amin Rakhsha, Goran Radanovic, Rati Devidze, Xiaojin Zhu, and Adish Singla.
Policy teaching in reinforcement learning via environment poisoning attacks.
Journal of Machine Learning Research, 22(210):1–45, 2021.

References XXXVII

- [100] Amin Rakhsha, Xuezhou Zhang, Xiaojin Zhu, and Adish Singla.

Reward poisoning in reinforcement learning: Attacks against unknown learners in unknown environments.

arXiv preprint arXiv:2102.08492, 2021.

- [101] Anshuka Rangi, Long Tran-Thanh, Haifeng Xu, and Massimo Franceschetti.

Saving stochastic bandits from poisoning attacks via limited data verification.

In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 8054–8061, 2022.

References XXXVIII

- [102] Anshuka Rangi, Haifeng Xu, Long Tran-Thanh, and Massimo Franceschetti.
Understanding the limits of poisoning attacks in episodic reinforcement learning.
In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 3394–3400. International Joint Conferences on Artificial Intelligence Organization, 7 2022.
Main Track.
- [103] Tummalapalli Sudhamsh Reddy, Vamsikrishna Gopikrishna, Gergely Zaruba, and Manfred Huber.
Inverse reinforcement learning for decentralized non-cooperative multiagent systems.
In *2012 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 1930–1935. IEEE, 2012.

References XXXIX

- [104] Martin Riedmiller, Thomas Gabel, Roland Hafner, and Sascha Lange.
Reinforcement learning for robot soccer.
Autonomous Robots, 27:55–73, 2009.
- [105] Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel.
Adversarial attacks on neural network policies, 2017.
- [106] Jernej Kos and Dawn Song.
Delving into adversarial attacks on deep policies, 2017.
- [107] Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua.
Safe, multi-agent, reinforcement learning for autonomous driving.
arXiv preprint arXiv:1610.03295, 2016.

References XLIII

[115] Yanchao Sun, Ruijie Zheng, Yongyuan Liang, and Furong Huang.

Who is the strongest enemy? towards optimal and efficient evasion attacks in deep rl.

arXiv preprint arXiv:2106.05087, 2021.

[116] Peter Szilágyi.

On the uniqueness of the optimal solution in linear programming.

Revue d'analyse numérique et de théorie de l'approximation, 35(2):225–244, 2006.

References XLIX

- [128] Fan Wu, Linyi Li, Zijian Huang, Yevgeniy Vorobeychik, Ding Zhao, and Bo Li.
Crop: Certifying robust policies for reinforcement learning through functional smoothing.
arXiv preprint arXiv:2106.09292, 2021.
- [129] Fan Wu, Linyi Li, Chejian Xu, Huan Zhang, Bhavya Kailkhura, Krishnaram Kenthapadi, Ding Zhao, and Bo Li.
Copa: Certifying robust policies for offline reinforcement learning against poisoning attacks.
arXiv preprint arXiv:2203.08398, 2022.

References L

[130] Young Wu, Jeremy McMahan, Xiaojin Zhu, and Qiaomin Xie.

Reward poisoning attacks on offline multi-agent reinforcement learning.

In *The Thirty-Seventh AAAI Conference on Artificial Intelligence (AAAI)*, 2023.

[131] Young Wu, Jeremy McMahan, Xiaojin Zhu, and Qiaomin Xie.

On faking a Nash equilibrium.

arXiv preprint arXiv:2306.08041, 2023.

References LI

[132] Young Wu, Jeremy McMahan, Xiaojin Zhu, and Qiaomin Xie.

Reward poisoning attacks on offline multi-agent reinforcement learning.

In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pages 10426–10434, 2023.

[133] Qiaomin Xie, Yudong Chen, Zhaoran Wang, and Zhuoran Yang.

Learning zero-sum simultaneous-move markov games using function approximation and correlated equilibrium.

In Conference on learning theory, pages 3674–3682. PMLR, 2020.

References LII

- [134] Yang Yang, Li Juntao, and Peng Lingling.
Multi-robot path planning based on a deep reinforcement learning dqn algorithm.
CAAI Transactions on Intelligence Technology, 5(3):177–183, 2020.
- [135] Haoqi Zhang and David C Parkes.
Value-based policy teaching with active indirect elicitation.
In *AAAI*, volume 8, pages 208–214, 2008.
- [136] Haoqi Zhang, David C Parkes, and Yiling Chen.
Policy teaching through reward function learning.
In *Proceedings of the 10th ACM conference on Electronic commerce*, pages 295–304, 2009.

References LIII

- [137] Xuezhou Zhang, Yuzhe Ma, Adish Singla, and Xiaojin Zhu. Adaptive reward-poisoning attacks against reinforcement learning.
In International Conference on Machine Learning, pages 11225–11234. PMLR, 2020.
- [138] Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui Hsieh. Robust deep reinforcement learning against adversarial perturbations on state observations.
Advances in Neural Information Processing Systems, 33:21024–21037, 2020.
- [139] Xuezhou Zhang, Yiding Chen, Jerry Zhu, and Wen Sun. Corruption-robust offline reinforcement learning.
arXiv preprint arXiv:2106.06630, 2021.

References LIV

- [140] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar.
Multi-agent reinforcement learning: A selective overview of theories and algorithms.
Handbook of Reinforcement Learning and Control, pages 321–384, 2021.
- [141] Xuezhou Zhang, Yiding Chen, Xiaojin Zhu, and Wen Sun.
Robust policy gradient against strong data corruption.
In International Conference on Machine Learning, pages 12391–12401. PMLR, 2021.
- [142] Stephan Zheng, Alexander Trott, Sunil Srinivasa, Nikhil Naik, Melvin Gruesbeck, David C Parkes, and Richard Socher.
The ai economist: Improving equality and productivity with ai-driven tax policies.
arXiv preprint arXiv:2004.13332, 2020.

References LV

- [143] Han Zhong, Wei Xiong, Jiyuan Tan, Liwei Wang, Tong Zhang, Zhaoran Wang, and Zhuoran Yang.
Pessimistic minimax value iteration: Provably efficient equilibrium learning from offline datasets.
arXiv preprint arXiv:2202.07511, 2022.
- [144] Shiliang Zuo.
Near optimal adversarial attack on ucb bandits.
arXiv preprint arXiv:2008.09312, 2020.
- [145] Avrim Blum and Yishay Mansour.
Learning, regret minimization, and equilibria.
Algorithmic Game Theory, 2007.