

Adapting to game trees in zero-sum imperfect information games



Côme Fiegel, Pierre Ménard, Tadashi Kozuno, Rémi Munos, Vianney Perchet and Michal Valko



Two-player Zero-sum IIG with Perfect Recall

\mathcal{S} : State space of size S , Horizon H

\mathcal{X} : Max-player's information set space of size X

\mathcal{A} : Max-player's action space of size A

\mathcal{Y} : Min-player's information set space of size Y

\mathcal{B} : Min-player's action space of size B

r_h, p_h : Reward/loss function and state-transition dynamics

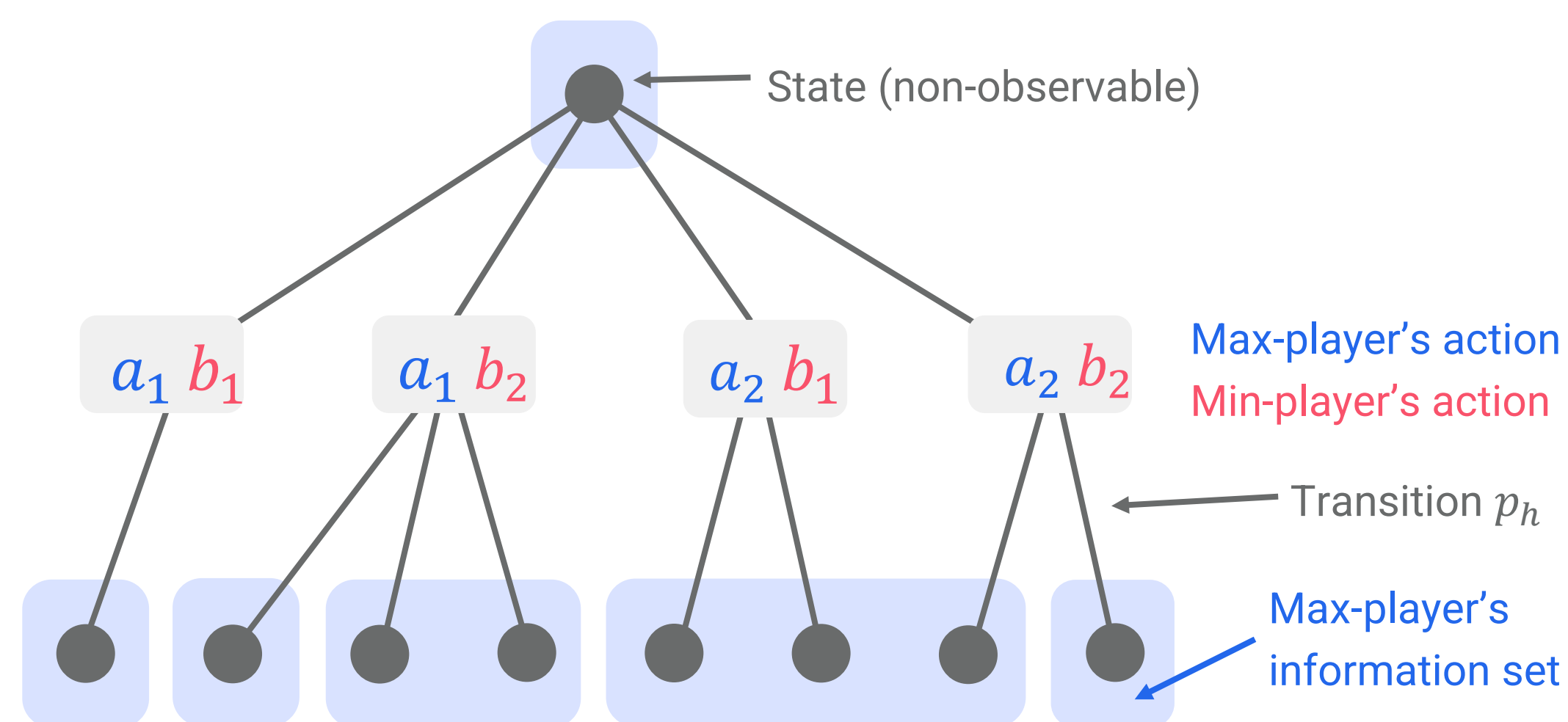


Fig 1. An IIG with $H = 2$, $\mathcal{A} = \{a_1, a_2\}$, and $\mathcal{B} = \{b_1, b_2\}$. Only max-player's information sets are shown.

Regret, Average Profile, and Nash Equilibrium

For a profile (μ, ν) , the expected return (of the max-player) is defined by

$$V^{\mu, \nu} = \mathbb{E}^{\mu, \nu} \left[\sum_{h=1}^H r_h(s_h, a_h, b_h) \right].$$

When a profile (μ, ν) satisfies the following, it is said to be an ε -NE:

$$\max_{\mu'} V^{\mu', \nu} - \min_{\nu} V^{\mu, \nu'} \leq \varepsilon.$$

For a sequence of profiles (μ^t, ν^t) , the regret of the max-player is

$$\mathcal{R}_{max}^T = \max_{\mu} \sum_{t=1}^T (V^{\mu, \nu^t} - V^{\mu^t, \nu^t}).$$

The time-averaged profile $(\bar{\mu}, \bar{\nu})$ is an ε -NE with:

$$\varepsilon = \frac{\mathcal{R}_{max}^T + \mathcal{R}_{min}^T}{T}.$$

The Problem

Find an approximation through self-play of an optimal strategy for a zero-sum imperfect information game only using trajectory feedback.

Our Contributions

- Propose two computationally efficient algorithms, by combining implicit exploration and follow the regularized leader.
- If applied by both players, the first has an optimal high-probability sample complexity of order $H(A\mathcal{X} + B\mathcal{Y})/\varepsilon^2$ requiring the knowledge of the structure.
- The second has a sample complexity of order $H^2(A\mathcal{X} + B\mathcal{Y})/\varepsilon^2$ without this knowledge, using an adaptive regularization.

Algorithm	Sample complexity	Structure-free
MCCFR (Farina et al., 2020; Bai et al., 2022)	$\tilde{O}(H^4(A\mathcal{X} + B\mathcal{Y})/\varepsilon^2)$	✗
IXOMD (Kozuno et al., 2021)	$\tilde{O}(H^2(A\mathcal{X}^2 + B\mathcal{Y}^2)/\varepsilon^2)$	✓
Balanced OMD (Bai et al., 2022)	$\tilde{O}(H^3(A\mathcal{X} + B\mathcal{Y})/\varepsilon^2)$	✗
Balanced FTRL (this paper)	$\tilde{O}(H(A\mathcal{X} + B\mathcal{Y})/\varepsilon^2)$	✗
Adaptive FTRL (this paper)	$\tilde{O}(H^2(A\mathcal{X} + B\mathcal{Y})/\varepsilon^2)$	✓
Lower bound (this paper)	$\tilde{O}(H(A\mathcal{X} + B\mathcal{Y})/\varepsilon^2)$	

Sample complexity for various algorithms. Structure-free means that the algorithm does not need to know the structure of the information set spaces in advance.

Algorithm 1 Adaptive FTRL for the max-player

1: **Input:**

Base learning rate η and IX bias γ

Uniform policy μ^0

$\mu_{1:h}(x, a)$ denotes the combined probability for the max-player of choosing actions that lead to (x, a)

2: **For** $t = 1$ to T :

For all h and $x_h \in \mathcal{X}_h$, compute learning rate:

$$\eta_h^t(x_h) \leftarrow \min_{x'_h, \geq x_h} \eta / (1 + \tilde{P}_h^{t-1}(x'_h))$$

For all $a_h \in \mathcal{A}(x_h)$, compute bias rate:

$$\gamma_h^t(x_h, a_h) \leftarrow \gamma / (1 + \tilde{P}_h^{t-1}(x_h, a_h))$$

Compute update:

$$\mu^t \leftarrow \operatorname{argmin}_{\mu} \sum_h \langle \mu_{1:h}, \tilde{L}_h^{t-1} \rangle + \mathcal{D}_{\eta^t}(\mu, \mu^0)$$

with $\mathcal{D}_{\eta}(\mu, \mu^0) = \sum_x \mu_{1:h}(x) \operatorname{KL}(\mu(\cdot|x), \mu^0(\cdot|x)) / \eta(x)$

For $h = 1$ to H :

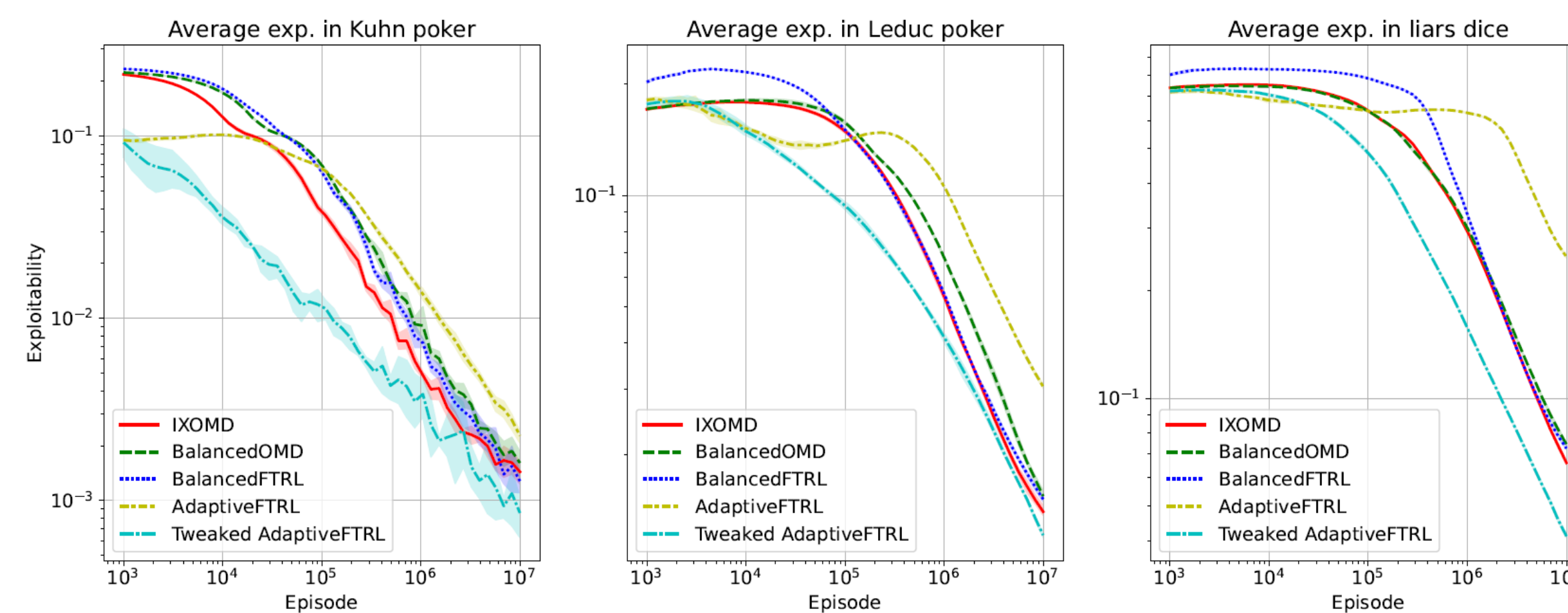
Observe information set x_h^t

Execute $a_h^t \sim \mu^t(\cdot|x_h^t)$ and **receive** reward r_h^t

$$\tilde{L}_h^t \leftarrow \tilde{L}_h^{t-1} + \mathbb{I}_{\{x_h^t, a_h^t\}} (1 - r_h^t) / (\mu_{1:h}^t + \gamma_h^t)$$

$$\tilde{P}_h^t \leftarrow \tilde{P}_h^{t-1} + \mathbb{I}_{\{x_h^t, a_h^t\}} / (\mu_{1:h}^t + \gamma_h^t)$$

Return average policy $\bar{\mu}$



Performances of various algorithms with respect to the number of episodes. The vertical axis corresponds to the smallest ε such that the output is an ε -NE.