

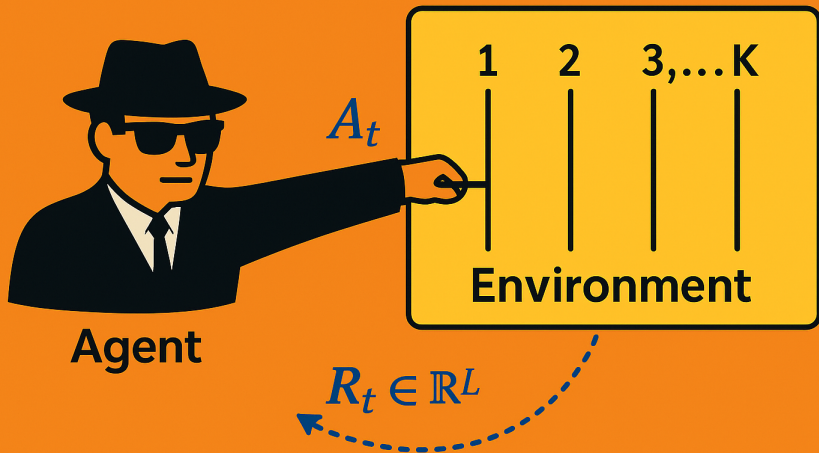


FraPPE: Fast and Efficient Preference-based Pure Exploration

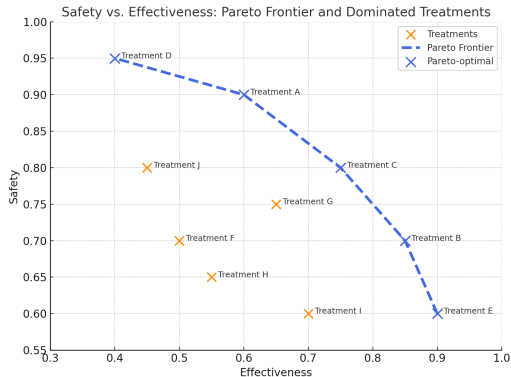
Udvas Das, Apurv Shukla, Debabrota Basu

NeurIPS, 2025

Multi-Objective Bandits



Pareto optimality: Conflicting objectives and need of preferences



Given a matrix of mean rewards of K -arms

$$M^{K \times L} \triangleq [M_1, M_2, \dots, M_K] \in \mathcal{M}$$

and a preference cone

$$\mathcal{C} \triangleq \{\mathbf{x} \in \mathbb{R}^L \mid W\mathbf{x} \geq 0\},$$

the pareto optimal policy set is:

$$\pi^* \in \Pi^P \triangleq \arg \max_{\pi \in \Delta_K} M\pi \text{ w.r.t } \mathcal{C}$$

Preference-based Pure Exploration (PrePEX)

Identify the set of Pareto optimal arms $\mathcal{P}^* \subset \{1, 2, \dots, K\}$ with probability at least $1 - \delta$, while keeping **the expected number of interactions** $\mathbb{E}[\tau_\delta]$ **as low as possible**.

Lower bound on sample complexity: Characteristic time

Theorem (Lower bound (Shukla and Basu, 2024))

Given a bandit instance $M \in \mathcal{M}$, a preference cone \mathcal{C} , the expected stopping time of any $(1 - \delta)$ -correct PrePEX algorithm satisfies

$$\mathbb{E}[\tau_\delta] \geq \mathcal{T}_{M,\mathcal{C}} \log \left(\frac{1}{2.4\delta} \right)$$

Optimal sampling ratio over arms

Most confusing instance

$$(\mathcal{T}_{M,\mathcal{C}})^{-1} \triangleq \sup_{\omega \in \Delta_K} \inf_{\substack{\pi \in \Delta_K \setminus \{\pi^*\} \\ \pi^* \in \Pi^P(M,\mathcal{C})}} \inf_{\tilde{M} \in \partial \Lambda(M)} \underbrace{\inf_{z \in \mathcal{C}^+} \sum_{k=1}^K \omega_k D_{\text{KL}} \left(z^\top M_k \parallel z^\top \tilde{M}_k \right)}_{\text{Preference that makes true and alt-instance most indistinguishable}}$$

Optimising over Pareto optimal policy space

Efficiently solving the lower bound for efficient PrePEX

Step 1: Preference optimisation

Preference cone is closed, convex, and compact \implies plug-in a cone programming solver

$$(\mathcal{T}_{M,\mathcal{C}})^{-1} \triangleq \sup_{\omega \in \Delta_K} \inf_{\substack{\pi \in \Delta_K \setminus \{\pi^*\} \\ \pi^* \in \Pi^P(M,\mathcal{C})}} \inf_{\tilde{M} \in \partial \Lambda(M)} \underbrace{\min_{z \in \mathcal{C}^+} \sum_{k=1}^K \omega_k D_{\text{KL}} \left(z^\top M_k \parallel z^\top \tilde{M}_k \right)}_{\text{Preference that makes true and alt-instance most indistinguishable}}$$

Optimal sampling ratio over arms

Most confusing instance

Optimising over Pareto optimal policy space

Efficiently solving the lower bound for efficient PrePEX

Step 2: Reducing policy set to $\mathcal{O}(\min\{K, L\})$ arms

- Π^P is spanned by p pure policies $\{\pi_i^*\}_{i=1}^p$ corresponding to p Pareto optimal arm.
- Number of neighbour of any Pareto optimal policy is $\mathcal{O}(\min\{K, L\})$.

Optimal sampling ratio over arms

Most confusing instance

$$(\mathcal{T}_{M,\mathcal{C}})^{-1} \triangleq \sup_{\omega \in \Delta_K} \underbrace{\min_{\substack{\pi_j \in \text{nbd}(\pi_i^*) \\ \pi_i^* \in \{\pi_i^*\}_{i=1}^p}} \inf_{\tilde{M} \in \partial \Lambda(M)} \underbrace{\min_{z \in \mathcal{C}^+} \sum_{k=1}^K \omega_k D_{\text{KL}} \left(z^\top M_k \parallel z^\top \tilde{M}_k \right)}_{\text{Preference that makes true and alt-instance most indistinguishable}}}_{\text{Optimising over Pareto optimal policy space}}$$

Efficiently solving the lower bound for efficient PrePEx

Step 3: Reducing the Alt-set to union of lines

- The solution always lies at the boundary of the Alt-set.
- The boundary of the Alt-set is union of $\mathcal{O}(KL)$ lines

$$\bar{\Lambda}_{ij}(M) \triangleq \left\{ \tilde{M} \in \mathcal{M} \setminus \{M\} : \exists \mathbf{y} \in \text{bd}(\mathcal{C}^\circ) \text{ such that } \tilde{M}(\pi_j - \pi_i^*) = \mathbf{y} \right\}$$

Optimal sampling ratio over arms

Most confusing instance

$$(\mathcal{T}_{M,\mathcal{C}})^{-1} \triangleq \sup_{\omega \in \Delta_K} \underbrace{\min_{\substack{\pi_j \in \text{nbd}(\pi_i^*) \\ \pi_i^* \in \{\pi_i^*\}_{i=1}^P}} \min_{\mathbf{y} \in \text{bd}(\mathcal{C}^\circ)} \underbrace{\min_{\mathbf{z} \in \mathcal{C}^+} \sum_{k=1}^K \omega_k D_{\text{KL}} \left(\mathbf{z}^\top M_k \parallel \mathbf{z}^\top \mathbf{y} f(\pi_i^*, \pi_j) \right)}_{\text{Preference that makes true and alt-instance most indistinguishable}}}_{\text{Optimising over Pareto optimal policy space}}$$

Reduces the complexity from $\mathcal{O}(K^L)$ of (Crepon et al., 2024) to $\mathcal{O}(L)$.

Efficiently solving the lower bound for efficient PrePEx

Step 4: Frank-Wolfe for allocation optimisation in $\mathcal{O}(1)$

- Simplex is convex and closed. So, apply

$$\omega_{t+1} \leftarrow \text{FRANKWOLFE}(\omega_t, r_t, \hat{M}_t, \mathcal{C})$$

- The function to be optimised has bounded gradient and curvature.

$$(\mathcal{T}_{M,\mathcal{C}})^{-1} \triangleq \underbrace{\max_{\omega \in \Delta_K} \min_{\substack{\pi_j \in \text{nbd}(\pi_i^*) \\ \pi_i^* \in \{\pi_i^*\}_{i=1}^p}}}_{\text{Optimal sampling ratio over arms}} \min_{\mathbf{y} \in \text{bd}(\mathcal{C}^\circ)} \underbrace{\min_{\mathbf{z} \in \mathcal{C}^+} \sum_{k=1}^K \omega_k D_{\text{KL}} \left(\mathbf{z}^\top M_k \parallel \mathbf{z}^\top \mathbf{y} f(\pi_i^*, \pi_j) \right)}_{\text{Most confusing instance}}$$

Preference that makes true and alt-instance most indistinguishable

Optimising over Pareto optimal policy space

We solve the lower bound with $\mathcal{O}(KL \min\{K, L\})$ complexity.

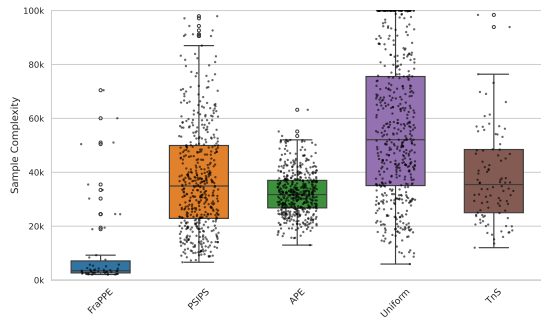
FraPPE: Frugal and Fast Preference-Based Pure Exploration

Algorithm FraPPE- Frugal and Fast Preference-Based Pure Exploration

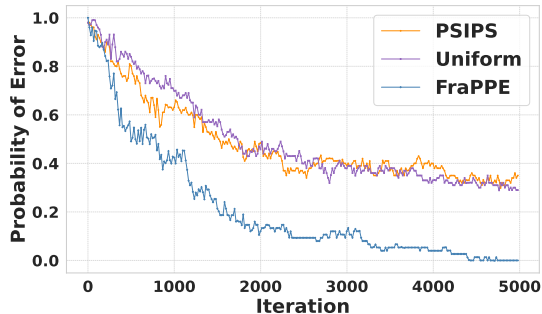
- 1: **Input:** Confidence level δ , sequence $\{r_t\}_{t \geq 1} = t^{-0.9}/K$
 - 2: **Initialise:** For $t \in [K]$, sample each arm once s.t. $\omega_K = (1/K, \dots, 1/K)$, estimate \hat{M}_K
 - 3: **for** $t > K$ **do**
 - 4: **Estimate Pareto Indices:** Compute \mathcal{P}_t from \hat{M}_t
 - 5: Compute the allocation policy with Frank-Wolfe: $\omega_{t+1} \leftarrow \text{FRANKWOLFE}(\omega_t, r_t, \hat{M}_t, \mathcal{C})$
 - 6: **C-tracking:** Play $A_t \in \arg \min_{a \in [K]} N_{a,t} - \sum_{s=1}^{t+1} \omega_{a,s}$ (ties broken arbitrarily)
 - 7: **Feedback:** Get $R_t \in \mathbb{R}^L$ and update $\hat{M}_t \rightarrow \hat{M}_{t+1}$
 - 8: **IF** $\min_{\pi_{i_t}^* \in \{\pi_i^*\}_{i=1}^p} \min_{\pi_j \in \text{nbd}(\pi_{i_t}^*)} \inf_{y \in \text{bd}(\mathcal{C}^\circ)} \min_{z \in \mathcal{C}^+} N_{t+1}^\top D_{\text{KL}} \left(z^\top \hat{M}_{t+1} \parallel z^\top y f(\pi_j, \pi_{i_t}^*) \right) > c(t+1, \delta)$
 break;
 - 9: **end for**
 - 10: **Recommend:** \mathcal{P}_t as Pareto optimal set
-

FraPPE achieves asymptotically optimal sample complexity and a per-iteration computational complexity dominated by Pareto set computation.

Empirical performance: Cov-Boost ($K = 20, L = 3, \delta = 0.01$)



$\sim 4 - 10\times$ reduction



Uniformly lower probability of error

The Parting Message

We resolve an extension of the open problem (Crepon et al., 2024) for solving PrePEX

We design

- **a computationally efficient (polynomial in both K and L) and**
- **statistically optimal PrePEX algorithm**
- **beyond Gaussian rewards and**
- **for arbitrary preference cones.**

What's Ahead?

To scale **FraPPE** to practical applications of PrePEX, e.g. aligning LLMs with RL under Human Feedback (RLHF) (Ji et al., 2023).

References I

- Crepon, E., Garivier, A., and M Koolen, W. (2024). Sequential learning of the Pareto front for multi-objective bandits. In Dasgupta, S., Mandt, S., and Li, Y., editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 3583–3591. PMLR.
- Ji, J., Liu, M., Dai, J., Pan, X., Zhang, C., Bian, C., Chen, B., Sun, R., Wang, Y., and Yang, Y. (2023). Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36:24678–24704.
- Shukla, A. and Basu, D. (2024). Preference-based pure exploration. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.