# Minimum Empirical Divergence for Sub-Gaussian Linear Bandits

Kapilan Balagopalan, Kwang-Sung Jun

kapilanbgp@arizona.edu, kjun@cs.arizona.edu

THE UNIVERSITY OF ARIZONA

## Preliminary

**Protocol:** For $t \le n$:
- Choose an arm $A_t \in \mathcal{A}_t \subset \mathbb{R}^d$ and receive the reward $Y_t$.

**Model:**
- Reward $Y_t = \langle \theta^*, A_t \rangle + \eta_t$, where $\theta^* \in \mathbb{R}^d$ is an unknown.
- Noise: $\eta_t$ is $\sigma_*^2$-sub-Gaussian.

**Goal:** Minimize cumulative regret,
- $\text{Reg}_n := \sum_{t=1}^{n} \langle a_t^*, \theta^* \rangle - \langle A_t, \theta^* \rangle$ where $a_t^* := \max_{a \in \mathcal{A}_t} \langle a, \theta^* \rangle$.

**Assumption 1:** For all $t \ge 1$, every arm $a \in \mathcal{A}_t$ satisfies $\|a\|_2 \le 1$ and for some constant $B$, $\Delta_{a,t} := \langle \theta^*, a_t^* \rangle - \langle \theta^*, a \rangle \le B$. Furthermore $\|\theta^*\|_2 = S_*$.

## Main results

**Theorem 1** (Instance-dependent bound). *Under Assumption 1, with $\delta_t = \frac{1}{t+1}$, Lin-MED satisfies, $\forall n \ge 1$,*

$$\mathbb{E}\,\text{Reg}_n = \left( \frac{1}{\Delta} d \log(n) \Big( \big( \sigma^2 d \log(n) + \lambda S^2 \big) \log(\log n) + \big( \sigma_*^2 d \log(n) + \lambda S_*^2 \big) H_{\max} \Big) \right).$$

Our algorithm achieves an instance dependent bound of $\hat{O}\!\left( \frac{1}{\Delta} d^2 (\log^2 n) \right)$. (Symbol $\hat{O}$ ignores $\log(\log(n))$ factor)

**Theorem 2** (Minimax Bound). *Under Assumption 1, with $\delta_t = \frac{1}{t+1}$, LinMED satisfies, $\forall n \ge 1$,*

$$\mathbb{E}\,\text{Reg}_n = \left( \sqrt{n} \left( \log^{\frac{1}{2}}(n) \Big( d\sigma \log(n) + \frac{\lambda S^2}{\sigma} \Big) + \frac{H_{\max}}{\sigma \log^{\frac{3}{2}}(n)} \Big( d\sigma_*^2 \log(n) + \lambda S_*^2 \Big) \right) \right).$$

$S$ and $\sigma^2$ are the guesses for $S_*$ and $\sigma_*^2$ respectively. Our algorithm achieves a near-optimal minimax bound ($\tilde{O}(d\sqrt{n})$) and a state-of-the art instance dependent bound ($\frac{1}{\Delta} d^2 \log^2 n$), even when $S_*$ and $\sigma_*^2$ are misspecified. (Many state-of-the-art algorithms including OFUL lacks an analysis when they are underspecified)

## Comparison

| Algorithms | Minimax regret | Instance dependent regret | Closed form probability | Probability assigned for all arms |
|---|---|---|---|---|
| OFUL(Abbasi-Yadkori et al., 2011) | $\tilde{O}(d\sqrt{n})$ | $O(\frac{d^2}{\Delta}\log^2 n)$ | N/A | ✗ |
| LinIMED( Bian and Tan Y.F., 2024) | $\tilde{O}(d\sqrt{n})$ | Unknown | N/A | ✗ |
| LinTS(Agrawal and Goyal, 2014) | $\tilde{O}(d^{\frac{3}{2}}\sqrt{n})$ | Unknown | ✗ | ✗ |
| RandUCB(Vaswanit et al., 2020) | $\tilde{O}(d\sqrt{n})$ | Unknown | ✗ | ✗ |
| SquareCB(Foster and Rakhlin, 2020) | $\tilde{O}(\sqrt{Kdn})$ | Unknown | ✓ | ✗ |
| E2D(Foster et al., 2023) | $\tilde{O}(d\sqrt{n})$ | Unknown | ✓ | ✗* |
| SpannerIGW(Zhu et al., 2022) | $\tilde{O}(d\sqrt{n})$ | $\Omega(\Delta\sqrt{n})$ | ✓ | ✗* |
| EXP2(Bubeck and Cesa-Bianchi, 2012) | $O(\sqrt{dn\log K})$ | $\Omega(\Delta\sqrt{n})$ | ✓ | ✓ |
| LinMED(ours) | $\tilde{O}(d\sqrt{n})$ | $\hat{O}(\frac{d^2}{\Delta}\log^2 n)$ | ✓ | ✓ |

Symbol '✗*' means that the algorithm can be modified to assign probability to all arms. Symbol $\tilde{O}$ ignores the $\log(n)$ factor. Symbol $\hat{O}$ ignores the $\log(\log(n))$ factor.
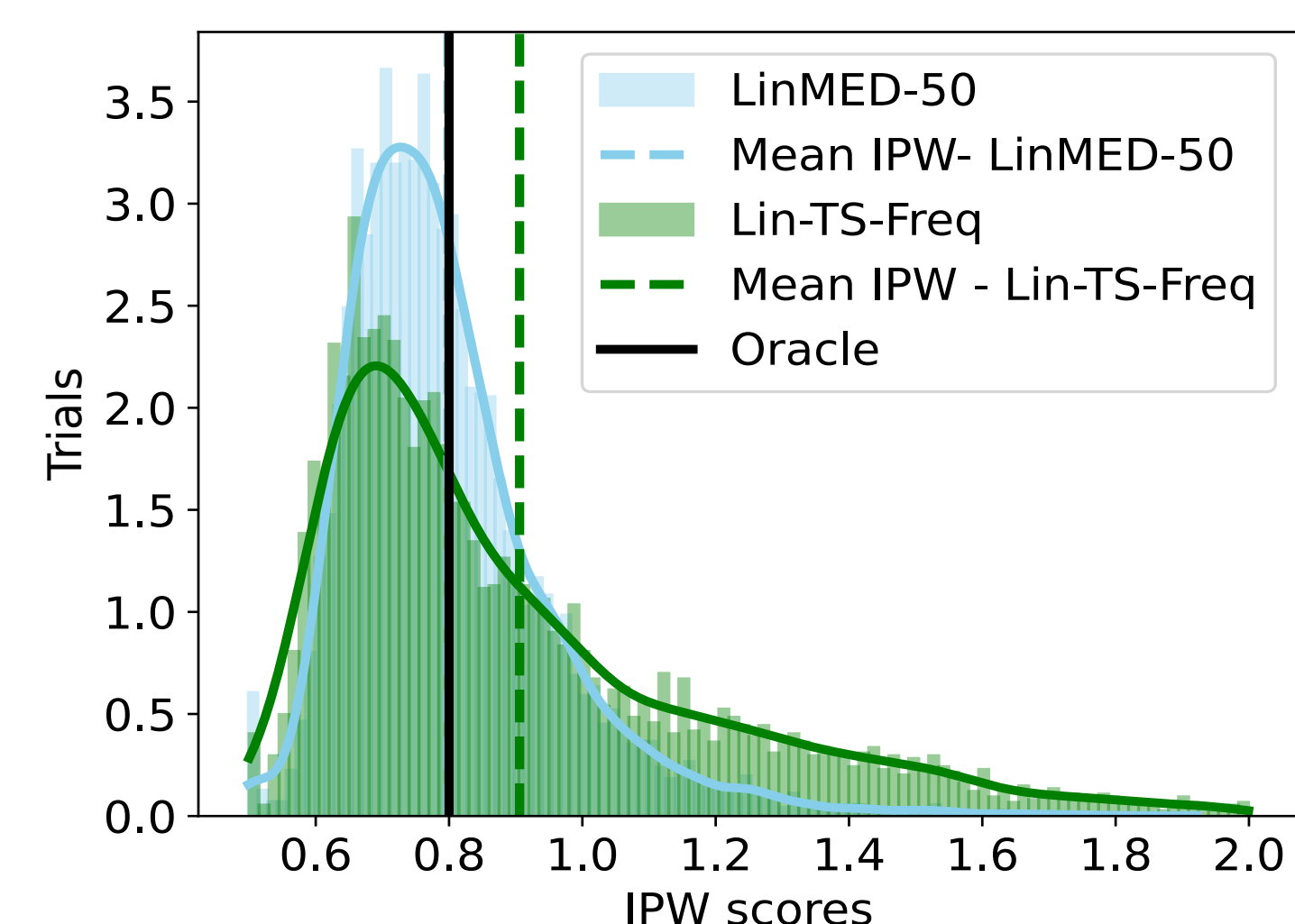
## Supplementary result

**Theorem 3.** *There exists a linear bandit problem for which the EXP2 and Spanner-IGW algorithms satisfy*

$$\mathbb{E}\,\text{Reg}_n \ge \Omega(\Delta\sqrt{n}).$$

EXP2 (Bubeck and Cesa-Bianchi, 2012) and SpannerIGW (Zhu et al., 2022) have polynomial instance dependent lower bound. LinMED achieves polylog instance dependent upper bound.

## *OPE-friendliness*



$$\text{IPW score} = \frac{1}{n} \sum_{t=1}^{n} \frac{\pi_t^{\text{target}}(A_t)}{p_t(A_t)} \cdot Y_t.$$

$$\left( \pi_t^{\text{target}}(A_t) = \frac{1}{|\mathcal{A}|} \right)$$
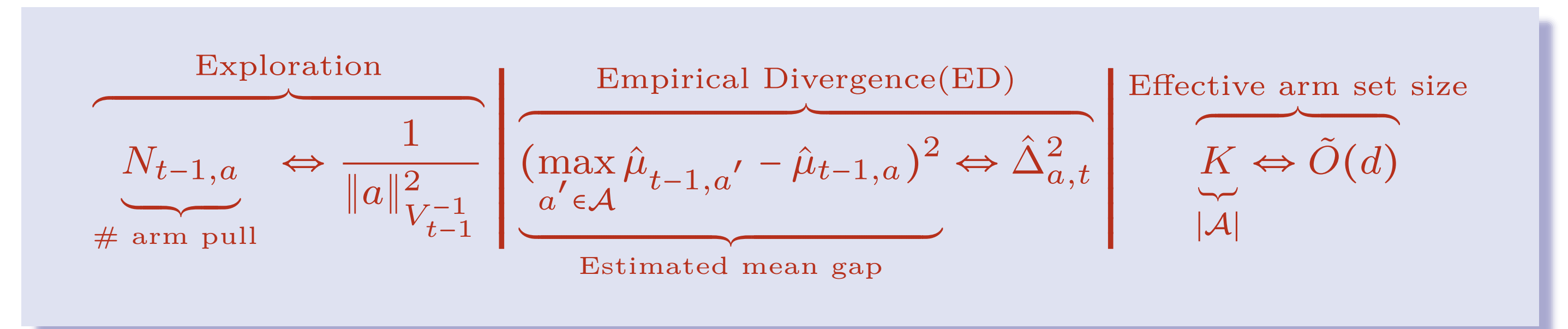
IPW scores (expected regret equivalent) of the uniform policy when the logging policy is LinMED and LinTS respectively. We used $10^3$ Monte Carlo samples to estimate the sampling probabilities of LinTS. Oracle denotes the expected reward of the uniform policy. LinTS shows a nontrivial amount of bias while LinMED is exactly aligned with the oracle.

$\mathcal{A} = \{a_1 = (1,0)^\top, a_2 = (0.6,0.8)^\top\}$, $\theta^* = (1,0)^\top$, LinTS's mean = 0.906, Oracle's mean ≈ LinMED's mean = 0.800

## Highlights and contributions

- LinMED: Linear extension for MED algorithm. (Bian and Jun, 2022, Honda and Takemura, 2011)
- Near-optimal minimax bound and logarithmic instance dependent bound even with noise misspecification.
- Polynomial instance dependent lower bounds for SpannerIGW (Zhu et al., 2022) and EXP2 (Bubeck and Cesa-Bianchi, 2012) ($\Omega(\Delta\sqrt{n})$), which are strictly worse than LinMED.
- Offline policy evaluation friendly algorithm: Our algorithm assigns a closed from probability for each arm, hence it can be used as a logging policy for offline performance evaluation of other policies.

## MED vs LinMED

$$p_{t,a} \propto \exp\left( -\frac{N_{t-1,a}}{2} \cdot \big( \max_{a' \in \mathcal{A}} \hat{\mu}_{t-1,a'} - \hat{\mu}_{t-1,a} \big)^2 \right). \qquad \text{(MED)}$$



## LinMED algorithm

---
**Algorithm 1** LinMED

---

**Input:** Regularization $\lambda$, failure rates $\{\delta_t\}_{t=0}^{\infty}$, optimal design fraction $\alpha_{\text{opt}}$, empirical best fraction $\alpha_{\text{emp}}$, $S$ (guess for $\|\theta^*\|_2$), and $\sigma^2$ (guess for $\sigma_*^2$)
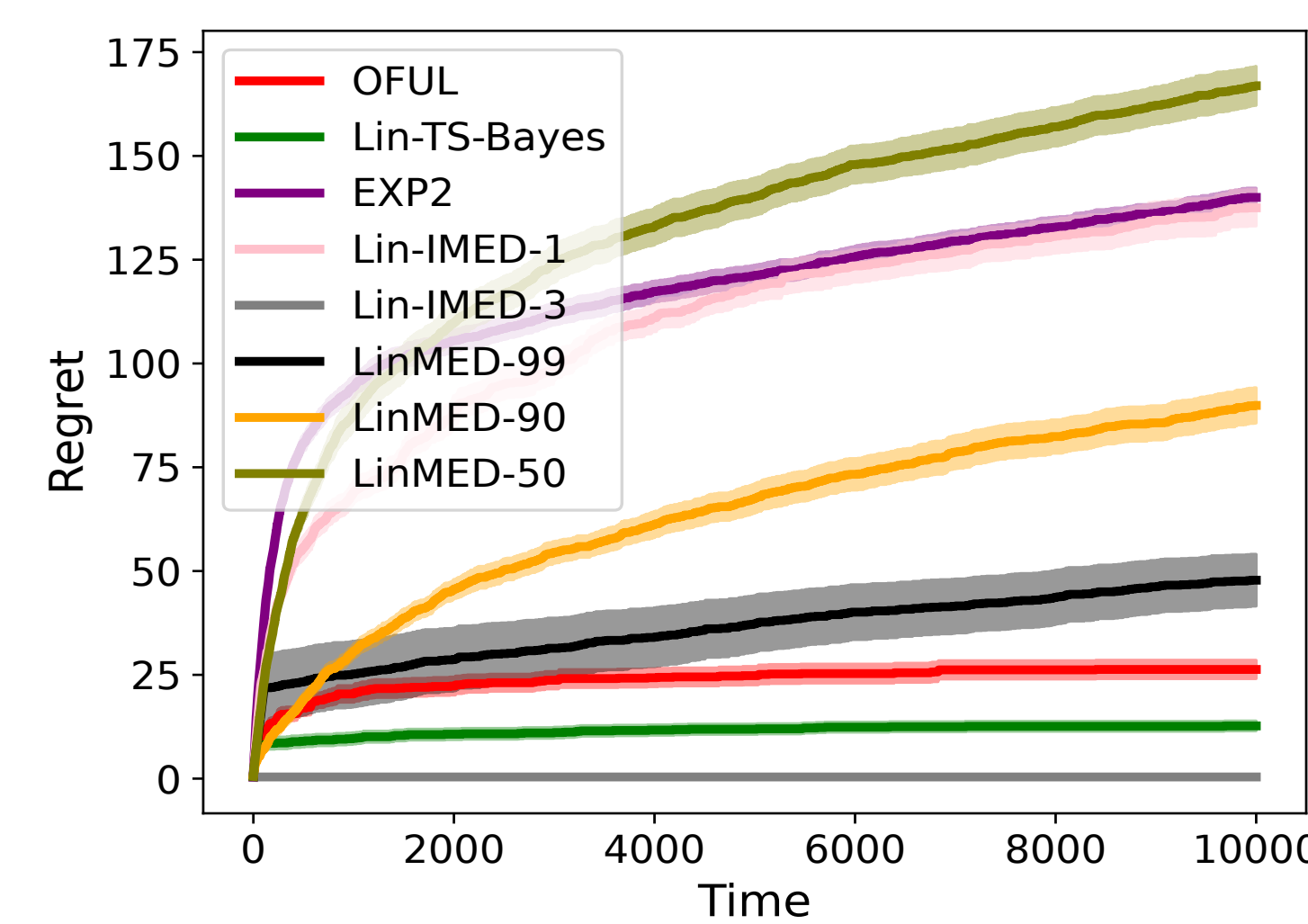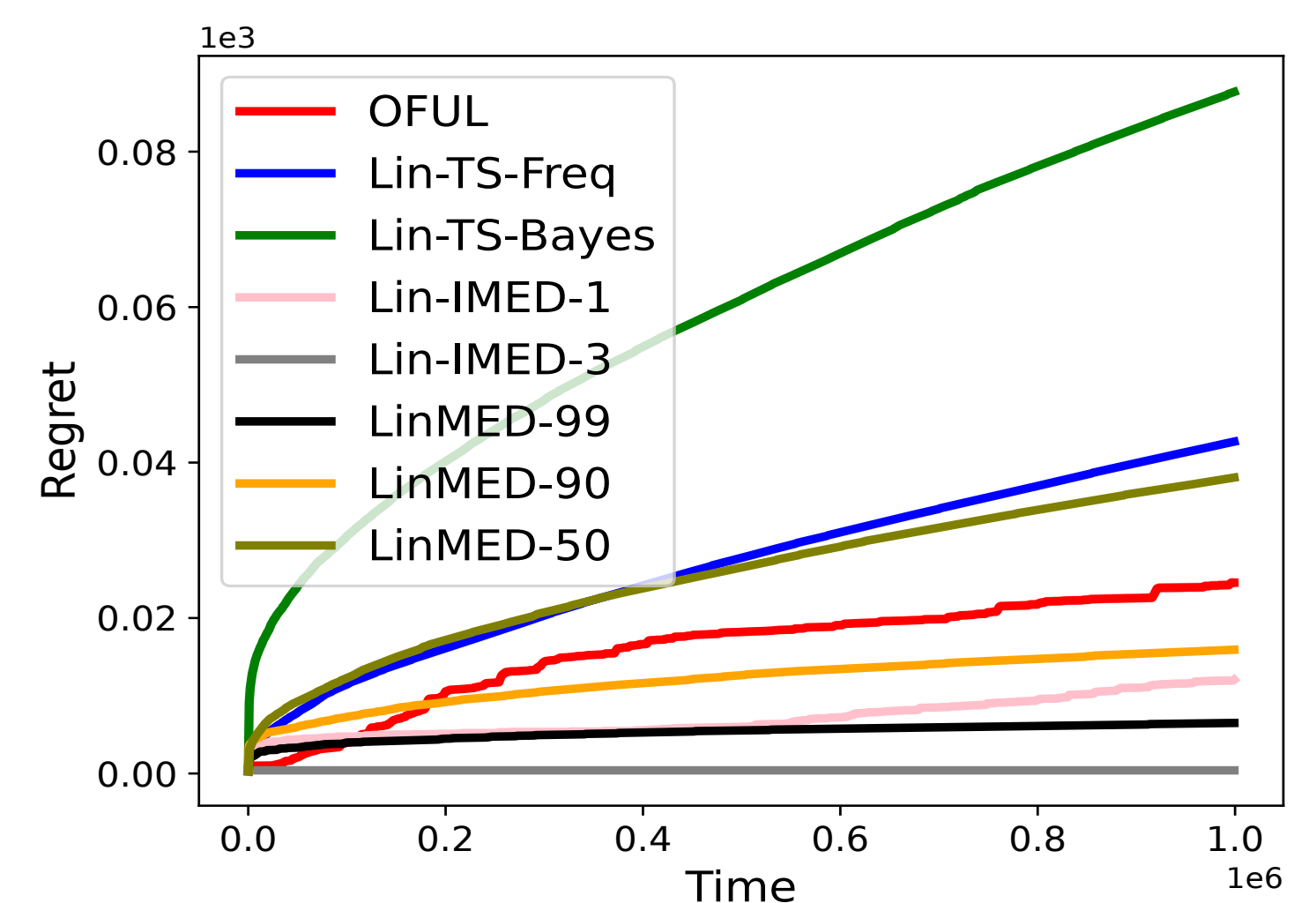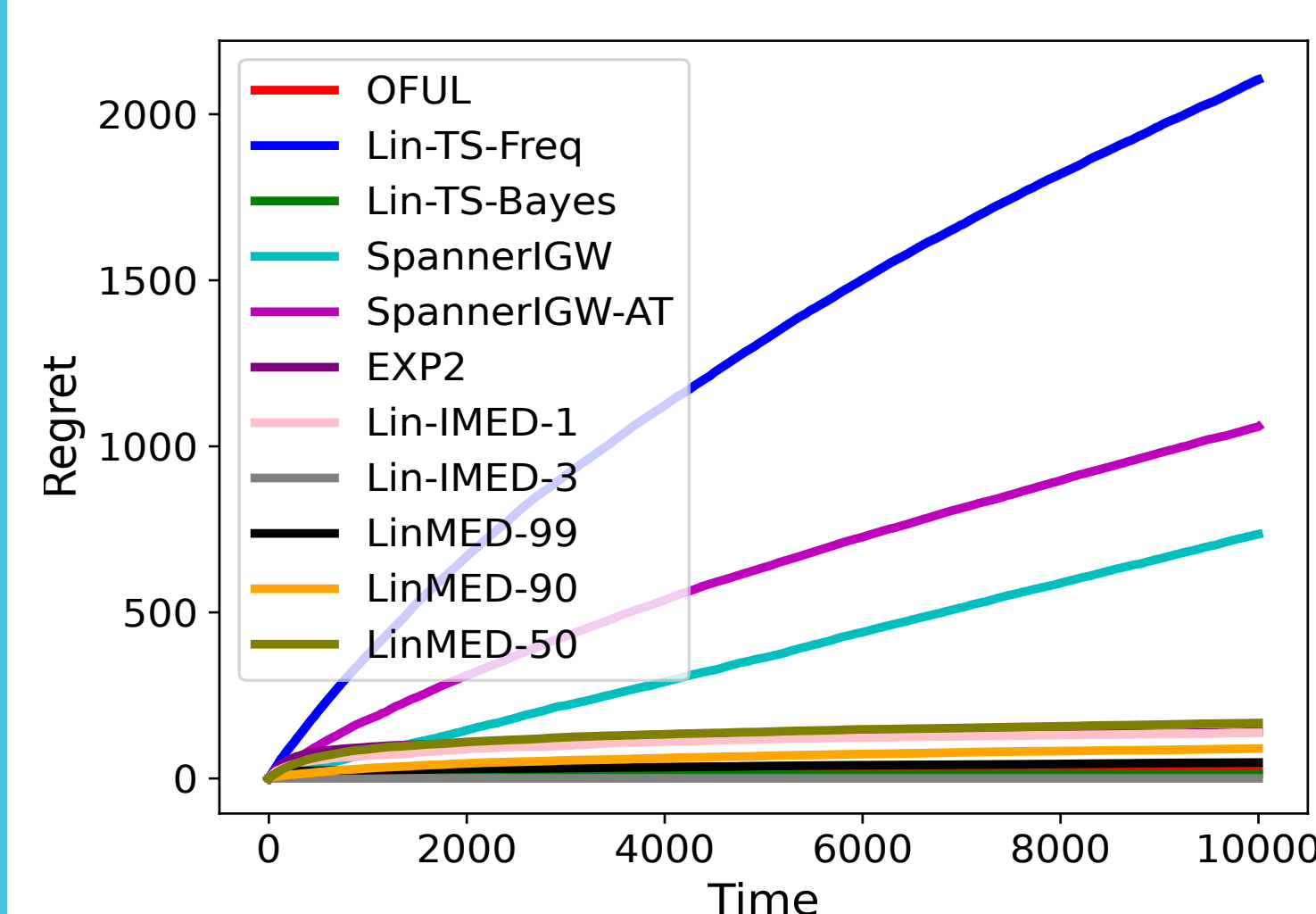
1: Initialize $\hat{\theta}_0 = 0$, $V_0 = \lambda I$.
2: **for** $t = 1, 2, \dots$ **do**
3:     Observe arm set $\mathcal{A}_t$.
4:     Estimate $\hat{a}_t = \max_{a' \in \mathcal{A}_t} \langle \hat{\theta}_{t-1}, a' \rangle$.
5:     Estimate $\hat{\Delta}_{a,t} := \langle \hat{\theta}_{t-1}, \hat{a}_t - a \rangle \quad \forall a \in \mathcal{A}_t$.
6:     Define $\forall a \in \mathcal{A}_t$

$$f_t(a) = \exp\left( -\frac{\hat{\Delta}_{a,t}^2}{\beta_{t-1}(\delta_{t-1}) \|\hat{a}_t - a\|_{V_{t-1}^{-1}}^2} \right)$$

$$\text{where we take } \frac{0}{0} = 0 \text{ and } \beta_t(\delta_t) := \left( \sigma \sqrt{\log\!\left( \frac{\det V_t}{\det V_0} \right) + 2\log\frac{1}{\delta_t}} + \sqrt{\lambda}S \right)^2$$

7:     Re-scale the arms: $\overline{\mathcal{A}}_{(t)} = \{ \sqrt{f_t(a)} \cdot a \mid a \in \mathcal{A}_t \}$.
8:     Compute $q_t^{\text{opt}} = \text{ApproxDesign}(\bar{\mathcal{A}}_t)$ such that $\|b\|_{V^{-1}(q_t^{\text{opt}})}^2 \le \tilde{O}(d), \forall b \in \bar{\mathcal{A}}_t$.
9:     Let $\forall a \in \mathcal{A}_t \quad q_t(a) = \alpha_{\text{opt}} \cdot q_t^{\text{opt}}(a) + \alpha_{\text{emp}} \cdot \mathbf{1}\{a = \hat{a}_t\} + (1 - \alpha_{\text{opt}} - \alpha_{\text{emp}}) \cdot \frac{1}{|\mathcal{A}_t|}$.
10:    Compute $p_t'(a) = \frac{q_t(a)f_t(a)}{\sum_{b \in \mathcal{A}_t} q_t(b)f_t(b)}$.
11:    Define $\mathcal{B}_t = \{a \in \mathcal{A}_t : \|a\|_{V_{t-1}^{-1}}^2 > 1\}$.
12:    **if** $|\mathcal{B}_t| > 0$ **then**
13:      $\forall a \in \mathcal{A}_t, \quad p_t(a) = \frac{1}{2}p_t'(a) + \frac{1}{2}\mathbf{1}\{a = B_t\}$ where $B_t \in \mathcal{B}_t$.
14:    **else**
15:      $\forall a \in \mathcal{A}_t \quad p_t(a) = p_t'(a)$.
16:    **end if**
17:    Take action $A_t \sim p_t$.
18:    Observe the reward $Y_t$ and update $V_t = V_{t-1} + A_t A_t^\top$ and $\hat{\theta}_t = V_t^{-1} \sum_{s=1}^{t} A_s Y_s$.
19: **end for**

## Experiments



**Large gap instance**
Model 1 $\leftarrow \theta^* = (1,0)$, $\eta_t \sim N(0, \sigma^2 = 1)$
$\mathcal{A} = \{(1,0), (0,1)\}$

**End of optimism instance**
Model 1, $\epsilon \in \{0.005, 0.01\}$
$\mathcal{A} = \{a_1 = (1,0), a_2 = (0,1), a_3 = (1-\epsilon, 2\epsilon)\}$

LinMED shows logarithmic growth for Large gap instance. Optimistic algorithms like OFUL and Thompson sampling fail under End of optimism experiments (Lattimore and Szepesvári, 2017).