

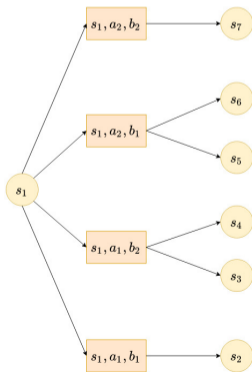
Adapting to game trees in zero-sum imperfect information games

Côme Fiegel¹, Pierre Ménard², Tadashi Kozuno³,
Rémi Munos⁴, Vianney Perchet^{1,5}, Michal Valko⁴

¹ENSAE Paris ²ENS Lyon ³Omron Sinic X ⁴DeepMind ⁵Criteo AI Lab

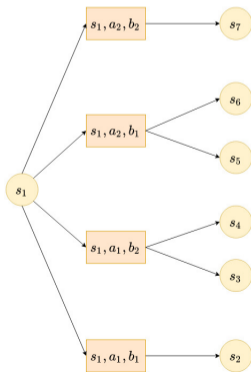
Zero-sum (two players) imperfect information games

- State space \mathcal{S} , initial state $s_1 \in \mathcal{S}$ and horizon $H > 0$
- At timestep $h \in [1..H]$, the two players take actions $\mathbf{a} \in \mathcal{A}$ and $\mathbf{b} \in \mathcal{B}$
- Reward $r_h(s, \mathbf{a}, \mathbf{b}) \in [0, 1]$ and transition to the next state $p_h(\cdot | s, \mathbf{a}, \mathbf{b})$



Zero-sum (two players) imperfect information games

- State space \mathcal{S} , initial state $s_1 \in \mathcal{S}$ and horizon $H > 0$
- At timestep $h \in [1..H]$, the two players take actions $\mathbf{a} \in \mathcal{A}$ and $\mathbf{b} \in \mathcal{B}$
- Reward $r_h(s, \mathbf{a}, \mathbf{b}) \in [0, 1]$ and transition to the next state $p_h(\cdot | s, \mathbf{a}, \mathbf{b})$

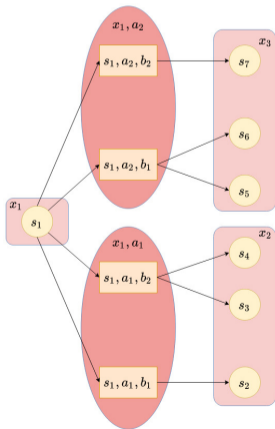


Assumptions:

Zero-sum: max-player receives r_h ,
min-player receives $-r_h$

Zero-sum (two players) imperfect information games

- State space \mathcal{S} , initial state $s_1 \in \mathcal{S}$ and horizon $H > 0$
- At timestep $h \in [1..H]$, the two players take actions $\mathbf{a} \in \mathcal{A}$ and $\mathbf{b} \in \mathcal{B}$
- Reward $r_h(s, \mathbf{a}, \mathbf{b}) \in [0, 1]$ and transition to the next state $p_h(\cdot | s, \mathbf{a}, \mathbf{b})$



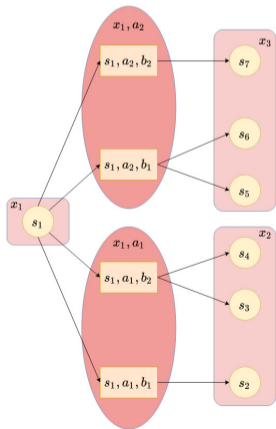
Assumptions:

Zero-sum: max-player receives r_h ,
min-player receives $-r_h$

Imperfect information: Players only observe information sets $\mathbf{x}(s) \in \mathcal{X}$ and $\mathbf{y}(s) \in \mathcal{Y}$

Zero-sum (two players) imperfect information games

- State space \mathcal{S} , initial state $s_1 \in \mathcal{S}$ and horizon $H > 0$
- At timestep $h \in [1..H]$, the two players take actions $\mathbf{a} \in \mathcal{A}$ and $\mathbf{b} \in \mathcal{B}$
- Reward $r_h(s, \mathbf{a}, \mathbf{b}) \in [0, 1]$ and transition to the next state $p_h(\cdot | s, \mathbf{a}, \mathbf{b})$



Assumptions:

Zero-sum: **max-player** receives r_h ,
min-player receives $-r_h$

Imperfect information: Players only observe information sets $\mathbf{x}(s) \in \mathcal{X}$ and $\mathbf{y}(s) \in \mathcal{Y}$

Perfect recall: Players do not forget past observations and actions

Approximate Nash equilibrium

Policies: Non-deterministic $\mu = (\mu(\cdot|x))_{x \in \mathcal{X}}$ and $\nu = (\nu(\cdot|y))_{y \in \mathcal{Y}}$

Approximate Nash equilibrium

Policies: Non-deterministic $\mu = (\mu(\cdot|x))_{x \in \mathcal{X}}$ and $\nu = (\nu(\cdot|y))_{y \in \mathcal{Y}}$

Value: $V^{\mu, \nu} = \mathbb{E}^{\mu, \nu} \left[\sum_{h=1}^H r_h \right]$

Approximate Nash equilibrium

Policies: Non-deterministic $\mu = (\mu(\cdot|x))_{x \in \mathcal{X}}$ and $\nu = (\nu(\cdot|y))_{y \in \mathcal{Y}}$

Value: $V^{\mu, \nu} = \mathbb{E}^{\mu, \nu} \left[\sum_{h=1}^H r_h \right]$

(μ, ν) is a **Nash equilibrium** if $\mu \in \operatorname{argmax} V^{\cdot, \nu}$ and $\nu \in \operatorname{argmin} V^{\mu, \cdot}$

Approximate Nash equilibrium

Policies: Non-deterministic $\mu = (\mu(\cdot|x))_{x \in \mathcal{X}}$ and $\nu = (\nu(\cdot|y))_{y \in \mathcal{Y}}$

Value: $V^{\mu, \nu} = \mathbb{E}^{\mu, \nu} \left[\sum_{h=1}^H r_h \right]$

(μ, ν) is a **Nash equilibrium** if $\mu \in \operatorname{argmax} V^{\cdot, \nu}$ and $\nu \in \operatorname{argmin} V^{\mu, \cdot}$

Objective \rightarrow Approximate a **Nash equilibrium**

Sequential learning

Interaction with the game: T episodes played using freely chosen profiles
 $(\mu^t, \nu^t)_{t \in [1..T]}$

Sequential learning

Interaction with the game: T episodes played using freely chosen profiles
 $(\mu^t, \nu^t)_{t \in [1..T]}$

Regrets: $\mathfrak{R}_{\max}^T = \max_{\mu} \sum_t [V^{\mu, \nu^t} - V^{\mu^t, \nu^t}]$, $\mathfrak{R}_{\min}^T = \max_{\nu} \sum_t [V^{\mu^t, \nu^t} - V^{\mu^t, \nu}]$

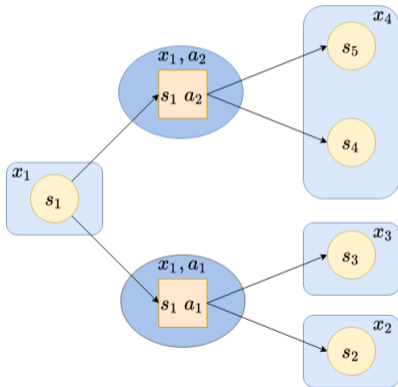
Sequential learning

Interaction with the game: T episodes played using freely chosen profiles $(\mu^t, \nu^t)_{t \in [1..T]}$

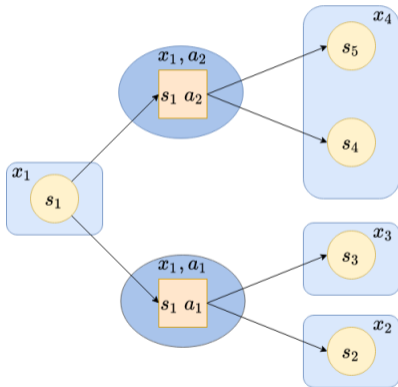
Regrets: $\mathfrak{R}_{\max}^T = \max_{\mu} \sum_t [V^{\mu, \nu^t} - V^{\mu^t, \nu^t}]$, $\mathfrak{R}_{\min}^T = \max_{\nu} \sum_t [V^{\mu^t, \nu^t} - V^{\mu^t, \nu}]$

Small **regrets** \iff average profile $(\bar{\mu}, \bar{\nu})$ approximates a **Nash equilibrium**

Max player's point of view



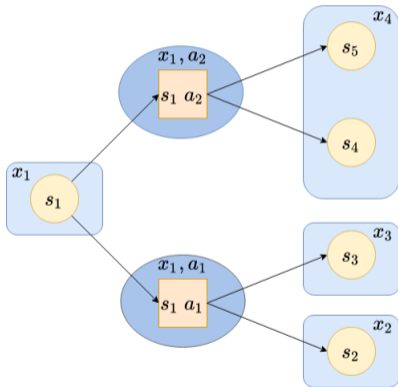
Max player's point of view



Advantage

Episodic MDP with a tree-structure

Max player's point of view



Advantage

Episodic MDP with a tree-structure

Difficulty

Adversarial transitions p^t that change between episodes

Back to regret minimization

Objective: minimize $\mathfrak{R}_{\max}^T = \max_{\mu} \sum_{t=1}^T [V^{\mu, \nu^t} - V^{\mu^t, \nu^t}]$

FTRL approach : $\mu^{t+1} = \operatorname{argmax}_{\mu} \sum_{k=1}^t \tilde{V}^{\mu, \nu^k} - \Psi(\mu)$ with

- \tilde{V}^{\cdot, ν^t} estimated value at episode t as a function of μ
- Ψ the regularizer

Back to regret minimization

Objective: minimize $\mathfrak{R}_{\max}^T = \max_{\mu} \sum_{t=1}^T [V^{\mu, \nu^t} - V^{\mu^t, \nu^t}]$

FTRL approach : $\mu^{t+1} = \operatorname{argmax}_{\mu} \sum_{k=1}^t \tilde{V}^{\mu, \nu^k} - \Psi(\mu)$ with

- \tilde{V}^{\cdot, ν^t} estimated value at episode t as a function of μ
- Ψ the regularizer

How to choose Ψ ?

Regularizer choice

First choice: **BalancedFTRL**

- $\Psi_p(\mu) = \text{negentropy}(\text{information set} | \mu, p) / \eta$
- Compute balanced transitions p^* and use Ψ_{p^*}
- $\mathfrak{R}_{\max}^T = \tilde{\mathcal{O}}\left(\sqrt{H|\mathcal{X}||\mathcal{A}|T}\right)$

Regularizer choice

First choice: **BalancedFTRL**

- $\Psi_p(\mu) = \text{negentropy}(\text{information set} | \mu, p) / \eta$
- Compute balanced transitions p^* and use Ψ_{p^*}
- $\mathfrak{R}_{\max}^T = \tilde{O}\left(\sqrt{H|\mathcal{X}||\mathcal{A}|T}\right)$

Impossible with initially unknown tree structure

Regularizer choice

First choice: **BalancedFTRL**

- $\Psi_p(\mu) = \text{negentropy}(\text{information set} | \mu, p) / \eta$
- Compute balanced transitions p^* and use Ψ_{p^*}
- $\mathfrak{R}_{\max}^T = \tilde{O}\left(\sqrt{H|\mathcal{X}||\mathcal{A}|T}\right)$

Impossible with initially unknown tree structure

Second choice: **AdaptiveFTRL**

- Estimate cumulative transitions \tilde{P}^t
- "Replace" p^* with \tilde{P}^t
- $\mathfrak{R}_{\max}^T = \tilde{O}\left(H\sqrt{|\mathcal{X}||\mathcal{A}|T}\right)$

Conclusion

Algorithm	Sample complexity	Structure-free
IXOMD	$\tilde{O}(H^2(\mathcal{X} ^2 \mathcal{A} + \mathcal{Y} ^2 \mathcal{B})/\epsilon^2)$	✓
BalancedOMD	$\tilde{O}(H^3(\mathcal{X} \mathcal{A} + \mathcal{Y} \mathcal{B})/\epsilon^2)$	✗
BalancedFTRL	$\tilde{O}(H(\mathcal{X} \mathcal{A} + \mathcal{Y} \mathcal{B})/\epsilon^2)$	✗
AdaptiveFTRL	$\tilde{O}(H^2(\mathcal{X} \mathcal{A} + \mathcal{Y} \mathcal{B})/\epsilon^2)$	✓
Lower bound	$\tilde{O}(H(\mathcal{X} \mathcal{A} + \mathcal{Y} \mathcal{B})/\epsilon^2)$	-

Contact: come.fiegel@gmail.com