

Closing the Computational-Statistical Gap in BAI for Combinatorial Semi-bandits



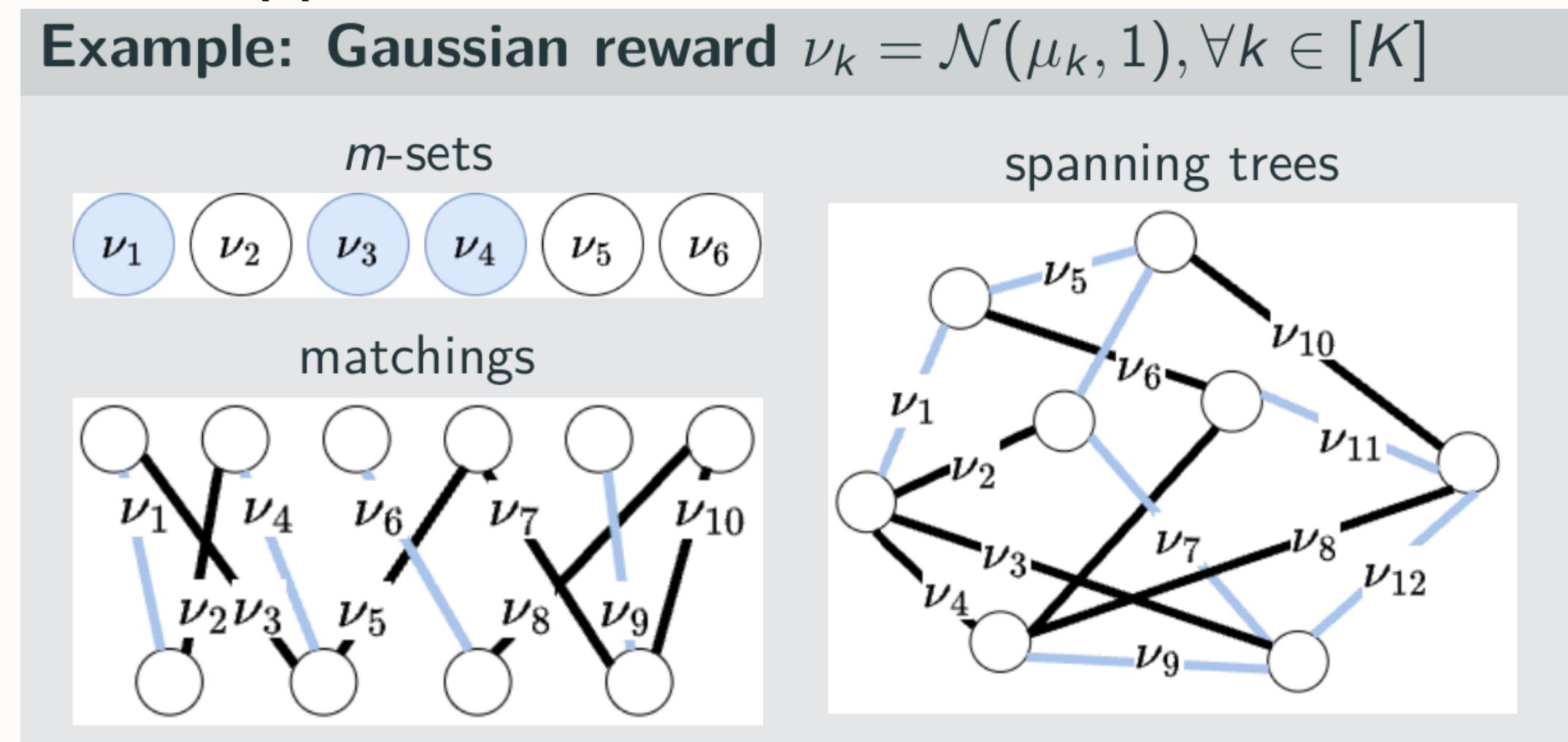
Ruo-Chun Tzeng¹, Po-An Wang¹, Alexandre Proutiere¹ Chi-Jen Lu²

¹EECS, KTH Royal Institute of Technology, Sweden
²Institute of Information Science, Academia Sinica, Taiwan



Combinatorial BAI with semi-bandit feedback

Input: K arms $(\nu_k)_{k \in [K]}$ with mean $\boldsymbol{\mu} \in \mathbb{R}^K$ and $\mathcal{X} \subseteq \{0, 1\}^K$.



Rule: At each round $t \in \mathbb{N}$, the learner pulls an action $\mathbf{x}(t) \in \mathcal{X}$ and observe $y_k(t) \sim \nu_k$ iff $x_k(t) = 1$, and returns her estimated best action $\hat{i} \in \mathcal{X}$ when she decides to terminate at round τ .

Goal: Design a δ -PAC learning strategy s.t. the best action $\mathbf{i}^*(\boldsymbol{\mu}) \in \operatorname{argmax}_{\mathbf{x}} \langle \boldsymbol{\mu}, \mathbf{x} \rangle$ is identified w.p. $\geq 1 - \delta$ and $\mathbb{P}_{\boldsymbol{\mu}}[\tau < \infty] = 1$ while minimizing $\mathbb{E}_{\boldsymbol{\mu}}[\tau]$.

Prior works: a computational-statistical gap

Any δ -PAC algorithm satisfies $\mathbb{E}_{\boldsymbol{\mu}}[\tau] \geq T^*(\boldsymbol{\mu}) \operatorname{kl}(\delta, 1 - \delta)$, where

$$T^*(\boldsymbol{\mu})^{-1} = \sup_{\boldsymbol{\omega} \in \Sigma} F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) \text{ with } F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \inf_{\lambda \in \operatorname{Alt}(\boldsymbol{\mu})} \sum_{k=1}^K \frac{\omega_k (\mu_k - \lambda_k)^2}{2}.$$

Track-and-Stop [6] is statistically optimal but requires to repeatedly solve $T^*(\hat{\boldsymbol{\mu}}(t-1))^{-1} \Rightarrow$ computationally inefficient.

FWS [8] at the FW-update round has to solve a potentially $\mathcal{O}(2^K)$ many convex programs \Rightarrow computationally inefficient

CombGame [7] is MCP-oracle efficient and statistically optimal \Rightarrow left open the design of an efficient MCP-oracle

Designing efficient MCP based on a structural observation

Let $f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu}) = \inf_{\lambda \in \mathbb{R}: \langle \mathbf{i}(\boldsymbol{\mu}) - \mathbf{x}, \boldsymbol{\lambda} \rangle < 0} \sum_{k=1}^K \frac{\omega_k (\mu_k - \lambda_k)^2}{2}$ s.t. $F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \min_{\mathbf{x} \neq \mathbf{i}^*(\boldsymbol{\mu})} f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu})$.

Property of $f_{\mathbf{x}}$ and its Lagrangian dual $g_{\boldsymbol{\omega}, \boldsymbol{\mu}}$:

$$f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu}) = \max_{\alpha \geq 0} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha) \quad (\text{known by [2]})$$

$$g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha) \text{ is linear in } \mathbf{x} \text{ and concave in } \alpha \quad (\text{our observation})$$

These properties $\Rightarrow F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \min_{\mathbf{x} \neq \mathbf{i}^*(\boldsymbol{\mu})} \max_{\alpha \geq 0} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha)$ as a two-player zero-sum game.

We not only want to estimate $F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$ but also the *equilibrium action* \mathbf{x}_e s.t.

$F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \max_{\alpha \geq 0} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}_e, \alpha)$. This rules out many existing results from applying.

- \mathbf{x}_e is required to solve $\max_{\boldsymbol{\omega} \in \Sigma} F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$ by the first-order methods
- Last-iterate convergence [1, 3] are mostly for saddle-point problems

Algorithm 1: (ϵ, θ) -MCP($\boldsymbol{\omega}, \boldsymbol{\mu}$)

for $n = 1, 2, \dots$ do

(Follow-the-Perturbed-Leader) $\mathcal{Z}_n \sim \exp(1)^K$ and $\eta_n = \frac{c_0}{\sqrt{n}}$

$$\mathbf{x}^{(n)} \in \operatorname{argmin}_{\mathbf{x} \neq \mathbf{i}^*(\boldsymbol{\mu})} \left(\sum_{m=1}^{n-1} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha^{(m)}) + \frac{\langle \mathcal{Z}_n, \mathbf{x} \rangle}{\eta_n} \right)$$

(Best-Response) $\alpha^{(n)} \in \operatorname{argmax}_{\alpha \geq 0} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}^{(n)}, \alpha)$

if $\sqrt{n} > \frac{c_{\theta}(1 + \epsilon)}{\epsilon \hat{F}}$, where $\begin{cases} \hat{F} = g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}^{(n^*)}, \alpha^{(n^*)}) \\ n^* \in \operatorname{argmin}_{m \leq n} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}^{(m)}, \alpha^{(m)}) \end{cases}$

then return $(\hat{F}, \mathbf{x}^{(n^*)})$;

end

(Theorem 1) Let $\epsilon, \theta \in (0, 1)$ and $(\boldsymbol{\omega}, \boldsymbol{\mu}) \in \Sigma_+ \times \Lambda$.

- $\mathbb{P}_{\boldsymbol{\mu}}[F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) \leq \hat{F} \leq (1 + \epsilon)F_{\boldsymbol{\mu}}(\boldsymbol{\omega})] \geq 1 - \theta$
- the number of \mathbf{i}^* -oracle calls: $\mathcal{O}\left(\frac{K^3 D^5 \ln K \ln \theta^{-1} \|\boldsymbol{\mu}\|_{\infty}^4 \|\boldsymbol{\omega}^{-1}\|_{\infty}^2}{\epsilon^2 F_{\boldsymbol{\mu}}(\boldsymbol{\omega})^2}\right)$

Our Perturbed Frank-Wolfe Sampling (P-FWS)

We use stochastic smoothing [5, 4] to overcome the nonsmoothness of $F_{\boldsymbol{\mu}}$ as: all we need is \mathbf{i}^* -oracle and its required gradient can be evaluated by envelope theorem [8].

The smoothed objective $\bar{F}_{\boldsymbol{\mu}, \eta}(\boldsymbol{\omega}) = \mathbb{E}_{\mathcal{Z} \sim \operatorname{Uniform}(\mathcal{B}_2)}[F_{\boldsymbol{\mu}}(\boldsymbol{\omega} + \eta \mathcal{Z})]$ satisfies:

- $\nabla \bar{F}_{\boldsymbol{\mu}, \eta}(\boldsymbol{\omega}) = \mathbb{E}_{\mathcal{Z} \sim \operatorname{Uniform}(\mathcal{B}_2)}[\nabla F_{\boldsymbol{\mu}}(\boldsymbol{\omega} + \eta \mathcal{Z})]$
- $\bar{F}_{\boldsymbol{\mu}, \eta}$ is $\frac{\ell K}{\eta}$ -smooth and $\bar{F}_{\boldsymbol{\mu}, \eta}(\boldsymbol{\omega}) \xrightarrow{\eta \downarrow 0} F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$

High-level design of P-FWS

Let \mathcal{X}_0 be a set s.t. $\forall k \in [K]$, there exists $\mathbf{x} \in \mathcal{X}_0$ s.t. $x_k = 1$.

P-FWS alternate between two phases:

- pull each $\mathbf{x} \in \mathcal{X}_0$ once (to avoid high cost and boundary cases)
- pull $\mathbf{x}(t) \in \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \nabla \bar{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t}(\hat{\boldsymbol{\omega}}(t-1)), \mathbf{x} \rangle$ (ideal FW update)

(Theorem 2) Let $\boldsymbol{\mu} \in \Lambda$ and $\delta \in (0, 1)$. P-FWS is δ -PAC, finishes in finite time, $\mathbb{P}_{\boldsymbol{\mu}}[\limsup_{\delta \rightarrow 0} \frac{\tau}{\ln \delta^{-1}} \leq T^*(\boldsymbol{\mu})] = 1$, $\mathbb{E}_{\boldsymbol{\mu}}[\tau]$ is bounded by Poly(K) in moderate-confidence regime and achieves the minimal in high-confidence regime, and the total number of \mathbf{i}^* -oracle calls is bounded by Poly(K).

References

- [1] K. Abe, K. Ariu, M. Sakamoto, K. Toyoshima, and A. Iwasaki. Last-iterate convergence with full-and noisy-information feedback in two-player zero-sum games. In *Proc. of AISTATS*, 2023.
- [2] L. Chen, A. Gupta, and J. Li. Pure exploration of multi-armed bandit under matroid constraints. In *Proc. of COLT*, 2016.
- [3] C. Daskalakis and I. Panageas. Last-iterate convergence: Zero-sum games and constrained min-max optimization. *Proc. of ITCS*, 2019.
- [4] J. C. Duchi, P. L. Bartlett, and M. J. Wainwright. Randomized smoothing for stochastic optimization. *SIAM Journal on Optimization*, 2012.
- [5] A. D. Flaxman, A. T. Kalai, and H. B. McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proc. of SODA*, 2005.
- [6] A. Garivier and E. Kaufmann. Optimal best arm identification with fixed confidence. In *Proc. of COLT*, 2016.
- [7] M. Jourdan, M. Mutnà, J. Kirschner, and A. Krause. Efficient pure exploration for combinatorial bandits with semi-bandit feedback. In *Proc. of ALT*, 2021.
- [8] P.-A. Wang, R.-C. Tzeng, and A. Proutiere. Fast pure exploration via frank-wolfe. In *Proc. of NeurIPS*, 2021.



digital futures

