

# Optimal Exploration is no harder than Thompson Sampling

Zhaoqi Li, Kevin Jamieson, Lalit Jain

University of Washington



## Motivation

market automation



optimal design

optimal recommendation

$\theta_*$ : unknown population preferences

**Goal: find the best action to maximize profit!**

However, large companies have millions of users and millions of items are listed

⇒ **Can we find the best action quickly on large scale?**

## Problem Statement

**Input:**  $X \subset \mathbb{R}^d$

**for**  $t = 1, 2, \dots$

1. Learner chooses  $x_t \in X$

2. Nature reveals  $y_t = \langle x_t, \theta_* \rangle + \epsilon_t \rightarrow N(0, 1)$

**Goal: identify**  $x_* := \arg \max_{x \in X} \langle x, \theta_* \rangle$  **as quickly as possible**

## Key Advantages of Our Algorithm

- Achieves **optimal** rate ⇒ find best action quickly
- Algorithm is **easy** to implement on **large scale**

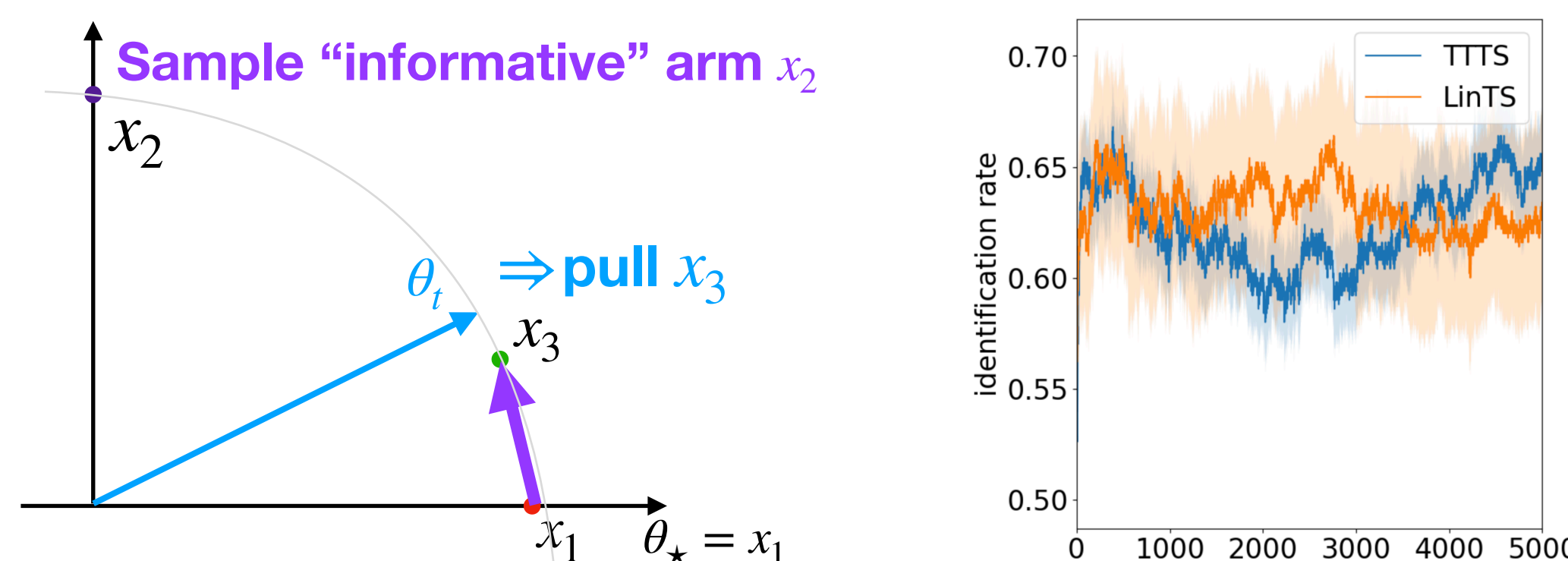
## Related Work

**Existing Optimal Approaches:** [Xu et al. 2018], [Fiez et al. 2019], [Degenne et al. 2020]

- Require **enumeration of X** or **complicated projections**

**Thompson Sampling:**

- Regret minimizing algorithm and **computationally easy**
- However **suboptimal** for BAI on certain instances



**Is there an algorithm as easy as TS yet still optimal?**

## Key Ideas behind Our Algorithm

The optimal allocation is implied by the lower bound  $\tau_*$   
⇒ solve the saddle point problem using online learning!

$$\tau_* := \max_{\lambda \in \Delta_X} \min_{\theta \in \Theta_{x_*}^c} \|\theta_* - \theta\|_{A(\lambda)}^2 = \max_{\lambda \in \Delta_X} \min_{p \in \Delta(\Theta_{x_*}^c)} \mathbb{E}_{\theta \sim p} [\|\theta_* - \theta\|_{A(\lambda)}^2]$$

Maintain a distribution over  $\Theta_{x_*}^c := \{\theta \in \Theta : x_* \neq \arg \max_{x \in X} x^\top \theta\}$

**Idea: Two player zero-sum game**

- Max-player: Exponential Weights on  $\lambda \in \Delta_X$
- Min-player: **Posterior Updates** on  $p \in \Delta(\Theta_{x_*}^c)$

$$p_{t+1, \theta} \propto \exp\left(-(\theta - \hat{\theta}_{t+1})^\top \left(\sum_{s=1}^t x_s x_s^\top\right) (\theta - \hat{\theta}_{t+1})\right) \\ \approx \exp\left(-\sum_{s=1}^t (x_s^\top (\theta - \hat{\theta}_s))^2\right) \propto p_{t, \theta} \exp\left(-(x_t^\top (\theta - \hat{\theta}_t))^2\right)$$

⇒ **exponential weights update!**

Also, **sampling from p is easy** since  $p_{t+1}$  is Gaussian posterior

By OCO argument,  $(\lambda_t, p_t) \rightarrow \max_{\lambda \in \Delta_X} \min_{p \in \Delta(\Theta_{x_*}^c)} \mathbb{E}_{\theta \sim p} [\|\theta_* - \theta\|_{A(\lambda)}^2] =: \tau_*$

## Our Algorithm

**PEPS:** Pure Exploration with Projection-Free Sampling

Input:  $X, T, \eta$

for  $t = 1, 2, \dots, T$

- compute the leader  $\hat{x}_t = \arg \max_{x \in X} x^\top \hat{\theta}_t$
- sample a challenger  $\theta_t \sim p_t$  whose best arm is not  $\hat{x}_t$
- sample  $x_t \sim \lambda_t$ , observe  $y_t$  pull arm x maximizing  $|x^\top (\theta_t - \hat{\theta}_t)|$
- Update  $\hat{\theta}_{t+1}$ ,  $\lambda_{t+1, x} \leftarrow \lambda_{t, x} e^{\eta(x^\top (\theta_t - \hat{\theta}_t))^2}$
- Update  $p_{t+1} = N(\hat{\theta}_{t+1}, (\sum_{s=1}^t x_s x_s^\top)^{-1})$  ← Posterior Update

Return  $\hat{x} = \arg \max_{x \in X} x^\top \theta, \theta \sim p_T$

## Performance Guarantee

**Theorem (Li, Jamieson, Jain)** Set  $\eta = O(1/\sqrt{T})$ , and assume  $\Theta$  is bounded. Then with probability 1

$$\lim_{T \rightarrow \infty} -\frac{1}{T} \log \mathbb{P}_{\theta \sim p_T} (\hat{x}_T \neq x_*) = \tau_* \quad \text{optimal rate!}$$

## Experimental Results

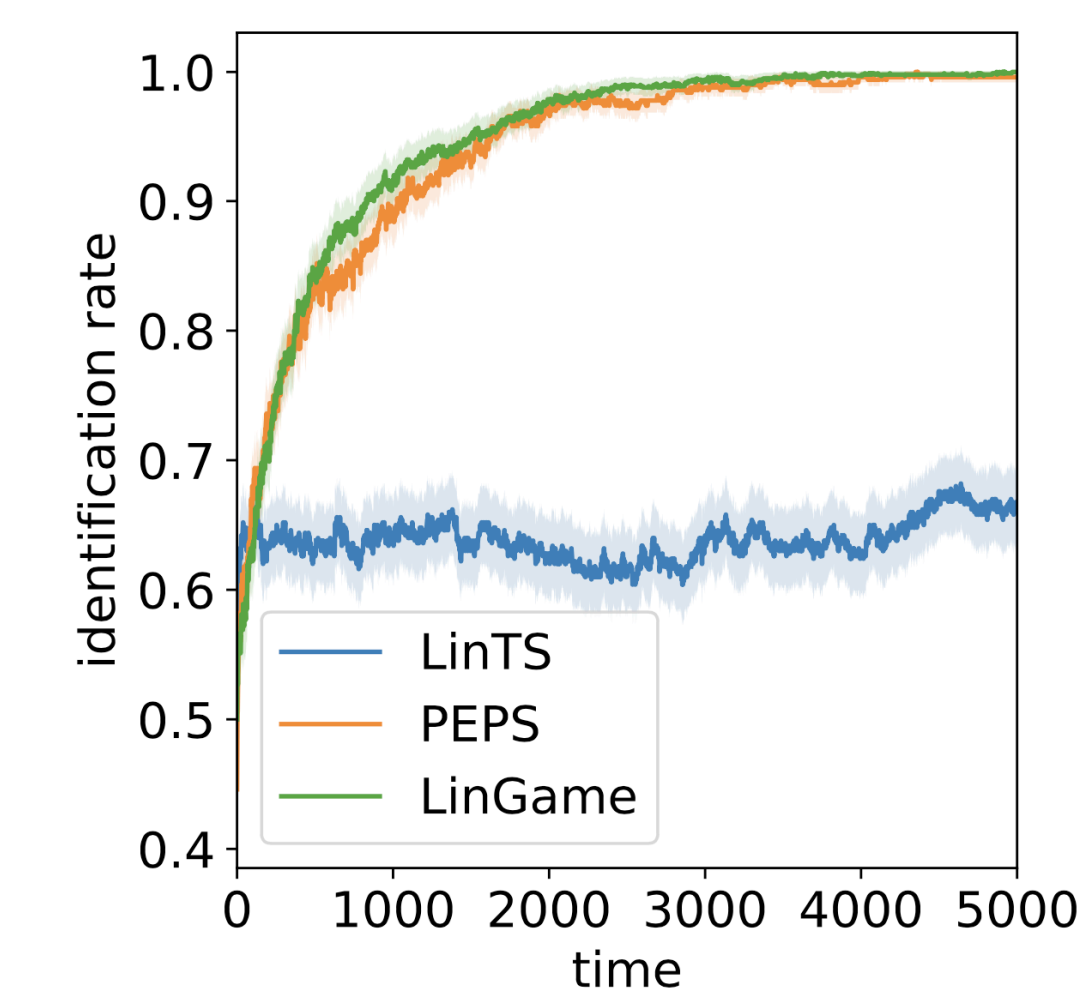


Figure: Hard Instance

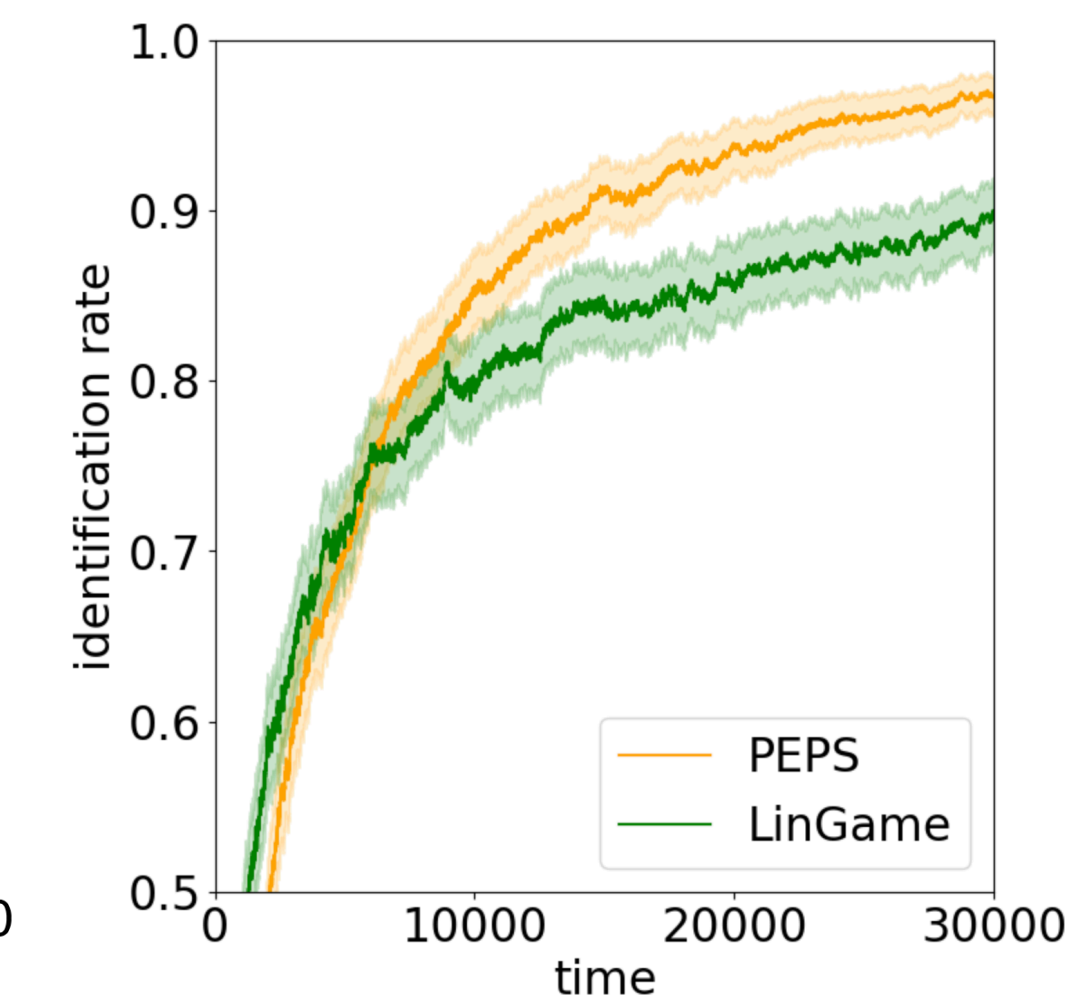


Figure: TopK Instance