

A Game-Theoretic Perspective on Risk-Sensitive Reinforcement Learning

Mathieu Godbout, Maxime Heuillet, Sharath Chandra,
Rupali Bhati, Audrey Durand



UNIVERSITÉ
LAVAL

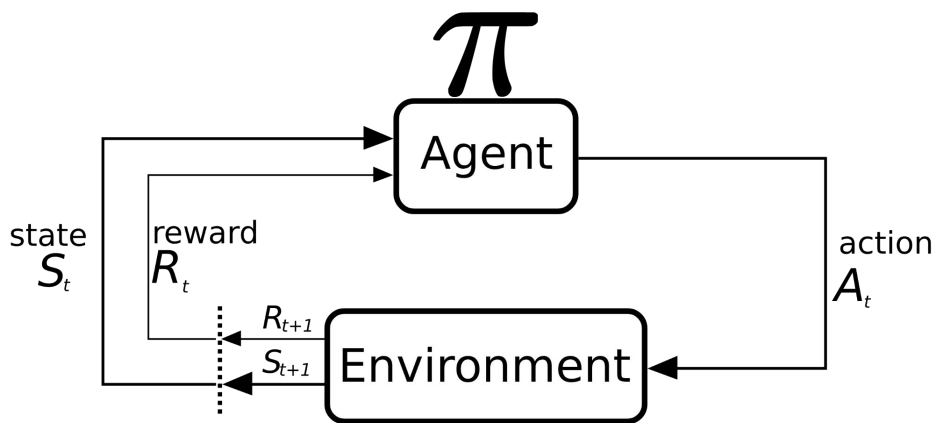


Institut
intelligence
et données

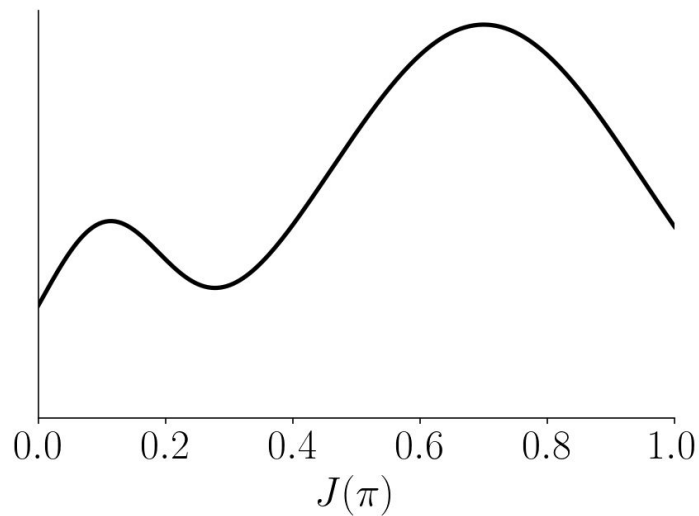


Mila

Reinforcement Learning (RL)

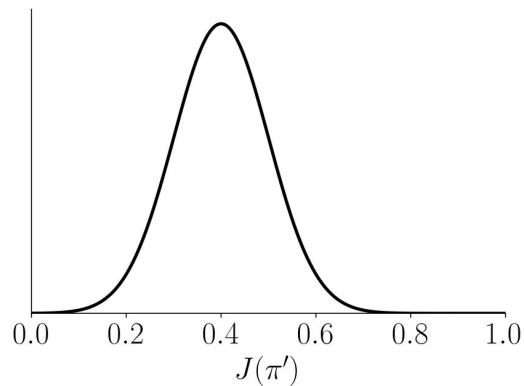
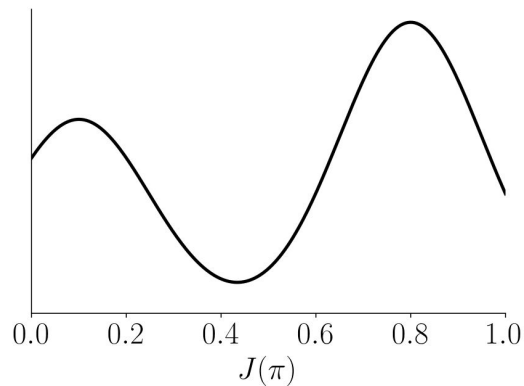


$$J(\pi) := \sum_{t=0}^{\infty} \gamma^t R_t$$



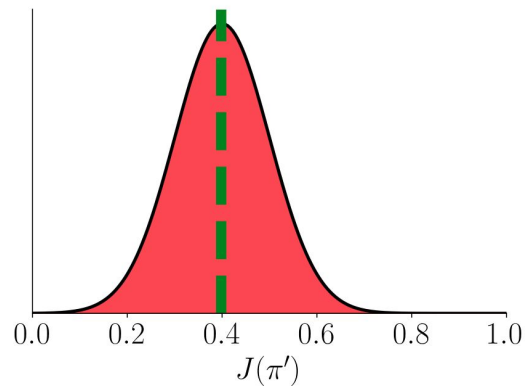
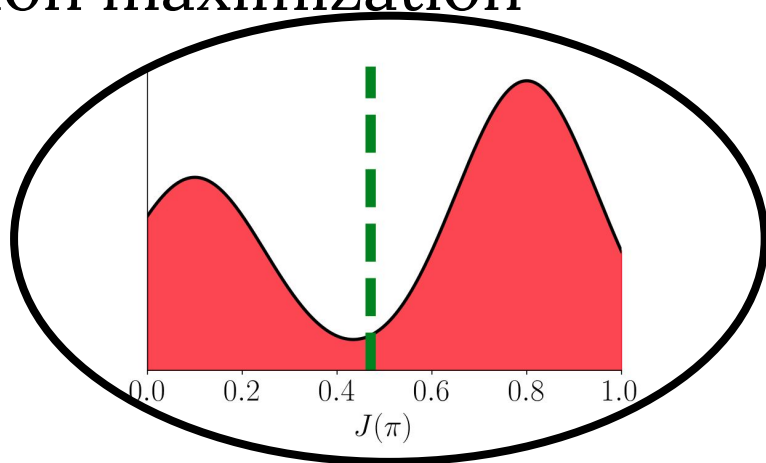
Classical RL objective: expectation maximization

$$\pi^{\star} = \arg \max_{\pi} \mathbb{E}[J(\pi)]$$

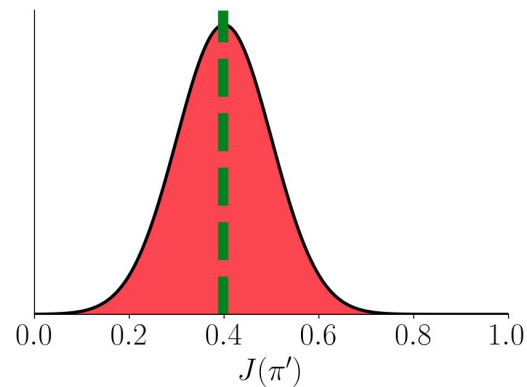
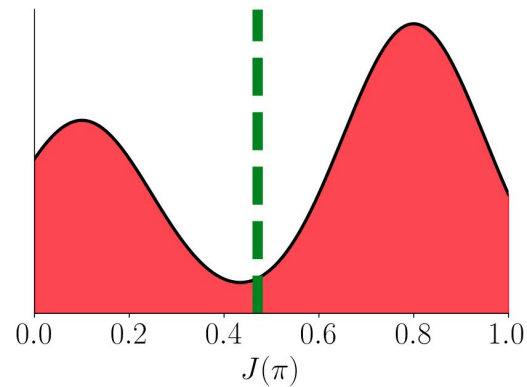
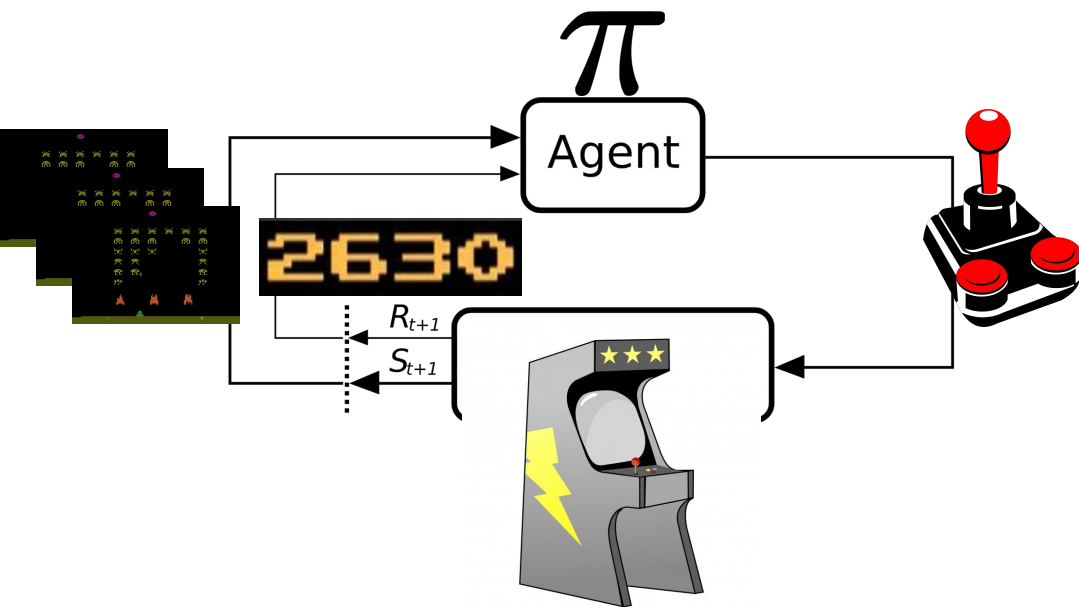


Classical RL objective: expectation maximization

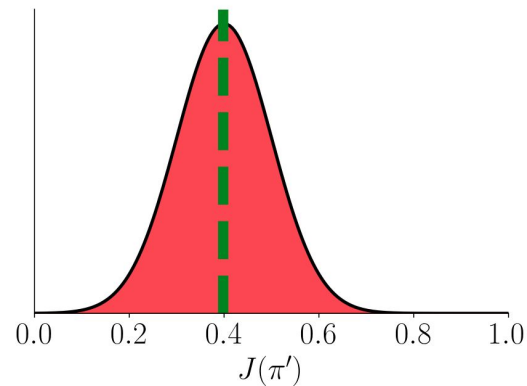
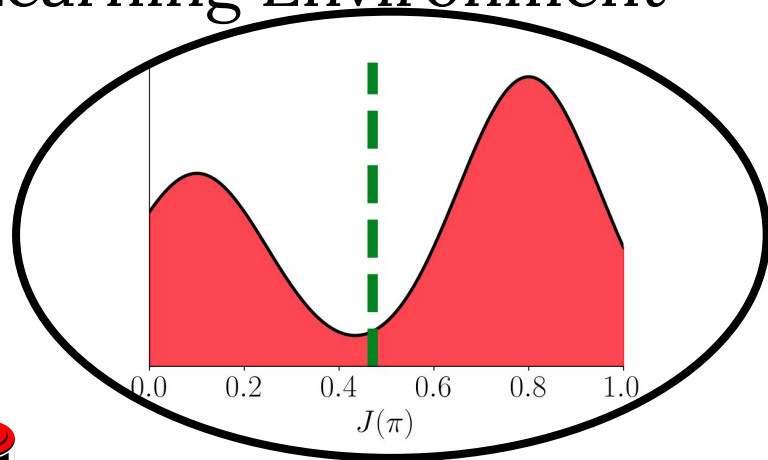
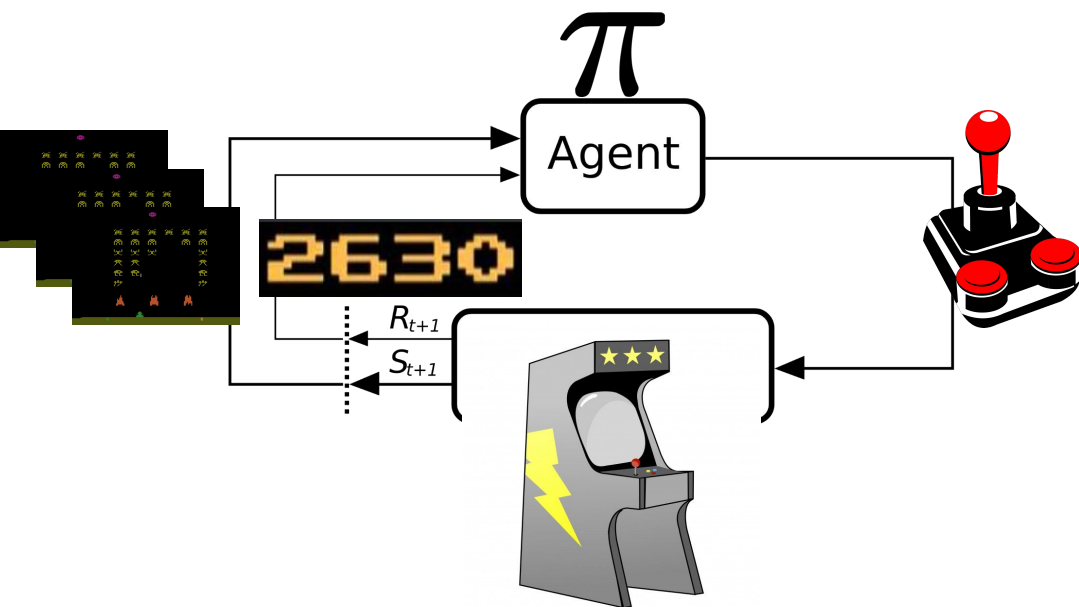
$$\pi^{\star} = \arg \max_{\pi} \mathbb{E}[J(\pi)]$$



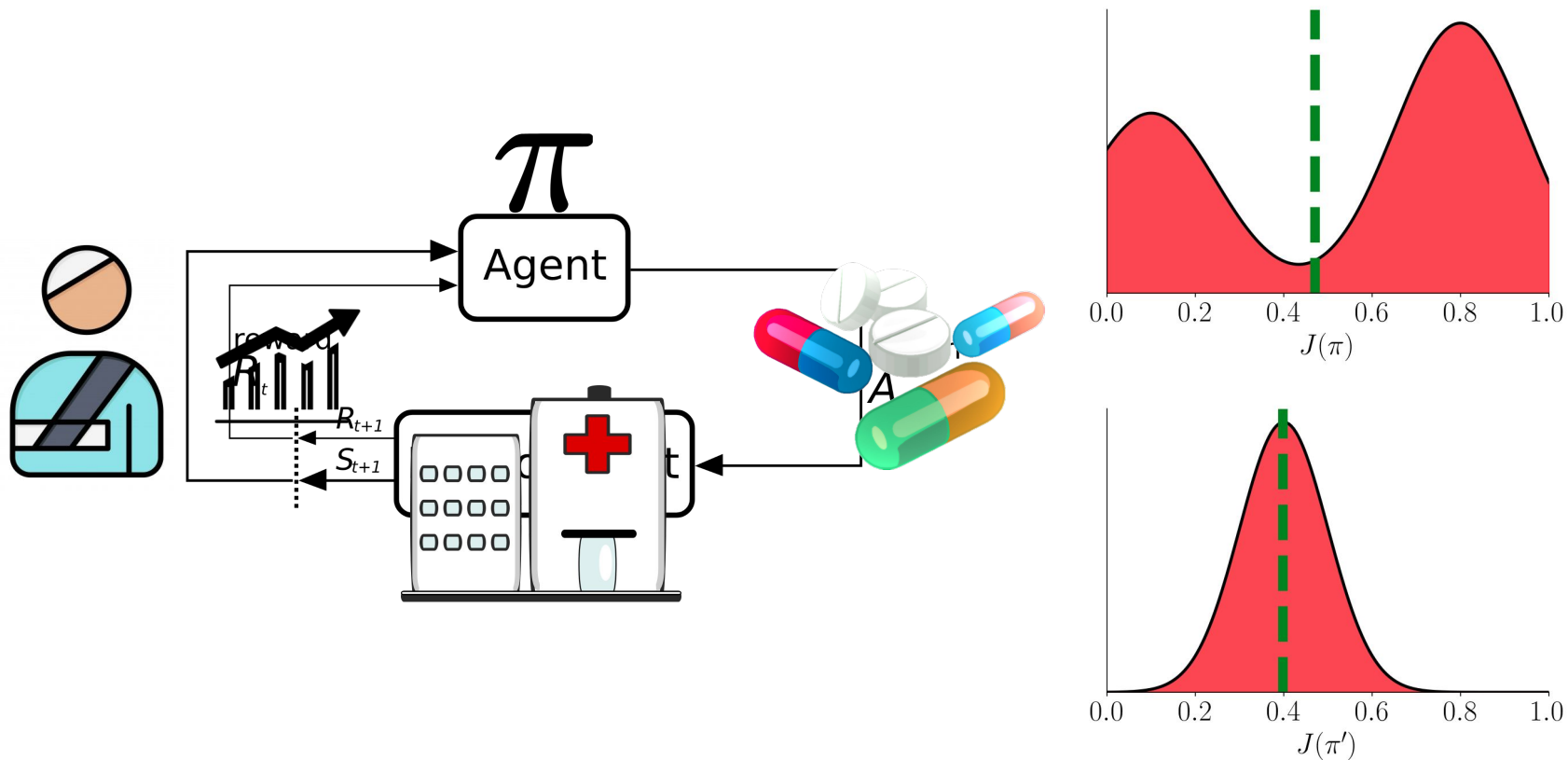
Good objective for the Arcade Learning Environment (Bellemare et al., 2013)



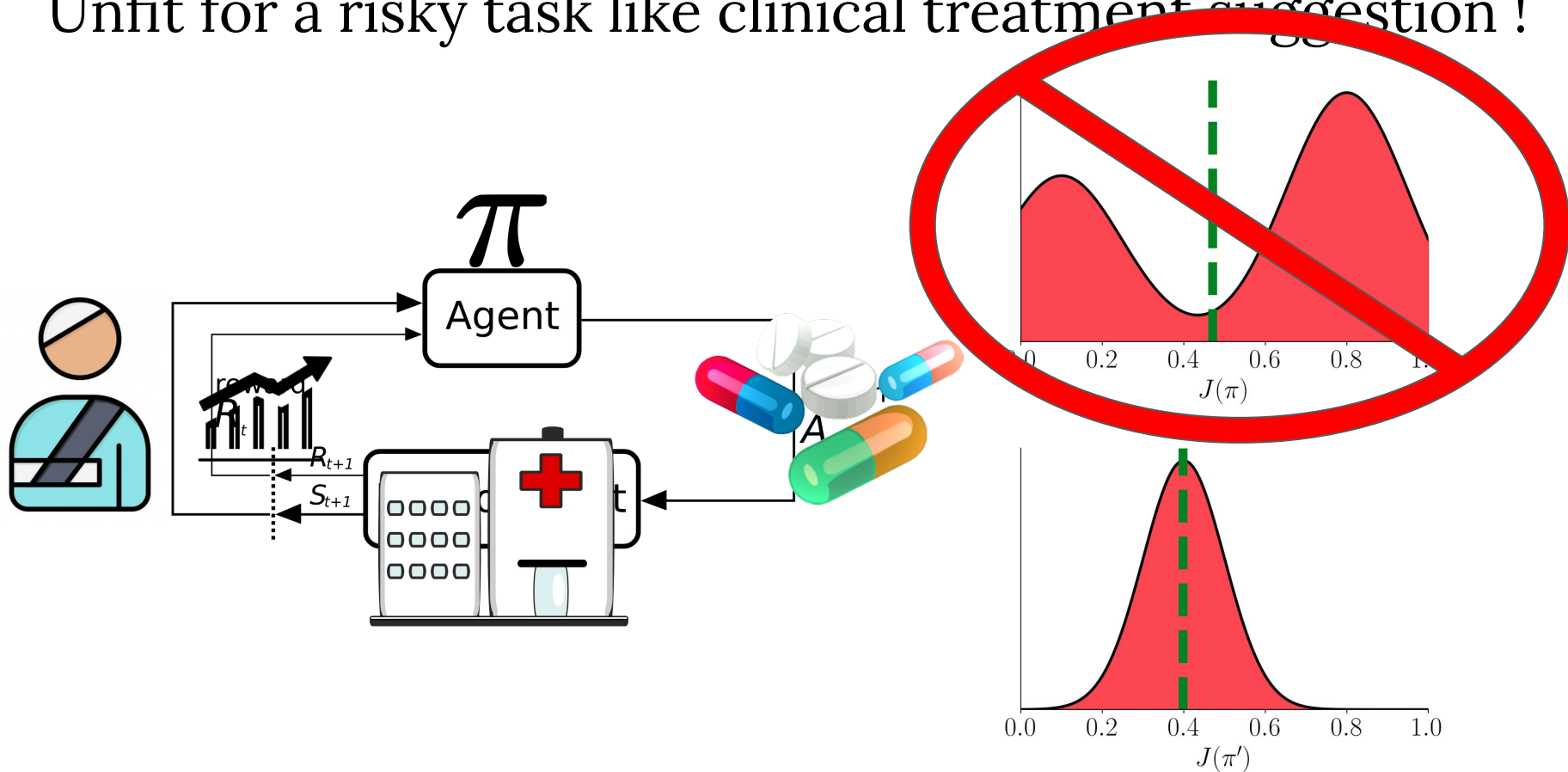
Good objective for the Arcade Learning Environment (Bellemare et al., 2013)



Unfit for a risky task like clinical treatment suggestion !



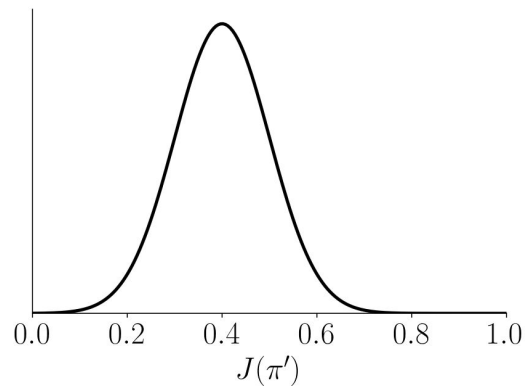
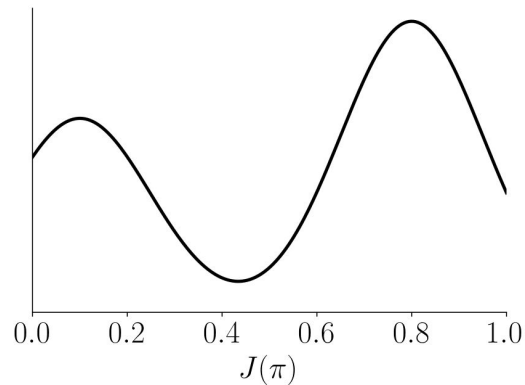
Unfit for a risky task like clinical treatment suggestion !



Conditional-value-at-risk RL: a risk-sensitive objective

$$\text{CVaR}_\alpha(Z) = \mathbb{E}[z \mid z \leq \text{VaR}_\alpha(Z)]$$

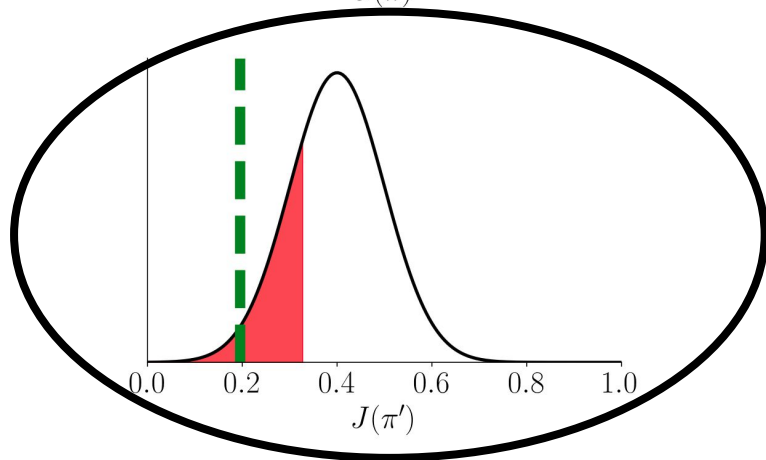
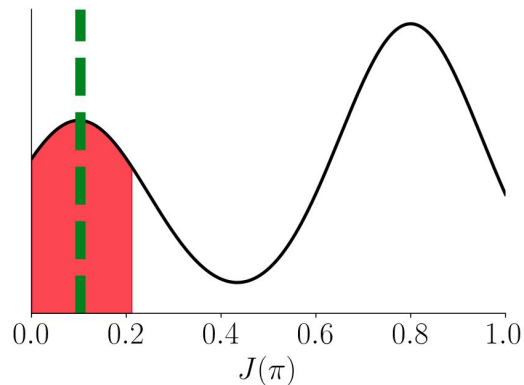
$$\pi^* = \arg \max_{\pi} \text{CVaR}_\alpha[J(\pi)]$$



Conditional-value-at-risk RL: a risk-sensitive objective

$$\text{CVaR}_\alpha(Z) = \mathbb{E}[z \mid z \leq \text{VaR}_\alpha(Z)]$$

$$\pi^* = \arg \max_{\pi} \text{CVaR}_\alpha[J(\pi)]$$

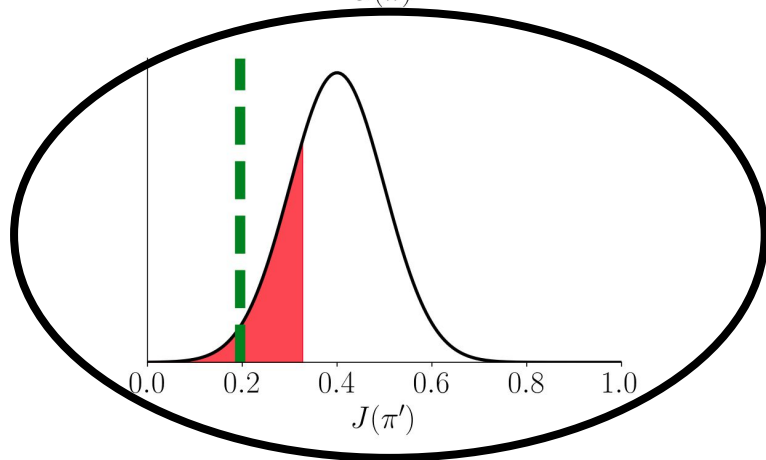
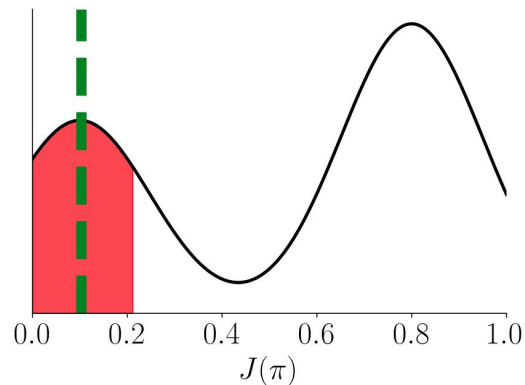


Conditional-value-at-risk RL: a risk-sensitive objective

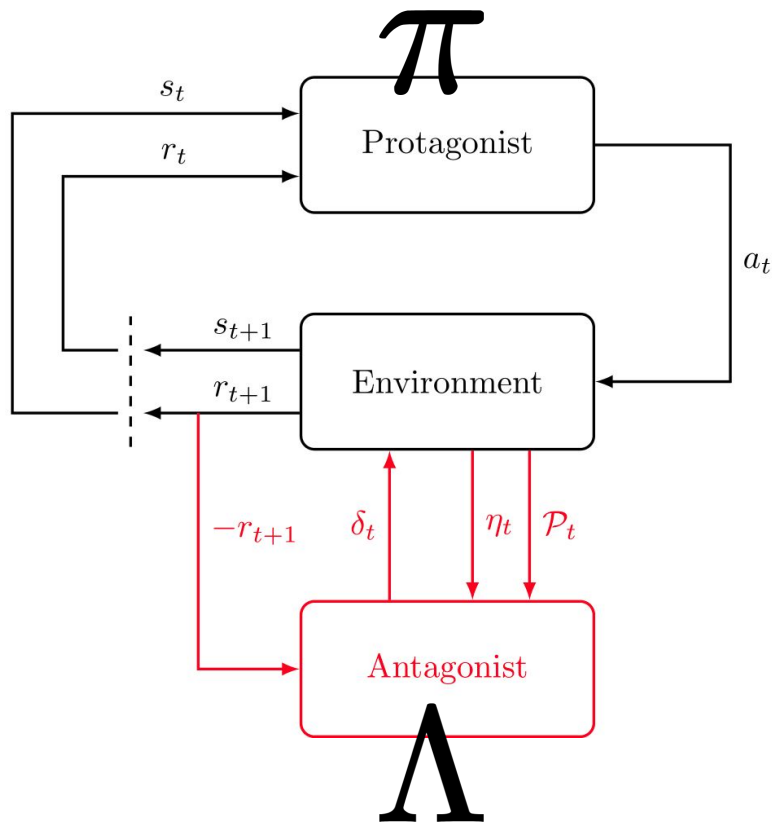
$$\text{CVaR}_\alpha(Z) = \mathbb{E}[z \mid z \leq \text{VaR}_\alpha(Z)]$$

$$\pi^* = \arg \max_{\pi} \text{CVaR}_\alpha[J(\pi)]$$

**All existing approaches require
distributional RL methods.**

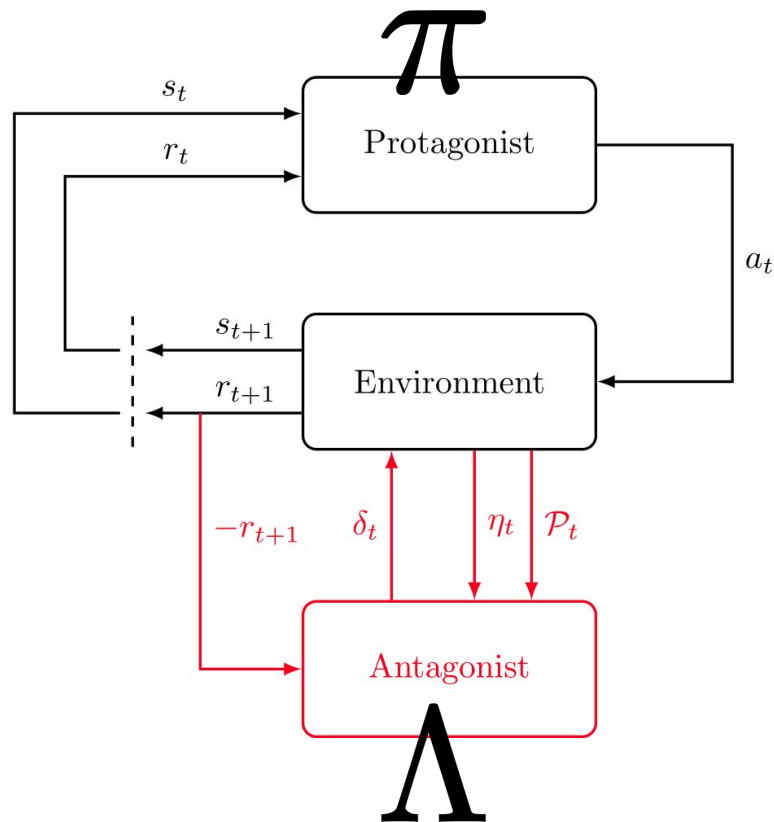


Game Structure



Antagonist produces perturbed next state transitions \hat{P}_t to minimize the protagonist's rewards

Game Structure

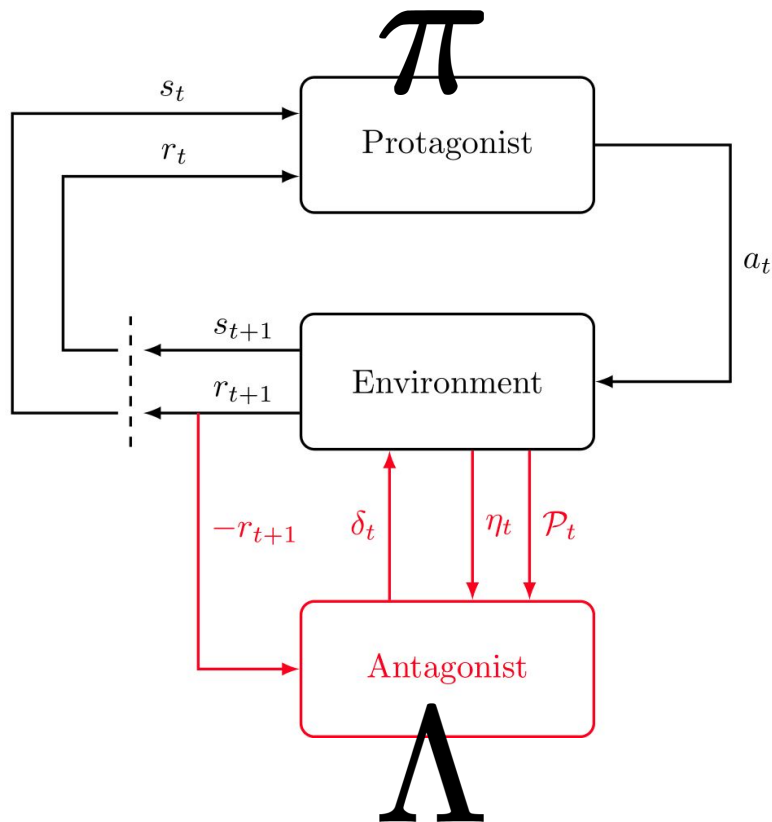


Antagonist produces perturbed next state transitions \hat{P}_t to minimize the protagonist's rewards

Perturbations are multiplicative

$$\hat{P}_t = P_t \circ \delta_t$$

Game Structure

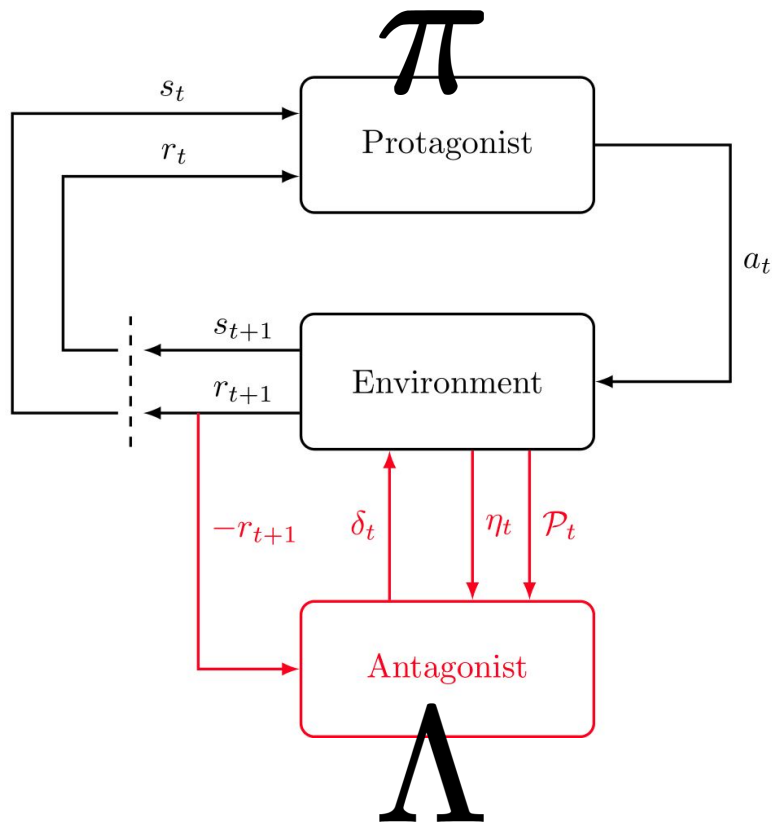


Antagonist produces perturbed next state transitions \hat{P}_t to minimize the protagonist's rewards

Perturbations are multiplicative

$$\hat{P}_t = P_t \circ \delta_t$$

Game Structure



Antagonist produces perturbed next state transitions \hat{P}_t to minimize the protagonist's rewards

Perturbations are multiplicative

$$\hat{P}_t = P_t \circ \delta_t$$

Antagonist has limited budget

$$\delta_0(s_1)\delta_1(s_2)\cdots\delta_{T-1}(s_T) \leq \eta$$

Retrieving CVaR RL optimal policies

Max-min objective:
$$\max_{\pi} \min_{\Lambda} \mathbb{E}[J^{\eta}(\pi, \Lambda)]$$

Retrieving CVaR RL optimal policies

Max-min objective:
$$\max_{\pi} \min_{\Lambda} \mathbb{E}[J^{\eta}(\pi, \Lambda)]$$

The solution is the equilibrium point $(\pi^{\star}, \Lambda^{\star})$, for which we have (Chow et al., 2015):

$$\pi^{\star} = \arg \max_{\pi} \text{CVaR}_{\frac{1}{\eta}} [J(\pi)]$$

Retrieving CVaR RL optimal policies

Max-min objective:
$$\max_{\pi} \min_{\Lambda} \mathbb{E}[J^{\eta}(\pi, \Lambda)]$$

The solution is the equilibrium point (π^*, Λ^*) , for which we have (Chow et al., 2015):

$$\pi^* = \arg \max_{\pi} \text{CVaR}_{\frac{1}{\eta}}[J(\pi)]$$

CARL: Game properties

The objective for both the agent and the adversary is to maximize their **expected** rewards.

CARL: Game properties

The objective for both the agent and the adversary is to maximize their **expected** rewards.

Risk tolerance is based on a single hyperparameter and is easy to interpret.

Stackelberg games for gradient updates

Updating each player naively is unstable due to the non-stationarity of games
(Fiez et al., 2019)

Stackelberg games for gradient updates

Updating each player naively is unstable due to the non-stationarity of games (Fiez et al., 2019)

Stackelberg game: a leader (π) takes for granted that its follower (Λ) is optimal with respect to itself.

$$\pi^* = \arg \max_{\pi} \left\{ \mathbb{E}[J^{\eta}(\pi, \Lambda')] \text{ s.t. } \Lambda' = \arg \max_{\Lambda} \mathbb{E}[J^{\eta}(\pi, \Lambda)] \right\}$$

$$\Lambda^* = \arg \max_{\Lambda} \mathbb{E}[J^{\eta}(\pi, \Lambda)]$$

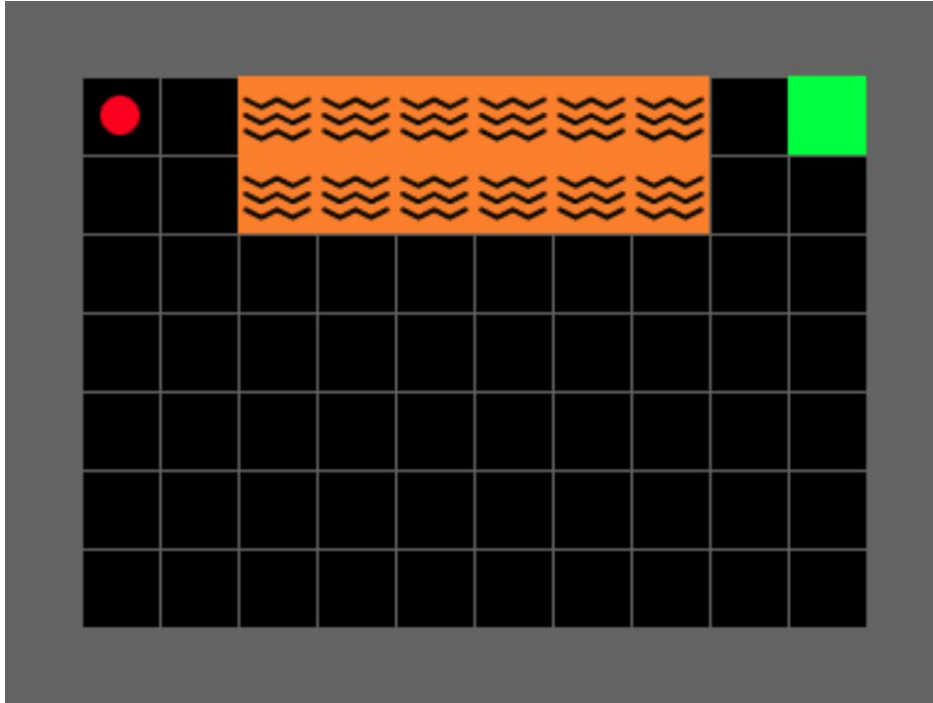
Practical Stackelberg-based algorithm

Algorithm 1: CVaR Adversarial Stackelberg Algorithm

Require: π_θ (protagonist), Λ_ω (antagonist), η (perturbation budget), K_{ant} (number of intermediate antagonist steps)

- 1: $N_{\text{updates}} = 0$
 - 2: **while** training not done **do**
 - 3: Get initial state s_t
 - 4: $\eta_\tau = \eta$ \triangleright Remaining antagonist budget
 - 5: **while** s_t not terminal **do**
 - 6: $a_t \sim \pi_\theta(s_t), \mathcal{P}_t = \mathcal{P}(s_t, a_t)$
 - 7: $\delta_t = \Lambda_\omega(\mathcal{P}_t, \eta_\tau)$
 - 8: $\hat{\mathcal{P}}_t = \mathcal{P}_t \circ \delta$
 - 9: $s_{t+1} \sim \hat{\mathcal{P}}_t, r_{t+1} \sim \mathcal{R}(s_{t+1})$
 - 10: $\eta_\tau = \frac{\eta_\tau}{\delta_t(s_{t+1})}$ \triangleright Update remaining budget
 - 11: **end while**
 - 12: Update θ or ω according to N_{updates} and K_{ant} .
 - 13: $N_{\text{updates}} = N_{\text{updates}} + 1$
 - 14: **end while**
-

Risky Gridworld: experimental setting

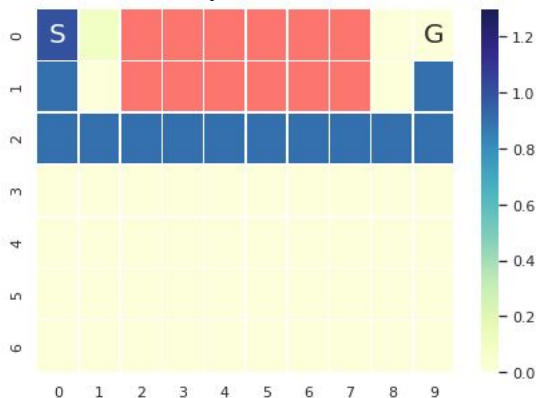


5 % chance that the environment executes a random action.

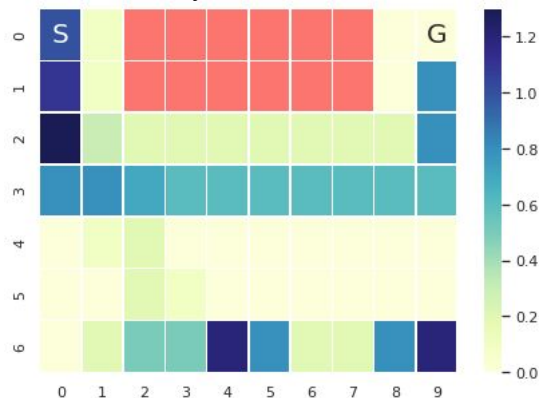
The agent's degree of caution is represented by its willingness to move lower on the grid to distance itself from the lava tiles.

Empirical results

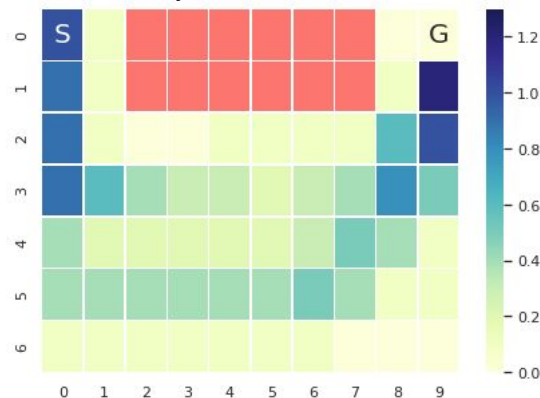
$\eta = 1$



$\eta = 25$



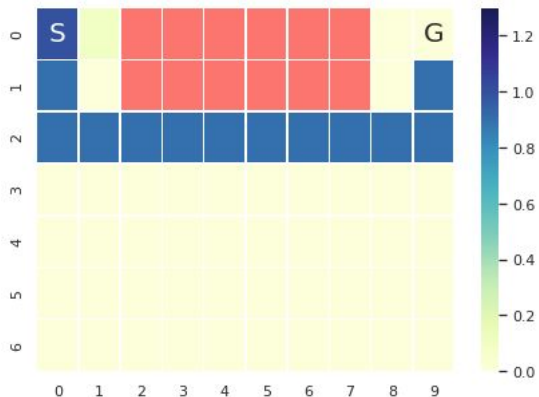
$\eta = 100$



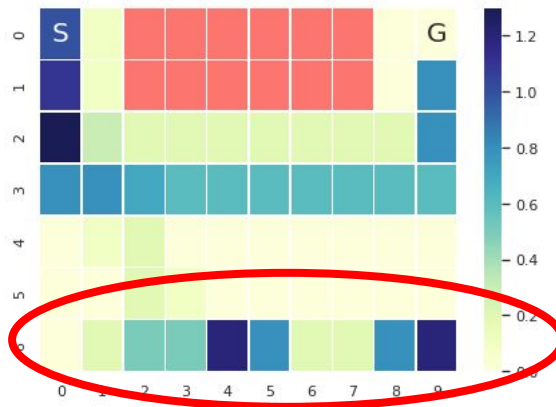
Increasing the adversary's budget leads to an increasingly cautious agent.

Empirical results

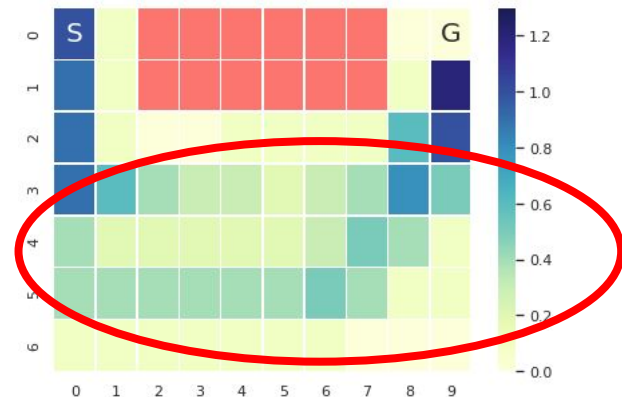
$\eta = 1$



$\eta = 25$



$\eta = 100$



Increasing the adversary's budget leads to an increasingly cautious agent.

There appears to be instability issues in the training procedure.

Conclusion

We proposed a new risk-sensitive RL method for the CVaR risk measure which does not require distributional RL algorithms.

Conclusion

We proposed a new risk-sensitive RL method for the CVaR risk measure which does not require distributional RL algorithms.

We estimate that our proposal can serve as a building block because it paves the way to incorporate results from the Game Theory literature to risk-sensitivity in RL.

References

Bellemare, M. G., Naddaf, Y., Veness, J., & Bowling, M. (2013). The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47, 253-279.

Chow, Y., Tamar, A., Mannor, S., & Pavone, M. (2015). Risk-sensitive and robust decision-making: a cvar optimization approach. arXiv preprint arXiv:1506.02188.

Fiez, T., Chasnov, B., & Ratliff, L. J. (2019). Convergence of learning dynamics in stackelberg games. arXiv preprint arXiv:1906.01217.