

Fast Pure Exploration via Frank-Wolfe

Po-An Wang¹, Ruo-Chun Tzeng², and Alexandre Proutiere¹

Conference on Neural Information Processing Systems, 2021

¹EECS, Division of Decision and Control System

²EECS, Division of Theoretical Computer Science
KTH Royal Institute of Technology



Pure exploration on structured bandits

Stochastic Multi-Armed Bandit (MAB)

K arms (K prob. distribution ν_1, \dots, ν_K), the mean of ν_k is μ_k



ν_1



ν_2



ν_3



ν_4



ν_5

Stochastic Multi-Armed Bandit (MAB)

K arms (K prob. distribution ν_1, \dots, ν_K), the mean of ν_k is μ_k



ν_1



ν_2



ν_3



ν_4



ν_5

In round t , an agent

1. pulls arm $A_t \in [K]$
2. receives the reward $X_{A_t}(t) \sim \nu_{A_t}$

Sequential sampling strategy: $A_t \in \mathcal{F}_t = \sigma[A_1, X_1, \dots, A_{t-1}, X_{t-1}]$

Pure exploration with fixed confidence

Goal: Identify a certain answer $i^*(\boldsymbol{\mu}) \in \mathcal{I}$

Example: Identify the best arm $i^*(\boldsymbol{\mu}) = \operatorname{argmax}_{k \in [K]} \mu_k$

A strategy consists of

- a sampling rule A_t (arm to explore)
- a stopping rule τ (time to stop)
- a \mathcal{F}_τ -measurable decision rule $\hat{i} \in \mathcal{I}$ (answer to return)

Pure exploration with fixed confidence

Goal: Identify a certain answer $i^*(\boldsymbol{\mu}) \in \mathcal{I}$

Example: Identify the best arm $i^*(\boldsymbol{\mu}) = \operatorname{argmax}_{k \in [K]} \mu_k$

A strategy consists of

- a sampling rule A_t (arm to explore)
- a stopping rule τ (time to stop)
- a \mathcal{F}_τ -measurable decision rule $\hat{i} \in \mathcal{I}$ (answer to return)

We wish to minimize $\mathbb{E}_\mu[\tau]$ subject to $\mathbb{P}_\mu[\hat{i} \neq i^*(\boldsymbol{\mu})] < \delta$



“Side information” is encoded by the **structure**

Popular structures: Unstructured, Linear, Lipschitz, Dueling,
Combinatorial, Unimodal, Monotone, Spectral and Cascading

“Side information” is encoded by the **structure**

Popular structures: Unstructured, Linear, Lipschitz, Dueling, Combinatorial, Unimodal, Monotone, Spectral and Cascading

Question 1. What is the sample complex gain achievable when exploiting the structure?

“Side information” is encoded by the **structure**

Popular structures: Unstructured, Linear, Lipschitz, Dueling, Combinatorial, Unimodal, Monotone, Spectral and Cascading

Question 1. What is the sample complex gain achievable when exploiting the structure?

Question 2. Can we devise a computational efficient algorithm achieving the promised gains for all structures?

Lower bound [GK16]

For any **good** strategy,

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau]}{\log(\frac{1}{\delta})} \geq T^*(\mu),$$

where $T^*(\mu)^{-1} = \sup_{\omega \in \Sigma} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{k=1}^K \omega_k d(\mu_k, \lambda_k)$

- Σ : $K - 1$ simplex
- $\text{Alt}(\mu) = \{\lambda \in \Lambda : i^*(\lambda) \neq i^*(\mu)\}$
- $d(\mu_k, \lambda_k)$: KL-divergent of arm-k reward distribution under λ and μ

Lower bound [GK16]

For any **good** strategy,

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau]}{\log(\frac{1}{\delta})} \geq T^*(\mu),$$

where $T^*(\mu)^{-1} = \sup_{\omega \in \Sigma} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{k=1}^K \omega_k d(\mu_k, \lambda_k)$

- Σ : $K - 1$ simplex
- $\text{Alt}(\mu) = \{\lambda \in \Lambda : i^*(\lambda) \neq i^*(\mu)\}$
- $d(\mu_k, \lambda_k)$: KL-divergent of arm- k reward distribution under λ and μ

\Rightarrow An optimal algorithm has a sampling strategy described by

$$\omega^*(\mu) = \operatorname{argmax}_{\omega \in \Sigma} F_{\mu}(\omega),$$

$$\text{where } F_{\mu}(\omega) = \inf_{\lambda \in \text{Alt}(\mu)} \sum_{k=1}^K \omega_k d(\mu_k, \lambda_k)$$



Frank-Wolfe based sampling (FWS)

Frank-Wolfe based sampling

- Devise a simple algorithm (FW-based) to track $\mathbf{x}(\mathbf{t}) \xrightarrow{t \rightarrow \infty} \omega^*(\boldsymbol{\mu})$

Frank-Wolfe based sampling

- Devise a simple algorithm (FW-based) to track $\mathbf{x}(\mathbf{t}) \xrightarrow{t \rightarrow \infty} \boldsymbol{\omega}^*(\boldsymbol{\mu})$
- **Envelope theorem** shows that $F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \min_{j \in J} f_j(\boldsymbol{\omega}, \boldsymbol{\mu})$, where J is a finite set and $f_j(\boldsymbol{\omega}, \boldsymbol{\mu})$ is smooth $\forall j \in J$ ($F_{\boldsymbol{\mu}}$ is **non-smooth**)

Frank-Wolfe based sampling

- Devise a simple algorithm (FW-based) to track $\mathbf{x}(\mathbf{t}) \xrightarrow{t \rightarrow \infty} \boldsymbol{\omega}^*(\boldsymbol{\mu})$
- **Envelope theorem** shows that $F_\mu(\boldsymbol{\omega}) = \min_{j \in J} f_j(\boldsymbol{\omega}, \boldsymbol{\mu})$, where J is a finite set and $f_j(\boldsymbol{\omega}, \boldsymbol{\mu})$ is smooth $\forall j \in J$ (F_μ is **non-smooth**)
- To deal with non-smoothness, define

$$H_{F_\mu}(\boldsymbol{\omega}, r) = \text{cov} \{ \nabla_{\boldsymbol{\omega}} f_j(\boldsymbol{\omega}, \boldsymbol{\mu}) : j \in \mathcal{J}, f_j(\boldsymbol{\omega}, \boldsymbol{\mu}) < F_\mu(\boldsymbol{\omega}) + r \}$$

Frank-Wolfe based sampling

- Devise a simple algorithm (FW-based) to track $\mathbf{x}(t) \xrightarrow{t \rightarrow \infty} \boldsymbol{\omega}^*(\boldsymbol{\mu})$
- **Envelope theorem** shows that $F_\mu(\boldsymbol{\omega}) = \min_{j \in J} f_j(\boldsymbol{\omega}, \boldsymbol{\mu})$, where J is a finite set and $f_j(\boldsymbol{\omega}, \boldsymbol{\mu})$ is smooth $\forall j \in J$ (F_μ is **non-smooth**)

- To deal with non-smoothness, define

$$H_{F_\mu}(\boldsymbol{\omega}, r) = \text{cov} \{ \nabla_{\boldsymbol{\omega}} f_j(\boldsymbol{\omega}, \boldsymbol{\mu}) : j \in \mathcal{J}, f_j(\boldsymbol{\omega}, \boldsymbol{\mu}) < F_\mu(\boldsymbol{\omega}) + r \}$$

- Update

$$\begin{cases} \mathbf{z}(t+1) \leftarrow \operatorname{argmax}_{\mathbf{z} \in \Sigma} \min_{h \in H_{F_\mu}(\mathbf{x}(t), r_t)} \langle \mathbf{z} - \mathbf{x}(t), h \rangle, \\ \mathbf{x}(t+1) \leftarrow \frac{t}{t+1} \mathbf{x}(t) + \frac{1}{t+1} \mathbf{z}(t+1) \end{cases}$$



Input: Confidence level δ , sequence $\{r_t\}_{t \geq 1}$

Initialization: Sample each arm once and update $\omega(K), \mathbf{x}(K) = (\frac{1}{K}, \dots, \frac{1}{K})$, and $\hat{\mu}(K)$
 $t \leftarrow K$

While $t F_{\hat{\mu}(t)}(\omega(t)) < \beta(\delta, t)$ \leftarrow **Stopping criteria** or $\hat{\mu}(t-1) \notin \Lambda$

IF $\sqrt{[t/K]} \in \mathbb{N}$ or $\hat{\mu}(t-1) \notin \Lambda$, (Forced exploration) $\mathbf{z}(t) \leftarrow (\frac{1}{K}, \dots, \frac{1}{K})$

Else, (FW update)

$$\mathbf{z}(t) \leftarrow \operatorname{argmax}_{\mathbf{z} \in \Sigma} \min_{h \in H_{F_{\hat{\mu}(t-1)}}(\mathbf{x}(t-1), r_t)} \langle \mathbf{z} - \mathbf{x}(t-1), h \rangle$$

$$\text{Update } \mathbf{x}(t) \leftarrow \frac{t-1}{t} \mathbf{x}(t-1) + \frac{1}{t} \mathbf{z}(t)$$

Sample $A_t \leftarrow \operatorname{argmax}_k \mathbf{x}_k(t) / \omega_k(t-1)$ (ties broken arbitrarily)

Update $\omega(t)$ and $\hat{\mu}(t)$

Output: $i^*(\hat{\mu}(t))$



Theoretical Results

Theorem

For most pure exploration problems in structured bandits, FWS satisfies:

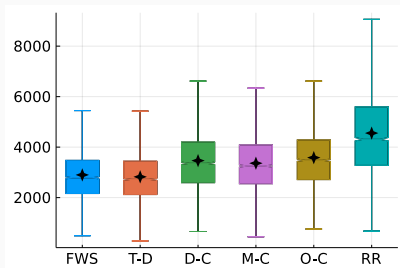
$$\mathbb{P}_{\mu}[\hat{i} \neq i^*(\mu)] < \delta \text{ and } \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau]}{\log(\frac{1}{\delta})} \leq T^*(\mu)$$

With further assumptions, we can provide **non-asymptotic** upper bound for $\mathbb{E}_{\mu}[\tau]$

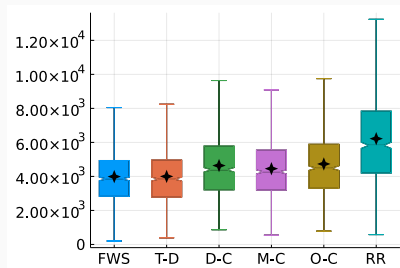
Numerical Results

Experiment (i) Unstructured bandits

Averaged sample complexity at $\delta = 0.01$



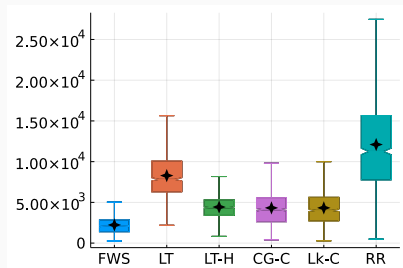
Bernoulli



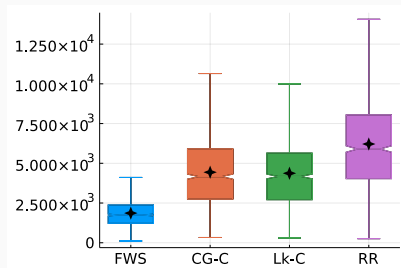
Gaussian

Experiment (ii) Linear bandits

Averaged sample complexity at $\delta = 0.01$



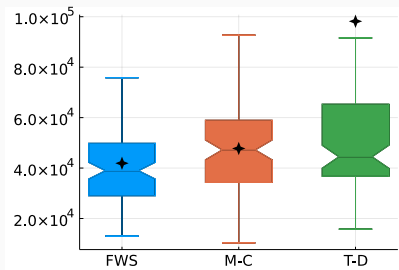
BAI



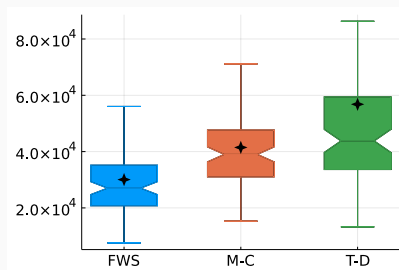
ThresholdingBandit

Experiment (iii) Lipschitz bandits

Averaged Sample complexity at $\delta = 0.01$



Experiment 1



Experiment 2

This is the first result for Lipschitz bandits in literatures



Related work and conclusion

Related works:

- LMA [Mén19]: Apply mirror ascent to update $\mathbf{x}(t)$
- Gamification [DMSV20, Sha21, JMKK21]: Use 2 player game to reach $\omega^*(\mu)$

Unclear to extend the above approaches to general structures

Related work and conclusion

Related works:

- LMA [Mén19]: Apply mirror ascent to update $\mathbf{x}(t)$
- Gamification [DMSV20, Sha21, JMKK21]: Use 2 player game to reach $\omega^*(\mu)$






Unclear to extend the above approaches to general structures

Conclusion:

- FWS is computationally and statistically efficient for general pure exploration problems
- Theoretically, FWS matches the instance-specific lower bounds
- Numerically, FWS outperforms all the other optimal algorithms in structured bandits



Reference

-  Rémy Degenne, Pierre Ménard, Xuedong Shang, and Michal Valko, *Gamification of pure exploration for linear bandits*, Proc. of ICML, 2020.
-  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.
-  Pierre Ménard, *Gradient ascent for active exploration in bandit problems*, arXiv (2019).
-  Xuedong Shang, *Linbai: Gamification of pure exploration for linear bandits*, <https://github.com/xuedong/LinBAI.jl>, 2021, [Online; accessed 09-May-2021].

