

Learning Nash Equilibria in Zero-Sum Stochastic Games via Entropy-Regularized Policy Approximation

Yue Guan Qifan Zhang Panagiotis Tsiotas
Georgia Institute of Technology

Background

Two-agent zero-sum stochastic game:

- A tuple $(\mathcal{S}, \mathcal{A}^{\text{pl}}, \mathcal{A}^{\text{op}}, T, \mathcal{R}, \gamma)$
- The Player maximizes; the Opponent minimizes

Policy π^{pl} (π^{op}): mapping from \mathcal{S} to \mathcal{A}^{pl} (\mathcal{A}^{op})

Q-function and value associated with $\pi = (\pi^{\text{pl}}, \pi^{\text{op}})$:

$$Q^\pi(s, a^{\text{pl}}, a^{\text{op}}) = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t^{\text{pl}}, A_t^{\text{op}}) \mid S_0 = s, A_0^{\text{pl}} = a^{\text{pl}}, A_0^{\text{op}} = a^{\text{op}} \right]$$

$$V^{\pi^{\text{pl}}, \pi^{\text{op}}}(s) = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t^{\text{pl}}, A_t^{\text{op}}) \mid S_0 = s \right]$$

Nash equilibrium:

- A coupled max-min optimization to find $(\pi^{\text{pl}*}, \pi^{\text{op}*})$:

$$V^{\pi^{\text{pl}*}, \pi^{\text{op}*}} = \max_{\pi^{\text{pl}}} \min_{\pi^{\text{op}}} V^{\pi^{\text{pl}}, \pi^{\text{op}}}$$

- Solved as Linear Programs at each state. **Expensive** to solve.

$$\begin{aligned} \max_v \quad & v \\ \text{subject to} \quad & 1^T - \pi^{\text{pl}}(s) \mathcal{Q}^\pi(s) \leq 0 \quad \text{subject to} \quad u1 - \mathcal{Q}^\pi(s) \pi^{\text{op}*T}(s) \leq 0 \\ & 1^T \pi^{\text{pl}}(s) = 1, \pi^{\text{pl}}(s) \geq 0 \quad & 1^T \pi^{\text{op}}(s) = 1, \pi^{\text{op}}(s) \geq 0 \end{aligned}$$

Shapley's method (minimax-Q) iterates between two operators:

$$(\pi_{\text{Nash}}^{\text{pl}}, \pi_{\text{Nash}}^{\text{op}}) = \Gamma_{\text{Nash}}(\mathcal{Q}); \quad \mathcal{Q} = \Gamma_1(\mathcal{Q}, \pi_{\text{Nash}}^{\text{pl}}, \pi_{\text{Nash}}^{\text{op}})$$

computes Nash based on Q-estimate updates Q-estimate based on computed Nash

Entropy-Regularized Policy Approximation

Fixed entropy regularization [1]

$$V^{\pi^{\text{pl}}, \pi^{\text{op}}}(s) = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(S_t, A_t^{\text{pl}}, A_t^{\text{op}}) \right]$$

regulated policy

$$\frac{1}{\beta^{\text{pl}}} \log \frac{\pi^{\text{pl}}(a_t^{\text{pl}} | s_t)}{\rho^{\text{pl}}(a_t^{\text{pl}} | s_t)} \rightarrow \frac{1}{\beta^{\text{op}}} \log \frac{\pi^{\text{op}}(a_t^{\text{op}} | s_t)}{\rho^{\text{op}}(a_t^{\text{op}} | s_t)}$$

fixed inverse temperature fixed reference policy information cost

Closed-form max-min soft solution under regularization

Marginalization:

$$Q_{\text{KL}}^{\text{pl}}(s, a^{\text{pl}}) = \frac{1}{\beta^{\text{op}}} \log \sum_{a^{\text{op}}} \rho^{\text{op}}(a^{\text{op}} | s) \exp(\beta^{\text{op}} Q_{\text{KL}}(s, a^{\text{pl}}, a^{\text{op}}))$$

Soft-Nash Policies:

$$\pi_{\text{KL}}^{\text{pl}}(a^{\text{pl}} | s) = \frac{1}{Z^{\text{pl}}(s)} \rho^{\text{op}}(a^{\text{op}} | s) \exp(\beta^{\text{op}} Q_{\text{KL}}^{\text{pl}}(s, a^{\text{pl}}, a^{\text{op}}))$$

Two soft operators: inverse temperature dependent

NOT original NE

$$(\pi_{\text{KL}}^{\text{pl}}, \pi_{\text{KL}}^{\text{op}}) = \Gamma_{\text{KL}}^{\beta}(\mathcal{Q}_{\text{KL}}; \rho); \quad \mathcal{Q}_{\text{KL}} = \mathbb{I}_2^{\beta}(\mathcal{Q}_{\text{KL}}; \rho)$$

computes soft Nash Updates Q_{KL} with reference ρ

Soft Nash Q2 Algorithm

SNQ2 learns two Q-values simultaneously:

- (1) Original Q-value Q
- (2) Entropy-regularized Q-value Q_{KL}

Slow Module:

- Learns standard Q-value and Nash policies
- Slow but produces Nash

Fast Module:

- Learns entropy-regularized Q-value and soft-optimal policies
- Fast but only an approximation of the Nash policies

- Coupling of the two Modules:

Use Nash policies from the **slow module** to update the priors used in the **fast module**

Use soft-policies from the **fast module** to update Q-values in the **slow module**

- Actively adapts entropy regularization

- Reduce inverse temperature β over time
- Update reference policies using Nash policies from original Q-estimate

- A dynamic schedule scheme is introduced to balance the two modules

- Observes the Q-difference between two updates
- Decides when to perform Nash prior updates and reduce inverse temperature

Challenges in convergence analysis:

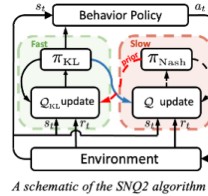
- With decreasing β , the operators used to update Q-value changes
- Standard fixed-point argument cannot be directly applied

Convergence Analysis

Theorem 1 Let (\mathcal{X}, ρ) be a complete metric space. Let $f^n: \mathcal{X} \rightarrow \mathcal{X}$ be a family of contraction operators such that for all $n = 1, 2, \dots$ there exists $d^n \in (0, 1)$, such that $\rho(f^n x, f^n y) \leq d^n \rho(x, y)$ for all $x, y \in \mathcal{X}$. Assume that $\lim_{n \rightarrow \infty} d^n = d \in (0, 1)$. Let $x \in \mathcal{X}$ be a starting point and let $x^n = f^n \dots f^1 x$ be the result of sequentially applying the operators f^1, \dots, f^n to x . If the sequence of operators $\{f^n\}_{n=1}^{\infty}$ convergence pointwise to f , then f is also a contraction mapping with contraction factor d . Furthermore, if x^* is the fixed point of f , then for every $x \in \mathcal{X}$, $\lim_{n \rightarrow \infty} x^n = x^*$.

The convergence of SNQ2 can be shown through the following argument:

- As β approaches zero, the update rule of SNQ2 converges to Shapley's method
- Per Theorem 1, SNQ2 converges to the fixed point of Shapley's method, which is the Nash Q-value of the game.



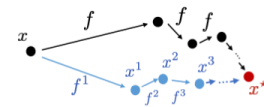
A schematic of the SNQ2 algorithm

Algorithm 1: SNQ2-Learning Algorithm

```

1 Inputs: Priors  $\rho$ , Learning rates  $\alpha$  and  $\eta$ ; initial prior update
   episode  $M = \Delta M_0$ ; Nash update frequency  $T$ ;
2 Set  $Q(s, a^{\text{pl}}, a^{\text{op}}) = Q_{\text{KL}}(s, a^{\text{pl}}, a^{\text{op}}) = 0$ ;
3 Set  $\beta^{\text{pl}}$  and  $\beta^{\text{op}}$  to some large values;
4 while  $Q$  not converged do
5   while episode  $i$  not end do
6     Compute  $\pi_{\text{KL}}(s_t) \leftarrow \Gamma_{\text{KL}}^{\beta}(\mathcal{Q}_{\text{KL}}, \rho)(s_t)$ ;
7     Collect transition  $(s_t, a_t^{\text{pl}}, a_t^{\text{op}}, r_t, s_{t+1})$  where
8        $a_t^{\text{pl}} \sim \pi_{\text{KL}}^{\text{pl}}(s_t)$ ,  $a_t^{\text{op}} \sim \pi_{\text{KL}}^{\text{op}}(s_t)$ ;
9     if  $t \bmod T = 0$  then
10      Compute  $V(s_{t+1}) =$ 
11         $\max_{\pi^{\text{pl}}} \min_{\pi^{\text{op}}} \sum_{a^{\text{pl}}, a^{\text{op}}} Q(s_{t+1}, a^{\text{pl}}, a^{\text{op}}) \pi^{\text{pl}}(a^{\text{pl}} | s_{t+1})$ ;
12      else
13        Compute
14         $V(s_{t+1}) = \pi_{\text{KL}}^{\text{pl}}(s_{t+1})^T \mathcal{Q}(s_{t+1}) \pi_{\text{KL}}^{\text{op}}(s_{t+1})$ ;
15      end
16      Update  $Q(s_t, a_t^{\text{pl}}, a_t^{\text{op}})$  with  $V(s_{t+1})$  via (13);
17      Update  $Q_{\text{KL}}(s_t, a_t^{\text{pl}}, a_t^{\text{op}})$  via (12);
18    end
19    if  $i = M$  then
20      Compute  $\pi_{\text{Nash}} \leftarrow \Gamma_{\text{Nash}} \mathcal{Q}$ ;
21      Update priors  $\rho_{\text{new}} \leftarrow \pi_{\text{Nash}}$ ;
22      Update schedule as in Algorithm 3:
23       $\Delta M, \beta_{\text{new}} = DS(\rho_{\text{new}}, \rho, \beta, \Delta M, \mathcal{Q})$ ;
24      Update next prior update schedule  $M \leftarrow \Delta M$ ;
25      Update priors  $\rho \leftarrow \rho_{\text{new}}$ ,  $\beta \leftarrow \beta_{\text{new}}$ ;
26      Decrease learning rates  $\alpha$  and  $\eta$ ;
27    end
28  return  $Q(s, a^{\text{pl}}, a^{\text{op}})$ .

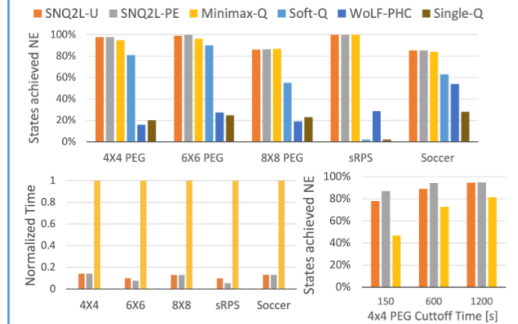
```



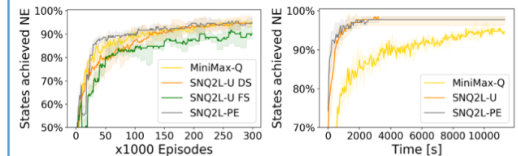
FULL PAPER

Numerical Experiments

Experiments are conducted in Pursuit-Evasion games (PEG), Sequential Rock-Paper-Scissor (sRPS) and Soccer games



- Without updating regularization, two-agent soft-Q [1] failed to converge to a Nash in sRPS.
- With updating regularization, SNQ2 achieves same level of convergence as Minimax-Q
- Significant reduction in learning time
- Warm starting (-PE) gives better convergence give the same cutoff time



- Similar episode-wise convergence trend as Minimax-Q
- Time-wise trend shows a significant speed up
- Warm starting (-PE) gives a better convergence trend comparing to uniform prior (-U)
- Dynamic scheduling (DS) improves episode-wise convergence speed

References

- [1] Grau-Moya, J., Leibfried, F., & Bou-Ammar, H. (2018). Balancing two-player stochastic games with soft q-learning. *arXiv preprint arXiv:1802.03216*.
- [2] Zhang, Q., Guan, Y., & Tsiotas, P. (2020). Learning Nash Equilibria in Zero-Sum Stochastic Games via Entropy-Regularized Policy Approximation. *arXiv preprint arXiv:2009.00162*.