

Closing the Computational-Statistical Gap in Best Arm Identification for Combinatorial Semi-bandits

Ruo-Chun Tzeng¹, Po-An Wang¹, Alexandre Proutiere¹, and Chi-Jen Lu²
Conference on Neural Information Processing Systems, 2023

¹EECS, KTH Royal Institute of Technology, Sweden

²Institute of Information Science, Academia Sinica, Taiwan



Computational-Statistical Gap in BAI for Combinatorial Semi-bandits

Motivation: combinatorial multi-armed bandits

- **(Network Routing)** Minimize the travel time from A to B when the latency of each individual link is unknown
 - Each **arm** is an edge in a graph
 - Each **action** is a path from A to B
 - Reward function is **linear**
 - Efficient algorithm exists for the offline problem
- **(Crowdsourcing)** Maximize the total number of correctly solved tasks by assigning tasks to workers
 - Each **arm** is an edge of a graph
 - Each **action** is a bipartite matching
 - Reward function is **linear**
 - Efficient algorithm exists for the offline problem



Problem: combinatorial BAI with fixed confidence

- **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 , the learner
 - pulls $\mathbf{x}(t) \in \mathcal{X}$ and observes $y_k(t) \sim \nu_k$ iff $x_k(t) = 1$
 - decides whether to stop and outputs $\hat{\mathbf{i}} \in \mathcal{X}$
 - let τ be the round it stops
- **Goal:** Identify $\mathbf{i}^*(\boldsymbol{\mu}) \in \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \boldsymbol{\mu} \rangle$ w.p. at least $1 - \delta$ and $\mathbb{P}_\mu[\tau < \infty] = 1$
- **Assumptions:** (i) $\mathbf{i}^*(\boldsymbol{\mu})$ is unique; (ii) for any $\mathbf{v} \in \mathbb{R}^K$, a best action $\mathbf{i}^*(\mathbf{v})$ can be found in polynomial time.



Problem: combinatorial BAI with fixed confidence

- **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 , the learner
 - pulls $\mathbf{x}(t) \in \mathcal{X}$ and observes $y_k(t) \sim \nu_k$ iff $x_k(t) = 1$
 - decides whether to stop and outputs $\hat{\mathbf{i}} \in \mathcal{X}$
 - let τ be the round it stops
- **Goal:** Identify $\mathbf{i}^*(\boldsymbol{\mu}) \in \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \boldsymbol{\mu} \rangle$ w.p. at least $1 - \delta$ and $\mathbb{P}_{\boldsymbol{\mu}}[\tau < \infty] = 1$
- **Assumptions:** (i) $\mathbf{i}^*(\boldsymbol{\mu})$ is unique; (ii) for any $\mathbf{v} \in \mathbb{R}^K$, a best action $\mathbf{i}^*(\mathbf{v})$ can be found in polynomial time.

Our contribution

We propose the first **computational efficient** and **statistical optimal** algorithm for this problem with **Gaussian** rewards.



Challenge in solving the lowerbound problem

Instance-specific sample complexity lower bound [GK16]

For any δ -PAC algorithm¹, $\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(\delta, 1 - \delta)$, where

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

- $\Sigma = \{\sum_{x \in \mathcal{X}} w_x \mathbf{x} : \mathbf{w} \in \Sigma_{|\mathcal{X}|}\}$: all possible arm allocations
- $\Lambda = \{\boldsymbol{\lambda} \in \mathbb{R}^K : |\mathbf{i}^*(\boldsymbol{\lambda})| = 1\}$: all possible parameters
- $\text{Alt}(\mu) = \{\boldsymbol{\lambda} \in \Lambda : \mathbf{i}^*(\boldsymbol{\lambda}) \neq \mathbf{i}^*(\mu)\}$: confusing parameters

¹Here we assume the arm- k reward distribution is $\nu_k = \mathcal{N}(\mu_k, 1)$



Challenge in solving the lowerbound problem

Instance-specific sample complexity lower bound [GK16]

For any δ -PAC algorithm¹, $\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(\delta, 1 - \delta)$, where

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

- $\Sigma = \{\sum_{x \in \mathcal{X}} w_x \mathbf{x} : \mathbf{w} \in \Sigma_{|\mathcal{X}|}\}$: all possible arm allocations
- $\Lambda = \{\boldsymbol{\lambda} \in \mathbb{R}^K : |\mathbf{i}^*(\boldsymbol{\lambda})| = 1\}$: all possible parameters
- $\text{Alt}(\mu) = \{\boldsymbol{\lambda} \in \Lambda : \mathbf{i}^*(\boldsymbol{\lambda}) \neq \mathbf{i}^*(\mu)\}$: confusing parameters

Each sampling strategy is represented by its arm allocation $\omega \in \Sigma$.

¹Here we assume the arm- k reward distribution is $\nu_k = \mathcal{N}(\mu_k, 1)$



Challenge in solving the lowerbound problem

Instance-specific sample complexity lower bound [GK16]

For any δ -PAC algorithm¹, $\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(\delta, 1 - \delta)$, where

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

- $\Sigma = \{\sum_{x \in \mathcal{X}} w_x \mathbf{x} : \mathbf{w} \in \Sigma_{|\mathcal{X}|}\}$: all possible arm allocations
- $\Lambda = \{\boldsymbol{\lambda} \in \mathbb{R}^K : |\mathbf{i}^*(\boldsymbol{\lambda})| = 1\}$: all possible parameters
- $\text{Alt}(\mu) = \{\boldsymbol{\lambda} \in \Lambda : \mathbf{i}^*(\boldsymbol{\lambda}) \neq \mathbf{i}^*(\mu)\}$: confusing parameters

The inner optimization measures the distance from μ to the *most confusing parameter* (MCP) with the best action different from $\mathbf{i}^*(\mu)$.
 \Rightarrow The **best** sampling strategy has the **largest** distance to the MCP.

¹Here we assume the arm- k reward distribution is $\nu_k = \mathcal{N}(\mu_k, 1)$



Challenge in solving the lowerbound problem

Instance-specific sample complexity lower bound [GK16]

For any δ -PAC algorithm¹, $\mathbb{E}_\mu[\tau] \geq T^*(\mu)k(\delta, 1 - \delta)$, where

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

A standard approach [GK16] achieving asymptotic optimality consists of:

- Chernoff stopping rule: $\tau = \inf\{t : tF_{\hat{\mu}(t)}(\hat{\omega}(t)) > \ln(\frac{t}{\delta}) + o(1)\}$
- Pull arms according to $\omega^*(\hat{\mu}(t)) = \operatorname{argmax}_{\omega \in \Sigma} F_{\hat{\mu}(t)}(\omega)$

Difficulty in determining the most confusing parameter (MCP)

The domain $\text{Alt}(\mu) = \{\lambda \in \Lambda : i^*(\lambda) \neq i^*(\mu)\}$ of $F_\mu(\omega)$

\Rightarrow The naive approach as to solve $|\mathcal{X}| - 1$ many convex programs by partitioning $\text{Alt}(\mu) = \cup_{x \neq i^*(\mu)} \{\lambda \in \Lambda : \langle i^*(\mu) - x, \lambda \rangle < 0\}$.

¹Here we assume the arm- k reward distribution is $\nu_k = \mathcal{N}(\mu_k, 1)$



Computational inefficiency in prior optimal algorithms

- Track-and-Stop [GK16] at each round t has to solve

$$\omega^*(\hat{\mu}(t)) \in \operatorname{argmax}_{\omega \in \Sigma} F_{\hat{\mu}(t)}(\omega), \text{ (computationally expensive)}$$

- FWS [WTP21] has to compute $f_x(\mathbf{w}(t), \hat{\mu}(t))$ of each $\mathbf{x} \neq \mathbf{i}^*(\hat{\mu}(t))$ to deal with the *nonsmoothness* of $F_{\hat{\mu}(t)}$
- CombGame [JMKK21] proposed a MCP-oracle efficient algorithm, but **no efficient MCP oracle exists** prior to our work

Our Perturbed Frank-Wolfe Sampling (P-FWS)

- P-FWS deals with $|\mathcal{X}| \leq 2^K$ actions by *stochastic smoothing*
- All P-FWS needs is the linear maximization \mathbf{i}^* oracle



**Our MCP Algorithm: a no-regret
algorithm for solving $F_\mu(\omega)$**

A crucial structural observation about $F_\mu(\omega)$

$$\text{Define } f_x(\omega, \mu) = \inf_{\lambda \in \mathbb{R}^K: \langle i^*(\mu) - x, \lambda \rangle < 0} \sum_{k=1}^K \frac{\omega_k (\mu_k - \lambda_k)^2}{2}.$$

Property of f_x and its Lagrangian dual $g_{\omega, \mu}$

$$f_x(\omega, \mu) = \max_{\alpha \geq 0} g_{\omega, \mu}(x, \alpha) \quad (\text{known by [CGL16]})$$

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

$$F_\mu(\omega) = \min_{x \neq i^*(\mu)} f_x(\omega, \mu) = \min_{x \neq i^*(\mu)} \max_{\alpha \geq 0} g_{\omega, \mu}(x, \alpha) \quad (1)$$



A crucial structural observation about $F_\mu(\omega)$

Define $f_x(\omega, \mu) = \inf_{\lambda \in \mathbb{R}^K: \langle i^*(\mu) - x, \lambda \rangle < 0} \sum_{k=1}^K \frac{\omega_k (\mu_k - \lambda_k)^2}{2}$.

Property of f_x and its Lagrangian dual $g_{\omega, \mu}$

$$f_x(\omega, \mu) = \max_{\alpha \geq 0} g_{\omega, \mu}(x, \alpha) \quad (\text{known by [CGL16]})$$

$g_{\omega, \mu}(x, \alpha)$ is **linear** in x and **concave** in α (our observation)

$$F_\mu(\omega) = \min_{x \neq i^*(\mu)} f_x(\omega, \mu) = \min_{x \neq i^*(\mu)} \max_{\alpha \geq 0} g_{\omega, \mu}(x, \alpha) \quad (1)$$

Requirement: Not only to estimate $F_\mu(\omega)$ but also the *equilibrium action* x_e s.t. $F_\mu(\omega) = \max_{\alpha \geq 0} g_{\omega, \mu}(x_e, \alpha)$.

- x_e is needed to solve $\max_{\omega \in \Sigma} F_\mu(\omega)$ by the first-order methods
- Existing results [DP19, LNP⁺21, APFS22, AAS⁺23] on last-iterate convergence are not applicable as they all consider convex domains



Our proposed MCP algorithm

Algorithm 1: (ϵ, θ) -MCP(ω, μ)

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

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

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

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

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

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

end



Our proposed MCP algorithm

Algorithm 1: (ϵ, θ) -MCP(ω, μ)

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

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

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

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

(Computational Cost Per Iteration)

- $\mathbf{x}^{(n)}$ can be computed by at most $D = \max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_0$ calls to $\mathbf{i}^*(\cdot)$
- $\alpha^{(n)}$ is evaluated in $\mathcal{O}(1)$



Our proposed MCP algorithm

The termination condition is designed based on Lemma 1

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

such that $\mathbb{P}\left[F_\mu(\omega) \leq \hat{F} \leq (1+\epsilon)F_\mu(\omega)\right] \geq 1 - \theta$ holds.

(Lemma 1) If Algorithm 1 runs for N iterations, then

$$\mathbb{P}\left[\underbrace{\frac{1}{N} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)})}_{\geq \min_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) = \hat{F}} - \underbrace{\frac{1}{N} \min_{\mathbf{x} \neq \mathbf{j}^*(\mu)} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)})}_{\leq \frac{1}{N} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}_e, \alpha^{(n)}) \leq F_\mu(\omega)} \leq \frac{c_\theta}{\sqrt{N}} \right] \geq 1 - \theta.$$



Our proposed MCP algorithm

Theorem 1 (MCP)

Let $(\omega, \mu) \in \Sigma_+ \times \Lambda$. The (ϵ, θ) -MCP(ω, μ) algorithm outputs $(\hat{F}, \hat{\mathbf{x}})$:

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

By envelop theorem [WTP21], we estimation (sub)gradient of $F_\mu(\omega)$ by

$$\nabla_\omega f_{\hat{\mathbf{x}}}(\omega, \mu) = \left(\frac{(\mu_k - \lambda_k^*)^2}{2} \right)_{k \in [K]},$$

where λ^* is the minimizer to the optimization problem of $f_{\hat{\mathbf{x}}}(\omega, \mu)$.



**Our P-FWS: the first poly-time
statistically optimal algorithm**

Solving $T^*(\mu)$ with stochastic smoothed objective

- The well-studied stochastic smoothing [FKM05, DBW12] takes the average value in a neighborhood of points:

$$\bar{F}_{\mu,\eta}(\omega) = \mathbb{E}_{\mathcal{Z} \sim \text{Uniform}(B_2)} [F_{\mu}(\omega + \eta\mathcal{Z})]$$

- F_{μ} is ℓ -Lipschitz and its smoothed objective satisfies:
 - $\bar{F}_{\mu,\eta}$ is $\frac{\ell K}{\eta}$ -smooth and $\bar{F}_{\mu,\eta}(\omega) \xrightarrow{\eta \downarrow 0} F_{\mu}(\omega)$
 - $\nabla \bar{F}_{\mu,\eta}(\omega) = \mathbb{E}_{\mathcal{Z} \sim \text{Uniform}(B_2)} [\nabla F_{\mu}(\omega + \eta\mathcal{Z})]$

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:

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



Theorem 2 (P-FWS)

Let $\mu \in \Lambda$ and $\delta \in (0, 1)$. P-FWS with proper parameters is δ -PAC and finishes in finite time;

- (i) $\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau]}{\ln(1/\delta)} = T^*(\mu)$;
- (ii) its $\mathbb{E}_{\mu}[\tau]$ being poly(K) in moderate confidence regime;
- (iii) the expected number of i^* upper bounded by poly(K).



Proof Sketch of Theorem 2

Define good events: $\mathcal{E}_t^{(1)}$ when $\hat{\mu}(t)$ is sufficiently close to μ , and $\mathcal{E}_t^{(2)}$ when $\mathbf{x}(t)$ is closed to the ideal FW-update.

- (Step 1)** By maximum theorem [FKV14], we derive uniform continuity of F_π and $\nabla \bar{F}_{\pi,\eta}$ in π
 \Rightarrow to simplify the analysis as if $\hat{\mu}(t) = \mu$ for $t \geq M$
- (Step 2)** Under $\mathcal{E}_t^{(1)} \cap \mathcal{E}_t^{(2)}$, we derive a recursive formula for the smoothed FW updates \Rightarrow the FW algorithm converges
- (Step 3)** $\mathbb{E}_\mu[\tau] \leq T_0(\delta) + \sum_{t \geq M} \mathbb{P}_\mu \left[(\mathcal{E}_t^{(1)} \cap \mathcal{E}_t^{(2)})^c \right]$, where
- $$\begin{cases} (\delta\text{-dependent}) \frac{T_0(\delta)}{\ln \delta^{-1}} \xrightarrow{\delta \rightarrow 0} T^*(\mu) \\ (\delta\text{-independent}) \sum_{t \geq M} \mathbb{P}_\mu \left[(\mathcal{E}_t^{(1)} \cap \mathcal{E}_t^{(2)})^c \right] \leq \text{poly}(K) \end{cases}$$



P-FWS: the first poly(K)-time optimal algorithm

Algorithm 1: P-FWS($\{\epsilon_t, \eta_t, n_t, \rho_t, \theta_t\}_t$)

Initialization: pull each $x \in \mathcal{X}_0$ four times and update estimates

for $t = 4|\mathcal{X}_0| + 1, \dots$ **do**

if $\sqrt{\frac{t}{|\mathcal{X}_0|}} \in \mathbb{N}$ **or** costly to estimate $F_{\hat{\mu}(t-1)}(\hat{\omega}(t-1))$ **then**
 | pull each $x \in \mathcal{X}_0$ once;

else

 | pull $\mathbf{x}(t) \in i^*(\nabla \tilde{F}_{\hat{\mu}(t-1), \eta_t, n_t}(\hat{\omega}(t-1)))$ and update estimates;

if not costly to estimate $F_{\hat{\mu}(t)}(\hat{\omega}(t))$ **then**

 | compute \hat{F}_t by $(\epsilon_t, \frac{\delta}{t^2})$ -MCP($\hat{\omega}(t), \hat{\mu}(t)$);

 | **return** $i^*(\hat{\mu}(t))$ **if** $t\hat{F}_t > (1 + \epsilon_t)\beta(t, \frac{4|\mathcal{X}_0|-1}{|\mathcal{X}_0|}\delta)$

end

$\nabla \tilde{F}_{\hat{\mu}(t-1), \eta_t, n_t}(\hat{\omega}(t-1))$ is a n_t -sample estimation to $\nabla \bar{F}_{\hat{\mu}(t-1), \eta_t}(\hat{\omega}(t-1))$



Preliminary Numerical Results

Empirical evaluation on \mathcal{X} as the set of spanning trees

All the experiments² are performed on a Macbook Air with 16 GB memory.

Table 1: Averaged sample complexity at $\delta = 0.1$ over 100 independent runs on a graph with $|\mathcal{X}| = 21\,025$ spanning trees.

Algorithm	Sample Complexity
P-FWS (ours)	1 176
CombGame [JMKK21] with our (ϵ, θ) -MCP	1 277

Table 2: Averaged sample complexity at $\delta = 0.1$ over 100 independent runs on a graph with $|\mathcal{X}| = 343\,385$ spanning trees.

Algorithm	Sample Complexity
P-FWS (ours)	1 501
CombGame [JMKK21] with our (ϵ, θ) -MCP	OOM

²Our code: <https://github.com/rctzeng/NeurIPS2023-PerturbedFWS>.



Conclusion and Future Works




Conclusion and open questions





- Our proposed P-FWS is the first algorithm to close the statistical-computational gap for combinatorial BAI by exploring the structural properties of the lowerbound problem.
- It remains largely unexplored whether one can close the computational-statistical gap for other tasks, such as

	reward distr.	comput. efficient	stat. optimal
comb. BAI semi-bandit	Gaussian Bernoulli	[TWPL23] open	[JMKK21, WTP21]
comb. BAI bandit	Gaussian Bernoulli	open	[WTP21]
regret min.	Gaussian Bernoulli	open	[CMP17]





Table 3: Computational-statistical gap in combinatorial semi-bandits.





-  Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, Kentaro Toyoshima, and Atsushi Iwasaki, *Last-iterate convergence with full-and noisy-information feedback in two-player zero-sum games*, Proc. of AISTATS, 2023.
-  Ioannis Anagnostides, Ioannis Panageas, Gabriele Farina, and Tuomas Sandholm, *On last-iterate convergence beyond zero-sum games*, Proc. of ICML, 2022.
-  Lijie Chen, Anupam Gupta, and Jian Li, *Pure exploration of multi-armed bandit under matroid constraints*, Proc. of COLT, 2016.

-  Richard Combes, Stefan Magureanu, and Alexandre Proutiere, *Minimal exploration in structured stochastic bandits*, Proc. of NeurIPS, 2017.
-  John C Duchi, Peter L Bartlett, and Martin J Wainwright, *Randomized smoothing for stochastic optimization*, SIAM Journal on Optimization (2012).
-  Constantinos Daskalakis and Ioannis Panageas, *Last-iterate convergence: Zero-sum games and constrained min-max optimization*, Proc. of ITCS (2019).
-  Abraham D Flaxman, Adam Tauman Kalai, and H Brendan McMahan, *Online convex optimization in the bandit setting: gradient descent without a gradient*, Proc. of SODA, 2005.



-  Eugene A Feinberg, Pavlo O Kasyanov, and Mark Voorneveld, *Berge's maximum theorem for noncompact image sets*, Journal of Mathematical Analysis and Applications (2014).
-  Aurélien Garivier and Emilie Kaufmann, *Optimal best arm identification with fixed confidence*, Proc. of COLT, 2016.
-  Marc Jourdan, Mojmír Mutný, Johannes Kirschner, and Andreas Krause, *Efficient pure exploration for combinatorial bandits with semi-bandit feedback*, Proc. of ALT, 2021.
-  Qi Lei, Sai Ganesh Nagarajan, Ioannis Panageas, et al., *Last iterate convergence in no-regret learning: constrained min-max optimization for convex-concave landscapes*, Proc. of AISTATS, 2021.

-  Ruo-Chun Tzeng, Po-An Wang, Alexandre Proutiere, and Chi-Jen Lu, *Closing the computational-statistical gap in best arm identification for combinatorial semi-bandits*, Proc. of NeurIPS, 2023.
-  Po-An Wang, Ruo-Chun Tzeng, and Alexandre Proutiere, *Fast pure exploration via frank-wolfe*, Proc. of NeurIPS, 2021.