

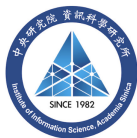
Closing the Computational-Statistical Gap in Best Arm Identification for Combinatorial Semi-bandits

Ruo-Chun Tzeng¹, **Po-An Wang**¹, Alexandre Proutiere¹, and Chi-Jen Lu²

May 2, 2026

¹EECS, KTH, Sweden

²IIS, AS, Taiwan



Introduction

1-page Summary

In combinatorial semi-bandits,...

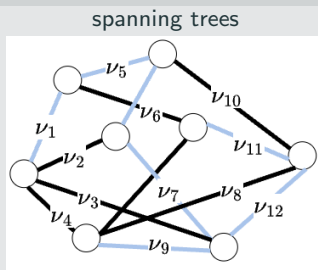
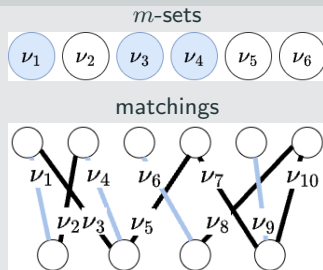
- most statistically efficient algorithms suffer **computational inefficiency** [Garivier and Kaufmann, 2016, Wang et al., 2021, Jourdan et al., 2021].
- most computationally efficient algorithms suffer **statistical inefficiency** [Katz-Samuels et al., 2020, Du et al., 2021].
- we fulfill the computational-statistical gap.



Combinatorial semi-bandits

- K arms (K probability distributions), say
 $\nu_1 = \mathcal{N}(\mu_1, 1), \nu_2 = \mathcal{N}(\mu_2, 1) \dots, \nu_K = \mathcal{N}(\mu_K, 1)$.
- Set of actions, $\mathcal{X} \subseteq \{0, 1\}^K$, with a combinatorial structure.

Examples for combinatorial structures



In round t , an agent

1. pulls an action $x(t) \in \mathcal{X}$
2. receives the reward $y_k(t) \sim \mathcal{N}(\mu_k, 1)$ if $x_k(t) = 1$

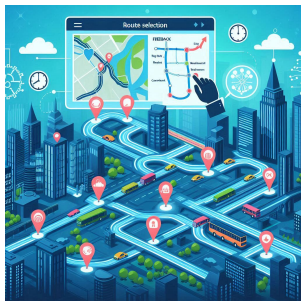
Best Action Identification with Fixed Confidence

Goal: identify the best action $i^*(\mu) \in \arg \max_{x \in \mathcal{X}} \langle x, \mu \rangle$.

A strategy consist of

- (sampling rule) $x(t) \in \mathcal{F}_t := \sigma(\mathbf{x}(1), \mathbf{y}(1), \dots, \mathbf{x}(t-1), \mathbf{y}(t-1))$ (arm to explore)
- (stopping rule) τ (round to stop)
- (decision rule) \mathcal{F}_τ -measurable $\hat{i} \in \mathcal{X}$ (guess of best action to return)

Wish to minimize $\mathbb{E}_\mu[\tau]$ subject to $\mathbb{P}_\mu[\hat{i} \neq i^*(\mu)] \leq \delta$.



In traffic management, authorities might select a set of routes for traffic flow optimization. Feedback on each route (e.g., congestion levels, travel times) helps improve future route selections.

Information-theoretic Lower Bound [Garivier and Kaufmann, 2016]

For any δ -PAC algorithm (stop in finite time a.e. and $\mathbb{P}_\mu[\hat{i} \neq i^*(\mu)] \leq \delta$),

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau]}{\ln(1/\delta)} \geq \left(\sup_{\omega \in \Sigma} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{k=1}^K \frac{\omega_k (\mu_k - \lambda_k)^2}{2} \right)^{-1}. \quad (1)$$

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

(1) \Rightarrow An optimal algorithm has a sampling strategy described by

$$\hat{\omega}(t) \xrightarrow{t \rightarrow \infty} \omega^*(\mu) := \arg \max_{\omega \in \Sigma} F_\mu(\omega),$$

where $\hat{\omega}(t)$ is the empirical allocation of arm draws up to round t , and

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

A Computational Challenge for Matching Lower Bound

A standard way [Garivier and Kaufmann, 2016] consists of:

- (sampling rule) pull an action by tracking $\arg \max_{\omega \in \Sigma} F_{\hat{\mu}(t)}(\omega)$
- (stopping rule) Generalized Likelihood Ratio Test (GLRT):

$$\tau = \inf \left\{ t : t F_{\hat{\mu}(t)}(\hat{\omega}(t)) > \ln \left(\frac{t}{\delta} \right) + o(1) \right\} \quad (2)$$

- (decision rule) return $\hat{i} \leftarrow i^*(\hat{\mu}(\tau))$

where $\hat{\mu}(t)$ is the estimate of μ , and $\hat{\omega}(t)$ is empirical pulled arm allocation up to t .

The value of $F_{\mu}(\omega) := \inf_{\lambda \in \text{Alt}(\mu)} \sum_{k=1}^K \frac{\omega_k(\mu_k - \lambda_k)^2}{2}$ is the most confusing parameter (MCP).

Difficulty in identifying MCP

Recall $\text{Alt}(\mu) = \{\lambda \in \Lambda : i^*(\lambda) \neq i^*(\mu)\}$

\Rightarrow Prior approach is to solve $|\mathcal{X}| - 1$ many convex programs by partitioning

$\text{Alt}(\mu) = \cup_{x \neq i^*(\mu)} \{\lambda \in \Lambda : \langle i^*(\mu) - x, \lambda \rangle < 0\}$

However, $|\mathcal{X}| = \mathcal{O}(2^K)!$

Other Statistically Optimal Sampling Rules

- FWS [Wang et al., 2021]: apply the first-order methods to maximize $F_{\mu}(\omega)$, where the gradient is derived by a MCP (Envelop theorem)
- CombGame [Jourdan et al., 2021]: need a MCP for chasing a saddle point.
- ACC [Qin and You, 2023]: need a MCP to get its best challenger.

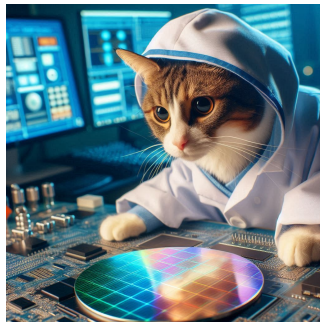


Figure 1: MCP in pure exploration resembles the role of the semi-conductor in modern technology.

Our Perturbed Frank-Wolfe Sampling (P-FWS)

- P-FWS uses perturbed Frank-Wolfe based algorithm to reach the optimal allocation
- All P-FWS needs are
 - (i) linear maximization (standard),
 - (ii) (ε, θ) -MCP, where $\varepsilon > 0, \theta \in (0, 1)$ (established in this work).

(ε, θ) -MCP

Return \hat{F} s.t.

$$\mathbb{P}[F_{\mu}(\omega) \leq \hat{F} \leq (1 + \varepsilon)F_{\mu}(\omega)] \geq 1 - \theta$$

How to get (ε, θ) -MCP?

Let Us Observe MCP

For each suboptimal action $\mathbf{x} \neq i^*(\boldsymbol{\mu})$, define

$$f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu}) := \inf_{\boldsymbol{\lambda} \in \mathbb{R}^K: \langle i^*(\boldsymbol{\mu}) - \mathbf{x}, \boldsymbol{\lambda} \rangle < 0} \sum_{k=1}^K \frac{\omega_k (\mu_k - \lambda_k)^2}{2}. \quad (3)$$

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

$$f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu}) = \max_{\alpha \geq 0} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha) \text{ (strong duality)}$$

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

$$\Rightarrow F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \min_{\mathbf{x} \neq i^*(\boldsymbol{\mu})} f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu}) = \min_{\mathbf{x} \neq i^*(\boldsymbol{\mu})} \max_{\alpha \geq 0} g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha) \quad (4)$$

(4) motivates us a design:

- **x-player** employs a regret minimization algorithm to **minimize** $g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha)$
- **α -player** employs a regret minimization algorithm to **maximize** $g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha)$

Two-players Game

However, the choices are limited since...

1. $\{x \neq i^*(\mu)\}$ is a discrete set which consist of $\mathcal{O}(2^K)$ elements
2. $\min \max = \max \min$ may not hold
3. *equilibrium action* x_e s.t. $F_\mu(\omega) = \max_{\alpha \geq 0} g_{\omega, \mu}(x_e, \alpha)$ need a maximizer in α
4. we want it to be computationally efficient, i.e. approximate x_e in polynomial time

Recall

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



Figure 2: Two-cats Game

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 $\eta_n = \frac{c_0}{\sqrt{n}}$

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

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

if $\sqrt{n} > \frac{c_\theta(1 + \varepsilon)}{\varepsilon \hat{F}}$, where $\hat{F} = \min_n g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)})$ **then return** \hat{F} ;

end

Main Results

Theorem 1

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

- $\mathbb{P} \left[F_\mu(\omega) \leq \hat{F} \leq (1 + \varepsilon)F_\mu(\omega) \right] \geq 1 - \theta$
- the number of calls to linear maximization ($i^*(\cdot)$) is polynomial

Theorem 2

Let $\mu \in \Lambda$ and $\delta \in (0, 1)$. P-FWS is δ -PAC, and

$$(i) \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau]}{\ln(1/\delta)} \leq \left(\sup_{\omega \in \Sigma} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{k=1}^K \frac{\omega_k (\mu_k - \lambda_k)^2}{2} \right)^{-1}$$

(ii) $\mathbb{E}_{\mu}[\tau] = \text{poly}(K)$ in moderate confidence regime

(iii) the expected number of calls for linear maximization upper bounded by $\text{poly}(K)$

Note: To our best knowledge, (ii) is even not yet established in unstructured BAI.



Figure 3: A cat wins three medals in *Olympicat*.

Conclusion and Future Works

Conclusion: P-FWS closes the statistical-computational gap for combinatorial BAI by exploiting the structural properties in lower bound, and it show strong numerical performance.





Table 1: 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 [Jourdan et al., 2021] with our (ϵ, θ) -MCP	OOM

Future works: It remains largely unexplored whether one can close the computational-statistical gap for other distributions.

	reward distr.	comput. & stat. efficient	stat. optimal
comb. BAI semi-bandit	Gaussian Bernoulli	[Tzeng et al., 2023] open	[Jourdan et al., 2021, Wang et al., 2021]
comb. BAI bandit	Gaussian Bernoulli	open	[Wang et al., 2021]
regret min.	Gaussian Bernoulli	open	[Combes et al., 2017]

Table 2: Computational-statistical gap in combinatorial semi-bandits.

-  Combes, R., Magureanu, S., and Proutiere, A. (2017).
Minimal exploration in structured stochastic bandits.
In *Proc. of NeurIPS*.
-  Du, Y., Kuroki, Y., and Chen, W. (2021).
Combinatorial pure exploration with full-bandit or partial linear feedback.
In *Proc. of AAAI*.
-  Garivier, A. and Kaufmann, E. (2016).
Optimal best arm identification with fixed confidence.
In *Proc. of COLT*.
-  Jourdan, M., Mutn y, M., Kirschner, J., and Krause, A. (2021).
Efficient pure exploration for combinatorial bandits with semi-bandit feedback.
In *Proc. of ALT*.

-  Katz-Samuels, J., Jain, L., Jamieson, K. G., et al. (2020).
An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits.
In Proc. of NeurIPS.
-  Qin, C. and You, W. (2023).
Dual-directed algorithm design for efficient pure exploration.
arXiv preprint arXiv:2310.19319.
-  Tzeng, R.-C., Wang, P.-A., Proutiere, A., and Lu, C.-J. (2023).
Closing the computational-statistical gap in best arm identification for combinatorial semi-bandits.
In Proc. of NeurIPS.
-  Wang, P.-A., Tzeng, R.-C., and Proutiere, A. (2021).
Fast pure exploration via frank-wolfe.
In Proc. of NeurIPS.

Typical Assumptions in Combinatorial BAI

Throughout this work, we assume

- (i) $\forall \boldsymbol{\mu} \in \Lambda$, the best action $i^*(\boldsymbol{\mu})$ is unique, where Λ denotes all possible parameters.
- (ii) $\forall k \in [K]$, arm- k reward distribution is $\nu_k = \mathcal{N}(\mu_k, 1)$.
- (iii) $\forall \boldsymbol{v} \in \mathbb{R}^K$, a maximizing action $i^*(\boldsymbol{v}) \in \arg \max_{\boldsymbol{x} \in \mathcal{X}} \langle \boldsymbol{x}, \boldsymbol{v} \rangle$ can be found in polynomial time.

Computational Cost of Our (ϵ, θ) -MCP

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)

- $\mathbf{x}^{(n)}$ is evaluated in $\max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_0$ calls to linear maximization $\mathbf{i}^*(\cdot)$
- $\alpha^{(n)}$ is evaluated in $\mathcal{O}(1)$

Termination Condition

The termination condition

$$\text{If } \sqrt{n} > \frac{c_\theta(1+\varepsilon)}{\varepsilon \hat{F}} \text{ where } \hat{F} = \min_n g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) \text{ then return } \hat{F} \quad (5)$$

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

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

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

Empirical Evaluation on \mathcal{X} as the Set of Spanning Trees

All the experiments¹ are performed on a Macbook Air with 16 GB memory.

Table 3: 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 [Jourdan et al., 2021] with our (ε, θ) -MCP	1 277

Table 4: 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 [Jourdan et al., 2021] with our (ε, θ) -MCP	OOM

¹Our code: <https://github.com/rctzeng/NeurIPS2023-PerturbedFWS>.