Improved Regret Bounds for Linear Bandits with Heavy-Tailed Rewards

Artin Tajdini University of Washington artin@cs.washington.edu Jonathan Scarlett
National University of Singapore
scarlett@comp.nus.edu.sg

Kevin Jamieson

University of Washington jamieson@cs.washington.edu

Abstract

We study stochastic linear bandits with heavy-tailed rewards, where the rewards have a finite $(1+\epsilon)$ -absolute central moment bounded by v for some $\epsilon \in (0,1]$. We improve both upper and lower bounds on the minimax regret compared to prior work. When $v = \mathcal{O}(1)$, the best prior known regret upper bound is $\tilde{\mathcal{O}}(dT^{\frac{1}{1+\epsilon}})$. While a lower with the same scaling has been given, it relies on a construction using $v = \mathcal{O}(d)$, and adapting the construction to the bounded-moment regime with $v = \mathcal{O}(1)$ yields only a $\Omega(d^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$ lower bound. This matches the known rate for multi-armed bandits and is generally loose for linear bandits, in particular being \sqrt{d} below the optimal rate in the finite-variance case ($\epsilon = 1$). We propose a new elimination-based algorithm guided by experimental design, which achieves regret $\tilde{\mathcal{O}}(d^{\frac{1+3\epsilon}{2(1+\epsilon)}}T^{\frac{1}{1+\epsilon}})$, thus improving the dependence on d for all $\epsilon \in (0,1)$ and recovering a known optimal result for $\epsilon=1$. We also establish a lower bound of $\Omega(d^{\frac{2\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$, which strictly improves upon the multi-armed bandit rate and highlights the hardness of heavy-tailed linear bandit problems. For finite action sets of size n, we derive upper and lower bounds of $\tilde{\mathcal{O}}(\sqrt{d}(\log n)^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$ and $\tilde{\Omega}(d^{\frac{\epsilon}{1+\epsilon}}(\log n)^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$, respectively. Finally, we provide action-set-dependent regret upper bounds, showing that (i) for some geometries, such as l_p -norm balls for $p \leq 1 + \epsilon$, we can further reduce the dependence on d, and (ii) for RKHS functions with the Matérn kernel we can attain sublinear regret for all $\epsilon \in (0, 1]$, thus substantially improving over the existing state-of-the-art.

1 Introduction

The stochastic linear bandit problem is a foundational setting of sequential decision-making under uncertainty, where the expected reward of each action is modeled as a linear function of known features. While most existing work assumes sub-Gaussian reward noise—enabling the use of concentration inequalities like Chernoff bounds—real-world noise often exhibits heavy tails, potentially with unbounded variance, violating these assumptions. Heavy-tailed noise naturally arises in diverse domains such as high-volatility asset returns in finance [CB00, Con01], conversion values in online advertising [CvdLG20, JSP21], cortical neural oscillations [RVHH15], and packet delays in communication networks [BTA02]. In such settings, reward distributions may be well-approximated by distributions such as Pareto, Student's t, or Weibull, all of which exhibit only polynomial tail decay.

The statistical literature has developed several robust estimation techniques for random variables with only bounded $(1 + \epsilon)$ -moments (for some $\epsilon \in (0, 1]$), such as median-of-means estimators [DLLO16, LM19b] and Catoni M-estimators [Cat12, BJL15] in the univariate case, as well as robust least squares [AC11, HS14, HW19] and adaptive Huber regression [SZF20] for multivariate settings.

Robustness to heavy tails was first introduced into sequential decision-making by [BCBL13] in the context of multi-armed bandits. Subsequent work including [MY16, SYKL18, XWWZ20] extended these ideas to linear bandits, where each action is represented by a feature vector and the reward includes heavy-tailed noise. Generalizing robust estimators from the univariate to the multivariate setting is nontrivial, and many works have focused on designing such estimators and integrating them into familiar algorithmic frameworks like UCB. However, the relative unfamiliarity of heavy-tailed noise can make it difficult to judge the tightness of the regret bounds. As we discuss later, this has led to some degree of misinterpretation of existing lower bounds, with key problems prematurely considered "solved" despite persistent, unrecognized gaps.

1.1 Problem Statement

We consider the problem of stochastic linear bandits with an action set $\mathcal{A} \subseteq \mathbb{R}^d$ and an unknown parameter $\theta^\star \in \mathbb{R}^d$. At each round $t=1,2,\ldots,T$, the learner chooses an action $x_t \in \mathcal{A}$ and observes the reward

$$y_t = \langle x_t, \theta^* \rangle + \eta_t,$$

where η_t are independent noise terms that satisfy $\mathbb{E}[\eta_t] = 0$ and $\mathbb{E}\big[|\eta_t|^{1+\epsilon}\big] \leq \upsilon$ for some $\epsilon \in (0,1]$ and finite $\upsilon > 0$. We adopt the standard assumption that the expected rewards and parameters are bounded, namely, $\sup_{x \in \mathcal{A}} |\langle x, \theta^\star \rangle| \leq 1$ and $\|\theta^\star\|_2 \leq 1$. Letting $x^\star \in \arg\max_{x \in \mathcal{A}} \langle x, \theta^\star \rangle$ be an optimal action, the cumulative expected regret after T rounds is

$$R_T = \sum_{t=1}^{T} (\langle x^{\star}, \theta^{\star} \rangle - \langle x_t, \theta^{\star} \rangle).$$

Given (A, ϵ, v) , the objective is to design a policy for sequentially selecting the points (i.e., x_t for t = 1, ..., T) in order to minimize R_T .

1.2 Contributions

We study the minimax regret of stochastic linear bandits under heavy-tailed noise and make several contributions that clarify and advance the current state of the art. Although valid lower bounds exist, we show that they have been misinterpreted as matching known upper bounds. After correcting this misconception, we provide improved upper and lower bounds in the following ways:

- Novel estimator and analysis: We introduce a new estimator inspired by [CJKS21] (who studied the finite-variance setting, $\epsilon=1$), adapted to the heavy-tailed setting ($\epsilon\in(0,1]$). Its analysis leads to an experimental design problem that accounts for the geometry induced by the heavy-tailed noise, which is potentially of independent interest beyond linear bandits.
- **Improved upper bounds:** We use this estimator within a phased elimination algorithm to obtain state-of-the-art regret bounds for both finite- and infinite-arm settings. Additionally, we derive a geometry-dependent regret bound that emerges naturally from the estimator's experimental design.
- Improved lower bounds: We establish novel minimax lower bounds under heavy-tailed noise that are the first to reveal a dimension-dependent gap between multi-armed and linear bandit settings (e.g., when the arms lie on the unit sphere). We provide such results for both the finite-arm and infinite-arm settings.

Table 1 summarizes our quantitative improvements over prior work, while Figure 1 illustrates the degree of improvement obtained and what gaps still remain.

In addition to these results for heavy-tailed linear bandits, we show that our algorithm permits the kernel trick, and that this leads to regret bounds for the Matérn kernel (with heavy-tailed noise) that significantly improve on the best existing bounds, in particular being sublinear for all $\epsilon \in (0,1]$. See Section 3.1 for a summary, and Appendix C for the details.

Table 1: Comparison of regret bounds (in the $\widetilde{O}(\cdot)$ or $\widetilde{\Omega}(\cdot)$ sense) with heavy-tailed rewards for the model $y_t = \langle x_t, \theta_* \rangle + \eta_t$ where $\mathbb{E}[\eta_t] = 0$, $\mathbb{E}[|\eta_t|^{1+\epsilon}] \leq 1$, $||\theta||_2 \leq 1$, $|\langle x, \theta \rangle| \leq 1$. The complexity measure $M(\mathcal{A})$ is defined in Theorem 3.

Paper	Setting	Regret Upper Bound	Regret Lower Bound
[SYKL18]	general	$dT^{rac{1}{1+\epsilon}}$	$d^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}}$ 1
[HZWY23]	$\mathbb{E}[\eta_t ^{1+\epsilon}] \le \upsilon_t^{1+\epsilon}$	V —	$a_{1+\epsilon} I_{1+\epsilon}$
[XWWZ20]	$ \mathcal{A} = n$	$\sqrt{d\log n}T^{\frac{1}{1+\epsilon}}$	$d^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}}$
[CG19]	$Mat\acute{ern}(\nu,d)$	$T^{\frac{2+\epsilon}{2(1+\epsilon)} + \frac{d}{2\nu + d}}$	$T^{\frac{\nu+d\epsilon}{\nu(1+\epsilon)+d\epsilon}}$
[BCBL13]	$MAB(\mathcal{A} = \Delta^d)$	$d^{rac{\epsilon}{1+\epsilon}}T^{rac{1}{1+\epsilon}}$	$d^{\frac{\epsilon}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}$
	\mathcal{A} -dependent	$M(\mathcal{A})^{\frac{1}{1+\epsilon}} \min(d, \log \mathcal{A})^{\frac{\epsilon}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}$ (Theorem 3)	
	general	$d^{rac{1+3\epsilon}{2(1+\epsilon)}}T^{rac{1}{1+\epsilon}}$ (Corollary 1)	$d^{rac{2\epsilon}{1+\epsilon}}T^{rac{1}{1+\epsilon}}$ (Theorem 1)
Our Work	$ \mathcal{A} = n$	$\sqrt{d}(\log n)^{rac{\epsilon}{1+\epsilon}}T^{rac{1}{1+\epsilon}}$ (Corollary 1)	$d^{\frac{\epsilon}{1+\epsilon}}(\log n)^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}}$ (Theorem 2)
	$Mat\'ern(\nu,d)$	$T^{1-rac{\epsilon}{1+\epsilon}rac{2 u}{2 u+d}}$ (Corollary 4)	
	$\mathrm{MAB}(\mathcal{A}=\Delta^d)$	$d^{rac{\epsilon}{1+\epsilon}}T^{rac{1}{1+\epsilon}}$ (Corollary 2)	

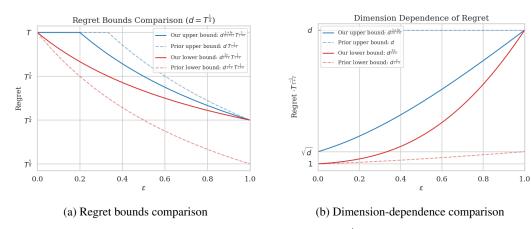


Figure 1: (a) Comparison of regret bounds across ϵ for $T=d^4$. (b) Scaling of the bounds in d.

1.3 Related Work

The first systematic study of heavy-tailed noise in bandits is due to [BCBL13], who replaced the empirical mean in UCB with robust mean estimators, and obtained a regret bound of $\widetilde{O}(n^{\frac{\epsilon}{1+\epsilon}}T^{1/(1+\epsilon)})$ with n arms, along with a matching lower bound. A sequence of follow-up works [YNKL18, LWHZ19, LYLO20, WS21, HDH22, CHDH25] refined these ideas and extended them to best-arm identification, adversarial, parameter-free, and Lipschitz settings. The first extension of heavy-tailed analysis from MAB to linear bandits is due to [MY16], who proposed truncation- and MoM-based algorithms and proved an $\widetilde{O}(dT^{\frac{2+\epsilon}{2(1+\epsilon)}})$ regret bound. Subsequently, [SYKL18, XWWZ20] improved the regret bounds for infinite and finite action sets, respectively (see Table 1). Huber-loss based estimators have emerged as another robustification strategy, for which [LS24, KK23, HZWY23, WZZZ25] provided moment-aware regret bounds. [ZHYW21] suggested median based estimators for symmetric error distributions without any bounded moments (e.g., Cauchy). Beyond linear bandits, [XWW+23] proved a similar $dT^{\frac{1}{1+\epsilon}}$ bound for generalized linear bandits, and [CG19] studied heavy-tailed kernel-based bandits, which we will cover in more detail in Appendix C. A summary of the best regret bounds of previous work and ours can be found in Table 1.

¹We refer to this as the multi-armed bandit (MAB) rate because it matches that of a MAB problem with d arms. Note that that the $dT^{\frac{1}{1+\epsilon}}$ lower bound from [SYKL18] was only proved for an instance with $\mathbb{E}[|\eta_t|^{1+\epsilon}] = O(d)$ rather than O(1); see Section 2 for further discussion.

2 Lower Bounds

Before describing our own lower bounds, we take a moment to clarify the state of lower bounds that exist in the literature, as there has been some apparent misinterpretation within the community. The regret lower bound construction presented in [SYKL18] leverages the reward distribution

$$y(x) = \begin{cases} (\frac{1}{\Delta})^{\frac{1}{\epsilon}} & \text{w.p. } \Delta^{\frac{1}{\epsilon}}\theta^{\top}x \\ 0 & \text{w.p. } 1 - \Delta^{\frac{1}{\epsilon}}\theta^{\top}x \end{cases}$$

under the choice $\Delta = \frac{1}{12} T^{-\frac{\epsilon}{1+\epsilon}}$, and with choices of θ and $\mathcal A$ that ensure $d\Delta \leq \theta^\top x \leq 2d\Delta$. A straightforward calculation shows that the reward distributions of this construction possesses a $(1+\epsilon)$ -absolute moment of $\Delta^{-1}(\theta^\top x) \geq d$ for all actions. Recall that in our problem statement we consider the $(1+\epsilon)$ -absolute moment to be a constant (that does not depend on the the dimension d or time horizon T). We can compare this with the canonical case of sub-Gaussian noise $(\epsilon=1)$ where it is assumed that the second moment is bounded by $\sigma^2 = \Omega(1)$, in which it is well-known that the optimal regret rate is on the order of $\sigma d\sqrt{T}$ [LS20]. If we were to set $\sigma^2 = \Theta(d)$, this would suggest a rate of $d^{3/2}\sqrt{T}$, but this only exceeds the usual $d\sqrt{T}$ because σ is artificially large. We stress that we are not claiming that the lower bound of [SYKL18] is in any way incorrect, and the authors even acknowledge that the bound on the moment scales with the dimension in the appendix of their work. We are simply pointing out that there has been some misinterpretation of the lower bound within the community.²

If we adjust the expected reward distributions such that $\Delta \leq \theta^{\top}x \leq 2\Delta$, so that the reward distribution maintains a constant $1+\epsilon$ absolute moment, the resulting regret lower bound turns out to scale as $d^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}}$, matching the known optimal lower bound for the Multi-Armed Bandit (MAB) setting with d arms. However, with a more precise analysis, we can prove a stronger lower bound on a similar instance (with modified parameters) having a constant $(1+\epsilon)$ -central moment of rewards, as we will see below.

2.1 Infinite Arm Set

Given the context above, we are ready to present our own lower bound that builds on the construction introduced by [SYKL18] but is specifically tailored to improving the d dependence.

Theorem 1. Fix the action set $\mathcal{A}=\{x\in[0,1]^{2d}:x_{2i-1}+x_{2i}=1\ \forall i\in[d]\}$. There exists a reward distribution with a $(1+\epsilon)$ -central moment bounded by 1 and a $\theta^*\in\mathbb{R}^{2d}$ with $\|\theta^*\|_2\leq 1$ and $\sup_{x\in\mathcal{A}}|x^\top\theta^*|\leq 1$, such that for $T\geq 4^{\frac{1+\epsilon}{\epsilon}}d^2$, the regret incurred is $\Omega(d^{\frac{2\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$.

Proof. For a parameter $\Delta \leq \frac{1}{4d}$ to be specified later, we let the reward distribution be a Bernoulli random variable defined as follows:

$$y(x) = \begin{cases} (\frac{1}{\gamma})^{\frac{1}{\epsilon}} & \text{w.p. } \gamma^{\frac{1}{\epsilon}}\theta^{\top}x \\ 0 & \text{w.p. } 1 - \gamma^{\frac{1}{\epsilon}}\theta^{\top}x \end{cases}$$

with $\gamma:=2d\Delta$. We consider parameter vectors θ lying in the set $\Theta:=\left\{\theta\in\{\Delta,2\Delta\}^{2d}:\theta_{2i-1}+\theta_{2i}=3\Delta\right\}$, from which the assumption $\Delta\leq\frac{1}{4d}$ readily implies $\|\theta\|_2\leq 1$ and $\sup_{x\in\mathcal{A}}|x^\top\theta^*|\leq 1$. For any $\theta\in\Theta$, the $(1+\epsilon)$ -raw moment of the reward distribution (and therefore the central moment, since the rewards are nonnegative) for each action is bounded by $\mathbb{E}[|y(x)|^{1+\epsilon}|x]=\gamma^{-(1+\epsilon)/\epsilon}\gamma^{1/\epsilon}\theta^\top x=\gamma^{-1}\theta^\top x\leq 1$, since $\gamma=2d\Delta$ and $\theta^\top x\leq 2d\Delta$.

Let $R_T(\mathcal{A}, \theta)$ be the cumulative regret for arm set \mathcal{A} and parameter θ , and let $\operatorname{ind}_i(\theta) := \arg \max_{b \in \{0,1\}} (\theta_{2i-1+b})$ for $\theta \in \Theta$, and write $x_t = (x_{t,1}, \dots, x_{t,d})$. We have

$$R_T(\mathcal{A}, \theta) = \sum_{t=1}^{T} \sum_{i=1}^{d} \left(\Delta - \Delta x_{t, 2i-1 + \operatorname{ind}_i(\theta)} \right) = \Delta \sum_{t=1}^{T} \sum_{i=1}^{d} \left(\frac{1}{2} - \frac{1}{2} (-1)^{\operatorname{ind}_i(\theta)} (x_{t, 2i-1} - x_{t, 2i}) \right)$$

²Previous works that indicate the minimax optimality of this bound (with respect to T and d) include [XWWZ20, XWW+23, HZWY23, WZZZ25].

 $^{^3}$ This is obtained by optimizing Δ for the adjusted regret $\Delta T(\frac{1}{4}-\frac{3}{2}\sqrt{d^{-1}\Delta^{\frac{1+\epsilon}{\epsilon}}T})$

$$\geq \frac{\Delta}{2} \sum_{i=1}^{d} \mathbb{E}_{\theta} \left[\sum_{t=1}^{T} \mathbb{I} \{ (-1)^{\operatorname{ind}_{i}(\theta)} (x_{t,2i-1} - x_{t,2i}) \leq 0 \} \right]$$

$$\geq \frac{\Delta T}{4} \sum_{i=1}^{d} \mathbb{P}_{\theta} \left[\sum_{t=1}^{T} \mathbb{I} \{ (-1)^{\operatorname{ind}_{i}(\theta)} (x_{t,2i-1} - x_{t,2i}) \leq 0 \} \geq \frac{T}{2} \right],$$

where the second equality follows by using $x_{t,2i-1} + x_{t,2i} = 1$ and checking the cases $\operatorname{ind}_i(\theta) = 0$ and $\operatorname{ind}_i(\theta) = 1$ separately.

For any $\theta \in \Theta, i \in [d]$, we define $\theta' \in \Theta$ with entries $\theta'_j = \begin{cases} 3\Delta - \theta_j & 2i - 1 \leq j \leq 2i \\ \theta_j & \text{otherwise} \end{cases}$, and let $p_{\theta,i} := \mathbb{P}_{\theta} \left[\sum_{t=1}^T \mathbb{I}\{(-1)^{\operatorname{ind}_i(\theta)}(x_{t,2i-1} - x_{t,2i}) \leq 0\} \geq \frac{T}{2} \right]$. We then have the following:

$$p_{\theta,i} + p_{\theta',i} \ge \frac{1}{2} \exp(-\text{KL}(\mathbb{P}_{\theta} || \mathbb{P}_{\theta'}))$$

$$= \frac{1}{2} \exp\left(-\mathbb{E}_{\theta} \left[\sum_{t=1}^{T} \text{KL}\left(\text{Ber}(\gamma^{\frac{1}{\epsilon}} \theta^{\top} x_{t}) || \text{Ber}(\gamma^{\frac{1}{\epsilon}} \theta'^{\top} x_{t})\right) \right] \right).$$
 (Chain rule)

Now we set $\Delta:=\frac{1}{2}d^{\frac{\epsilon-1}{1+\epsilon}}T^{-\frac{\epsilon}{1+\epsilon}}$. Note that since $T\geq 4^{\frac{1+\epsilon}{\epsilon}}d^2$, the above-mentioned condition $\Delta\leq \frac{1}{4d}$ holds, ensuring the Bernoulli parameter is in [0,1]. Under this choice of Δ , we have

$$\mathrm{KL}\left(\mathrm{Ber}(\gamma^{\frac{1}{\epsilon}}\theta^{\top}x_{t})\|\mathrm{Ber}(\gamma^{\frac{1}{\epsilon}}\theta'^{\top}x_{t})\right) \leq \frac{2^{\frac{2}{\epsilon}}4\Delta^{\frac{2}{\epsilon}}d^{\frac{2}{\epsilon}}\Delta^{2}}{2^{\frac{1}{\epsilon}}\Delta^{\frac{1+\epsilon}{\epsilon}}d^{\frac{1+\epsilon}{\epsilon}} \cdot \frac{1}{2}} = 2^{\frac{1}{\epsilon}}8\Delta^{\frac{1+\epsilon}{\epsilon}}d^{\frac{1-\epsilon}{\epsilon}} = 4T^{-1},$$

where in the first inequality we used $\mathrm{KL}(\mathrm{Ber}(p)\|\mathrm{Ber}(q)) \leq \frac{(p-q)^2}{q(1-q)};$ we get $|p-q| \leq 2\gamma^{\frac{1}{\epsilon}}\Delta = 2(2d\Delta)^{\frac{1}{\epsilon}}\Delta$ because θ and θ' differ only via a single swap of $(\Delta, 2\Delta)$ by $(2\Delta, \Delta)$, $q \geq \gamma^{\frac{1}{\epsilon}}\Delta d = (2d\Delta)^{\frac{1}{\epsilon}}\Delta d$ by construction, and $1-q \geq 1-\gamma^{\frac{1}{\epsilon}}2d\Delta \geq \frac{1}{2}$ via $\Delta \leq \frac{1}{4d}$.

Combining the preceding display equations gives $p_{\theta,i}+p_{\theta',i}\geq \frac{1}{2}\exp(-4)$, and averaging over all (θ,θ') (with $\theta'\neq \theta$) and summing over i, we obtain $\frac{1}{|\Theta|}\sum_{\theta\in\Theta}\sum_{i=1}^d p_{\theta,i}\geq \frac{1}{4}d\exp(-4)$. Hence, there exists $\theta^*\in\Theta$ such that $\sum_{i=1}^d p_{\theta^*,i}\geq \frac{1}{4}d\exp(-4)$, and substituting into our earlier lower bound on R_T gives $R_T(\mathcal{A},\theta^*)\geq \frac{1}{16}\exp(-4)\Delta dT=\frac{1}{32}\exp(-4)d^{\frac{2\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}}$.

The setting in Theorem 1 is not the only one that gives regret $\Omega(d^{\frac{2\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$. In fact, the same lower bound turns out to hold for the unit ball action set with a slight change in reward distribution to avoid large KL divergences when $\theta^{\top}x$ is small. The details are given in Appendix B.

2.2 Finite Arm Set

The best known lower bound for finite arm sets matches the MAB lower bound of $d^{\frac{e}{1+e}}T^{\frac{1}{1+e}}$ with d arms (see [XWWZ20] and the summary in Table 1). We provide the first n-dependent lower bound (where $n:=|\mathcal{A}|$) by combining ideas from the MAB lower bound construction for m arms with the construction used in Theorem 1 for dimension $\frac{d}{m}$, where $m^{\frac{d}{m}}\approx n$. When $n=2^{\mathcal{O}(d)}$ or $n=T^{\mathcal{O}(d)}$, which arises naturally when finely quantizing in each dimension, our lower bound matches the infinite arm case (in the $\widetilde{\Omega}(\cdot)$ sense) as one might expect.

Theorem 2. For each $n \in [d, 2^{\lfloor \frac{d}{4} \rfloor}]$, there exists an action set \mathcal{A} with $|\mathcal{A}| \leq n$, a reward distribution with a $(1+\epsilon)$ -central moment bounded by 1, and a $\theta^* \in \mathbb{R}^d$ with $\|\theta^*\|_2 \leq 1$ and $\sup_{x \in \mathcal{A}} |x^\top \theta^*| \leq 1$, such that for $T \geq 4^{\frac{1+\epsilon}{\epsilon}} d^{\frac{1+\epsilon}{\epsilon}}$, the regret incurred is $\Omega(T^{\frac{1}{1+\epsilon}} d^{\frac{\epsilon}{1+\epsilon}} (\frac{\log n}{\log d})^{\frac{\epsilon}{1+\epsilon}})$.

Proof Sketch. We outline the main proof steps here, and defer the full details to Appendix E.

Consider $\log(\cdot)$ with base 2, and define m to be the smallest integer such that $\frac{m}{\log m} \geq \frac{d}{\log n}$. From the assumption $n \in [d, 2^{\lfloor \frac{d}{4} \rfloor}]$ we can readily verify that d > 4 and $m \in [4, d]$. For convenience, we assume that d is a multiple of m, since otherwise we can form the construction of the lower bound

with $d' = d - (d \mod m)$ and pad the action vectors with zeros. Letting $d_i := (i-1)m$, we define the action set and the parameter set as follows for some $\Delta \leq \frac{1}{4d}$ to be specified later:

$$\mathcal{A} := \left\{ a \in \{0, 1\}^d : \sum_{j=d_i+1}^{d_{i+1}} a_j = 1, \ \forall i \in [d/m] \right\}$$

$$\theta^* \in \Theta := \left\{ \theta \in \{\Delta, 2\Delta\}^d : \sum_{j=d_i+1}^{d_{i+1}} \theta_j = (m+1)\Delta, \ \forall i \in [d/m] \right\}.$$

In simple terms, the d-dimensional vectors are arranged in d/m groups of size m; each block in $a \in \mathcal{A}$ has a single entry of 1 (with 0 elsewhere), and each block in θ^* has a single entry of 2Δ (with Δ elsewhere). The condition $\Delta \leq \frac{1}{4d}$ readily implies $\|\theta^*\|_2 \leq 1$ and $x^\top \theta^* \leq 1$ as required. Moreover, we have $|\mathcal{A}| = m^{\frac{d}{m}}$, and thus $\log |\mathcal{A}| = \frac{d}{m} \log m \leq \log n$ by the definition of m.

Similar to Theorem 1, we let the rewards distribution be $y(x) = \begin{cases} (\frac{1}{\gamma})^{\frac{1}{\epsilon}} & \text{w.p. } \gamma^{\frac{1}{\epsilon}} \theta^{\top} x \\ 0 & \text{w.p. } 1 - \gamma^{\frac{1}{\epsilon}} \theta^{\top} x \end{cases}$ with

 $\gamma := 2 \frac{d}{m} \Delta$. The choices of \mathcal{A} and Θ give $\theta^\top x \leq 2 \Delta \frac{d}{m}$, so by the same reasoning as in Theorem 1, the $(1+\epsilon)$ -moment of the reward distribution is bounded by 1.

Let $\operatorname{ind}_i(x) := \operatorname{arg} \max_{b \in [m]} (x_{d_i+b})$ for any $x \in \mathcal{A} \cup \Theta$, and define $T_{i,b} := |\{t : x_{t,d_i+b} = 1\}|$. Let $\operatorname{ind}_i(x) := \arg\max_{b \in [m]} (x_{d_i+b})$ for any $x \in \mathcal{A} \cup \Theta$, and define $T_{i,b} := |\{t : x_{t,d_i+b} = 1\}|$. Moreover, define t_U to be a random integer drawn uniformly from [T], which immediately implies that $\mathbb{P}_{\theta}[x_{t_U,d_i+b} = 1] = \frac{\mathbb{E}_{\theta}[T_{i,b}]}{T}$. Then we can rewrite regret as $R_T(\mathcal{A},\theta) = \Delta T \sum_{i=1}^{d/m} (1 - \mathbb{P}_{\theta}[x_{t_U,d_i+\operatorname{ind}_i(\theta)} = 1])$. For any $\theta \in \Theta$ and $i \in [\frac{d}{m}]$, and any $b \in [m]$, we define $\theta^{(b)} \in \Theta$ to have entries given by $\theta_j^{(b)} = \begin{cases} \Delta + \Delta \mathbb{I}\{j = d_i + b\} & j \in [d_i + 1, d_{i+1}] \\ \theta_j & \text{otherwise} \end{cases}$; and define the base parameter $\theta^{(0)}$ with entries $\theta_j^{(0)} = \begin{cases} \Delta & j \in [d_i + 1, d_{i+1}] \\ \theta_j & \text{otherwise} \end{cases}$. Note that $\theta^{(\operatorname{ind}_i(\theta))} = \theta$. Moreover, similar to

Theorem 1, the KL divergence of reward distribution of $\theta^{(0)}$ and $\theta^{(b)}$ problem instances at action x_t can be bounded by $2^{\frac{1+\epsilon}{\epsilon}}\Delta^{\frac{1+\epsilon}{\epsilon}}\left(\frac{d}{m}\right)^{\frac{1-\epsilon}{\epsilon}}\mathbb{I}\{x_{t,d_i+b}=1\}$. We set $\Delta:=\frac{1}{8}\left(\frac{d}{m}\right)^{\frac{\epsilon-1}{1+\epsilon}}\left(\frac{T}{m}\right)^{\frac{-\epsilon}{1+\epsilon}}$, from which the condition $T\geq 4^{\frac{1+\epsilon}{\epsilon}}d^{\frac{1+\epsilon}{\epsilon}}$ readily yields $\gamma^{\frac{1}{\epsilon}}(\frac{2d}{m}\Delta)\leq \frac{1}{2}$.

Next, using Pinsker's inequality along with averaging over $b \in m$, we can show that $\frac{1}{m} \sum_b \mathbb{P}_{\theta^{(b)}}[x_{t,d_i+b}=1] \leq \frac{1}{m}+\frac{1}{2}$. Averaging over all $\theta \in \Theta$, summing over $i \in [d/m]$, and recalling that $m \geq 4$, we obtain

$$\frac{1}{|\Theta|} \sum_{\theta \in \Theta} \sum_{i=1}^{d/m} \left(1 - \mathbb{P}_{\theta}[x_{t,d_i + \mathrm{ind}_i(\theta)} = 1] \right) \ge \frac{d}{m} \left(1 - \frac{1}{m} - \frac{1}{2} \right) \ge \frac{d}{4m}.$$

Hence, there exists $\theta^* \in \Theta$ such that $\sum_{i=1}^{d/m} \left(1 - \mathbb{P}_{\theta^*}[x_{t,d_i + \mathrm{ind}_i(\theta^*)} = 1]\right) \geq \frac{d}{4m}$, substituting into our earlier lower bound on R_T along with our choice of Δ , we obtain

$$R_T(\mathcal{A}, \theta^*) \ge \frac{d}{4m} \Delta T = \frac{1}{32} d^{\frac{\epsilon}{1+\epsilon}} \left(\frac{d}{m}\right)^{\frac{\epsilon}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}.$$

Since $f(x) = \frac{x}{\log x}$ is increasing for $x \ge e$, and $m \in [4,d]$, the definition of m gives the $\frac{d}{\log n} > \frac{m-1}{\log(m-1)} > \frac{m-1}{\log m} \ge \frac{m-1}{\log d}$. Rearranging the terms and bounding m, we obtain $\frac{d}{m} \ge \frac{\log n}{2\log d}$, which

Proposed Algorithm and Upper Bounds

In this section, we propose a phased elimination—style algorithm called MED-PE that achieves the best known minimax regret upper bound for linear bandits with noise that has bounded $(1+\epsilon)$ moments. In each phase ℓ , the algorithm operates as follows:

Algorithm 1 Moment-based Experimental Design Phased Elimination (MED-PE)

$$\begin{split} & \text{Input: } \mathcal{A}, \gamma > 0 \,, \beta \geq 0, \epsilon \in (0,1], \, v, \, T, \, \text{robust mean estimator } \widehat{\mu}(S,\delta) \\ & \text{Initialization } \ell \leftarrow 1, \, t \leftarrow 0, \, \mathcal{A}_1 \leftarrow \mathcal{A} \\ & \text{while } t < T \, \text{ and } |\mathcal{A}_\ell| > 1 \, \text{ do} \\ & M_{1+\epsilon}(\lambda; \mathcal{A}_\ell, \gamma, \beta) \leftarrow \max_{a \in \mathcal{A}_\ell} \mathbb{E}_{x \sim \lambda} \Big[\big| a^\top A^{(\gamma)}(\lambda)^{-1} x \big|^{1+\epsilon} \Big] + \beta^{1+\epsilon} \|a\|_{A^{(\gamma)}(\lambda)^{-1}}^{1+\epsilon} \\ & \qquad \qquad (A^{(\gamma)}(\lambda) := \gamma I + \mathbb{E}_{x \sim \lambda}[xx^\top]) \\ & \lambda_\ell^* \leftarrow \underset{\lambda \in \Delta_{\mathcal{A}_\ell}}{\operatorname{argmin}} \, M_{1+\epsilon}(\lambda; \mathcal{A}_\ell, \gamma, \beta) \\ & \varepsilon_\ell \leftarrow 2^{-\ell}, \tau_\ell \leftarrow 32^{\frac{1+\epsilon}{\epsilon}} (1+v)^{\frac{1}{\epsilon}} \varepsilon_\ell^{-\frac{1+\epsilon}{\epsilon}} M_{1+\epsilon}(\lambda_\ell^*; \mathcal{A}_\ell, \gamma, \beta)^{\frac{1}{\epsilon}} \log(2\ell^2 |\mathcal{A}_\ell| T) \\ & \text{for } s \leftarrow 1 \, \text{to } \tau_\ell \, \text{do} \\ & \big[\text{Draw } x_s \sim \lambda_\ell^*, \, \text{observe reward } y_s \\ & W^{(a)} \leftarrow \widehat{\mu} \left(\big\{ a^\top A^{(\gamma)}(\lambda_\ell^*)^{-1} x_s \, y_s \big\}_{s=1}^{\tau_\ell}, \, \frac{1}{2\ell^2 T |\mathcal{A}_\ell|} \right) \\ & \widehat{\theta}_\ell \leftarrow \underset{\alpha \in \mathcal{A}_\ell}{\operatorname{arg min }} \underset{\alpha \in \mathcal{A}_\ell}{\max} |\theta^\top a - W^{(a)}| \\ & \mathcal{A}_{\ell+1} \leftarrow \big\{ a \in \mathcal{A}_\ell : \widehat{\theta}_\ell^\top a \geq \underset{\alpha' \in \mathcal{A}_\ell}{\max} \, \widehat{\theta}_\ell^\top a' \, - \, 4\varepsilon_\ell \big\}, \, \ell \leftarrow \ell+1, \, t \leftarrow t + \tau_\ell \end{split}$$

- 1. Design a sampling distribution over the currently active arms that minimizes the $(1+\epsilon)$ -absolute moment of a certain estimator of θ^* in the worst-case direction among all active arms (see Lemma 1), along with a suitable regularization term.
- 2. Pull a budgeted number of samples (scaled by $2^{\ell \cdot \frac{1+\epsilon}{\epsilon}}$) from that distribution, and estimate the reward for each active arm separately using a robust mean estimator.
- 3. Fit a parameter $\hat{\theta}$ that minimizes the maximum distance of $\hat{\theta}^{\top}a$ to the estimated reward of a over all active arms.
- 4. Eliminate suboptimal arms from the active set.

This process is repeated with progressively tighter accuracy until the time horizon is reached or a single arm remains. In the latter case, the remaining arm is pulled for all remaining rounds.

MED-PE is a generalization of Robust Inverse Propensity Score estimator in [CJKS21] which assumes a bounded variance for the rewards. We first find an experimental design that minimizes the $(1+\epsilon)$ -absolute moment of $a^{\top}A^{(\gamma)}(\lambda)^{-1}x$, with suitable regularization, for all a that are active (and therefore the confidence interval of the robust estimator). Note that $A^{(\gamma)}(\lambda)^{-1}x_sy_s$ (with $A^{(\gamma)}(\lambda) := \gamma I + \mathbb{E}_{x \sim \lambda}[xx^{\top}]$) can be interpreted as a single-sample regularized least squares estimator, which is then robustified through a robust mean estimation subroutine $\hat{\mu}$ for each arm. The overall accuracy guarantee of this estimator turns out to depend directly on $M_{1+\epsilon}(\lambda; \mathcal{A}, \gamma, \beta)$ (see Lemma 1 below), which is why we seek to minimize this quantity in our design λ_{ℓ}^* . Moreover, we include a regularization term for design optimization to mitigate the estimator's bias, as $A^{(\gamma)}(\lambda)^{-1}x_sy_s$ is biased for $\gamma \neq 0$.

Any robust mean estimator such as truncated (trimmed) mean, median-of-means, or Catoni's M estimator [LM19a, Cat12], can be used as the subroutine $\hat{\mu}$ of MED-PE. We adopt the truncated mean for concreteness and simplicity. The following lemma provides our main confidence interval for our regression estimator.

Lemma 1. Consider $(x_i, y_i)_{i=1}^n$, where $x_i \sim \lambda(\mathcal{A})$ are i.i.d. vectors from distribution λ over \mathcal{A} , and suppose that $y_i = \langle \theta^*, x_i \rangle + \eta_i$, where η_i are independent zero-mean noise terms such that $\mathbb{E}[|\eta_i|^{1+\epsilon}] \leq v$, and $\max_{a \in \mathcal{A}} |\langle \theta^*, a \rangle| \leq 1$. The estimator $\widehat{\theta}(\gamma)$ with a robust mean estimator $\widehat{\mu}$ as a subroutine is defined as follows:

$$\widehat{\theta}(\gamma) := \arg\min_{\theta} \max_{a \in \mathcal{A}} \left| \theta^{\top} a - \widehat{\mu} \left(\{ a^{\top} A^{(\gamma)}(\lambda)^{-1} x_i y_i \}_{i=1}^n, \frac{\delta}{|\mathcal{A}|} \right) \right|,$$

where $A^{(\gamma)}(\lambda) := \gamma I + \mathbb{E}_{x \sim \lambda}[xx^{\top}]$. For any $\beta \geq 0$, $\widehat{\theta}(\gamma)$ with the truncated empirical mean $\widehat{\mu}(\{X_i\}_{i=1}^n, \delta) := \frac{1}{n} \sum X_i \mathbb{I}\{|X_i| \leq \left(\frac{vt}{\log(\delta^{-1})}\right)^{\frac{1}{1+\epsilon}}\}$ as a subroutine, satisfies the following with probability at least $1 - \delta$:

$$\max_{a \in \mathcal{A}} |\langle \widehat{\theta} - \theta^*, a \rangle| \le \left(2\gamma^{1/2} \|\theta^*\|_2 \beta^{-1} + 32(1+\upsilon)^{\frac{1}{1+\epsilon}} \left(\frac{\log(|\mathcal{A}|/\delta)}{n} \right)^{\frac{\epsilon}{1+\epsilon}} \right) M_{1+\epsilon}(\lambda; \mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}},$$

where
$$M_{1+\epsilon}(\lambda; \mathcal{A}, \gamma, \beta) := \max_{a \in \mathcal{A}} \mathbb{E}_{x \sim \lambda} \left[\left| a^{\top} A^{(\gamma)}(\lambda)^{-1} x \right|^{1+\epsilon} \right] + \beta^{1+\epsilon} \|a\|_{A^{(\gamma)}(\lambda)^{-1}}^{1+\epsilon}.$$

Proof Sketch. In order to use the robust mean estimator guaranties, we bound the $(1+\epsilon)$ -absolute moment of our samples $a^{\top}A^{(\gamma)}(\lambda)^{-1}xy$ for $x\sim\lambda$. Using the boundedness of the expected rewards and the $(1+\epsilon)$ -absolute moment of the noise η , we show that the moment is bounded by $4(1+\upsilon)M_{1+\epsilon}(\lambda;\mathcal{A},\gamma,\beta)$. Moreover, the expected reward estimator for arm a (denoted by $W^{(a)}$) is biased if $\gamma>0$, and we can bound the bias as follows:

$$\left| \langle \theta^*, a \rangle - \mathbb{E}[W^{(a)}] \right| \le \sqrt{\gamma} \beta^{-1} \|\theta^*\|_2 M_{1+\epsilon}(\lambda; \mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}}.$$

Using the triangle inequality and the union bound then gives the desired result. The detailed proof is given in Appendix A. \Box

The following theorem states our general action set dependent regret bound for MED-PE.

Theorem 3. For any linear bandit problem with finite action set $A \subseteq \mathbb{R}^d$, define

$$M_{1+\epsilon}^*(\mathcal{A}, \gamma, \beta) := \max_{\mathcal{V} \subseteq \mathcal{A}} \min_{\lambda \in \Delta^{\mathcal{V}}} M_{1+\epsilon}(\lambda; \mathcal{V}, \gamma, \beta).$$

If $\mathbb{E}[|\eta_t|^{1+\epsilon}] \leq v$, $\|\theta^*\|_2 \leq b$, and $\sup_{x \in \mathcal{A}} |a^\top \theta^*| \leq 1$, then MED-PE with the truncated empirical mean estimator (Lemma 1) and $\gamma = T^{-\frac{2\epsilon}{1+\epsilon}}$ achieves regret bounded by

$$R_T \le \left(C_0 \beta^{-1} b + C_1 (1+v)^{\frac{1}{1+\epsilon}} \log(|\mathcal{A}| T \log^2 T)^{\frac{\epsilon}{1+\epsilon}} \right) M_{1+\epsilon}^* (\mathcal{A}, T^{\frac{-2\epsilon}{1+\epsilon}}, \beta)^{\frac{1}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}$$

for some constants C_0 and C_1 .

Proof Sketch. Using Lemma 1, with probability at least $1 - (2\ell^2 T)^{-1}$, we have

$$\max_{a \in \mathcal{A}_{\ell}} |a^{\top} \theta^* - a^{\top} \widehat{\theta}_{\ell}| \le \epsilon_{\ell} + 2\gamma^{1/2} b \beta^{-1} M_{1+\epsilon}^* (\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}}.$$

Therefore, in the phases where ϵ_ℓ is large compared to $\gamma^{1/2}\beta^{-1}M_{1+\epsilon}^*(\mathcal{A},\gamma,\beta)^{\frac{1}{1+\epsilon}}$, suboptimal arms are eliminated, and no optimal arm is eliminated with high probability. In the phases where ϵ_ℓ is smaller, each arm pull incurs regret $\widetilde{\mathcal{O}}(\gamma^{1/2}\beta^{-1}M_{1+\epsilon}^*(\mathcal{A},\gamma,\beta)^{\frac{1}{1+\epsilon}})$. Setting $\gamma=T^{\frac{-2\epsilon}{1+\epsilon}}$, balances the two regret terms, and leads to the final regret bound. The detailed proof is given in Appendix A. \square

Remark 1. If A is not finite, we can cover the domain with $T^{O(d)}$ elements in A, such that the expected reward of each arm can be approximated by one of the covered elements with T^{-1} error, and therefore the bound of Theorem 3 can be written as

$$R_T \le \left(C_0 \beta^{-1} b + C_1' (1+v)^{\frac{1}{1+\epsilon}} d^{\frac{\epsilon}{1+\epsilon}} \log(T^2 \log^2 T)^{\frac{\epsilon}{1+\epsilon}}\right) M_{1+\epsilon}^* (\mathcal{A}, T^{\frac{-2\epsilon}{1+\epsilon}}, \beta)^{\frac{1}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}.$$

The quantity $M_{1+\epsilon}^*$ in Theorem 3 may be difficult to characterize precisely in general, but the following lemma gives a universal upper bound.

Lemma 2. For any action set A and $\epsilon \in (0,1]$, setting $\gamma = T^{\frac{-2\epsilon}{1+\epsilon}}$ and $\beta = 1$, we have

$$M_{1+\epsilon}^*(\mathcal{A}, T^{-\frac{2\epsilon}{1+\epsilon}}, 1) \le d^{\frac{1+\epsilon}{2}}.$$

Moreover, a design λ with $M_{1+\epsilon}(\lambda;\mathcal{A},T^{\frac{-2\epsilon}{1+\epsilon}},1)=O(d^{\frac{1+\epsilon}{2}})$ can be found with $O(d\log\log d)$ time.

Proof. We upper bound the first term in the objective function as follows:

$$\mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}x\big|^{1+\epsilon}\Big] \leq \mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}x\big|^{2}\Big]^{\frac{1+\epsilon}{2}} \qquad \text{(Jensen's inequality)}$$

$$= \mathbb{E}\big[a^{\top}A^{(\gamma)}(\lambda)^{-1}xx^{\top}A^{(\gamma)}(\lambda)^{-1}a\big]^{\frac{1+\epsilon}{2}}$$

$$\leq \|a\|_{A^{(\gamma)}(\lambda)^{-1}}^{1+\epsilon}. \qquad (\mathbb{E}[xx^{\top}] = \sum_{x} \lambda(x)xx^{\top} \leq A^{(\gamma)}(\lambda))$$

Hence, the minimization of $M_{1+\epsilon}$ is upper bounded in terms of a minimization of $\max_a \|a\|_{A(\lambda)^{-1}}$. This is equivalent to G-optimal design which is well-studied and the following is known (e.g., see [LS20, Chapter 21]): (i) The problem is convex and its optimal value is at most \sqrt{d} ; (ii) There are efficient algorithms such as Frank–Wolfe that can find a design having $\max_a \|a\|_{A(\lambda)^{-1}} = O(\sqrt{d})$ with $O(d \log \log d)$ iterations.

Combining Theorem 3 and Lemma 2, we obtain the following.

Corollary 1. For any action set A, MED-PE achieves regret $\widetilde{\mathcal{O}}(d^{\frac{1+3\epsilon}{2(1+\epsilon)}}T^{\frac{1}{1+\epsilon}})$. Moreover, for a finite action set with |A|=n, the regret bound is lowered to $\widetilde{\mathcal{O}}(\sqrt{d}T^{\frac{1}{1+\epsilon}}(\log n)^{\frac{\epsilon}{1+\epsilon}})$.

Computational complexity. By Lemma 2, a design over general action sets can be computed efficiently. The truncated sample-mean estimator can also be computed in linear time. Moreover, the minimum-distance estimator for $\hat{\theta}$ is obtained by solving a linear optimization problem and is therefore computable in polynomial time; in the infinite-dimensional case this is handled via a dual formulation (see Appendix C). The dominant per-round cost is the linear pass over the active arms to update estimates and apply elimination tests, which is standard for finite-arm algorithms.

The bound in Corollary 1 is the worst-case regret over all possible action sets A. However, based on geometry of the action set, we can achieve tighter regret bounds, as we see below.

3.1 Special Cases of the Action Set

Simplex. When \mathcal{A} is the simplex, the problem is essentially one of multi-armed bandits with d arms. Consider λ being uniform over canonical basis; then $A(\lambda) = \frac{1}{d}I$, and for each $a \in \mathcal{A}$, we have

$$\mathbb{E}_{x \sim \lambda}[|a^{\top} A^{-1} x|^{1+\epsilon}] = \mathbb{E}_{x \sim \lambda}[|da^{\top} x|^{1+\epsilon}] = d^{1+\epsilon} \sum_{i=1}^{d} d^{-1} |a^{\top} e_i|^{1+\epsilon} = d^{\epsilon} \sum_{i=1}^{d} |a_i|^{1+\epsilon} \le d^{\epsilon}.$$

Since one of the canonical basis vectors (or its negation) must be optimal when \mathcal{A} is the simplex, we can simply restrict to this subset of 2d actions, giving the following corollary, which recovers the well-known scaling for heavy-tailed MAB [BCBL13].

Corollary 2. For the simplex action set $\mathcal{A}=\Delta^d$, if the assumptions of Theorem 3 hold, then MED-PE, with parameters $\gamma=T^{\frac{-2\epsilon}{1+\epsilon}}, \beta=d^{\frac{\epsilon-1}{2}}$ achieves regret $\widetilde{\mathcal{O}}(d^{\frac{\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$.

 l_p -norm ball with radius r for $p \le 1 + \epsilon$. Similarly to the simplex, if we define λ to be uniform over $\{r\mathbf{e}_i\}_{i=1}^d$, then $A(\lambda) = \frac{r^2}{d}I$ for any $v \in \mathcal{B}(\|\cdot\|_p, r)$, and we have

$$\mathbb{E}_{x \sim \lambda}[|a^{\top}A^{-1}x|^{1+\epsilon}] = \mathbb{E}_{x \sim \lambda}\left[\left|\frac{d}{r^2}a^{\top}x\right|^{1+\epsilon}\right] = d^{\epsilon}\sum_{i=1}^{d}\left|\frac{a_i}{r}\right|^{1+\epsilon} \le d^{\epsilon}\sum_{i=1}^{d}\left|\frac{a_i}{r}\right|^{p} \le d^{\epsilon},$$

where the last inequality is by the definition of the l_p -norm ball.

Corollary 3. For the action set $\mathcal{A}=\{x:\|x\|_p\leq r\}$ with $p\leq 1+\epsilon$, if the assumptions of Theorem 3 hold, then MED-PE, with parameters $\gamma=T^{\frac{-2\epsilon}{1+\epsilon}}, \beta=d^{\frac{\epsilon-1}{2}}$, has regret of $\widetilde{\mathcal{O}}(d^{\frac{2\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$.

Matérn Kernels. Our algorithm does not require the action features to lie in a finite-dimensional space, as long as the design and the estimator $a^{\top}A^{(\gamma)}(\lambda)^{-1}x$ can be computed efficiently. In particular, following the approach of [CJKS21], our method extends naturally to kernel bandits, where the reward function belongs to a reproducing kernel Hilbert space (RKHS) associated with a

kernel K satisfying $K(x,y) = \phi(x)^{\top}\phi(y)$ for some (possibly infinite-dimensional) feature map ϕ . Since our focus is on linear bandits, we defer a full description of the kernel setting to Appendix C, where we also establish the following corollary (stated informally here, with the formal version deferred to Appendix C).

Corollary 4. (Informal) For the kernel bandit problem with domain $[0,1]^d$ for a constant value of d, under the Matérn kernel with smoothness parameter $\nu > 0$, the kernelized version of MED-PE (with suitably-chosen parameters) achieves regret $\widetilde{\mathcal{O}}(T^{1-\frac{\epsilon}{1+\epsilon}\cdot\frac{2\nu}{2\nu+d}})$.

While this does not match the known lower bound (except when $\epsilon=1$ or in the limit as $\epsilon\to 0$), it significantly improves over the best existing upper bound [CG19], which is only sublinear in T for a relatively narrow range of (ϵ,d,ν) . In contrast, our bound is sublinear in T for all such choices.

4 Conclusion

In this paper, we revisited stochastic linear bandits with heavy-tailed rewards and substantially narrowed the gap between known minimax lower and upper regret bounds in both the infinite- and finite-action settings. Our new regression estimator, guided by geometry-aware experimental design, yields improved instance-dependent guarantees that leverage the structure of the action set. Since our geometry-dependent bounds recover the $d^{\frac{2\epsilon}{1+\epsilon}}$ dimension dependence that also appears in our minimax lower bound, it is plausible that this gives the correct minimax rate for general action sets. Closing the remaining gap to establish true minimax-optimal rates for all moment parameters, and precisely characterizing the action-set-dependent complexity term under different geometries, remain promising directions for future work.

Acknowledgement

This work was supported by the Singapore National Research Foundation (NRF) under its AI Visiting Professorship programme, and by NSF Award TRIPODS 202323.

References

- [AC11] Jean-Yves Audibert and Olivier Catoni. Robust linear least squares regression. *The Annals of Statistics*, 39(5), October 2011.
- [BCBL13] Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *IEEE Transactions on Information Theory*, 59(11):7711–7717, 2013.
 - [BJL15] Christian Brownlees, Emilien Joly, and Gábor Lugosi. Empirical risk minimization for heavy-tailed losses. *The Annals of Statistics*, 43(6), December 2015.
 - [BTA02] François Baccelli, Gérard H. Taché, and Etienne Altman. Flow complexity and heavy-tailed delays in packet networks. *Performance Evaluation*, 49(1–4):427–449, 2002.
 - [Cat12] Olivier Catoni. Challenging the Empirical Mean and Empirical Variance: A Deviation Study, volume 1906 of Lecture Notes in Mathematics. Springer, 2012.
 - [CB00] Rama Cont and Jean-Philippe Bouchaud. Herd behavior and aggregate fluctuations in financial markets. *Macroeconomic Dynamics*, 4(2):170–196, 2000.
 - [CG19] Sayak Ray Chowdhury and Aditya Gopalan. Bayesian optimization under heavy-tailed payoffs. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- [CHDH25] Yu Chen, Jiatai Huang, Yan Dai, and Longbo Huang. uniINF: Best-of-both-worlds algorithm for parameter-free heavy-tailed MABs. In *International Conference on Learning Representations (ICLR)*, 2025.
- [CJKS21] Romain Camilleri, Kevin Jamieson, and Julian Katz-Samuels. High-dimensional experimental design and kernel bandits. In *International Conference on Machine Learning (ICML)*, pages 1227–1237. PMLR, 2021.

- [Con01] R. Cont. Empirical properties of asset returns: Stylized facts and statistical issues. *Quantitative Finance*, 1(2):223–236, 2001.
- [CvdLG20] Yong-Kyu Choi, Erko van der Laan, and Omar Ghattas. Modeling heavy-tailed conversion values in real-time bidding. In *ACM International Conference on Web Search and Data Mining (WSDM)*, pages 870–878, 2020.
- [DLLO16] Luc Devroye, Matthieu Lerasle, Gabor Lugosi, and Roberto I. Oliveira. Sub-Gaussian mean estimators. *The Annals of Statistics*, 44(6):2695 2725, 2016.
- [HDH22] Jiatai Huang, Yan Dai, and Longbo Huang. Adaptive best-of-both-worlds algorithm for heavy-tailed multi-armed bandits. In *International Conference on Machine Learning (ICML)*, volume 162, pages 9173–9200. PMLR, 17–23 Jul 2022.
 - [HS14] Daniel Hsu and Sivan Sabato. Heavy-tailed regression with a generalized median-of-means. In *International Conference on Machine Learning (ICML)*, volume 32, pages 37–45. PMLR, 22–24 Jun 2014.
- [HW19] Qiyang Han and Jon A. Wellner. Convergence rates of least squares regression estimators with heavy-tailed errors. *The Annals of Statistics*, 47(4):2286 2319, 2019.
- [HZWY23] Jiayi Huang, Han Zhong, Liwei Wang, and Lin Yang. Tackling heavy-tailed rewards in reinforcement learning with function approximation: Minimax optimal and instancedependent regret bounds. In Conference on Neural Information Processing Systems (NeurIPS), 2023.
 - [JSP21] Saravanan Jebarajakirthy, Paurav Shukla, and Prashant Palvia. Heavy-tailed distributions in online ad response: A marketing analytics perspective. *Journal of Business Research*, 124:818–830, 2021.
 - [KK23] Minhyun Kang and Gi-Soo Kim. Heavy-tailed linear bandit with Huber regression. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 216, pages 1027–1036. PMLR, 31 Jul–04 Aug 2023.
 - [LM19a] Gábor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190, 2019.
 - [LM19b] Gábor Lugosi and Shahar Mendelson. Sub-Gaussian estimators of the mean of a random vector. *The Annals of Statistics*, 47(2):783 794, 2019.
 - [LS20] Tor Lattimore and Csaba Szepesvári. Bandit Algorithms. Cambridge University Press, 2020.
 - [LS24] Xiang Li and Qiang Sun. Variance-aware decision making with linear function approximation under heavy-tailed rewards. *Transactions on Machine Learning Research*, 2024.
- [LWHZ19] Shiyin Lu, Guanghui Wang, Yao Hu, and Lijun Zhang. Optimal algorithms for Lipschitz bandits with heavy-tailed rewards. In *International Conference on Machine Learning (ICML)*, volume 97, pages 4154–4163. PMLR, 09–15 Jun 2019.
- [LYLO20] Kyungjae Lee, Hongjun Yang, Sungbin Lim, and Songhwai Oh. Optimal algorithms for stochastic multi-armed bandits with heavy tailed rewards. In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 33, pages 8452–8462, 2020.
 - [MY16] Andres Muñoz Medina and Scott Yang. No-regret algorithms for heavy-tailed linear bandits. In *International Conference on Machine Learning (ICML)*, pages 1642–1650, 2016.
- [RVHH15] James A. Roberts, László A. E. Varnai, Brenton H. Houghton, and David Hughes. Heavy-tailed distributions in the amplitude of neural oscillations. *Journal of Neuro-science*, 35(19):7313–7323, 2015.

- [Sas15] Igal Sason. An improved reverse pinsker inequality for probability distributions on a finite set. *CoRR*, abs/1503.03417, 2015.
- [SBC17] Jonathan Scarlett, Ilija Bogunovic, and Volkan Cevher. Lower bounds on regret for noisy Gaussian process bandit optimization. In Conference on Learning Theory (COLT). 2017.
- [SYKL18] Han Shao, Xiaotian Yu, Irwin King, and Michael R Lyu. Almost optimal algorithms for linear stochastic bandits with heavy-tailed payoffs. In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
 - [SZF20] Qiang Sun, Wen-Xin Zhou, and Jianqing Fan. Adaptive Huber regression. *Journal of the American Statistical Association*, 115(529):254–265, 2020.
- [VBJ⁺21] Sattar Vakili, Nacime Bouziani, Sepehr Jalali, Alberto Bernacchia, and Da-shan Shiu. Optimal order simple regret for Gaussian process bandits. *Conference on Neural Information Processing Systems (NeurIPS)*, 34:21202–21215, 2021.
- [VKP21] Sattar Vakili, Kia Khezeli, and Victor Picheny. On information gain and regret bounds in Gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2021.
- [WS21] Lai Wei and Vaibhav Srivastava. Minimax policy for heavy-tailed bandits. *IEEE Control Systems Letters*, 5(4):1423–1428, 2021.
- [WZZZ25] Jing Wang, Yu-Jie Zhang, Peng Zhao, and Zhi-Hua Zhou. Heavy-tailed linear bandits: Huber regression with one-pass update. *arXiv preprint arXiv:2503.00419*, 2025.
- [XWW⁺23] Bo Xue, Yimu Wang, Yuanyu Wan, Jinfeng Yi, and Lijun Zhang. Efficient algorithms for generalized linear bandits with heavy-tailed rewards. In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 36, pages 70880–70891, 2023.
- [XWWZ20] Bo Xue, Guanghui Wang, Yimu Wang, and Lijun Zhang. Nearly optimal regret for stochastic linear bandits with heavy-tailed payoffs. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2936–2942, 7 2020.
- [YNKL18] Xiaotian Yu, Yichuan Nevmyvaka, Irwin King, and Michael R. Lyu. Pure exploration of multi-armed bandits with heavy-tailed payoffs. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018.
- [ZHYW21] Han Zhong, Jiayi Huang, Lin Yang, and Liwei Wang. Breaking the moments condition barrier: No-regret algorithm for bandits with super heavy-tailed payoffs. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The main claims made in the abstract and introduction accurately reflect the paper's contributions and scope. The claims are validated by detailed proofs.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper discuss the limitations of the results.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The paper provides detailed assumptions and proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: This is a theoretical paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: This paper does not include experiments.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: This paper does not include experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: This paper does not include experiments.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: This paper does not include experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research conducted in the paper conforms, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This is a theoretical work. There is no societal impact of the work performed.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: This paper does not use existing assets.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

• If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- · Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- · For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: This research does not involve LLMs as any important, original, or non-standard components.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Upper Bound Proofs

A.1 Proof of Lemma 1 (Confidence Interval)

We first state a well known guarantee of the truncated mean estimator.

Lemma 3. (Lemma 1 of [BCBL13]) Let X_1, \ldots, X_n be i.i.d. random variables such that $\mathbb{E}[|X_i|^{1+\epsilon}] \leq u$ for some $\epsilon \in (0,1]$. Then the truncated empirical mean estimator $\widehat{\mu}(\{X_i\}_{i=1}^n, \delta) := \frac{1}{n} \sum_{i=1}^n X_i \mathbb{I}\{|X_i| \leq \left(\frac{ut}{\log(\delta^{-1})}\right)^{\frac{1}{1+\epsilon}}\}$ satisfies with probability at least $1-\delta$ that

$$|\widehat{\mu}(\{X_i\}_{i=1}^n, \delta) - \mu| \le 4u^{\frac{1}{1+\epsilon}} \left(\frac{\log(\delta^{-1})}{n}\right)^{\frac{\epsilon}{1+\epsilon}}.$$

Let $W^{(a)} := \widehat{\mu}\left(\{a^{\top}A^{(\gamma)}(\lambda)^{-1}x_i\,y_i\}_{i=1}^n, \frac{\delta}{|\mathcal{A}|}\right)$. We first observe that

$$\begin{split} \max_{a \in \mathcal{A}} |a^{\top} \widehat{\theta}(\gamma) - a^{\top} \theta^*| &= \max_{a \in \mathcal{A}} |a^{\top} \widehat{\theta}(\gamma) - W^{(a)} + W^{(a)} - a^{\top} \theta^*| \\ &\leq \max_{a \in \mathcal{A}} |a^{\top} \widehat{\theta}(\gamma) - W^{(a)}| + \max_{a \in \mathcal{A}} |W^{(a)} - a^{\top} \theta^*| \\ &= \min_{\theta} \max_{a \in \mathcal{A}} |a^T \theta - W^{(a)}| + \max_{a \in \mathcal{A}} |W^{(a)} - a^{\top} \theta^*| \\ &\leq 2 \max_{a \in \mathcal{A}} |W^{(a)} - a^{\top} \theta^*|. \end{split} \tag{def. } \widehat{\theta}(\gamma))$$

For fixed a, we bound the $(1 + \epsilon)$ -moment of $a^{\top}A^{(\gamma)}(\lambda)^{-1}xy$, where $x \sim \lambda$ and $y = x^{\top}\theta^* + \eta$, as follows:

$$\begin{split} \mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}xy\big|^{1+\epsilon}\Big] &= \mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}x(x^{\top}\theta^* + \eta)\big|^{1+\epsilon}\Big] \\ &\leq 2^{1+\epsilon}\mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}x(x^{\top}\theta^*)\big|^{1+\epsilon}\Big] + 2^{1+\epsilon}\mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}x\big|^{1+\epsilon}|\eta|^{1+\epsilon}\Big] \\ &\qquad \qquad (|a+b| \leq 2\max\{|a|,|b|\}) \\ &\leq 4\mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}x\big|^{1+\epsilon}\Big] + 4v\mathbb{E}\Big[\big|a^{\top}A^{(\gamma)}(\lambda)^{-1}x\big|^{1+\epsilon}\Big] \\ &\qquad \qquad (|x^{\top}\theta^*| \leq 1 \text{ and } \mathbb{E}[|\eta|^{1+\epsilon}] \leq v) \\ &\leq 4(1+v)M_{1+\epsilon}(\lambda;\mathcal{A},\gamma,\beta). \end{split}$$

Using this moment bound and Lemma 3, for any a, we have with probability at least $1 - \frac{\delta}{|\mathcal{A}|}$ that

$$|W^{(a)} - \mathbb{E}[W^{(a)}]| \le 16(1+\upsilon)^{\frac{1}{1+\epsilon}} M_{1+\epsilon}(\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}} \left(\frac{\log(\delta^{-1}|\mathcal{A}|)}{n}\right)^{\frac{\epsilon}{1+\epsilon}}.$$

Moreover, we have

$$\begin{split} |a^{\top}\theta^* - \mathbb{E}[W^{(a)}]| &= |\langle \theta^*, a \rangle - \mathbb{E}[a^{\top}A^{(\gamma)}(\lambda)^{-1}xx^{\top}\theta^*]| & (\text{def. } W^{(a)}) \\ &= |\langle \theta^*, a \rangle - a^{\top}A^{(\gamma)}(\lambda)^{-1}A(\lambda)\theta^*| & (\text{where } A(\lambda) = \mathbb{E}[xx^T]) \\ &= |\langle \theta^*, a \rangle - a^{\top}A^{(\gamma)}(\lambda)^{-1}\left(A^{(\gamma)}(\lambda) - \gamma I\right)\theta^*| & (A(\lambda) = A^{(\gamma)}(\lambda) - \gamma I) \\ &= \gamma |a^{\top}A^{(\gamma)}(\lambda)^{-1}\theta^*| & \\ &= \gamma |a^{\top}(A(\lambda) + \gamma I)^{-1/2}(A(\lambda) + \gamma I)^{-1/2}\theta^*| & \\ &\leq \gamma \|a\|_{A^{(\gamma)}(\lambda)^{-1}}\gamma^{-1/2}\|\theta^*\|_{(I+\gamma^{-1}A(\lambda))^{-1}} & (\text{Cauchy-Schwarz}) \\ &\leq \gamma^{1/2}\|a\|_{A^{(\gamma)}(\lambda)^{-1}}\|\theta^*\|_2 & (I+\gamma^{-1}A(\lambda) \succeq I) \\ &\leq \gamma^{1/2}\beta^{-1}\|\theta^*\|_2 M_{1+\epsilon}(\lambda; \mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}}. & (\text{def. } M_{1+\epsilon}) \end{split}$$

Putting the two inequalities together, and using the union bound completes the proof.

A.2 Proof of Theorem 3 (Regret Bound for MED-PE)

Using Lemma 1 for action set A_{ℓ} , we have with probability of at least $1 - \frac{1}{2\ell^2T}$,

$$\max_{a \in \mathcal{A}_{\ell}} |a^{\top} \theta^* - a^{\top} \widehat{\theta}_{\ell}| \leq \left(2\gamma^{1/2} \|\theta^*\|_{2} \beta^{-1} + 32(1+\upsilon)^{\frac{1}{1+\epsilon}} \left(\frac{\log(2l^2T|\mathcal{A}_{\ell}|)}{\tau_{\ell}} \right)^{\frac{\epsilon}{1+\epsilon}} \right) M_{1+\epsilon} (\lambda_{\ell}^*; \mathcal{A}_{\ell}, \gamma, \beta)^{\frac{1}{1+\epsilon}} \\
\leq 2\gamma^{1/2} b \beta^{-1} M_{1+\epsilon} (\lambda_{\ell}^*; \mathcal{A}_{\ell}, \gamma, \beta)^{\frac{1}{1+\epsilon}} + \epsilon_{\ell} \qquad \text{(choice of } \tau_{\ell} \text{ in Algorithm 1)} \\
\leq 2\gamma^{1/2} b \beta^{-1} M_{1+\epsilon}^* (\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}} + \epsilon_{\ell} \qquad \text{(def. } M_{1+\epsilon}^*)$$

Now we define the event $\mathcal{E} := \bigcap_{\ell=1}^{\infty} \bigcap_{x \in A_{\ell}} \mathcal{E}_{x,l}(\mathcal{A}_{\ell})$, where

$$\mathcal{E}_{x,l}(\mathcal{V}) := \left\{ |x^{\top} \widehat{\theta}_{\ell}(\mathcal{V}) - x^{\top} \theta^*| \le \epsilon_{\ell} + 2b\gamma^{1/2} \beta^{-1} M_{1+\epsilon}^* (\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}} \right\},\,$$

with $\widehat{\theta}_{\ell}(\cdot)$ corresponding to $\widehat{\theta}_{\ell}$ in Algorithm 1 with an explicit dependence on the action subset. Then, we have

$$\mathbb{P}\left(\bigcup_{\ell=1}^{\infty}\bigcup_{x\in\mathcal{A}_{\ell}}\left\{\mathcal{E}_{x,\ell}^{c}(\mathcal{A}_{\ell})\right\}\right) \leq \sum_{\ell=1}^{\infty}\mathbb{P}\left(\bigcup_{x\in\mathcal{A}_{\ell}}\left\{\mathcal{E}_{x,\ell}^{c}(\mathcal{A}_{\ell})\right\}\right) \\
= \sum_{\ell=1}^{\infty}\sum_{\mathcal{V}\subseteq\mathcal{A}}\mathbb{P}\left(\bigcup_{x\in\mathcal{V}}\left\{\mathcal{E}_{x,\ell}^{c}(\mathcal{V})\right\} \middle| \mathcal{A}_{\ell} = \mathcal{V}\right)\mathbb{P}(\mathcal{A}_{\ell} = \mathcal{V}) \\
\leq \sum_{\ell=1}^{\infty}\sum_{\mathcal{V}\subseteq\mathcal{A}}\frac{1}{2\ell^{2}T}\mathbb{P}(\mathcal{A}_{\ell} = \mathcal{V}) \leq \frac{1}{T}, \quad \text{(union bound and } \sum_{\ell=1}^{\infty}\frac{1}{\ell^{2}} < 2)$$

As $\mathbb{E}[R_T \mathbf{1}_{\mathcal{E}^c}] = \mathbb{E}[R_T | \mathcal{E}^c] \mathbb{P}[\mathcal{E}^c] \leq (\sup_{x,x'} x'^\top \theta^* - x^\top \theta^*) T \frac{1}{T} \leq 2$, for the rest of the proof we assume event \mathcal{E} .

Let $x^* = \operatorname{argmax}_{x \in \mathcal{A}} x^{\top} \theta^*$; then, for every ℓ such that $2\epsilon_{\ell} \geq 4b\gamma^{1/2}\beta^{-1}M_{1+\epsilon}(\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}}$ and any $x \in \mathcal{A}_{\ell}$, we have

$$\begin{split} x^{\top}\widehat{\theta}_{\ell} - x^{*\top}\widehat{\theta}_{\ell} &= (x^{\top}\widehat{\theta}_{\ell} - x^{\top}\theta^{*}) + (x^{\top}\theta^{*} - x^{*\top}\theta^{*}) + (x^{*\top}\theta^{*} - x^{*\top}\widehat{\theta}_{\ell}) \\ &\leq 2\epsilon_{\ell} + 4b\gamma^{1/2}\beta^{-1}M_{1+\epsilon}^{*}(\mathcal{A},\gamma,\beta)^{\frac{1}{1+\epsilon}} & (\text{def. } \mathcal{E} \text{ and def. } x^{*}) \\ &\leq 4\epsilon_{\ell}. & (\text{assumption on } \epsilon_{\ell}) \end{split}$$

Therefore, recalling the elimination rule in Algorithm 1, we have by induction that $x^* \in \mathcal{A}_{\ell+1}$. We also claim that all suboptimal actions of gap more than $8\epsilon_\ell = 16\epsilon_{\ell+1}$ are eliminated at the end of epoch ℓ . To see this, let $x' \in \mathcal{A}_\ell$ be such an action, and observe that

$$\begin{split} \max_{x \in \mathcal{A}_{\ell}} \left(x'^{\top} \widehat{\theta}_{\ell} - x^{\top} \widehat{\theta}_{\ell} \right) &\geq x^{*\top} \widehat{\theta}_{\ell} - x^{\top} \widehat{\theta}_{\ell} \\ &\geq x^{*\top} \theta^{*} - x^{\top} \theta^{*} - 2\epsilon_{\ell} - 4b\gamma^{1/2} \beta^{-1} M_{1+\epsilon}^{*} (\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}} \quad \text{(shown above)} \\ &\geq x^{*\top} \theta^{*} - x^{\top} \theta^{*} - 4\epsilon_{\ell} \\ &> 4\epsilon_{\ell}. \end{split} \tag{assumption on } \epsilon_{\ell})$$

In summary, the above arguments show that when $2\epsilon_\ell \geq 4b\gamma^{1/2}\beta^{-1}M_{1+\epsilon}^*(\mathcal{A},\gamma,\beta)^{\frac{1}{1+\epsilon}}$, the regret incurred in epoch $\ell+1$ is at most $16\epsilon_{\ell+1}$. Since $\mathcal{A}_{\ell+1} \subseteq \mathcal{A}_\ell$, this also implies that even when ℓ increases beyond such a point, we still incur regret at most $32b\gamma^{1/2}\beta^{-1}M_{1+\epsilon}^*(\mathcal{A},\gamma,\beta)^{\frac{1}{1+\epsilon}}$.

Finally, we can upper bound the regret as follows:

$$\mathbb{E}[R_T] \leq \sum_{\ell} \tau_{\ell} \left(\sup_{x \in \mathcal{A}_{\ell}} x^{*T} \theta^* - x^T \theta^* \right)$$

$$\leq \sum_{\ell} \tau_{\ell} \max \{ 16\epsilon_{\ell}, 32b\gamma^{1/2}\beta^{-1} M_{1+\epsilon}^* (\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}} \} \qquad \text{(shown above)}$$

$$\leq \sum_{\ell} 16\tau_{\ell} \epsilon_{\ell} + T\zeta \qquad (\zeta := 32b\gamma^{1/2}\beta^{-1} M_{1+\epsilon} (\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}})$$

$$\leq \sum_{\ell \,:\, 16\epsilon_{\ell} \geq \omega} 16\epsilon_{\ell}\tau_{\ell} + T\omega + T\zeta \qquad \qquad \text{(for any } \omega \geq 0)$$

$$\leq \sum_{\ell \,:\, 16\epsilon_{\ell} \geq \omega} 16\epsilon_{\ell} 32^{\frac{1+\epsilon}{\epsilon}} (1+v)^{\frac{1}{\epsilon}} \varepsilon_{\ell}^{-\frac{1+\epsilon}{\epsilon}} M_{1+\epsilon}^{*} (\mathcal{A}, \gamma, \beta)^{\frac{1}{\epsilon}} \log(2l^{2}|\mathcal{A}|T) + T(\omega + \zeta)$$

$$(\text{def. } \tau_{\ell} \text{ in Alg. 1})$$

$$\leq \sum_{\ell \,:\, 16\epsilon_{\ell} \geq \omega} C_{1}' (1+v)^{\frac{1}{\epsilon}} \varepsilon_{\ell}^{-\frac{1}{\epsilon}} M_{1+\epsilon}^{*} (\mathcal{A}, \gamma, \beta)^{\frac{1}{\epsilon}} \log(2l^{2}|\mathcal{A}|T) + T(\omega + \zeta)$$

$$(\text{for some constant } C_{1}')$$

$$\leq C_{1}(1+v)^{\frac{1}{1+\epsilon}} M_{1+\epsilon}^{*} (\mathcal{A}, \gamma, \beta)^{\frac{1}{1+\epsilon}} \log(2|\mathcal{A}|T \log_{2}^{2}T)^{\frac{\epsilon}{1+\epsilon}} T^{\frac{1}{1+\epsilon}} + T\zeta$$

$$(\omega := M_{1+\epsilon}^{*} (\cdot)^{\frac{1}{1+\epsilon}} \log(2|\mathcal{A}|T \log_{2}^{2}T)^{\frac{\epsilon}{1+\epsilon}} T^{\frac{1}{1+\epsilon}} \text{ and } \ell \leq \log_{2}T; \text{ see below)}$$

$$\leq \left(C_{0}\beta^{-1}b + C_{1}(1+v)^{\frac{1}{1+\epsilon}} \log(|\mathcal{A}|T \log_{2}^{2}T)^{\frac{\epsilon}{1+\epsilon}}\right) M_{1+\epsilon}^{*} (\mathcal{A}, T^{\frac{-2\epsilon}{1+\epsilon}}, \beta)^{\frac{1}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}.$$

$$(\text{def. } \zeta \text{ and } \gamma = T^{\frac{-2\epsilon}{1+\epsilon}})$$

In more detail, the second-last step upper bounds $\sum_{\ell:16\epsilon_\ell\geq\omega}\epsilon_\ell^{-\frac{1}{\epsilon}}$ by a constant times its largest possible term $\omega^{-\frac{1}{\epsilon}}$, since $\{\epsilon_\ell\}_{\ell\geq 1}$ is exponentially decreasing. Since the choice of ω contains $(M_{1+\epsilon}^*)^{\frac{1}{1+\epsilon}}$, the overall $M_{1+\epsilon}^*$ dependence simplifies as $\left(\frac{M_{1+\epsilon}^*}{(M_{1+\epsilon}^*)^{\frac{1}{1+\epsilon}}}\right)^{\frac{1}{\epsilon}}=\left((M_{1+\epsilon}^*)^{\frac{\epsilon}{1+\epsilon}}\right)^{\frac{1}{\epsilon}}=(M_{1+\epsilon}^*)^{\frac{1}{1+\epsilon}}$.

B Unit Ball Lower Bound

In this appendix, we prove the following lower bound for the case that the action set is the unit ball. **Theorem 4.** Let the action set be $\mathcal{A}=\{x\in\mathbb{R}^d:\|x\|_2\leq 1\}$, and the $(1+\epsilon)$ -absolute moment of the error distribution be bounded by 1. Then, for any algorithm, there exists $\theta^*\in\mathbb{R}^d$ such that $\sup_{x\in\mathcal{A}}|x^\top\theta^*|\leq 1$, and such that for $T\geq d^2$, the regret incurred is $\Omega(d^{\frac{2\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}})$.

Since the KL divergence between Bernoulli random variables Ber(p) and Ber(q) goes to infinity as $p \to 0$, and $\theta^\top x$ can be zero for unit ball, we cannot use the same reward distribution as before. However, we can overcome this by shifting all probabilities and adding -1 to the support of the reward random variable. Specifically, we set the error distribution to be:

$$y(x) = \begin{cases} \left(\frac{1}{\gamma}\right)^{\frac{1}{\epsilon}} & w.p. \ \gamma^{\frac{1}{\epsilon}}(\theta^{\top}x + 2\sqrt{d}\Delta) \\ 0 & w.p. \ 1 - \gamma^{\frac{1}{\epsilon}}(\theta^{\top}x + 2\sqrt{d}\Delta) - 2\sqrt{d}\Delta \\ -1 & w.p. \ 2\sqrt{d}\Delta \end{cases}$$

with $\gamma:=24\sqrt{d}\Delta$ and Δ to be specified later. For any $\theta\in\{\pm\Delta\}^d$, the absolute value of rewards are bounded by $\sum_{i=1}^d\frac{1}{\sqrt{d}}\Delta=\sqrt{d}\Delta$. Then, assuming $\Delta\leq\frac{1}{24\sqrt{d}}$, we have $|\theta^\top x|\leq\sqrt{d}\Delta\leq\frac{1}{8}$ and $\|\theta\|_2\leq 1$ as well as $\gamma\leq 1$, and the $(1+\epsilon)$ -central absolute moment is bounded by:

$$\begin{split} &\mathbb{E}[|y(x) - \theta^\top x|^{1+\epsilon} \mid x] \\ & \leq |\gamma^{-\frac{1}{\epsilon}} - \theta^\top x|^{1+\epsilon} (\theta^\top x + 2\sqrt{d}\Delta) + |\theta^\top x|^{1+\epsilon} + |-1 - \theta^\top x|^{1+\epsilon} 2\sqrt{d}\Delta \qquad (\gamma \leq 1) \\ & \leq 2^{1+\epsilon} \gamma^{-1} 3\sqrt{d}\Delta + (\sqrt{d}\Delta)^{1+\epsilon} + 2\sqrt{d}\Delta (\sqrt{d}\Delta + 1)^{1+\epsilon} \qquad (|\theta^\top x| \leq \sqrt{d}\Delta \leq 1 \text{ and } \gamma^{-\frac{1}{\epsilon}} \geq 1) \\ & \leq \frac{2^{1+\epsilon}}{8} + \left(\frac{1}{24}\right)^{1+\epsilon} + \frac{1}{12}\left(\frac{9}{24}\right)^{1+\epsilon} < 1. \qquad (\text{def. } \gamma, \Delta \leq \frac{1}{24\sqrt{d}}, \text{ and } \epsilon \in (0,1]) \end{split}$$

Defining $T_i := T \wedge \min(s : \sum_{t=1}^s x_{t,i}^2 \ge \frac{T}{d})$, we have

$$\begin{split} R_T(\mathcal{A}, \theta) &= \Delta \mathbb{E}_{\theta} \left[\sum_{t=1}^T \sum_{i=1}^d \left(\frac{1}{\sqrt{d}} - x_{t,i} \text{sign}(\theta_i) \right) \right] \\ &\geq \frac{\Delta \sqrt{d}}{2} \mathbb{E}_{\theta} \left[\sum_{t=1}^T \sum_{i=1}^d \left(\frac{1}{\sqrt{d}} - x_{t,i} \text{sign}(\theta_i) \right)^2 \right] \\ & \text{(by expanding the square and applying } \|x_t\|_2^2 \leq 1) \end{split}$$

$$\geq \frac{\Delta \sqrt{d}}{2} \sum_{i=1}^d \mathbb{E}_{\theta} \left[\sum_{t=1}^{T_i} \left(\frac{1}{\sqrt{d}} - x_{t,i} \mathrm{sign}(\theta_i) \right)^2 \right].$$

Now we define $U_i(b) := \sum_{t=1}^{T_i} \left(\frac{1}{\sqrt{d}} - x_{t,i}b\right)^2$, which gives

$$U_i(1) \le 2\sum_{t=1}^{T_i} \frac{1}{d} + 2\sum_{t=1}^{T_i} x_{t,i}^2 \le \frac{4T}{d} + 2.$$

Then, for any $\theta, \theta' \in \{\pm \Delta\}^d$ that only differ in *i*-th element, we have

$$\begin{split} \mathbb{E}_{\theta}[U_{i}(1)] &\geq \mathbb{E}_{\theta'}[U_{i}(1)] - \left(\frac{4T}{d} + 2\right) \sqrt{\frac{1}{2}} \mathrm{KL}(\mathbb{P}_{\theta} \| \mathbb{P}_{\theta'}) & \text{(Pinsker's inequality)} \\ &\geq \mathbb{E}_{\theta'}[U_{i}(1)] - \left(\frac{4T}{d} + 2\right) \sqrt{\frac{1}{2}} \mathbb{E}_{\theta} \left[\sum_{t=1}^{T_{i}} \mathrm{KL}(y_{\theta}(x_{t}) \| y_{\theta'}(x_{t})) \right] & \text{(Chain rule)} \\ &\geq \mathbb{E}_{\theta'}[U_{i}(1)] - \left(\frac{4T}{d} + 2\right) \sqrt{\frac{1}{2}} \mathbb{E}_{\theta} \left[\sum_{t=1}^{T_{i}} 24^{\frac{1}{\epsilon}} 8\sqrt{d}^{\frac{1-\epsilon}{\epsilon}} \Delta^{\frac{1+\epsilon}{\epsilon}} x_{t,i}^{2} \right] & \text{(Inverse Pinsker's inequality; see below)} \\ &\geq \mathbb{E}_{\theta'}[U_{i}(1)] - 24^{\frac{1}{2\epsilon}} 2\Delta^{\frac{1+\epsilon}{2\epsilon}} \sqrt{d}^{\frac{1-\epsilon}{2\epsilon}} \left(\frac{4T}{d} + 2\right) \sqrt{\mathbb{E}_{\theta} \left[\sum_{t=1}^{T_{i}} x_{t,i}^{2} \right]} \\ &\geq \mathbb{E}_{\theta'}[U_{i}(1)] - 24^{\frac{1}{2\epsilon}} 12\sqrt{2}\Delta^{\frac{1+\epsilon}{2\epsilon}} \sqrt{d}^{\frac{1-\epsilon}{2\epsilon}} \frac{T}{d} \sqrt{\frac{T}{d}}. & (d \leq T, \sum_{t=1}^{T_{i}} x_{t,i}^{2} \leq \frac{T}{d} + 1) \end{split}$$

Note that the version of the chain rule with a random stopping time can be found in [LS20, Exercise 15.7]. We detail the step using inverse Pinsker's inequality ([Sas15]) as follows:

$$KL(y_{\theta}(x_t)||y_{\theta'}(x_t)) \leq \frac{2}{\min_{a \in \{\gamma^{-\frac{1}{\epsilon}}, 0, -1\}} \mathbb{P}[y_{\theta'}(x_t) = a]} \sup_{a} |\mathbb{P}[y_{\theta}(x_t) = a] - \mathbb{P}[y_{\theta'}(x_t) = a]|^2$$

$$\leq \frac{2}{\gamma^{\frac{1}{\epsilon}} \sqrt{d\Delta}} (\gamma^{\frac{1}{\epsilon}} 2\Delta x_{t,i})^2$$

$$\leq 24^{\frac{1}{\epsilon}} 8\sqrt{d}^{\frac{1}{\epsilon}-1} \Delta^{\frac{1}{\epsilon}+1} x_{t,i}^2. \qquad (\gamma = 24\sqrt{d}\Delta)$$

Using the above lower bound on $\mathbb{E}_{\theta}[U_i(1)]$, and setting $\Delta := 24^{\frac{-1}{1+\epsilon}} d^{\frac{3\epsilon-1}{2(1+\epsilon)}} (288T)^{\frac{-\epsilon}{1+\epsilon}}$ (noting $288 = (12\sqrt{2})^2$), we have the following:

$$\mathbb{E}_{\theta}[U_{i}(1)] + \mathbb{E}_{\theta'}[U_{i}(-1)] \geq \mathbb{E}_{\theta'}[U_{i}(1) + U_{i}(-1)] - 24^{\frac{1}{2\epsilon}}12\sqrt{2}\Delta^{\frac{1+\epsilon}{2\epsilon}}\sqrt{d}^{\frac{1-\epsilon}{2\epsilon}}\frac{T}{d}\sqrt{\frac{T}{d}}$$

$$= 2\mathbb{E}_{\theta'}\left[\frac{T_{i}}{d} + \sum_{t=1}^{T_{i}}x_{t,i}^{2}\right] - 24^{\frac{1}{2\epsilon}}12\sqrt{2}\Delta^{\frac{1+\epsilon}{2\epsilon}}\sqrt{d}^{\frac{1-\epsilon}{2\epsilon}}\frac{T}{d}\sqrt{\frac{T}{d}}$$

$$\geq \frac{2T}{d} - \frac{T}{d} = \frac{T}{d}. \qquad (T_{i} \geq 0, \text{ def. } T_{i}, \text{ choice of } \Delta)$$

Note also that $\Delta \leq \frac{1}{24\sqrt{d}}$ (as required earlier) since $T \geq d^2$. We now combine the preceding equation with our earlier lower bound on R_T . By averaging overall $\theta \in \{\pm \Delta\}^d$, we conclude that there exists some θ^* such that

$$\begin{split} R_T(\mathcal{A}, \theta^*) &\geq \frac{\Delta \sqrt{d}}{2} \frac{1}{2^d} \sum_{\theta \in \{-\Delta, \Delta\}^d} R_T(\mathcal{A}, \theta) \\ &\geq \frac{\Delta \sqrt{d}}{4} \sum_{i=1}^d \sum_{\theta_i \in \{-\Delta, \Delta\}} \mathbb{E}_{\theta}[U_i(\mathtt{sign}(\theta_i))]. \qquad (R_T \text{ bound and } \sum_{\{\theta_j\}_{j \neq i}} 1 = 2^{d-1}) \end{split}$$

$$\geq \frac{1}{4}T\sqrt{d}\Delta \qquad (\mathbb{E}_{\theta}[U_{i}(1)] + \mathbb{E}_{\theta'}[U_{i}(-1)] \geq \frac{T}{d})$$

$$\geq \frac{1}{4 \cdot 24 \cdot 12\sqrt{2}}d^{\frac{2\epsilon}{1+\epsilon}}T^{\frac{1}{1+\epsilon}}. \qquad (\text{choice of } \Delta, \epsilon \in [0,1])$$

C Extension to Kernel Bandits

C.1 Problem Setup

We consider an unknown reward function $f: \mathcal{A} \to \mathbb{R}$ lying in the reproducing kernel Hilbert space (RKHS) \mathcal{H} associated with a given kernel K, i.e., $f(x) = \langle f, K(x, \cdot) \rangle_K$. Similar to the linear bandit setting, we assume that $\max_{x \in \mathcal{A}} |f(x)| \leq 1$ and $||f||_K \leq b$ for some b > 0.

At each round t = 1, 2, ..., T, the learner chooses an action $x_t \in \mathcal{A} \subseteq [0, 1]^d$ and observes the reward

$$y_t = f(x_t) + \eta_t,$$

where η_t are independent noise terms that satisfy $\mathbb{E}[\eta_t] = 0$ and $\mathbb{E}\big[|\eta_t|^{1+\epsilon}\big] \leq \upsilon$ for some $\epsilon \in (0,1]$ and finite $\upsilon > 0$. Letting $x^\star \in \arg\max_{x \in [0,1]^d} f(x)$ be an optimal action, the cumulative expected regret after T rounds is

$$R_T = \sum_{t=1}^{T} (f(x^*) - f(x_t)).$$

Given (A, ϵ, v) , the objective is to design a policy for sequentially selecting the points (i.e., x_t for t = 1, ..., T) in order to minimize R_T . We focus on the Matérn kernel, defined as follows:

$$K_{\nu,l}(x,x') := \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\|x - x'\|_2 \sqrt{2\nu}}{l} \right)^{\nu} B_{\nu} \left(\frac{\|x - x'\|_2 \sqrt{2\nu}}{l} \right),$$

where Γ is the Gamma function, B_{ν} is the modified Bessel function, and (ν, l) are parameters corresponding to smoothness and lengthscale.

We focus on the case that \mathcal{A} is a *finite* subset of $[0,1]^d$, but it is well known (e.g., see [VBJ⁺21, Assumption 4]) that the resulting regret bounds extend to the continuous domain $[0,1]^d$ via a discretization argument with with $\log |\mathcal{A}| = O(\log T)$.

C.2 Proof of Corollary 4

We state a more precise version of Corollary 4 as follows.

Theorem 5. For any unknown reward function $f: A \to \mathbb{R}$ lying in the RKHS of the Matérn kernel with parameters (ν, l) , for some finite set $A \subseteq [0, 1]^d$, assuming that $\max_{x \in A} |f(x)| \le 1$ and $||f||_K \le b$ for some b > 0, we have

$$M^*(A, T^{\frac{-2\epsilon}{1+\epsilon}}, 1) \le CT^{\epsilon \cdot \frac{d}{2\nu+d}},$$

for some constant C, and Algorithm 1 achieves regret of

$$R_T(f, \mathcal{A}) \le \left(C_0'b + C_1'(1+v)^{\frac{1}{1+\epsilon}}\log(|\mathcal{A}|T\log^2 T)^{\frac{\epsilon}{1+\epsilon}}\right)T^{1-\frac{\epsilon}{1+\epsilon}\frac{2\nu}{2\nu+d}},$$

for some constants C'_0, C'_1 . Note that the constants may depend on the kernel parameters (ν, l) and the dimension d.

We now proceed with the proof. We first argue that Algorithm 1 and Theorem 3 can still be applied (with x replacing a and f(x) replacing $a^{T}\theta^{*}$) in the kernel setting. The reasoning is the same as the case $\epsilon = 1$ handled in [CJKS21], so we keep the details brief.

Recall that for any kernel K, there exists a (possibly infinite dimensional) feature map $\phi: \mathcal{A} \to \mathcal{H}$ such that $K(x,x') = \phi(x)^{\top}\phi(x')$. For any $\lambda \in \Delta_{\mathcal{A}}$, we define $k_{\lambda}(\cdot) \in \mathbb{R}^{|\mathcal{A}|}$ such that for $\psi \in \mathcal{H}$, $k_{\lambda}(\psi)_i := \sqrt{\lambda_i}\phi(x_i)^{\top}\psi$, and $K_{\lambda} \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{A}|}$ such that $(K_{\lambda})_{i,j} := \sqrt{\lambda_i}\sqrt{\lambda_j}K(x_i,x_j)$. Then similar to [CJKS21, Lemma 2], we have for any $\psi, \rho \in \mathcal{H}$ that

$$\psi^{\top} A^{(\gamma)}(\lambda)^{-1} \rho = \gamma^{-1} \psi^{\top} \rho - \gamma^{-1} k_{\lambda}(\psi) (K_{\lambda} + I_{|\mathcal{A}|})^{-1} k_{\lambda}(\rho).$$

Then the gradient for the experimental design problem $\inf_{\lambda \in \Delta_{\mathcal{V}}} \max_{v \in \mathcal{V}} \|\phi(v)\|_{A^{(\gamma)}(\lambda)^{-1}}$ (which is an upper bound for our experimental design objective $M_{1+\epsilon}(\lambda; \mathcal{V}, \gamma, 1)$ by the proof of Lemma 2) can be computed efficiently. Moreover, Theorem 3 still holds because the the kernel setup can be viewed as a linear setup in an infinite-dimensional feature space (after applying the feature map ϕ to the action set), and our analysis does not use the finiteness of the dimension.

Given Theorem 3, the main remaining step is to upper bound $M_{1+\epsilon}^*$. To do so, we use the well-known polynomial eigenvalue decay of the Matérn kernel. Specifically, the j-th eigenvalue φ_j satisfies $\varphi_j \leq \mathcal{O}(j^{-\kappa})$ with $\kappa = \frac{2\nu + d}{d}$ (e.g., see [VBJ⁺21]). We let $\lambda_D^* \in \arg\max_{\lambda \in \Delta_{\mathcal{A}}} \log\det\left(A^{(\gamma)}(\lambda)\right)$, and proceed as follows:

$$\begin{split} M_{1+\epsilon}^*(\mathcal{A},\gamma,1)^{\frac{2}{1+\epsilon}} &\leq \max_{\mathcal{V}\in\mathcal{A}}\inf_{\lambda\in\Delta_{\mathcal{V}}} 2^{\frac{2}{1+\epsilon}}\max_{v\in\mathcal{V}}\|\phi(v)\|_{A^{(\gamma)}(\lambda)^{-1}}^2 \qquad \text{(shown in the proof of Lemma 2)} \\ &\leq 4\mathrm{Tr}\left(A(\lambda_D^*)(A(\lambda_D^*)+\gamma I)^{-1}\right) \qquad \qquad \text{[CJKS21, Lemma 3]} \\ &= 4\mathrm{Tr}\left(K_{\lambda_D^*}(K_{\lambda_D^*}+\gamma I)^{-1}\right) \\ &= 4\sum_{j=1}^{|\mathcal{A}|}\frac{\varphi_j}{\varphi_j+\gamma} \qquad \qquad \text{(for some constant } c\geq 1 \text{ dependent on } l,\nu,d) \\ &\leq 4c\sum_{j\leq\gamma^{-\frac{1}{\kappa}}}\frac{j^{-\kappa}}{j^{-\kappa}+\gamma} + 4c\sum_{j>\gamma^{-\frac{1}{\kappa}}}\frac{j^{-\kappa}}{j^{-\kappa}+\gamma} \qquad \qquad \text{($c\geq 1$)} \\ &\leq 4c\gamma^{-1/\kappa} + 4c\sum_{j>\gamma^{-\frac{1}{\kappa}}}\frac{j^{-\kappa}}{\gamma} \qquad \qquad \text{(dropping terms in denominators)} \\ &\leq 4c\gamma^{-\frac{1}{\kappa}} + 4c(\gamma^{-\frac{1}{\kappa}})^{1-\kappa}\frac{1}{(\kappa-1)\gamma} \qquad \qquad \text{(bounding sum by integral; $\kappa>1$)} \\ &= 4c\gamma^{-\frac{1}{\kappa}}\left(1+\frac{1}{\kappa-1}\right) \\ &= 4c\frac{2\nu+d}{2\nu}T^{\frac{2\epsilon}{1+\epsilon}\frac{d}{2\nu+d}}. \qquad \qquad (\gamma=T^{\frac{-2\epsilon}{1+\epsilon}}\text{ and $\kappa=\frac{2\nu+d}{d}$)} \end{split}$$

Taking the square root on both sides gives $M_{1+\epsilon}^*(\mathcal{A},\gamma,1)^{\frac{1}{1+\epsilon}}=\widetilde{\mathcal{O}}\big(T^{\frac{\epsilon}{1+\epsilon}}\frac{d}{2\nu+d}\big)$, and multiplying by $\widetilde{\mathcal{O}}(T^{\frac{1}{1+\epsilon}})=\widetilde{\mathcal{O}}(T^{1-\frac{\epsilon}{1+\epsilon}})$ from the regret bound in Theorem 3 gives $\widetilde{\mathcal{O}}(T^{1-\frac{\epsilon}{1+\epsilon}}\cdot\frac{2\nu}{2\nu+d})$ regret as claimed in Corollary 4. By the same reasoning but keeping track of the logarithmic terms, we obtain the regret bound stated in Theorem 5.

C.3 Comparisons of Bounds

Comparison to existing lower bound. In Figure 2, we compare our regret upper bound to the lower bound of $\Omega(T^{\frac{\nu+d\epsilon}{\nu(1+\epsilon)+d\epsilon}})$ proved in [CG19]. We see that the upper and lower bounds coincide in certain limits and extreme cases:

- As $\nu/d \to \infty$, the regret approaches $T^{\frac{1}{1+\epsilon}}$ scaling, which matches the regret of linear heavy-tailed bandits in constant dimension.
- As $\nu/d \to 0$ and/or $\epsilon \to 0$, the regret approaches trivial linear scaling in T.
- When $\epsilon=1$, the regret scales as $\widetilde{\Theta}(T^{\frac{\nu+d}{2\nu+d}})$, which matches the optimal scaling for the sub-Gaussian noise setting [SBC17]. As we discussed earlier, this finite-variance setting was already handled in [CJKS21].

For finite ν/d and fixed $\epsilon \in (0,1)$, we observe from Figure 2 that gaps still remain between the upper and lower bounds, but they are typically small, especially when ν/d is not too small.

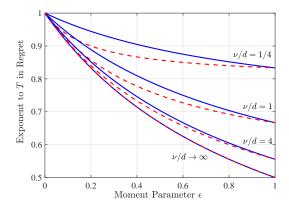


Figure 2: Comparison of our regret upper bound (solid) and the lower bound of [CG19] (dashed). We plot the exponent c such that the regret bound has dependence T^c , with the 4 pairs of curves corresponding to $\nu/d \in \{0.25, 1, 4\}$ and $\nu/d \to \infty$.

Comparison to existing upper bound. In [CG19], a regret upper bound of $\widetilde{\mathcal{O}}(\gamma_T T^{\frac{2+\epsilon}{2(1+\epsilon)}})$ was established, where γ_T is an *information gain* term that satisfies $\gamma_T = \widetilde{\mathcal{O}}(T^{\frac{d}{2\nu+d}})$ for the Matérn kernel [VKP21]. We did not plot this upper bound in Figure 2, because its high degree of suboptimality is easier to describe textually:

- For $\nu/d = 1/4$ and $\nu/d = 1$, their bound exceeds the trivial $\mathcal{O}(T)$ bound for all $\epsilon \in (0,1]$.
- For $\nu/d=4$, their bound still exceeds $\mathcal{O}(T)$ for $\epsilon \lesssim 0.28$, and is highly suboptimal for larger ϵ .
- As $\nu/d \to \infty$, the γ_T term becomes insignificant and their bound simplifies to $\widetilde{\mathcal{O}}(T^{\frac{2+\epsilon}{2(1+\epsilon)}})$, which is never better than $\widetilde{\mathcal{O}}(T^{3/4})$ (achieved when $\epsilon = 1$).
- A further weakness when $\epsilon=1$ is that the optimal γ_T dependence should be $\sqrt{\gamma_T}$ rather than linear in γ_T [SBC17, CJKS21].

For the *squared exponential kernel*, which has exponentially decaying eigenvalues rather than polynomial, these weaknesses were overcome in [CG19] using kernel approximation techniques, to obtain an optimal $\widetilde{\mathcal{O}}(T^{\frac{1}{1+\epsilon}})$ regret bound. Our main contribution above is to establish a new state of the art for the Matérn kernel, which is significantly more versatile in being able to model both highly smooth (high ν) and less smooth (small ν) functions.

D Numerical Experiments

In this section, we perform a simple proof-of-concept experiment to demonstrate that our algorithm can outperform existing methods as the ambient dimension increases. However, we emphasize that our main contributions are theoretical, and we leave detailed experimental studies for future work.

We conduct experiments with horizon T=100,000 and action set size N=2d. The true parameter is $\theta^\star=\frac{1}{\sqrt{d}}\mathbf{1}$ (so $\|\theta^\star\|_2=1$), and the action set is the subset of the normalized hypercube given by the signed coordinate directions $\mathcal{A}=\{\pm e_i\}_{i=1}^d$. Rewards follow $r_t=x_t^\top\theta^\star+\eta_t$ with heavy-tailed noise $\eta_t\sim \mathrm{ParetoII}(\alpha=2,\sigma=1)-\mathbb{E}[\mathrm{ParetoII}(\alpha=2,\sigma=1)]$ (centered to zero mean). The proposed algorithm is instantiated with input $\epsilon=0.5$ and evaluated against the Confidence Region with Truncated Mean (CRTM) Algorithm in [XWW+23]. Performance is measured via cumulative pseudo-regret $\sum_{t=1}^T \left(x_t^{\star\top}\theta^\star-x_t^{\top}\theta^\star\right)$, aggregated over 10 independent repetitions (with identical arm sets and independent noise).

Results. As shown in Figure 3, for all $d \ge 40$ Algorithm 1 achieves comparable or lower mean regret than CRTM, and notably, the gap widens as d increases. Both procedures remain sublinear in T in this controlled setting; however, the regret of Algorithm 1 grows more slowly with d, consistent with the guarantees in Theorem 3.

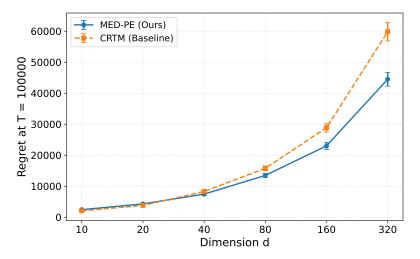


Figure 3: Regret vs dimension d with time horizon T = 100,000, and N = 2d arms

E Proof of Theorem 2 (Finite-Arm Lower Bound)

Consider $\log(\cdot)$ with base 2, and define m to be the smallest integer such that $\frac{m}{\log m} \geq \frac{d}{\log n}$. From the assumption $n \in [d, 2^{\lfloor \frac{d}{4} \rfloor}]$ we can readily verify that d > 4 and $m \in [4, d]$. For convenience, we assume that d is a multiple of m, since otherwise we can form the construction of the lower bound with $d' = d - (d \mod m)$ and pad the action vectors with zeros. Letting $d_i := (i-1)m$, we define the action set and the parameter set as follows for some Δ to be specified later:

$$\mathcal{A} := \left\{ a \in \{0, 1\}^d : \sum_{j=d_i+1}^{d_{i+1}} a_j = 1, \ \forall i \in [d/m] \right\}$$

$$\theta^* \in \Theta := \left\{ \theta \in \{\Delta, 2\Delta\}^d : \sum_{j=d_i+1}^{d_{i+1}} \theta_j = (m+1)\Delta, \ \forall i \in [d/m] \right\}.$$

In simple terms, the d-dimensional vectors are arranged in d/m groups of size m; each block in $a \in \mathcal{A}$ has a single entry of 1 (with 0 elsewhere), and each block in θ^* has a single entry of 2Δ (with Δ elsewhere). Observe that if $\Delta \leq \min(\frac{m}{4d}, \frac{1}{4\sqrt{d}})$, then $\|\theta^*\|_2 \leq 1$ and $x^\top \theta^* \leq 1$ as required.

Moreover, we have $|\mathcal{A}| = m^{\frac{d}{m}}$, and thus $\log |\mathcal{A}| = \frac{d}{m} \log m \le \log n$ by the definition of m.

Similar to Theorem 1, we let the reward distribution be

$$y(x) = \begin{cases} (\frac{1}{\gamma})^{\frac{1}{\epsilon}} & \text{w.p. } \gamma^{\frac{1}{\epsilon}} \theta^{\top} x \\ 0 & \text{w.p. } 1 - \gamma^{\frac{1}{\epsilon}} \theta^{\top} x \end{cases}$$

with $\gamma:=2\Delta\frac{d}{m}$. The choices of \mathcal{A} and Θ give $\theta^{\top}x\leq 2\Delta\frac{d}{m}$, so by the same reasoning as in Theorem 1, the $(1+\epsilon)$ -moment of the reward distribution is bounded by 1.

Let $\operatorname{ind}_i(x) := \operatorname{arg\,max}_{b \in [m]}(x_{d_i+b})$ for fixed $x \in \mathcal{A} \cup \Theta$, and define $T_{i,b} := |\{t : x_{t,d_i+b} = 1\}|$. Moreover, define t_U to be a random integer drawn uniformly from [T], which immediately implies that $\mathbb{P}_{\theta}[x_{t_U,d_i+b} = 1] = \frac{\mathbb{E}_{\theta}[T_{i,b}]}{T}$. Then,

$$R_T(\mathcal{A}, \theta) = \sum_{t=1}^{T} \sum_{i=1}^{d/m} \left(\Delta - \Delta \mathbb{I}\{ \text{ind}_i(\theta) = \text{ind}_i(x_t) \} \right)$$
$$= \Delta \sum_{i=1}^{d/m} \left(T - \mathbb{E}_{\theta} \left[T_{i, \text{ind}_i(\theta)} \right] \right)$$

$$= \Delta T \sum_{i=1}^{d/m} \left(1 - \mathbb{P}_{\theta}[x_{t_{\mathrm{U}}, d_i + \mathrm{ind}_i(\theta)} = 1] \right).$$

For fixed $\theta \in \Theta$ and $i \in [\frac{d}{m}]$, and any $b \in [m]$, we define $\theta^{(b)} \in \Theta$ to have entries given by $\theta_j^{(b)} = \begin{cases} \Delta + \Delta \mathbb{I}\{j = d_i + b\} & j \in [d_i + 1, d_{i+1}] \\ \theta_j & \text{otherwise} \end{cases}; \text{ and define the base parameter } \theta^{(0)} \text{ with entries otherwise}$ $\theta_j^{(0)} = \begin{cases} \Delta & j \in [d_i + 1, d_{i+1}] \\ \theta_j & \text{otherwise} \end{cases}. \text{ Note that } \theta^{(\text{ind}_i(\theta))} = \theta, \text{ and that the dependence of } \theta^{(b)} \text{ on } i \text{ is left implicit.}$

Then, for $b \in [m]$, we have

$$\begin{split} \mathbb{P}_{\theta^{(b)}}[x_{t,d_i+b} = 1] &\leq \mathbb{P}_{\theta^{(0)}}[x_{t,d_i+b} = 1] + \sqrt{\frac{1}{2} \text{KL}(\mathbb{P}_{\theta^{(0)}} \| \mathbb{P}_{\theta^{(b)}})} \qquad \text{(Pinsker's Inequality)} \\ &= \mathbb{P}_{\theta^{(0)}}[x_{t,d_i+b} = 1] + \sqrt{\frac{1}{2} \mathbb{E}_{\theta^{(0)}} \left[\sum_{t=1}^{T} \text{KL}\left(\text{Ber}(\gamma^{\frac{1}{\epsilon}} \theta^{(0)}^{\top} x_t) \| \text{Ber}(\gamma^{\frac{1}{\epsilon}} \theta^{(b)}^{\top} x_t) \right) \right]}. \end{split}$$

$$\text{(Chain rule)}$$

Similarly to the proof of Theorem 1, applying $\mathrm{KL}(\mathrm{Ber}(p)\|\mathrm{Ber}(q)) \leq \frac{(p-q)^2}{q(1-q)}$ along with $\Delta d/m \leq \theta^\top x \leq 2\Delta d/m$ and $|(\theta^{(0)} - \theta^{(b)})^\top x| \leq \Delta$ gives

$$\begin{split} \mathrm{KL}\left(\mathrm{Ber}(\gamma^{\frac{1}{\epsilon}}\theta^{(0)}{}^{\top}x_{t})\|\mathrm{Ber}(\gamma^{\frac{1}{\epsilon}}\theta^{(b)}{}^{\top}x_{t})\right) &\leq \frac{2(\gamma^{\frac{1}{\epsilon}}(\theta^{(0)}-\theta^{(b)})^{\top}x_{t})^{2}}{\gamma^{\frac{1}{\epsilon}}\theta^{(b)}{}^{\top}x_{t}} \\ &\leq \frac{2^{\frac{2+\epsilon}{\epsilon}}\Delta^{\frac{2}{\epsilon}}(\frac{d}{m})^{\frac{2}{\epsilon}}\Delta^{2}\mathbb{I}\{x_{t,d_{i}+b}=1\}}{2^{\frac{1+\epsilon}{\epsilon}}\Delta^{\frac{1+\epsilon}{\epsilon}}\left(\frac{d}{m}\right)^{\frac{1-\epsilon}{\epsilon}}} = 2^{\frac{1+\epsilon}{\epsilon}}\Delta^{\frac{1+\epsilon}{\epsilon}}\left(\frac{d}{m}\right)^{\frac{1-\epsilon}{\epsilon}}\mathbb{I}\{x_{t,d_{i}+b}=1\}. \end{split}$$

We set $\Delta:=\frac{1}{8}\left(\frac{d}{m}\right)^{\frac{\epsilon-1}{1+\epsilon}}\left(\frac{T}{m}\right)^{\frac{-\epsilon}{1+\epsilon}}.$ We claim that under this choice, the condition $T\geq 4^{\frac{1+\epsilon}{\epsilon}}d^{\frac{1+\epsilon}{\epsilon}}$ implies $\Delta\leq \min(\frac{m}{4d},\frac{1}{4\sqrt{d}})$, as we required earlier. To see this, we rewrite $\Delta=\frac{1}{8}d^{\frac{\epsilon-1}{1+\epsilon}}m^{\frac{1}{1+\epsilon}}T^{-\frac{\epsilon}{1+\epsilon}}$ and substitute the bound on T to obtain $\Delta\leq\frac{1}{32}d^{\frac{\epsilon-1}{1+\epsilon}}m^{\frac{1}{1+\epsilon}}d^{-1}.$ Dividing both sides by m gives $\frac{\Delta}{m}\leq\frac{1}{32d},$ whereas applying $m\leq d$ gives $\Delta\leq\frac{1}{32}d^{-\frac{1}{1+\epsilon}}\leq\frac{1}{32\sqrt{d}}.$

Combining the preceding two display equations and averaging over all $b \in m$, we have

$$\begin{split} \frac{1}{m} \sum_{b} \mathbb{P}_{\theta^{(b)}}[x_{t,d_i+b} = 1] &\leq \frac{1}{m} + \frac{1}{m} \sum_{b} \sqrt{2^{\frac{1+\epsilon}{\epsilon}} \Delta^{\frac{1+\epsilon}{\epsilon}} \left(\frac{d}{m}\right)^{\frac{1-\epsilon}{\epsilon}}} \mathbb{E}_{\theta^{(0)}}[T_{i,b}] \\ &\leq \frac{1}{m} + \sqrt{2^{\frac{1+\epsilon}{\epsilon}} \frac{1}{m} \Delta^{\frac{1+\epsilon}{\epsilon}} \left(\frac{d}{m}\right)^{\frac{1-\epsilon}{\epsilon}}} \sum_{b} \mathbb{E}_{\theta^{(0)}}[T_{i,b}] &\leq \frac{1}{m} + \frac{1}{2}. \\ & (\text{Jensen}, \sum_{b} T_{i,b} = T \text{ & choice of } \Delta) \end{split}$$

Averaging over all $\theta \in \Theta$, summing over $i \in [d/m]$, and recalling that $m \ge 4$, we obtain

$$\frac{1}{|\Theta|} \sum_{\theta \in \Theta} \sum_{i=1}^{d/m} \left(1 - \mathbb{P}_{\theta}[x_{t,d_i + \operatorname{ind}_i(\theta)} = 1] \right) \ge \frac{d}{m} \left(1 - \frac{1}{m} - \frac{1}{2} \right) \ge \frac{d}{4m}.$$

Hence, there exists $\theta^* \in \Theta$ such that $\sum_{i=1}^{d/m} \left(1 - \mathbb{P}_{\theta^*}[x_{t,d_i + \operatorname{ind}_i(\theta^*)} = 1]\right) \geq \frac{d}{4m}$. Substituting into our earlier lower bound on R_T and again using our choice of Δ , we obtain

$$R_T(\mathcal{A}, \theta^*) \ge \frac{d}{4m} \Delta T = \frac{1}{32} d^{\frac{\epsilon}{1+\epsilon}} \left(\frac{d}{m}\right)^{\frac{\epsilon}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}.$$

Since $f(x) = \frac{x}{\log x}$ is increasing for $x \ge e$, and $m \in [4, d]$, the definition of m gives the following:

$$\frac{d}{\log n} > \frac{m-1}{\log (m-1)} > \frac{m-1}{\log m} \ge \frac{m-1}{\log d}.$$

Rearranging the above, we obtain $\frac{d}{m} > \frac{\log n}{\log d} \left(1 - \frac{1}{m}\right) \ge \frac{\log n}{2 \log d}$, completing the proof.