

# On the Power of Adaptivity for $\varepsilon$ -Best Arm Identification in Linear Bandits

**Arnab Maiti**

University of Washington

ARNABM2@UW.EDU

**Yunbei Xu**

National University of Singapore

YUNBEI@NUS.EDU.SG

**Kevin Jamieson**

University of Washington

JAMIESON@CS.WASHINGTON.EDU

**Editors:** Steve Hanneke and Tor Lattimore

## Abstract

We study the minimax sample complexity of  $\varepsilon$ -best arm identification in linear bandits, a classical pure-exploration problem. Given a compact action set  $\mathcal{X}$  that spans  $\mathbb{R}^d$  and an unknown reward vector  $\theta \in \mathbb{R}^d$ , the goal is to output an arm  $\hat{x} \in \mathcal{X}$  such that  $\langle \hat{x}, \theta \rangle \geq \max_{x \in \mathcal{X}} \langle x, \theta \rangle - \varepsilon$  with probability at least  $1 - \delta$ , using as few samples as possible. Our aim is to better understand the power and limitations of adaptivity in this setting.

We begin with non-adaptive algorithms. We present a non-adaptive fixed-design method with sample complexity  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$ , where  $w(\mathcal{X})$  is a Gaussian width term dependent on  $\mathcal{X}$ , and we prove a matching lower bound  $\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$  for all non-adaptive fixed-design methods. Moreover,  $w(\mathcal{X}) \leq \mathcal{O}(d)$  for general  $\mathcal{X}$ , which is tight for sets such as the unit  $\ell_2$  ball, and  $w(\mathcal{X}) \leq \mathcal{O}(\sqrt{d \log |\mathcal{X}|})$  when  $\mathcal{X}$  is finite, which is tight for the canonical basis  $\{e_1, \dots, e_d\}$ .

We then turn to adaptive sampling. For any finite action set  $\mathcal{X}$ , we prove the existence of an adaptive algorithm with sample complexity  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{d \log(|\mathcal{X}|/d)}{\varepsilon^2}\right)$  via a generalization of Median Elimination, which is known to yield a  $\log d$  improvement for the canonical basis. This raises a structural question: beyond the canonical basis, are there structured action sets for which adaptivity yields only logarithmic-factor improvements over the optimal non-adaptive rate? We answer in the affirmative for several natural action sets, namely the hypercube, the  $\ell_2$  ball,  $m$ -sets, and multi-task multi-armed bandits.

Finally, we show that logarithmic improvements are not the whole story. To our knowledge, we provide the first construction of an action set  $\mathcal{X}$  for which adaptivity yields a *polynomial-factor improvement* over every non-adaptive algorithm. A key ingredient behind this separation is an  $\ell_2$ -norm estimation subroutine: we design an adaptive algorithm that uses  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples from the unit  $\ell_2$  ball in  $\mathbb{R}^d$  and outputs an estimate  $\hat{r}$  satisfying  $|\hat{r} - \|\theta\|_2| \leq \varepsilon$  with probability at least  $1 - \delta$ , where  $\theta$  is the unknown reward vector.

Taken together, these results illustrate when adaptivity can offer only modest savings and when it can enable genuine polynomial gains, sharpening our understanding of the role of adaptivity and geometry in pure exploration and experimental design.

**Keywords:** linear bandits,  $\varepsilon$ -best arm identification, adaptive vs non-adaptive designs, minimax sample complexity, polynomial gap, action set geometry, experimental design and pure exploration

## 1. Introduction

Experimental design is a classical and widely studied topic dating back to Fisher’s foundational work (Fisher, 1935), with applications spanning clinical trials, A/B testing, engineering, and scientific discovery. A major part of its appeal is that it allows one to commit in advance to a fixed set of measurements to collect, that is, a non-adaptive design; in many linear settings, there is a rich theory for constructing such designs (Pukelsheim, 2006). Such non-adaptive designs also offer practical advantages over adaptive approaches, including computational simplicity and the ability to parallelize or batch data collection.

Moreover, in certain linear settings, such non-adaptive approaches are known to be near-minimax optimal (Arias-Castro et al., 2012; Even-Dar et al., 2002; Fan et al., 2025). Motivated by these advantages, we revisit a fundamental sequential decision problem in a linear setting, namely pure exploration in linear bandits, and ask how much adaptivity can help beyond the best non-adaptive design. In particular, we study the minimax complexity of non-adaptive designs and ask whether the structure of the action set confines adaptive sampling to logarithmic gains, or allows polynomial improvements over any non-adaptive design.

In this paper, we study the  $\varepsilon$ -best arm identification problem (i.e., pure exploration) in linear bandits. We are given a compact arm set  $\mathcal{X} \subset \mathbb{R}^d$  with  $\text{span}(\mathcal{X}) = \mathbb{R}^d$  and an unknown reward vector  $\theta \in \mathbb{R}^d$ . At each round  $t = 1, 2, \dots$ , an algorithm selects an arm  $x_t \in \mathcal{X}$  and observes  $y_t = \langle x_t, \theta \rangle + \eta_t$ , where  $\eta_t \sim \mathcal{N}(0, 1)$ <sup>1</sup>. Given  $\varepsilon \in (0, 1)$  and  $\delta \in (0, 1)$ , the objective is to design an algorithm, consisting of a sampling rule and a stopping time, such that upon stopping it outputs  $\hat{x} \in \mathcal{X}$  satisfying  $\mathbb{P}(\max_{x \in \mathcal{X}} \langle x - \hat{x}, \theta \rangle \leq \varepsilon) \geq 1 - \delta$ . We call an algorithm non-adaptive if its sampling rule does not depend on the observed rewards  $y_t$ , and adaptive otherwise. Our goal is to characterize the power and limitations of adaptivity in this problem setting.

When  $\mathcal{X} = \{e_1, e_2, \dots, e_d\}$ , the problem reduces to the classical stochastic multi-armed bandit setting. A simple non-adaptive algorithm samples each arm  $\mathcal{O}\left(\frac{\log(d/\delta)}{\varepsilon^2}\right)$  times and outputs the arm with the largest empirical mean. A union bound shows that, with probability at least  $1 - \delta$ , the returned arm is  $\varepsilon$ -best. One way to extend this approach to linear bandits is via a fixed design. In optimal experimental design, a design distribution  $\lambda$  is a probability distribution over  $\mathcal{X}$ . Using a finitely supported  $\lambda$ , one can construct a fixed design, that is, a deterministic sequence of arms, in various ways. Typically, the number of times an arm  $x \in \mathcal{X}$  appears in the sequence is chosen proportional to  $\lambda(x)$ , though other constructions are also possible. The sampling rule then pulls the arms in this predetermined order and observes the corresponding rewards. Using these observations, one typically forms an unbiased estimator  $\hat{\theta}$  and outputs the empirically best arm. We refer to such procedures as non-adaptive fixed-design algorithms.

One can view the non-adaptive multi-armed bandit algorithm discussed above as a fixed-design method induced by the uniform design distribution over the arms. This procedure is already near-optimal for multi-armed bandits: adaptive algorithms such as Median Elimination (Even-Dar et al., 2002) achieve the optimal minimax sample complexity  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$ , improving over the optimal non-adaptive rate only by a logarithmic factor. This motivates the following question for the more general linear bandit setting:

*What is the minimax sample complexity of non-adaptive fixed-design algorithms for  $\varepsilon$ -best arm identification in linear bandits? Moreover, can adaptivity improve over the optimal non-adaptive rate only by logarithmic factors, as in multi-armed bandits, or are polynomial-factor improvements possible?*

## 1.1. Notations

Before we address the question raised above, we would like to mention a few notations. The unit  $\ell_2$  ball refers to the set  $\mathbb{B}_d := \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$ . If  $\lambda$  is a distribution over the set  $\mathcal{X}$ , we say  $\lambda \in \Delta_{\mathcal{X}}$  and will denote  $A(\lambda; \mathcal{X}) := \mathbb{E}_{x \sim \lambda} [xx^\top]$ . We define the gaussian width term  $w(\mathcal{X})$  depending on  $\mathcal{X}$  as

$$w(\mathcal{X}) := \inf_{\lambda \in \Delta_{\mathcal{X}}} \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in \mathcal{X}} \langle x, A(\lambda, \mathcal{X})^{-1/2} \eta \rangle \right]. \quad (1)$$

For any matrix  $W \in \mathbb{R}^{d \times d}$  and a vector  $x \in \mathbb{R}^d$ , we will denote  $\|x\|_W^2 := x^\top W x$ . For any positive scalar  $x$  define  $(x)_+ = \max\{1, x\}$  as the maximum of  $x$  and one, *not* zero. Define  $x_* = \arg \max_{x \in \mathcal{X}} \langle x, \theta \rangle$ .

In this paper, we focus on  $(\varepsilon, \delta)$ -PAC algorithms that, for any  $\varepsilon \in (0, 1]$  and  $\delta \in (0, 1)$ , use  $\frac{H_1 \log(1/\delta) + H_2}{\varepsilon^2}$  samples, where  $0 < H_1 \leq H_2$ , and output  $\hat{x} \in \mathcal{X}$  such that  $\mathbb{P}(\langle x_* - \hat{x}, \theta \rangle \leq \varepsilon) \geq 1 - \delta$ .

1. All our algorithmic results extend to independent mean-zero subgaussian noise with bounded variance proxy. We restrict attention to Gaussian noise for simplicity of presentation.

## 1.2. Our Contributions and Techniques

We now address the question posed above by summarizing our main contributions and techniques. We first present a non-adaptive algorithm based on a fixed design with sample complexity  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$ , where  $w(\mathcal{X})$  is the Gaussian width term defined by (1) depending on  $\mathcal{X}$ . Our fixed design has length  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$  and is constructed using two design distributions,  $\lambda_1$  and  $\lambda_2$ . Here,  $\lambda_1$  minimizes  $\max_{x \in \mathcal{X}} \|x\|_{A(\lambda; \mathcal{X})}^2$ , while  $\lambda_2$  minimizes  $\mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} [\max_{x \in \mathcal{X}} \langle x, A(\lambda; \mathcal{X})^{-1/2} \eta \rangle]$ . Given the resulting samples, we form the unbiased least-squares estimator  $\hat{\theta}$  and output an arm  $\hat{x} \in \arg \max_{x \in \mathcal{X}} \langle x, \hat{\theta} \rangle$ . Using the Borell-TIS inequality, we show that  $\hat{x}$  is an  $\varepsilon$ -best arm with probability at least  $1 - \delta$ .

We then prove a matching minimax lower bound of  $\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$  for any non-adaptive fixed-design algorithm. We first show a lower bound of  $\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  for any algorithm (possibly adaptive). For the  $\Omega\left(\frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$  lower bound, our proof reduces  $\varepsilon$ -best arm identification to bounding the minimax simple regret after  $T$  rounds. We then select a Gaussian prior over the reward vector  $\theta$  tailored to the fixed design  $(x_t)_{t \in [T]}$ , more precisely to the design matrix  $\sum_{t=1}^T x_t x_t^\top$ , and use it to derive a lower bound on the simple regret.

We can make our algorithm adaptive as follows. We partition the arm set into  $d$  regions and attempt to identify a candidate good arm from each region. We then run Median Elimination on these  $d$  candidate arms to obtain an  $\varepsilon$ -best arm. In certain cases, this yields a  $\log d$  improvement, since the analysis can be localized to the region containing the best arm and the Borell–TIS inequality is applied only over that region.

This observation also motivates us to ask whether there are structured action sets beyond the multi-armed bandit case for which adaptivity yields at most logarithmic factors improvement over our non-adaptive algorithm. Towards that we study several structured action sets of interest, namely the hypercube, the unit  $\ell_2$ -ball,  $m$ -sets, and multi-task multi-armed bandits (MAB), and characterize the corresponding Gaussian width terms  $w(\mathcal{X})$ . For each of these sets, we show, by establishing an adaptive lower bound, that adaptive algorithms can improve over our non-adaptive fixed-design sample complexity by at most logarithmic factors only. Our main technical contributions in this part are the lower bound arguments for  $m$ -sets and multi-task MAB. We begin by constructing a family of hard instances that is oblivious to the algorithm: each instance specifies a reward vector whose coordinates are chosen as  $\varepsilon$  times coordinate-specific scaling factor. We then show that, unless the algorithm collects sufficiently many samples, it must incur simple regret  $\Omega(\varepsilon)$ . To analyze the simple regret, we localize the KL-divergence-based argument. Concretely, we restrict attention to a structured subset of instances in which a portion of the reward vector is fixed and any two instances differ in exactly two coordinates within the remaining portion. This reduction enables a KL-divergence analysis closely analogous to the classical multi-armed bandit setting. Finally, a global averaging argument lifts the localized bound to the full family, yielding simple regret  $\Omega(\varepsilon)$  and, consequently, a minimax sample complexity lower bound.

The above result leads us to our main question: does adaptivity in linear bandits improve over non-adaptive approaches by only logarithmic factors? We answer this in the negative by constructing an action set  $\mathcal{X}$  for which adaptivity yields a *polynomial-factor* improvement over all non-adaptive algorithms. We take  $\mathcal{X}$  to be the union of  $k = \text{poly}(d)$  unit  $\ell_2$ -balls, each contained in a distinct  $d$ -dimensional subspace. A naive non-adaptive approach allocates on the order of  $d^2/\varepsilon^2$  samples to each ball to form an unbiased estimator and identify an  $\varepsilon$ -best arm, leading to total sample complexity  $kd^2/\varepsilon^2$ . In contrast, an adaptive strategy can save a factor of  $d$  by first identifying a ball that contains a near-optimal arm and then allocating additional samples to that ball to recover a near-optimal arm using just on the order of  $\frac{kd+d^2}{\varepsilon^2}$  samples.

A key step in identifying such a “good” unit ball is an  $\ell_2$ -norm estimation subroutine, which aims to estimate the best value over the unit ball, namely the  $\ell_2$ -norm of the reward vector. Specifically, we design

an adaptive algorithm that takes  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples from the unit  $\ell_2$ -ball in  $\mathbb{R}^d$ , observes noisy rewards of the form  $\langle x, \theta \rangle + \eta$  for an unknown vector  $\theta \in \mathbb{R}^d$ , and outputs an estimate  $\hat{r}$  such that, with probability at least  $1 - \delta$ ,  $\hat{r} \in [\|\theta\|_2 - \varepsilon, \|\theta\|_2 + \varepsilon]$ . A key technical component of our algorithm is to sample multiple Rademacher unit vectors uniformly at random and take multiple observations along each sampled direction. We then construct an  $\ell_2$ -norm estimator by squaring the empirical mean reward along each direction and aggregating these squared estimates. This squaring step leads to a subexponential-tail analysis, which is the main non-trivial part of our proof required to establish the stated sample complexity.

By constructing such a set  $\mathcal{X}$ , we highlight a key technical insight: accurately estimating the best value over a region  $\mathcal{Z} \subset \mathcal{X}$ , namely  $\max_{x \in \mathcal{Z}} \langle x, \theta \rangle$ , plays a crucial role in achieving improved performance. In the special case where  $\mathcal{Z}$  is a unit  $\ell_2$ -ball, this task reduces to  $\ell_2$ -norm estimation. More generally, this observation suggests that region-wise value estimation is a fundamental ingredient in the design of near-optimal policies in related settings.

### 1.3. Related Works

Our work is closely related to a line of research on instance-dependent sample complexity for best arm identification (Soare et al., 2014; Karnin, 2016; Xu et al., 2018; Tao et al., 2018; Fiez et al., 2019; Jedra and Proutiere, 2020; Katz-Samuels et al., 2020; Degenne et al., 2020). A key distinction between instance-dependent and minimax guarantees is the role of  $\theta$ . In the instance-dependent regime, one seeks sample complexity bounds of the form  $\mathbf{H}_1 \log(1/\delta) + \mathbf{H}_2$ , where  $\mathbf{H}_1$  and  $\mathbf{H}_2$  depend on the arm set  $\mathcal{X}$  and the reward vector  $\theta$ . In contrast, minimax characterizations aim for bounds of the same form with  $\mathbf{H}_1$  and  $\mathbf{H}_2$  depending on the arm set  $\mathcal{X}$  and the accuracy parameter  $\varepsilon$ , but not on  $\theta$ , since the guarantee is worst-case over  $\theta$ . While several instance-dependent algorithmic ideas extend to the minimax setting, lower bounds are less transferable since they typically hinge on simple-regret arguments, and minimax complexity characterizations remain relatively sparse (Even-Dar et al., 2002; Shamir, 2015; Chen et al., 2024).

One approach to designing algorithms in the instance-dependent regime is through experimental design. Early works such as Soare et al. (2014) and Karnin (2016) used G-optimal designs. Later, Fiez et al. (2019) proposed an algorithm based on sequential experimental design that is near-optimal for the first term  $\mathbf{H}_1$ , but can incur a large second term  $\mathbf{H}_2$ . This motivated Katz-Samuels et al. (2020) to design an algorithm whose design distribution minimizes certain Gaussian width terms, thereby reducing  $\mathbf{H}_2$ . Such a Gaussian-width based approach to experimental design is also used in a related regret minimization setting by Wagenmaker et al. (2021). For a detailed discussion of experimental design and its connection to linear bandits, we refer the reader to Lattimore and Szepesvári (2020).

Our work is also related to the literature on norm estimation and value estimation. Kong et al. (2020) studied optimal policy value estimation in contextual bandits. In a different direction, Cai and Low (2011) studied estimating  $\frac{1}{n} \sum_i |\theta_i|$  from an observation  $Y \sim \mathcal{N}(\theta, I_n)$ , and Collier et al. (2020) later generalized this to estimating  $\frac{1}{n} \sum_i |\theta_i|^\gamma$  for  $\gamma > 0$ . Han et al. (2020) studied  $L_r$ -norm estimation in Gaussian white noise models. More recently, Cleanthous et al. (2025) studied adaptive estimation of the  $L_2$ -norm of a probability density on  $\mathbb{R}^d$ .

## 2. Results on the Power of Adaptivity

In this section, we analyze the power of adaptivity for  $\varepsilon$ -best arm identification in linear bandits. In Section 2.1, we present a non-adaptive fixed design algorithm with sample complexity  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$ , where  $w(\mathcal{X})$  is a Gaussian width term depending on  $\mathcal{X}$ . We then prove a matching minimax lower bound of  $\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$  for non-adaptive fixed-design algorithms in Section 2.2. In Section 2.3, we study the Gaussian width term  $w(\mathcal{X})$  and characterize it for the structured sets of interest: the hypercube, the unit

$\ell_2$ -ball,  $m$ -sets, and multi-task multi-armed bandits. For each of these sets, we show in Section 2.4 that adaptivity improves upon the minimax sample complexity of our non-adaptive algorithm by at most logarithmic factors only. Finally, in Section 2.5, we construct a set  $\mathcal{X}$  for which an adaptive algorithm yields a polynomial-factor improvement over all non-adaptive algorithms.

## 2.1. Non-Adaptive Fixed Design Algorithm

In this section, we present a non-adaptive fixed design algorithm that uses  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$  samples and returns an  $\varepsilon$ -best arm with probability at least  $1 - \delta$ . Our analysis relies on several standard technical results, which are stated in Appendix B.1.

Let  $\lambda_1$  be a distribution over  $\mathcal{X}$  that minimizes the quantity  $\mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} [\max_{x \in \mathcal{X}} \langle x, A(\lambda; \mathcal{X})^{-1/2} \eta \rangle]$ . Similarly, let  $\lambda_2$  be a distribution over  $\mathcal{X}$  that minimizes  $\max_{x \in \mathcal{X}} \|x\|_{A(\lambda; \mathcal{X})^{-1}}^2$ . Assume that both  $\lambda_1$  and  $\lambda_2$  have finite support (such minimizers always exist due to Carathéodory's theorem). Finally, define  $\lambda_0$  to be the mixture distribution that samples from  $\lambda_1$  with probability  $1/2$  and from  $\lambda_2$  with probability  $1/2$ .

Since  $\frac{1}{2}A(\lambda_2; \mathcal{X}) \preceq A(\lambda_0; \mathcal{X})$ , it follows that for every  $x \in \mathcal{X}$ ,  $x^\top A(\lambda_0; \mathcal{X})^{-1} x \leq 2 x^\top A(\lambda_2; \mathcal{X})^{-1} x$ . Because  $\lambda_2$  is  $G$ -optimal, we obtain  $\|x\|_{A(\lambda_0; \mathcal{X})^{-1}}^2 \leq 2d$  for all  $x \in \mathcal{X}$ .

Similarly, since  $\frac{1}{2}A(\lambda_1; \mathcal{X}) \preceq A(\lambda_0; \mathcal{X})$ , the Sudakov–Fernique inequality yields

$$\mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in \mathcal{X}} \langle x, A(\lambda_0; \mathcal{X})^{-1/2} \eta \rangle \right] \leq \sqrt{2} \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in \mathcal{X}} \langle x, A(\lambda_1; \mathcal{X})^{-1/2} \eta \rangle \right] = \sqrt{2} w(\mathcal{X}).$$

Now consider a fixed design  $x_1, x_2, \dots, x_T \in \mathcal{X}$  such that  $\tau(A_T) \leq 2 \tau(A(\lambda_0; \mathcal{X}))$ , where

$$A_T := \frac{1}{T} \sum_{i=1}^T x_i x_i^\top \text{ and } \tau(A) := \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in \mathcal{X}} x^\top A^{-1/2} \eta \right]^2 + 2 \max_{x \in \mathcal{X}} \|x\|_{A^{-1}}^2 \log(2/\delta).$$

Such a fixed design exists for any  $T \geq 180d$  (Katz-Samuels et al., 2020; Allen-Zhu et al., 2021). By the bounds established above, it follows that  $\tau(A_T) \leq 4 w(\mathcal{X})^2 + 8d \log(2/\delta)$ .

Let  $y_t = \langle x_t, \theta \rangle + \eta_t$  denote the noisy rewards under the fixed design, where  $\eta_t \sim \mathcal{N}(0, 1)$  are i.i.d. Define the least-squares estimator  $\hat{\theta} = (\sum_{t=1}^T x_t x_t^\top)^{-1} \sum_{t=1}^T x_t y_t$ . Note that  $\hat{\theta}$  is distributionally equivalent to  $\theta + (\sum_{t=1}^T x_t x_t^\top)^{-1/2} \eta$ , where  $\eta \sim \mathcal{N}(0, I_d)$ . Consequently, applying the Borell–TIS inequality to the Gaussian process  $V_x = x^\top (\hat{\theta} - \theta)$ , we obtain the following with probability at least  $1 - \delta$ :

$$\left| \max_{x \in \mathcal{X}} \langle x, \hat{\theta} - \theta \rangle \right| \leq \frac{1}{\sqrt{T}} \cdot \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in \mathcal{X}} x^\top A_T^{-1/2} \eta \right] + \frac{1}{\sqrt{T}} \cdot \sqrt{2 \max_{x \in \mathcal{X}} \|x\|_{A_T^{-1}}^2 \log(2/\delta)} \leq \sqrt{\frac{2 \tau(A_T)}{T}}.$$

We select  $\hat{x} \in \arg \max_{x \in \mathcal{X}} \langle x, \hat{\theta} \rangle$  and output it as our candidate  $\varepsilon$ -best arm. If we set  $T = 360 \left( \frac{d}{\varepsilon^2} \log(2/\delta) + \frac{w(\mathcal{X})^2}{\varepsilon^2} \right)$ , then with probability at least  $1 - \delta$  we have  $\left| \max_{x \in \mathcal{X}} \langle x, \hat{\theta} - \theta \rangle \right| \leq \varepsilon/2$ . It follows that, with probability at least  $1 - \delta$ ,  $\langle \hat{x}, \theta \rangle \geq \max_{x \in \mathcal{X}} \langle x, \theta \rangle - \varepsilon$ .

Hence, we have the following theorem.

**Theorem 1** *There exists a non-adaptive fixed-design algorithm that first selects a deterministic sequence of arms  $x_1, x_2, \dots, x_T \in \mathcal{X}$  with  $T = 360 \left( \frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2} \right)$ , then observes the corresponding noisy rewards  $y_t = \langle x_t, \theta \rangle + \eta_t$ , where  $\eta_t \sim \mathcal{N}(0, 1)$  are i.i.d., and outputs an arm  $\hat{x} \in \mathcal{X}$  such that, with probability at least  $1 - \delta$ ,  $\max_{x \in \mathcal{X}} \langle x - \hat{x}, \theta \rangle \leq \varepsilon$ .*

## 2.2. Lower Bound for Non-Adaptive Fixed Design Algorithms

In this section, we prove a lower bound of  $\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$  on the sample complexity of any non-adaptive algorithm that commits to a fixed design  $x_1, x_2, \dots, x_T \in \mathcal{X}$ , then observes the corresponding noisy rewards  $y_1, y_2, \dots, y_T$ , and finally outputs an arm  $\hat{x} \in \mathcal{X}$  that is  $\varepsilon$ -best with probability at least  $1 - \delta$ .

We begin by stating the following theorem, which follows from a hypothesis-testing reduction and standard KL-divergence arguments. We refer the reader to Appendix B.2 for the complete proof.

**Theorem 2** *Any algorithm (possibly adaptive) that outputs  $\hat{x} \in \mathcal{X}$  satisfying  $\mathbb{P}(\max_{x \in \mathcal{X}} \langle x - \hat{x}, \theta \rangle \leq \varepsilon) \geq 1 - \delta$  for all  $\theta \in \mathbb{R}^d$  must, for some worst-case  $\theta$ , use  $\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples.*

We next present our proof idea for establishing a lower bound of  $\Omega\left(\frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$ . Towards that we look at a simple regret minimization problem. Consider a deterministic sequence  $(x_t)_{t=1}^T$  with  $x_t \in \mathcal{X}$ . Recall that we observe

$$y_t = \langle x_t, \theta \rangle + \eta_t, \quad \eta_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1), \quad t = 1, \dots, T.$$

Define the matrix

$$A := \sum_{t=1}^T x_t x_t^\top = X^\top X,$$

where  $X \in \mathbb{R}^{T \times d}$  has rows  $x_t^\top$ . Assume that  $A$  is invertible (the other case is considered in Appendix B.4). An algorithm  $\mathcal{A}$  observes the full sequence  $\mathcal{D} = \{(x_t, y_t)\}_{t=1}^T$  and outputs  $\hat{x}_{\mathcal{A}} \in \mathcal{X}$ . The *simple regret* is

$$r(\hat{x}_{\mathcal{A}}, \theta) := \max_{x \in \mathcal{X}} \langle x - \hat{x}_{\mathcal{A}}, \theta \rangle.$$

Fix a prior  $\theta \sim \mathcal{N}(0, \tau^2 A^{-1})$  independent of the noise  $\eta = (\eta_1, \dots, \eta_T) \sim \mathcal{N}(0, I_T)$  and (if the algorithm is randomized) an internal random seed  $U$ , also independent of  $(\theta, \eta)$ . Let  $y = X\theta + \eta$ ,  $\mathcal{D} = \{(x_t, y_t)\}_{t=1}^T$ , and let the (possibly randomized) non-interactive rule output  $\hat{x}_{\mathcal{A}} = \hat{x}_{\mathcal{A}}(\mathcal{D}, U) \in \mathcal{X}$ . We study the minimax expected regret

$$\mathfrak{R}(A; \mathcal{X}) := \inf_{\mathcal{A}} \mathbb{E}[r(\hat{x}_{\mathcal{A}}, \theta)],$$

where the infimum is over all (possibly randomized) non-interactive procedures  $\mathcal{A}$  using  $\mathcal{D}$ . Throughout this proof,  $\mathbb{E}[\cdot]$  denotes expectation over all randomness in  $(\theta, \eta)$  and in the algorithm (if randomized), unless stated otherwise.

We define the gaussian width term  $w(\mathcal{X}; A) := \mathbb{E}_{g \sim \mathcal{N}(0, I_d)} [\sup_{x \in \mathcal{X}} \langle x, A^{-1/2} g \rangle]$ . If we establish that  $\mathfrak{R}(A; \mathcal{X}) \geq c_0 w(\mathcal{X}; A)$  for some absolute constant  $c_0$ , then combining this with the inequality  $w(\mathcal{X}; A) \geq w(\mathcal{X})/\sqrt{T}$  and the straightforward calculations in Appendix B.3 yields the following theorem.

**Theorem 3** *Consider any non-adaptive fixed-design algorithm that first selects a deterministic sequence of arms  $x_1, x_2, \dots, x_T \in \mathcal{X}$ , then observes the corresponding noisy rewards  $y_t = \langle x_t, \theta \rangle + \eta_t$ , where  $\eta_t \sim \mathcal{N}(0, 1)$  are i.i.d., and outputs an arm  $\hat{x} \in \mathcal{X}$  such that  $\max_{x \in \mathcal{X}} \langle x - \hat{x}, \theta \rangle \leq \varepsilon$  with probability at least  $1 - \delta$ . Then  $T \geq \Omega\left(\frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$ .*

For the remainder of this section, we lower bound  $\mathfrak{R}(A; \mathcal{X})$ . Fix any (possibly randomized) non-interactive rule  $\hat{x}_{\mathcal{A}}(\mathcal{D}, U) \in \mathcal{X}$ . By the definition of simple regret,

$$\mathbb{E}[r(\hat{x}_{\mathcal{A}}, \theta)] = \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, \theta \rangle\right] - \mathbb{E}[\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), \theta \rangle].$$

Applying the tower rule and conditioning on  $(\mathcal{D}, U)$ , we have

$$\mathbb{E}[\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), \theta \rangle] = \mathbb{E}[\mathbb{E}[\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), \theta \rangle \mid \mathcal{D}, U]].$$

Since  $\hat{x}_{\mathcal{A}}(\mathcal{D}, U)$  is measurable with respect to  $(\mathcal{D}, U)$  and  $U$  is independent of  $\theta$ , we can pull it outside the inner expectation, which gives

$$\mathbb{E}[\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), \theta \rangle] = \mathbb{E}[\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), \mathbb{E}[\theta \mid \mathcal{D}] \rangle].$$

Under the Gaussian prior  $\theta \sim \mathcal{N}(0, \tau^2 A^{-1})$  and likelihood  $y \mid \theta \sim \mathcal{N}(X\theta, I_T)$ , the posterior mean is  $\mathbb{E}[\theta \mid \mathcal{D}] = \frac{\tau^2}{1+\tau^2} A^{-1} X^\top y$ . Substituting this expression yields

$$\mathbb{E}[r(\hat{x}_{\mathcal{A}}, \theta)] = \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, \theta \rangle\right] - \frac{\tau^2}{1+\tau^2} \mathbb{E}\left[\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), A^{-1} X^\top y \rangle\right].$$

Using  $y = X\theta + \eta$  and  $A = X^\top X$ , we obtain

$$\mathbb{E}[r(\hat{x}_{\mathcal{A}}, \theta)] = \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, \theta \rangle\right] - \frac{\tau^2}{1+\tau^2} \mathbb{E}\left[\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), \theta + A^{-1} X^\top \eta \rangle\right].$$

Next, for any vector  $v$ , we have  $\langle \hat{x}_{\mathcal{A}}(\mathcal{D}, U), v \rangle \leq \max_{x \in \mathcal{X}} \langle x, v \rangle$ . Applying this bound gives

$$\mathbb{E}[r(\hat{x}_{\mathcal{A}}, \theta)] \geq \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, \theta \rangle\right] - \frac{\tau^2}{1+\tau^2} \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, \theta + A^{-1} X^\top \eta \rangle\right].$$

As  $\max_{x \in \mathcal{X}} \langle x, \theta + A^{-1} X^\top \eta \rangle \leq \max_{x \in \mathcal{X}} \langle x, \theta \rangle + \max_{x \in \mathcal{X}} \langle x, A^{-1} X^\top \eta \rangle$ , we get

$$\mathbb{E}[r(\hat{x}_{\mathcal{A}}, \theta)] \geq \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, \theta \rangle\right] - \frac{\tau^2}{1+\tau^2} \left( \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, \theta \rangle\right] + \mathbb{E}\left[\max_{x \in \mathcal{X}} \langle x, A^{-1} X^\top \eta \rangle\right] \right).$$

Since  $\theta \sim \mathcal{N}(0, \tau^2 A^{-1})$ , we have  $\theta \stackrel{d}{=} \tau A^{-1/2} \xi$  for  $\xi \sim \mathcal{N}(0, I_d)$ , and thus  $\mathbb{E}[\max_{x \in \mathcal{X}} \langle x, \theta \rangle] = \tau \mathbb{E}_{\xi \sim \mathcal{N}(0, I_d)}[\max_{x \in \mathcal{X}} \langle x, A^{-1/2} \xi \rangle]$ . Similarly, since  $\eta \sim \mathcal{N}(0, I_T)$  and  $A = X^\top X$ , we have  $A^{-1} X^\top \eta \stackrel{d}{=} \mathcal{N}(0, A^{-1})$ , and hence  $A^{-1} X^\top \eta \stackrel{d}{=} A^{-1/2} \xi$  for  $\xi \sim \mathcal{N}(0, I_d)$ . Therefore,  $\mathbb{E}[\max_{x \in \mathcal{X}} \langle x, A^{-1} X^\top \eta \rangle] = \mathbb{E}_{\xi \sim \mathcal{N}(0, I_d)}[\max_{x \in \mathcal{X}} \langle x, A^{-1/2} \xi \rangle]$ . Substituting these two identities yields

$$\mathbb{E}[r(\hat{x}_{\mathcal{A}}, \theta)] \geq \frac{\tau(1-\tau)}{1+\tau^2} \mathbb{E}_{\xi \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in \mathcal{X}} \langle x, A^{-1/2} \xi \rangle \right].$$

Finally, the function  $\frac{\tau(1-\tau)}{1+\tau^2}$  is maximized at  $\tau = \sqrt{2} - 1$ , which yields a universal constant of approximately 0.207.

### 2.3. Properties of Gaussian Width

In this section, we study basic properties of the term  $w(\mathcal{X})$ . First, we present the following proposition.

**Proposition 4** *For any  $\mathcal{X} \subset \mathbb{R}^d$ , we have  $w(\mathcal{X}) \geq \Omega(\sqrt{d \log d})$  and  $w(\mathcal{X}) \leq O(d)$ . Moreover, if  $\mathcal{X}$  is finite, then  $w(\mathcal{X}) \leq O(\sqrt{d \log |\mathcal{X}|})$ .*

These bounds are tight for suitable choices of  $\mathcal{X}$ . For instance, in the classical multi-armed bandit setting,  $w(\mathcal{X}) = \Theta(\sqrt{d \log d})$  (as discussed in the introduction). In contrast, for structured sets such as the hypercubes  $\{-1, +1\}^d$  and  $\{0, 1\}^d$ , as well as the unit  $\ell_2$  ball in  $\mathbb{R}^d$ , one can verify that  $w(\mathcal{X}) = \Theta(d)$ . A more nuanced example is the class of  $m$ -sets, defined by  $\mathcal{X} := \{x \in \{0, 1\}^d : \|x\|_1 = m\}$ , for which one can show  $w(\mathcal{X}) \geq \Omega(\sqrt{md})$ ; this matches the finite-set upper bound of  $O(\sqrt{d \log |\mathcal{X}|})$  up to logarithmic factors. We provide formal proofs of these claims and Proposition 4 in Appendix F.

These observations raise a natural question: does there exist a set  $\mathcal{X} \subset \mathbb{R}^d$  of dimension  $\Theta(d)$  for which  $w(\mathcal{X})$  is polynomially smaller than what the upper bounds in Proposition 4 might suggest? We answer this question in the affirmative. To this end, we introduce the multi-task multi-armed bandit problem, a well-known generalization of the classical multi-armed bandit setting that has been studied in prior work (Cesa-Bianchi and Lugosi, 2012; Cohen et al., 2017; Maiti et al., 2025; Fan et al., 2025).

Informally, in the multi-task multi-armed bandit (MAB) problem, we are given  $m$  bandit problems, where the  $i$ -th problem contains  $d_i \geq 2$  arms. In each round, the learner selects one arm from every problem simultaneously and observes as feedback the sum of the corresponding rewards plus Gaussian noise.

We now formalize the multi-task MAB problem. Let  $d = \sum_{i=1}^m d_i$ , and define  $d_{1:i} = \sum_{j=1}^i d_j$  with the convention  $d_{1:0} = 0$ . The arm set  $\mathcal{X}$  is defined as follows:

$$\mathcal{X} = \left\{ x \in \{0, 1\}^d : \forall j \in [m] \sum_{i=d_{1:j-1}+1}^{d_{1:j}} x_i = 1 \right\}.$$

For this set  $\mathcal{X}$ , we show that its dimension is  $d - m + 1 = \Theta(d)$  and that its Gaussian width satisfies  $w(\mathcal{X}) \leq O(\sum_{i=1}^m \sqrt{d_i \log d_i})$ ; the proof appears in Appendix G. Now suppose that  $d_i = 2$  for all  $i \in [m-1]$  and  $d_m = m^2$ . Then  $\dim(\mathcal{X}) = \Theta(m^2)$  and  $\min\{d, \sqrt{d \log |\mathcal{X}|\}\} = \Theta(m^{3/2})$ . On the other hand, the Gaussian width can be upper bounded as:

$$w(\mathcal{X}) \leq O\left(\sum_{i=1}^m \sqrt{d_i \log d_i}\right) = O\left(m + \sqrt{m^2 \log m}\right) = O\left(m\sqrt{\log m}\right),$$

which is polynomially smaller than  $\Theta(m^{3/2})$ . We summarize the above discussion in the following theorem.

**Theorem 5** *There exists a finite set  $\mathcal{X} \subset \mathbb{R}^d$  such that  $\dim(\mathcal{X}) = \Theta(d)$  and  $w(\mathcal{X}) \leq O(\sqrt{d \log d})$ , while  $\sqrt{d \log |\mathcal{X}|} \geq \Omega(d^{3/4})$ .*

This example shows that  $w(\mathcal{X})$  can be genuinely non-trivial, and in particular need not scale as  $\Theta(d)$  nor as  $\Theta(\sqrt{d \log |\mathcal{X}|})$  even when  $\mathcal{X}$  is finite.

## 2.4. Structured Sets $\mathcal{X}$ for Which Adaptivity Yields Only Logarithmic-Factor Improvements

For classical multi-armed bandits, Median Elimination, an adaptive algorithm, improves upon the non-adaptive approach described in the introduction by a factor of  $\log d$ . In the same spirit, we can make our algorithm from Section 2.1 adaptive as follows. We partition the arm set into  $d$  regions and attempt to identify a candidate good arm from each region. We then run Median Elimination on these  $d$  candidate arms to obtain an  $\varepsilon$ -best arm. The analysis of this adaptive approach can be localized to the region containing the best arm, and the Borell–TIS inequality is applied only over that region. This can lead to a  $\log d$  improvement in certain regimes. In particular, for a finite set  $\mathcal{X} \subset \mathbb{R}^d$ , our adaptive algorithm satisfies the following guarantee.

**Theorem 6** *There exists an  $(\varepsilon, \delta)$ -PAC adaptive algorithm for  $\varepsilon$ -best arm identification with sample complexity  $O\left(\frac{d \log((|\mathcal{X}|/d)_+/\delta)}{\varepsilon^2}\right)$  for any finite set  $\mathcal{X} \subset \mathbb{R}^d$ .*

Table 1: Sample complexity bounds for different structured action sets for which the adaptivity only yields logarithmic improvement.

Set	Non-adaptive upper bound	Adaptive lower bound
unit $\ell_2$ ball, $\{-1, 1\}^d, \{0, 1\}^d$	$O\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{d^2}{\varepsilon^2}\right)$	$\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{d^2}{\varepsilon^2}\right)$
$m$ -sets ( $m \leq d/21$ )	$O\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{md \log(d/m)}{\varepsilon^2}\right)$	$\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{md}{\varepsilon^2}\right)$
Multi-task MAB	$O\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{\left(\sum_{j=1}^m \sqrt{d_j \log d_j}\right)^2}{\varepsilon^2}\right)$	$\Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{\left(\sum_{j=1}^m \sqrt{d_j}\right)^2}{\varepsilon^2}\right)$

Consequently, for those finite sets  $\mathcal{X} \subset \mathbb{R}^d$  for which a lower bound of  $\Omega\left(\frac{d \log(|\mathcal{X}|/\delta)}{\varepsilon^2}\right)$  is known to hold for  $(\varepsilon, \delta)$ -PAC non-adaptive algorithms, our adaptive approach improves the sample complexity to  $O\left(\frac{d \log(|\mathcal{X}|/d + \delta)}{\varepsilon^2}\right)$ . We refer the reader to Appendix C for formal details. This raises a natural question of whether there are structured settings beyond multi-armed bandits in which adaptivity yields at most logarithmic improvements, or even no improvement at all. We answer this question in the affirmative for the structured sets studied in the previous section, namely the unit  $\ell_2$ -ball, the  $\{-1, +1\}^d$  and  $\{0, 1\}^d$  hypercubes,  $m$ -sets, and multi-task MAB.

For these structured sets, the optimal minimax sample complexity takes the form  $\mathbf{H}_1 \log(1/\delta) + \mathbf{H}_2$ . Combining the results of Sections 2.1 and 2.2, we conclude that  $\mathbf{H}_1 = \Theta\left(\frac{d}{\varepsilon^2}\right)$ . The more interesting term is  $\mathbf{H}_2$ , which can dominate when  $\delta$  is an absolute constant. We now establish lower bounds on  $\mathbf{H}_2$  for the various structured sets; these bounds are also summarized in Table 1.

For the unit  $\ell_2$  ball, the simple-regret lower bound of Chen et al. (2024), which is based on a Gaussian-prior construction, together with the straightforward calculations in Appendix I, implies that  $\mathbf{H}_2 \geq \Omega\left(\frac{d^2}{\varepsilon^2}\right)$ .

For the finite sets  $\{-1, +1\}^d$  hypercube,  $\{0, 1\}^d$  hypercube,  $m$ -sets, and multi-task MAB, we prove lower bounds via a different approach, based on constructing a finite family of hard instances rather than using a continuous prior. To illustrate this approach, we provide high-level intuition for the hard-instance construction in the multi-task MAB setting. We defer the formal proofs of the lower bounds to Appendix H.

**Intuition:** It is well known that in the  $K$ -armed bandit problem, if the total number of samples satisfies  $T \leq \frac{cK}{\varepsilon^2}$  for a sufficiently small universal constant  $c$ , then no algorithm can identify an  $\varepsilon$ -best arm with constant success probability. A standard hard family consists of spiked instances indexed by  $i_\star \in [K]$ , where the unique optimal arm has mean  $\mu_{i_\star} = 10\varepsilon$  and all other arms have mean 0, together with an alternative instance in which all arms have mean 0. Under the uniform distribution over these hard instances, a KL-divergence based argument implies that when  $T \leq \frac{cK}{\varepsilon^2}$ , any deterministic algorithm outputs  $\hat{i}$  with  $\mathbb{E}[\mu_{\hat{i}}] \leq 7\varepsilon$ , and Markov's inequality yields constant-probability failure; Yao's minimax principle extends the conclusion to randomized algorithms.

For the multi-task setting with  $m$  MAB problems, where problem  $j$  has  $d_j$  arms, we take the cartesian-product hard family obtained by choosing one arm  $i_j$  per problem  $j$  and setting its mean to  $10\varepsilon_j$  while all other arms have mean 0, yielding  $\prod_{j=1}^m d_j$  spiked instances, with the corresponding alternative instances formed by zeroing out a single problem  $j$  while keeping the others fixed. Conditioning on the spiked choices in problems  $k \neq j$ , the  $j$ -th problem reduces to the standard  $d_j$ -armed problem, so if  $T \leq \frac{cd_j}{\varepsilon_j^2}$  then  $\mathbb{E}[\mu_{\hat{i}_j}] \leq 7\varepsilon_j$ . Summing over  $j$  gives  $\mathbb{E}\left[\sum_{j \in [m]} \mu_{\hat{i}_j}\right] \leq 7 \sum_{j \in [m]} \varepsilon_j$ , while choosing the spiked arm from each problem leads to a total value  $10 \sum_{j \in [m]} \varepsilon_j$ , and Markov's inequality again yields constant-probability

failure to achieve additive error  $\sum_{j \in [m]} \varepsilon_j$ , with extension to randomized algorithms by Yao's minimax principle.

Now set  $\varepsilon_j := \frac{\varepsilon \sqrt{d_j}}{\sum_{s \in [m]} \sqrt{d_s}}$ . Then  $\frac{d_j}{\varepsilon_j^2} = \frac{(\sum_{s \in [m]} \sqrt{d_s})^2}{\varepsilon^2}$  for all  $j \in [m]$ , and  $\sum_{j \in [m]} \varepsilon_j = \varepsilon \cdot \frac{\sum_{j \in [m]} \sqrt{d_j}}{\sum_{s \in [m]} \sqrt{d_s}} = \varepsilon$ . Hence, if an algorithm takes at most  $T \leq \frac{c(\sum_{s \in [m]} \sqrt{d_s})^2}{\varepsilon^2}$  samples for a sufficiently small constant  $c$ , then it fails to identify an  $\varepsilon$ -best arm with high constant probability.

## 2.5. A Structured Set $\mathcal{X}$ for Which Adaptivity Yields a Polynomial-Factor Improvement

In the previous section, we showed that for several well-known structured sets, adaptivity can yield only logarithmic-factor improvements. This raises an important question: can adaptivity lead to polynomial-factor improvements? We answer this question in the affirmative by constructing an action set  $\mathcal{X}$  for which adaptivity yields a polynomial-factor improvement.

Consider positive integers  $k, d$  with  $d \leq k \leq d^2$ . For each  $i \in [k]$ , define

$$\mathcal{X}_i := \left\{ x \in \mathbb{R}^{kd} : \text{supp}(x) \subseteq \{(i-1)d+1, \dots, id\}, \|x\|_2 \leq 1 \right\},$$

and let  $\mathcal{X} := \bigcup_{i=1}^k \mathcal{X}_i$ . We first claim that, for the set  $\mathcal{X}$ , any non-adaptive algorithm has minimax sample complexity at least  $\Omega\left(\frac{kd \log(1/\delta)}{\varepsilon^2} + \frac{kd^2}{\varepsilon^2}\right)$  for the  $\varepsilon$ -best arm identification problem. The term  $\Omega\left(\frac{kd \log(1/\delta)}{\varepsilon^2}\right)$  follows directly from Theorem 3 together with the fact that  $\text{span}(\mathcal{X})$  has dimension  $kd$ . The  $\Omega\left(\frac{kd^2}{\varepsilon^2}\right)$  term arises from the block structure of  $\mathcal{X}$ . Each set  $\mathcal{X}_i$  is a unit  $\ell_2$ -ball in a  $d$ -dimensional coordinate subspace, and a non-adaptive algorithm does not know in advance which block contains the optimal arm. As a result, the sampling budget must be spread across all  $k$  blocks, requiring on the order of  $\Omega\left(\frac{d^2}{\varepsilon^2}\right)$  samples per block to ensure that, regardless of which  $\mathcal{X}_i$  contains the best arm, the algorithm can identify an  $\varepsilon$ -best arm within that block upon termination. We prove this claim formally in Appendix J.

We now discuss how an adaptive approach can improve upon the above non-adaptive lower bound. A natural strategy is to first identify an index  $i$  such that the best arm in  $\mathcal{X}_i$  is  $\varepsilon/2$ -close to the best arm in  $\mathcal{X}$ , and then find an  $\varepsilon/2$ -best arm within  $\mathcal{X}_i$ . To identify such an index  $i$ , we can proceed as follows: for each  $j \in [k]$ , estimate the quantity  $\max_{x \in \mathcal{X}_j} \langle x, \theta \rangle$  to within additive error  $\varepsilon/4$ , and output the index with the largest estimated value. After selecting this “good” set  $\mathcal{X}_i$ , we allocate additional samples to  $\mathcal{X}_i$  in order to identify an  $\varepsilon/2$ -best arm in that set.

Let  $\theta^{(i)} := (\theta_{(i-1)d+1}, \theta_{(i-1)d+2}, \dots, \theta_{id})^\top$ . Then  $\max_{x \in \mathcal{X}_i} \langle x, \theta \rangle = \|\theta^{(i)}\|_2$ . This motivates us to consider the following  $\ell_2$ -norm estimation problem: given an arm set  $\mathbb{B}_d := \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$ , estimate the  $\ell_2$ -norm of the reward vector using as few samples as possible while observing the corresponding noisy rewards.

Assume that we can solve this  $\ell_2$ -norm estimation problem with an algorithm whose minimax sample complexity is  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$ . Then, using  $\mathcal{O}\left(\frac{kd \log(k/\delta)}{\varepsilon^2}\right)$  samples and a union bound, we can find an index  $i \in [k]$  such that, with probability at least  $1 - \delta/2$ ,  $\max_{x \in \mathcal{X}_i} \langle x, \theta \rangle \geq \max_{x \in \mathcal{X}} \langle x, \theta \rangle - \varepsilon/2$ . Next, using our non-adaptive algorithm, we can find an  $\varepsilon/2$ -best arm in  $\mathcal{X}_i$  with probability at least  $1 - \delta/2$  using an additional  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{d^2}{\varepsilon^2}\right)$  samples. Since  $k \geq d$ , this yields an  $\varepsilon$ -best arm in  $\mathcal{X}$  with probability at least  $1 - \delta$  using at most  $\mathcal{O}\left(\frac{kd \log(1/\delta)}{\varepsilon^2} + \frac{kd \log k}{\varepsilon^2}\right)$  samples, yielding a polynomial improvement over all non-adaptive algorithms. We summarize the above discussion in the following theorem.

**Theorem 7** *Fix positive integers  $k$  and  $d$  satisfying  $d \leq k \leq d^2$ . There exists a set  $\mathcal{X} \subset \mathbb{R}^{kd}$  such that any non-adaptive  $(\varepsilon, \delta)$ -PAC algorithm for  $\varepsilon$ -best arm identification requires  $\Omega\left(\frac{kd \log(1/\delta)}{\varepsilon^2} + \frac{kd^2}{\varepsilon^2}\right)$  samples,*

while there exists an adaptive  $(\varepsilon, \delta)$ -PAC algorithm with sample complexity  $O\left(\frac{kd \log(1/\delta)}{\varepsilon^2} + \frac{kd \log k}{\varepsilon^2}\right)$ , yielding a polynomial improvement over all non-adaptive approaches.

The only remaining task is to design an algorithm that solves the  $\ell_2$ -norm estimation problem using at most  $O\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples, which we do in the next section.

### 3. An Algorithm for $\ell_2$ -Norm Estimation

In this section, we address the  $\ell_2$ -norm estimation problem, a key ingredient underlying the polynomial gap between adaptive and non-adaptive algorithms in Theorem 7. Recall that we are given the arm set  $\mathbb{B}_d := \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$  and wish to estimate the  $\ell_2$ -norm of an unknown reward vector  $\theta \in \mathbb{R}^d$  using as few samples  $x_t \in \mathbb{B}_d$  as possible, while observing noisy rewards  $y_t = \langle x_t, \theta \rangle + \eta_t$ , where  $\eta_t \sim \mathcal{N}(0, 1)$ . We solve this problem in three steps. First, we test whether  $\lambda_0 \varepsilon \leq \|\theta\|_2 \leq \lambda_1 \sqrt{d}$  for some positive constants  $\lambda_0$  and  $\lambda_1$ . If so, we obtain a constant-factor (within 2) multiplicative estimate of  $\|\theta\|_2$ . Finally, we use this coarse estimate to refine our estimate of  $\|\theta\|_2$  to within an additive error of  $\varepsilon$ . The final step is the most technically involved part of the analysis, and we outline its high-level ideas here. We formalize the guarantee of the  $\ell_2$ -norm estimation procedure in the following theorem, and defer the complete description of the procedure and its analysis to Appendix D.

**Theorem 8** *Let  $\theta \in \mathbb{R}^d$  be an unknown reward vector. At each round, a learner selects an action  $x_t \in \mathbb{B}_d$  and observes  $y_t = \langle x_t, \theta \rangle + \eta_t$ , where  $\eta_t \sim \mathcal{N}(0, 1)$  are i.i.d. There exists an adaptive algorithm that uses  $O\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  such observations and returns an estimate  $\hat{r}$  such that  $|\hat{r} - r| \leq \varepsilon$  with probability at least  $1 - \delta$ , where  $r := \|\theta\|_2$ .*

Begin by defining  $r := \|\theta\|_2$  and assuming that we are given a value  $\varepsilon < r_0 < 2\sqrt{d}$  such that  $\frac{r}{2} < r_0 \leq 2r$ . We now describe an algorithm that takes  $O\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples from  $\mathbb{B}_d$  and outputs an estimate  $\hat{r}$  such that with probability  $1 - \delta/2$ , we have  $|\hat{r} - r| \leq \varepsilon$ .

---

#### Algorithm 1: $\varepsilon$ -additive error $\ell_2$ -norm estimation algorithm

---

**Input:**  $\delta \in (0, 1), \varepsilon > 0, s = \frac{c_0 d}{r_0^2}, K = c_1 r_0^2 \cdot \varepsilon^{-2} \log(4/\delta)$  where  $c_0, c_1$  are large absolute constants.

1 **for**  $k = 1, \dots, K$  **do**

2     Draw a Rademacher unit vector  $x^{(k)} = \frac{(\varepsilon_1, \dots, \varepsilon_d)}{\sqrt{d}}$  where  $\varepsilon_i \stackrel{iid}{\sim} \{\pm 1\}$ .

3     Take  $s$  samples along this direction and observe:

4      $y_{k,\ell} = \mu_k + \eta_{k,\ell}$ , where  $\mu_k = \langle x^{(k)}, \theta \rangle$ ,  $\eta_{k,\ell} \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ ,  $\ell = 1, \dots, s$ .

4     Let  $\bar{y}_k := \frac{1}{s} \sum_{\ell=1}^s y_{k,\ell}$  and define  $Z_k := d \left( \bar{y}_k^2 - \frac{1}{s} \right)$ .

5 **end**

6 Define  $\bar{Z} := \frac{1}{K} \sum_{k=1}^K Z_k$  and output  $\hat{r} := \begin{cases} 0 & \text{if } \bar{Z} < 0 \\ \sqrt{\bar{Z}} & \text{otherwise} \end{cases}$

---

We now begin with some high-level analysis. Recall  $x^{(k)} = \frac{(\varepsilon_1, \dots, \varepsilon_d)}{\sqrt{d}}$  where  $\varepsilon_i \stackrel{iid}{\sim} \{\pm 1\}$  and  $\mu_k = \langle x^{(k)}, \theta \rangle$ . Then we have the following:

$$\mathbb{E}[\mu_k^2] = \theta^\top \mathbb{E} \left[ x^{(k)} (x^{(k)})^\top \right] \theta = \theta^\top \left( \frac{1}{d} I_d \right) \theta = \frac{r^2}{d}.$$

Now conditioning on  $x^{(k)}$  (hence  $\mu_k$ ), we have

$$\bar{y}_k \sim \mathcal{N} \left( \mu_k, \frac{1}{s} \right) \Rightarrow \mathbb{E}[\bar{y}_k^2 \mid \mu_k] = \mu_k^2 + \frac{1}{s} \Rightarrow \mathbb{E}[Z_k \mid \mu_k] = d\mu_k^2.$$

We will now aim to bound

$$\bar{Z} - r^2 = \underbrace{\left( \frac{1}{K} \sum_{k=1}^K d\mu_k^2 - r^2 \right)}_{\text{Term 1}} + \underbrace{\left( \frac{1}{K} \sum_{k=1}^K (Z_k - d\mu_k^2) \right)}_{\text{Term 2}}.$$

**Bounding Term 1:** Let  $X_k := d\mu_k^2 - r^2$ . With the help of Hanson–Wright inequality, we can show that  $X_k$  is sub-exponential and show the following for some absolute constant  $c$ :

$$\Pr \left( \left| \frac{1}{K} \sum_{k=1}^K X_k \right| > t \right) \leq 2 \exp \left( -c \cdot K \min \left( \frac{t^2}{r^4}, \frac{t}{r^2} \right) \right). \quad (2)$$

Choosing  $t = \frac{r\varepsilon}{4}$ , and  $K = c_1 r_0^2 \varepsilon^{-2} \log \frac{4}{\delta}$  for some large constant  $c_1$ , we get:

$$\left| \frac{1}{K} \sum_{k=1}^K d\mu_k^2 - r^2 \right| \leq \frac{r\varepsilon}{4} \quad \text{with probability} \geq 1 - \frac{\delta}{4}.$$

**Bounding Term 2:** Let  $W_k := Z_k - d\mu_k^2 = d(\bar{y}_k^2 - \frac{1}{s} - \mu_k^2)$ . Conditioning on  $\mu_k$ , we have:

$$\sqrt{d} \cdot \bar{y}_k \mid \mu_k \sim \mathcal{N}(\sqrt{d} \cdot \mu_k, \sigma^2), \quad \text{with } \sigma^2 = \frac{d}{s}.$$

Hence  $W_k \mid \mu_k$  is sub-exponential and therefore conditioning on  $\mu_{1:K} := \mu_1, \dots, \mu_K$ , we get:

$$\Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \mid \mu_{1:K} \right) \leq 2 \exp \left( -cK \min \left( \frac{t^2}{\bar{V}}, \frac{t}{b} \right) \right).$$

where  $\bar{V} := 8 \left( \frac{d^2}{s^2} + \frac{d^2}{s} \cdot \frac{1}{K} \sum_{k=1}^K \mu_k^2 \right)$  and  $b = 4d/s$ . Using Eq. (2), we have  $\bar{V} \leq \tau := 8 \left( \frac{d^2}{s^2} + \frac{3dr^2}{2s} \right)$  with probability at least  $1 - \frac{\delta}{8}$ . Hence, substituting the values of  $\tau$  and  $b$ , choosing  $t = \frac{r\varepsilon}{4}$  and  $K = c_1 r_0^2 \varepsilon^{-2} \log \frac{4}{\delta}$  for a sufficiently large constant  $c_1$ , and using  $r/2 < r_0 \leq 2r$ , we obtain:

$$\Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \right) \leq \frac{\delta}{8} + \Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \mid \bar{V} \leq \tau \right) \leq \frac{\delta}{8} + 2 \exp \left( -cK \min \left( \frac{t^2}{\tau}, \frac{t}{b} \right) \right) \leq \frac{\delta}{4}.$$

**Final guarantee:** Combining the above results, we have, with probability at least  $1 - \frac{\delta}{2}$ ,  $|\bar{Z} - r^2| \leq \frac{r\varepsilon}{2}$ . Let us now assume that  $|\bar{Z} - r^2| \leq \frac{r\varepsilon}{2}$ . As  $\bar{Z} \geq 0$  we have  $\hat{r} = \sqrt{\bar{Z}}$ , which implies the following:

$$|\hat{r} - r| = \left| \sqrt{\bar{Z}} - r \right| = \left| \bar{Z} - r^2 \right| / \left( \sqrt{\bar{Z}} + r \right) \leq \frac{\frac{r\varepsilon}{2}}{r} = \frac{\varepsilon}{2} < \varepsilon.$$

#### 4. Conclusion and Future Works

In this paper, we studied the power of adaptivity for pure exploration in linear bandits. We established matching upper and lower bounds on the minimax sample complexity of non-adaptive fixed-design algorithms. We then identified structured action sets for which adaptivity yields at most a logarithmic improvement over our non-adaptive approach, and we constructed a structured action set for which adaptive sampling attains a polynomial improvement over any non-adaptive approach. To obtain the polynomial separation, we develop an  $\ell_2$ -norm estimation procedure that uses  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples. As a consequence, we also obtain

a polynomial separation between the sample complexity of identifying an  $\varepsilon$ -best arm and that of estimating the optimal value  $\max_{x \in \mathcal{X}} \langle x, \theta \rangle$  when  $\mathcal{X}$  is the unit  $\ell_2$  ball.

Our work raises several open questions. Is there an intrinsic and interpretable property of  $\mathcal{X}$  that determines whether adaptivity can yield at most logarithmic improvements or instead enables polynomial improvements over non-adaptive approaches? Is the minimax sample complexity of estimating  $\max_{x \in \mathcal{X}} \langle x, \theta \rangle$  always  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  for all action sets  $\mathcal{X} \subset \mathbb{R}^d$ ? Can our value-estimation-based adaptive approach be extended to other structured sets to obtain polynomial improvements, and can our adaptive lower-bound techniques be generalized beyond the classes considered in this paper?

Finally, we hope these results help connect various lines of work: for experimental design, they suggest that pure exploration in linear bandits *can* exhibit polynomial advantages from adaptivity under suitable action-set geometry; for reinforcement learning through a PAC-learning lens, they suggest that the structure of the feature mapping in linear function approximation (see, e.g., [Bradtke and Barto \(1996\)](#); [Jin et al. \(2020\)](#)) may fundamentally determine when adaptive exploration across episodes is beneficial.

## Acknowledgements

KJ and AM were supported in part by NSF 2141511, 2023239, and a Singapore AI Visiting Professorship award. YX was supported by NUS A-0010008-00-00.

## References

- Zeyuan Allen-Zhu, Yuanzhi Li, Aarti Singh, and Yining Wang. Near-optimal discrete optimization for experimental design: A regret minimization approach. *Mathematical Programming*, 186(1):439–478, 2021.
- Ery Arias-Castro, Emmanuel J Candes, and Mark A Davenport. On the fundamental limits of adaptive sensing. *IEEE Transactions on Information Theory*, 59(1):472–481, 2012.
- Keith Ball et al. An elementary introduction to modern convex geometry. *Flavors of geometry*, 31(1-58): 26, 1997.
- Steven J Bradtke and Andrew G Barto. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1):33–57, 1996.
- T Tony Cai and Mark G Low. Testing composite hypotheses, hermite polynomials and optimal estimation of a nonsmooth functional. *The Annals of Statistics*, 39(244):1012–1041, 2011.
- Nicolo Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.
- Fan Chen, Dylan J Foster, Yanjun Han, Jian Qian, Alexander Rakhlin, and Yunbei Xu. Assouad, fano, and le cam with interaction: A unifying lower bound framework and characterization for bandit learnability. *Advances in Neural Information Processing Systems*, 37:75585–75641, 2024.
- G Cleanthous, AG Georgiadis, and OV Lepski. Adaptive estimation of the  $\ell_2$ -norm of a probability density and related topics ii. upper bounds via the oracle approach. *The Annals of Statistics*, 53(3):1275–1297, 2025.
- Alon Cohen, Tamir Hazan, and Tomer Koren. Tight bounds for bandit combinatorial optimization. In *Conference on Learning Theory*, pages 629–642. PMLR, 2017.

- Olivier Collier, Laëtítia Comminges, and Alexandre B Tsybakov. On estimation of nonsmooth functionals of sparse normal means. *Bernoulli*, 2020.
- Rémy Degenne, Pierre Ménard, Xuedong Shang, and Michal Valko. Gamification of pure exploration for linear bandits. In *International Conference on Machine Learning*, pages 2432–2442. PMLR, 2020.
- Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Pac bounds for multi-armed bandit and markov decision processes. In *International Conference on Computational Learning Theory*, pages 255–270. Springer, 2002.
- Zhiyuan Fan, Arnab Maiti, Kevin Jamieson, Lillian J Ratliff, and Gabriele Farina. On the universal near optimality of hedge in combinatorial settings. *arXiv preprint arXiv:2510.17099*, 2025.
- Tanner Fiez, Lalit Jain, Kevin G Jamieson, and Lillian Ratliff. Sequential experimental design for transductive linear bandits. *Advances in neural information processing systems*, 32, 2019.
- Ronald A. Fisher. *The Design of Experiments*. 1935.
- John A Gubner. The gamma function and stirling’s formula, 2021.
- YanJun Han, Jiantao Jiao, and Rajarshi Mukherjee. On estimation of  $l_r$ -norms in gaussian white noise models. *Probability Theory and Related Fields*, 177(3):1243–1294, 2020.
- Yassir Jedra and Alexandre Proutiere. Optimal best-arm identification in linear bandits. *Advances in Neural Information Processing Systems*, 33:10007–10017, 2020.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on learning theory*, pages 2137–2143. PMLR, 2020.
- Zohar S Karnin. Verification based solution for structured mab problems. *Advances in Neural Information Processing Systems*, 29, 2016.
- Julian Katz-Samuels, Lalit Jain, Kevin G Jamieson, et al. An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits. *Advances in Neural Information Processing Systems*, 33:10371–10382, 2020.
- Weihao Kong, Emma Brunskill, and Gregory Valiant. Sublinear optimal policy value estimation in contextual bandits. In *International conference on artificial intelligence and statistics*, pages 4377–4387. PMLR, 2020.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Arnab Maiti, Zhiyuan Fan, Kevin Jamieson, Lillian J Ratliff, and Gabriele Farina. Efficient near-optimal algorithm for online shortest paths in directed acyclic graphs with bandit feedback against adaptive adversaries. *arXiv preprint arXiv:2504.00461*, 2025.
- Friedrich Pukelsheim. *Optimal design of experiments*. SIAM, 2006.
- Ohad Shamir. On the complexity of bandit linear optimization. In *Conference on Learning Theory*, pages 1523–1551. PMLR, 2015.
- Marta Soare, Alessandro Lazaric, and Rémi Munos. Best-arm identification in linear bandits. *Advances in neural information processing systems*, 27, 2014.

Chao Tao, Saúl Blanco, and Yuan Zhou. Best arm identification in linear bandits with linear dimension dependency. In *International Conference on Machine Learning*, pages 4877–4886. PMLR, 2018.

Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.

Andrew Wagenmaker, Julian Katz-Samuels, and Kevin Jamieson. Experimental design for regret minimization in linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 3088–3096. PMLR, 2021.

Liyuan Xu, Junya Honda, and Masashi Sugiyama. A fully adaptive algorithm for pure exploration in linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 843–851. PMLR, 2018.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Notations . . . . .	2
1.2	Our Contributions and Techniques . . . . .	3
1.3	Related Works . . . . .	4
<b>2</b>	<b>Results on the Power of Adaptivity</b>	<b>4</b>
2.1	Non-Adaptive Fixed Design Algorithm . . . . .	5
2.2	Lower Bound for Non-Adaptive Fixed Design Algorithms . . . . .	6
2.3	Properties of Gaussian Width . . . . .	7
2.4	Structured Sets $\mathcal{X}$ for Which Adaptivity Yields Only Logarithmic-Factor Improvements . . . . .	8
2.5	A Structured Set $\mathcal{X}$ for Which Adaptivity Yields a Polynomial-Factor Improvement . . . . .	10
<b>3</b>	<b>An Algorithm for <math>\ell_2</math>-Norm Estimation</b>	<b>11</b>
<b>4</b>	<b>Conclusion and Future Works</b>	<b>12</b>
<b>A</b>	<b>Technical Lemmas</b>	<b>17</b>
<b>B</b>	<b>Non-Adaptive Fixed Design Bounds</b>	<b>20</b>
B.1	Technical Lemmas for Non-Adaptive Fixed-Design Algorithm . . . . .	20
B.2	Adaptive Lower Bound . . . . .	22
B.3	Gaussian Width Lower Bound . . . . .	26
B.4	Singular Case of the Gaussian Width Lower Bound . . . . .	26
<b>C</b>	<b>Adaptive Version of Our Non-Adaptive Algorithm.</b>	<b>27</b>
<b>D</b>	<b><math>\ell_2</math> Norm Estimation Algorithm</b>	<b>28</b>
D.1	Estimating with the Help of a Multiplicative Estimate $r_0$ . . . . .	28
D.2	Estimating $r$ up to a Constant Factor . . . . .	32
D.3	Estimation in the Large-Norm Regime . . . . .	36
D.4	Meta Algorithm for $\ell_2$ -Norm Estimation . . . . .	41
<b>E</b>	<b>Tail to Sub-Exponential Technical Lemma</b>	<b>41</b>
<b>F</b>	<b>Gaussian Width Properties</b>	<b>45</b>
<b>G</b>	<b>Gaussian Width Calculations for Multi-Task MAB</b>	<b>48</b>
<b>H</b>	<b>Lower Bounds for Structured Sets</b>	<b>51</b>
H.1	Lower Bound for Multi-task MAB . . . . .	51
H.1.1	$d_j < 100$ Case . . . . .	54
H.2	Hypercubes . . . . .	54
H.3	$m$ -Sets Lower Bound . . . . .	55
<b>I</b>	<b>Unit Ball Lower Bound</b>	<b>57</b>
<b>J</b>	<b>Polynomial Separation Instance’s Non-Adaptive Lower Bound</b>	<b>58</b>

## Appendix A. Technical Lemmas

**Lemma 9** *If  $X \sim \mathcal{N}(\mu, \sigma^2)$  and  $Y := X^2 - \mathbb{E}[X^2]$ , then for all  $|\lambda| \leq \frac{1}{4\sigma^2}$ , we have:*

$$\log \mathbb{E}(e^{\lambda Y}) \leq 4(\sigma^4 + \mu^2 \sigma^2) \lambda^2$$

**Proof:**

Let  $X \sim \mathcal{N}(\mu, \sigma^2)$ . For  $|\lambda| \leq \frac{1}{4\sigma^2}$ , we compute:

$$\mathbb{E} \left[ e^{\lambda X^2} \right] = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp \left( \lambda x^2 - \frac{(x - \mu)^2}{2\sigma^2} \right) dx$$

Now, complete the square in the exponent:

$$\begin{aligned} \lambda x^2 - \frac{(x - \mu)^2}{2\sigma^2} &= -\frac{1}{2\sigma^2} [(1 - 2\sigma^2\lambda)x^2 - 2\mu x + \mu^2] \\ &= -\frac{1 - 2\sigma^2\lambda}{2\sigma^2} \left( x - \frac{\mu}{1 - 2\sigma^2\lambda} \right)^2 + \frac{\mu^2\lambda}{1 - 2\sigma^2\lambda} \end{aligned}$$

Hence:

$$\mathbb{E} \left[ e^{\lambda X^2} \right] = \frac{e^{\mu^2\lambda/(1-2\sigma^2\lambda)}}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp \left( -\frac{(1 - 2\sigma^2\lambda)}{2\sigma^2} \left( x - \frac{\mu}{1 - 2\sigma^2\lambda} \right)^2 \right) dx$$

This integral is Gaussian:

$$\begin{aligned} &= \frac{e^{\mu^2\lambda/(1-2\sigma^2\lambda)}}{\sqrt{2\pi\sigma^2}} \cdot \sqrt{\frac{2\pi\sigma^2}{1 - 2\sigma^2\lambda}} \\ &= (1 - 2\sigma^2\lambda)^{-1/2} \exp \left( \frac{\mu^2\lambda}{1 - 2\sigma^2\lambda} \right) \end{aligned}$$

As

$$\lambda \mathbb{E}[X^2] = \lambda(\mu^2 + \sigma^2)$$

$$\log \mathbb{E}[e^{\lambda Y}] = -\frac{1}{2} \log(1 - 2\sigma^2\lambda) + \frac{\mu^2\lambda}{1 - 2\sigma^2\lambda} - \lambda(\mu^2 + \sigma^2)$$

Let  $\nu := 2\sigma^2\lambda$ . As  $\lambda \leq \frac{1}{4\sigma^2}$ , we have  $|\nu| \leq \frac{1}{2}$ .

$$\Rightarrow \log \mathbb{E}[e^{\lambda Y}] = \left( -\frac{1}{2} \log(1 - \nu) - \frac{\nu}{2} \right) + \mu^2\lambda \left( \frac{1}{1 - \nu} - 1 \right)$$

For any  $\nu \in [-\frac{1}{2}, \frac{1}{2}]$ , we have:

1.

$$-\frac{1}{2} \log(1 - \nu) - \frac{\nu}{2} \leq \nu^2$$

2.

$$\left| \frac{1}{1 - \nu} - 1 \right| = \left| \frac{\nu}{1 - \nu} \right| \leq \frac{|\nu|}{1 - |\nu|} \leq 2|\nu|$$

Hence we have:

$$\log \mathbb{E}[e^{\lambda Y}] \leq \nu^2 + \mu^2|\lambda||2\nu| = 4\sigma^4\lambda^2 + 4\mu^2\sigma^2\lambda^2$$

**Lemma 10** *If  $U_1, \dots, U_K$  are independent mean-zero, sub-exponential random variables with parameters  $(V_1, b), \dots, (V_K, b)$  respectively, then we have the following:*

$$\Pr \left( \left| \frac{1}{K} \sum_{s=1}^K U_s \right| > t \right) \leq 2 \exp \left( -cK \min \left\{ \frac{t^2}{\bar{V}}, \frac{t}{b} \right\} \right)$$

where  $\bar{V} = \frac{1}{K} \sum_{k=1}^K V_k$  and  $c > 0$  is some absolute constant.

**Proof** Let  $U$  be a mean-zero sub-exponential random variables with parameters  $(V, b)$ . First, we show the following:

$$\Pr(|U| > t) \leq 2 \exp \left( -c \min \left\{ \frac{t^2}{V}, \frac{t}{b} \right\} \right)$$

For any  $t \geq 0$  and any  $\lambda \in [0, \frac{1}{b}]$ ,

$$\begin{aligned} \Pr[U \geq t] &= \Pr \left[ e^{\lambda U} \geq e^{\lambda t} \right] \\ &\leq e^{-\lambda t} \mathbb{E}[e^{\lambda U}] \\ &\leq \exp \left( -\lambda t + \frac{\lambda^2 V}{2} \right) \end{aligned}$$

Now we minimize over  $\lambda \in [0, \frac{1}{b}]$ . Observe that the unconstrained minimizer is:

$$\lambda^* = \frac{t}{V}$$

If  $t \leq \frac{V}{b}$ , then  $\lambda^* \leq \frac{1}{b}$ , so we get:

$$\Pr(U \geq t) \leq \exp \left( -\frac{t^2}{2V} \right)$$

If  $t > \frac{V}{b}$ , the minimum occurs at  $\lambda = \frac{1}{b}$ :

$$\Pr(U \geq t) \leq \exp \left( -\frac{t}{b} + \frac{V}{2b^2} \right) \leq \exp \left( -\frac{1}{2} \cdot \frac{t}{b} \right)$$

If  $U_1, \dots, U_K$  are independent mean-zero, sub-exponential random variables with parameters

$$(V_s, b) \quad (\text{same } b \text{ for all}),$$

then for any  $|\lambda| \leq \frac{K}{b}$ , we have:

$$\begin{aligned} \mathbb{E} \left[ \exp \left( \frac{\lambda}{K} \sum_{s=1}^K U_s \right) \right] &= \prod_{s=1}^K \mathbb{E} \left[ e^{\frac{\lambda}{K} U_s} \right] \leq \exp \left( \frac{\lambda^2}{2K^2} \sum_{s=1}^K V_s \right) \\ \Rightarrow \bar{U} := \frac{1}{K} \sum_{s=1}^K U_s &\text{ is sub-exponential with parameters } \left( \frac{1}{K^2} \sum_{s=1}^K V_s, \frac{b}{K} \right) \\ \Rightarrow \Pr \left( \left| \frac{1}{K} \sum_{s=1}^K U_s \right| > t \right) &\leq 2 \exp \left( -cK \min \left\{ \frac{t^2}{\bar{V}}, \frac{t}{b} \right\} \right) \end{aligned}$$

where

$$\bar{V} = \frac{1}{K} \sum_{k=1}^K V_k$$

■

**Lemma 11 (Hanson–Wright inequality)** *Let  $X = (X_1, \dots, X_n) \in \mathbb{R}^n$  be a random vector with independent components  $X_i$  satisfying:*

$$\mathbb{E}[X_i] = 0, \quad \|X_i\|_{\psi_2} \leq K$$

*Let  $A$  be an  $n \times n$  matrix. Then for every  $t \geq 0$ ,*

$$\Pr \left( \left| X^\top A X - \mathbb{E}[X^\top A X] \right| > t \right) \leq 2 \exp \left( -c \min \left( \frac{t^2}{K^4 \|A\|_F^2}, \frac{t}{K^2 \|A\|} \right) \right)$$

where

$$\|X\|_{\psi_2} = \sup_{p \geq 1} p^{-1/2} (\mathbb{E}|X|^p)^{1/p}$$

**Lemma 12 (Chain Rule)** *Let  $f(x_1, x_2, \dots, x_n)$  and  $g(x_1, x_2, \dots, x_n)$  be two joint PDFs for a tuple of random variables  $(X_i)_{i \in [n]}$ . Let the sample space be  $\Omega = \mathbb{R}^n$ . Then we have the following:*

$$KL(f, g) = \int_{\omega \in \Omega} f(\omega) \left( KL(f(X_1), g(X_1)) + \sum_{i=2}^n KL(f(X_i|X_{-i} = \omega_{-i}), g(X_i|X_{-i} = \omega_{-i})) \right) d\omega$$

where  $X_{-i} = (X_1, \dots, X_{i-1})$ ,  $\omega_{-i} = (\omega_1, \dots, \omega_{i-1})$ .

**Proof**

$$\begin{aligned} KL(f, g) &= \int_{\omega \in \Omega} f(\omega) \log \left( \frac{f(\omega)}{g(\omega)} \right) d\omega \\ &= \int_{\omega \in \Omega} f(\omega) \log \left( \frac{f(\omega_1) \prod_{i=2}^n f(\omega_i|\omega_{-i})}{g(\omega_1) \prod_{i=2}^n g(\omega_i|\omega_{-i})} \right) d\omega \\ &= \int_{\omega \in \Omega} f(\omega) \left( \log \left( \frac{f(\omega_1)}{g(\omega_1)} \right) + \sum_{i=2}^n \log \left( \frac{f(\omega_i|\omega_{-i})}{g(\omega_i|\omega_{-i})} \right) \right) d\omega \\ &= \int_{\omega \in \Omega} f(\omega) \log \left( \frac{f(\omega_1)}{g(\omega_1)} \right) d\omega + \sum_{i=2}^n \int_{\omega \in \Omega} f(\omega) \log \left( \frac{f(\omega_i|\omega_{-i})}{g(\omega_i|\omega_{-i})} \right) d\omega \\ &= \int_{\omega_1 \in \mathbb{R}} f(\omega_1) \log \left( \frac{f(\omega_1)}{g(\omega_1)} \right) d\omega_1 + \sum_{i=2}^n \int_{\omega_{-i} \in \mathbb{R}^{i-1}} f(\omega_{-i}) \int_{\omega_i \in \mathbb{R}} f(\omega_i|\omega_{-i}) \log \left( \frac{f(\omega_i|\omega_{-i})}{g(\omega_i|\omega_{-i})} \right) d\omega_i d\omega_{-i} \\ &= KL(f(X_1), g(X_1)) + \sum_{i=2}^n \int_{\omega_{-i} \in \mathbb{R}^{i-1}} f(\omega_{-i}) KL(f(X_i|X_{-i} = \omega_{-i}), g(X_i|X_{-i} = \omega_{-i})) d\omega_{-i} \end{aligned}$$

$$\begin{aligned}
 &= \int_{\omega \in \Omega} f(\omega) KL(f(X_1), g(X_1)) d\omega + \sum_{i=2}^n \int_{\omega \in \Omega} f(\omega) KL(f(X_i|X_{-i} = \omega_{-i}), g(X_i|X_{-i} = \omega_{-i})) d\omega \\
 &= \int_{\omega \in \Omega} f(\omega) \left( KL(f(X_1), g(X_1)) + \sum_{i=2}^n KL(f(X_i|X_{-i} = \omega_{-i}), g(X_i|X_{-i} = \omega_{-i})) \right) d\omega
 \end{aligned}$$

■

## Appendix B. Non-Adaptive Fixed Design Bounds

### B.1. Technical Lemmas for Non-Adaptive Fixed-Design Algorithm

In this section, we state various technical lemmas and explain how they are used in Section 2.1.

**Lemma 13 (Carathéodory's theorem)** *Let  $S \subset \mathbb{R}^p$  and let  $x \in \text{conv}(S)$ , where  $\text{conv}(S)$  denotes the convex hull of  $S$ . Then there exist points  $s_0, \dots, s_p \in S$  and coefficients  $\alpha_0, \dots, \alpha_p \geq 0$  with  $\sum_{i=0}^p \alpha_i = 1$  such that*

$$x = \sum_{i=0}^p \alpha_i s_i.$$

*In particular, every point in the convex hull of  $S$  can be expressed as a convex combination of at most  $p + 1$  points of  $S$ .*

Since,  $\mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)}[\max_{x \in \mathcal{X}} \langle x, A(\lambda; \mathcal{X})^{-1/2} \eta \rangle]$  depends on  $\lambda$  through  $A(\lambda; \mathcal{X}) = \mathbb{E}_{x \sim \lambda}[xx^\top]$ , applying Carathéodory's theorem on  $S = \{xx^\top : x \in \mathcal{X}\}$  yields a finite supported minimizer  $\lambda_1$ .

**Lemma 14 (G-optimal design)** *Consider a compact set  $\mathcal{X} \subset \mathbb{R}^d$  such that it spans  $\mathbb{R}^d$ . There exists a distribution  $\lambda_*$  over  $\mathcal{X}$  such that  $\max_{x \in \mathcal{X}} \|x\|_{A(\lambda_*; \mathcal{X})^{-1}}^2 = d$  and  $|\text{supp}(\lambda_*)| \leq d(d + 1)/2$ .*

Applying the above lemma yields a finite supported minimizer  $\lambda_2$ .

**Lemma 15 (Sudakov–Fernique inequality)** *Let  $(X_t)_{t \in T}$  and  $(Y_t)_{t \in T}$  be two mean-zero Gaussian processes. Assume that for all  $t, s \in T$ , we have*

$$\mathbb{E} \left[ (X_t - X_s)^2 \right] \leq \mathbb{E} \left[ (Y_t - Y_s)^2 \right].$$

Then

$$\mathbb{E} \left[ \sup_{t \in T} X_t \right] \leq \mathbb{E} \left[ \sup_{t \in T} Y_t \right].$$

Define  $X_x = \langle x, A(\lambda_0; \mathcal{X})^{-1/2} \eta \rangle$  and  $Y_x = \sqrt{2} \langle x, A(\lambda_1; \mathcal{X})^{-1/2} \eta \rangle$  with  $\eta \sim \mathcal{N}(0, I_d)$ . Then for any  $s, t \in \mathcal{X}$ ,  $X_s - X_t = \langle s - t, A(\lambda_0; \mathcal{X})^{-1/2} \eta \rangle$ , so by  $\mathbb{E}[\eta \eta^\top] = I_d$  we have  $\mathbb{E}[(X_s - X_t)^2] = (s - t)^\top A(\lambda_0; \mathcal{X})^{-1} (s - t)$  (and similarly  $\mathbb{E}[(Y_s - Y_t)^2] = 2(s - t)^\top A(\lambda_1; \mathcal{X})^{-1} (s - t)$ ). Thus the increment condition of Sudakov–Fernique holds, yielding  $\mathbb{E} \sup_{x \in \mathcal{X}} X_x \leq \mathbb{E} \sup_{x \in \mathcal{X}} Y_x$ .

**Lemma 16 (Linear image of a standard Gaussian)** *Let  $A \in \mathbb{R}^{m \times n}$  be a fixed matrix and let  $\eta \sim \mathcal{N}(0, I_n)$ . Then  $A\eta$  is distributionally equivalent to a mean-zero Gaussian vector in  $\mathbb{R}^m$  with covariance matrix  $AA^\top$ . That is  $A\eta \stackrel{d}{=} \xi$  where  $\xi \sim \mathcal{N}(0, AA^\top)$ .*

We used the above lemma to conclude that  $\widehat{\theta}$  is distributionally equivalent to  $\theta + (\sum_{t=1}^T x_t x_t^\top)^{-1/2} \eta$  where  $\eta \sim \mathcal{N}(0, I_d)$ .

**Lemma 17 (Rounding lemma (Katz-Samuels et al., 2020; Allen-Zhu et al., 2021))** *Let  $\mathcal{Z} \subset \mathbb{R}^d$  be finite with  $m := |\mathcal{Z}|$ , and let  $N \in \mathbb{N}$ . Define*

$$S_N := \left\{ \kappa \in \mathbb{N}^m : \sum_{x \in \mathcal{Z}} \kappa_x \leq N \right\}, \quad C_N := \left\{ \pi \in [0, N]^m : \sum_{x \in \mathcal{Z}} \pi_x \leq N \right\}.$$

For any weight vector  $v \in \mathbb{R}_+^m$ , define the (design) matrix

$$A(v) := \sum_{x \in \mathcal{Z}} v_x x x^\top \in \mathbb{S}_d^+.$$

Let  $F : \mathbb{S}_d^+ \rightarrow \mathbb{R}$  satisfy: (i) if  $A \preceq B$  then  $F(A) \geq F(B)$ , and (ii) for all  $A \in \mathbb{S}_d^+$  and  $t \in (0, 1)$ ,  $F(tA) = t^{-1} F(A)$ . Fix  $\epsilon \in (0, 1/6]$ . If  $m \geq N \geq 5d/\epsilon^2$ , then for every  $\pi \in C_N$  there exists an algorithm running in  $\widetilde{O}(md^2)$  time that outputs  $\kappa \in S_N$  such that

$$F(A(\kappa)) \leq (1 + 6\epsilon) F(A(\pi)).$$

Moreover, for any set  $V \subset \mathbb{R}^d$ , the functions  $F_V, G_V : \mathbb{S}_d^+ \rightarrow \mathbb{R}$  defined by

$$F_V(A) := \mathbb{E}_{\eta \sim \mathcal{N}(0, I)} \left[ \max_{v \in V} v^\top A^{-1/2} \eta \right], \quad G_V(A) := \max_{v \in V} v^\top A v$$

satisfy conditions (i) and (ii).

For any  $T \geq 180d$ , we used the above lemma to show the existence of a fixed design  $x_1, x_2, \dots, x_T \in \mathcal{X}$  such that  $\tau(A_T) \leq 2\tau(A(\lambda_0; \mathcal{X}))$  where  $A_T := \frac{1}{T} \sum_{i=1}^T x_i x_i^\top$  and  $\tau(A) := \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in \mathcal{X}} x^\top A^{-1/2} \eta \right]^2 + 2 \max_{x \in \mathcal{X}} \|x\|_{A^{-1}}^2 \log(2/\delta)$ . In the lemma, we set  $\mathcal{Z}$  as the finite support of  $\lambda_0$ ,  $N = T$ , and  $\pi$  as  $\pi_x = T \cdot \lambda_0(x)$ .

**Lemma 18 (Borell-TIS inequality)** *Let  $S \subset \mathbb{R}^d$  be bounded. Let  $(V_s)_{s \in S}$  be a Gaussian process such that*

$$\mathbb{E}[V_s] = 0 \quad \text{for all } s \in S.$$

Define

$$\sigma^2 = \sup_{s \in S} \mathbb{E}[V_s^2].$$

Then, for all  $u > 0$ ,

$$\mathbb{P} \left( \left| \sup_{s \in S} V_s - \mathbb{E} \sup_{s \in S} V_s \right| \geq u \right) \leq 2 \exp \left( -\frac{u^2}{2\sigma^2} \right).$$

We applied Borell-TIS inequality on the gaussian process  $V_x = x^\top (\widehat{\theta} - \theta)$ . We then used the facts that  $V_x$  is distributionally equivalent to  $\langle x, (\sum_{t=1}^T x_t x_t^\top)^{-1/2} \eta \rangle$  where  $\eta \sim \mathcal{N}(0, I_d)$  and  $\sup_{x \in \mathcal{X}} \mathbb{E}[V_x^2] = \max_{x \in \mathcal{X}} x^\top (\sum_{t=1}^T x_t x_t^\top)^{-1} x$ .

## B.2. Adaptive Lower Bound

**Theorem 19** Assume  $\mathcal{X} \subset \mathbb{R}^d$  is compact,  $\text{span}(\mathcal{X}) = \mathbb{R}^d$ , and  $d \geq 2$ . Consider the Gaussian linear observation model

$$y_t = \langle x_t, \theta \rangle + \eta_t, \quad \eta_t \sim \mathcal{N}(0, 1) \text{ i.i.d.},$$

where the (possibly adaptive) actions satisfy  $x_t \in \mathcal{X}$ . Any algorithm that outputs  $\hat{x} \in \mathcal{X}$  satisfying, for all  $\theta \in \mathbb{R}^d$ ,

$$\langle \hat{x}, \theta \rangle \geq \max_{x \in \mathcal{X}} \langle x, \theta \rangle - \varepsilon \quad \text{with probability at least } 1 - \delta,$$

must use

$$n = \Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$$

samples for some worst-case  $\theta$  and  $\delta \in (0, 1/16)$ .

**Proof** Let

$$B := \text{conv}(\mathcal{X} \cup (-\mathcal{X})).$$

Let  $E$  be the Löwner ellipsoid of  $B$ , the minimum-volume ellipsoid containing  $B$ . Since  $B$  is centrally symmetric and full-dimensional,  $E$  is an origin-centered ellipsoid and there exists an invertible linear map  $A$  such that

$$A(E) = \mathbb{B}_2^d \quad \text{and hence} \quad A(B) \subseteq \mathbb{B}_2^d.$$

Define

$$K := A(B), \quad \mathcal{X}' := A(\mathcal{X}),$$

so  $\mathcal{X}' \subseteq K \subseteq \mathbb{B}_2^d$ .

We work in the transformed instance  $(\mathcal{X}', \theta')$  where  $\theta' := (A^{-1})^\top \theta$ . Indeed,  $\langle Ax, \theta' \rangle = \langle x, \theta \rangle$ , so any lower bound for  $\mathcal{X}'$  implies the same for  $\mathcal{X}$ . Hence it suffices to prove the statement for  $\mathcal{X}'$ , and we now drop primes and write  $\mathcal{X} \subseteq \mathbb{B}_2^d$  with

$$K = \text{conv}(\mathcal{X} \cup (-\mathcal{X})) \subseteq \mathbb{B}_2^d,$$

such that  $\mathbb{B}_2^d$  is the Löwner ellipsoid of  $K$ .

Since  $\mathbb{B}_2^d$  is the minimum-volume ellipsoid containing  $K$ , By applying John's theorem on the polar  $K^\circ$  (see Theorem 3.1 from [Ball et al. \(1997\)](#)), there exist points  $u_1, \dots, u_m \in K \cap \mathbb{S}^{d-1}$  and weights  $c_1, \dots, c_m > 0$  such that

$$\sum_{j=1}^m c_j u_j u_j^\top = I_d, \quad \text{and hence} \quad \sum_{j=1}^m c_j = \text{tr}(I_d) = d. \quad (3)$$

We now note that these  $u_j$  actually belong to  $\mathcal{X} \cup (-\mathcal{X})$ . Fix  $j$ . Because  $u_j \in \mathbb{S}^{d-1}$  and  $K \subseteq \mathbb{B}_2^d$ , we have

$$\max_{x \in K} \langle x, u_j \rangle \leq \max_{x \in K} \|x\|_2 \|u_j\|_2 = 1.$$

On the other hand,  $u_j \in K$  implies  $\max_{x \in K} \langle x, u_j \rangle \geq \langle u_j, u_j \rangle = 1$ . Hence

$$\max_{x \in K} \langle x, u_j \rangle = 1.$$

Since  $K = \text{conv}(\mathcal{X} \cup (-\mathcal{X}))$  and the objective is linear, the maximum over  $K$  equals the maximum over  $\mathcal{X} \cup (-\mathcal{X})$ . Thus there exists  $v \in \mathcal{X} \cup (-\mathcal{X})$  such that  $\langle v, u_j \rangle = 1$ . By Cauchy-Schwarz,  $\langle v, u_j \rangle \leq \|v\|_2 \|u_j\|_2 \leq 1$ , so equality forces  $v = u_j$ . Therefore  $u_j \in \mathcal{X} \cup (-\mathcal{X})$ .

Define

$$v_j := \begin{cases} u_j, & u_j \in \mathcal{X}, \\ -u_j, & u_j \in -\mathcal{X}, \end{cases}$$

so that  $v_j \in \mathcal{X} \cap \mathbb{S}^{d-1}$  for every  $j$ . Since  $v_j v_j^\top = u_j u_j^\top$ , the same weights  $c_1, \dots, c_m$  satisfy

$$\sum_{j=1}^m c_j v_j v_j^\top = I_d, \quad \text{and} \quad \sum_{j=1}^m c_j = d. \quad (4)$$

Let  $\mathcal{D}$  be the distribution on  $[m]$  such that  $\mathbb{P}_{J \sim \mathcal{D}}(J = j) = c_j/d$ . Then (4) implies

$$\mathbb{E}[v_J v_J^\top] = \frac{1}{d} I_d. \quad (5)$$

We first prove the following lemma.

**Lemma 20** *Any algorithm that outputs  $\hat{x} \in \mathcal{X} \subseteq \mathbb{B}_2^d$  satisfying, for all  $\theta \in \mathbb{R}^d$ ,*

$$\langle \hat{x}, \theta \rangle \geq \max_{x \in \mathcal{X}} \langle x, \theta \rangle - \varepsilon \quad \text{with probability at least } 1 - \delta,$$

*must use*

$$n \geq \frac{2 - \sqrt{2}}{32 \varepsilon^2} \log\left(\frac{1}{4\delta}\right) = \Omega\left(\frac{\log(1/\delta)}{\varepsilon^2}\right)$$

*samples for some worst-case  $\theta$  and  $\delta \in (0, 1/16)$ .*

**Proof** Recall that  $J \sim \mathcal{D}$  is a random index with  $\mathbb{P}_{J \sim \mathcal{D}}(J = j) = c_j/d$ . Hence, we have:

$$\mathbb{E}[v_J v_J^\top] = \frac{1}{d} I_d.$$

Fix  $v_1$ . Then, we have:

$$\mathbb{E}[(v_1^\top v_J)^2] = v_1^\top \mathbb{E}[v_J v_J^\top] v_1 = \frac{1}{d} \|v_1\|_2^2 = \frac{1}{d}.$$

Hence there exists  $k \in [m]$  such that  $(v_1^\top v_k)^2 \leq 1/d$ , and thus  $v_1^\top v_k \leq 1/\sqrt{d} \leq 1/\sqrt{2}$  (since  $d \geq 2$ ). Therefore

$$\|v_1 - v_k\|_2^2 = 2 - 2v_1^\top v_k \geq 2 - \sqrt{2}. \quad (6)$$

Set  $a := v_1$ ,  $b := v_k$ , and denote  $\rho := \|a - b\|_2$ . Then (6) gives  $\rho^2 \geq 2 - \sqrt{2}$ .

Let  $u := (a - b)/\|a - b\|_2 \in \mathbb{S}^{d-1}$  so that  $\langle a, u \rangle - \langle b, u \rangle = \rho$ . Define

$$\theta_+ := \frac{4\varepsilon}{\rho} u, \quad \theta_- := -\frac{4\varepsilon}{\rho} u.$$

Consider an algorithm that uses  $n$  samples and outputs  $\hat{x}$  under each of the hypothesis above while satisfying the conditions in the theorem statement. Let  $p_{\max} := \max_{x \in \mathcal{X}} \langle x, u \rangle$  and  $p_{\min} := \min_{x \in \mathcal{X}} \langle x, u \rangle$ . Note that  $p_{\max} \geq \langle a, u \rangle$  and  $p_{\min} \leq \langle b, u \rangle$ .

Under  $\theta_+$ , the optimal value is  $\max_{x \in \mathcal{X}} \langle x, \theta_+ \rangle = (4\varepsilon/\rho) p_{\max}$ . If  $\hat{x}$  is  $\varepsilon$ -optimal for  $\theta_+$ , then

$$\frac{4\varepsilon}{\rho} \langle \hat{x}, u \rangle = \langle \hat{x}, \theta_+ \rangle \geq \frac{4\varepsilon}{\rho} p_{\max} - \varepsilon \quad \Rightarrow \quad \langle \hat{x}, u \rangle \geq p_{\max} - \frac{\rho}{4} \geq \langle a, u \rangle - \frac{\rho}{4}.$$

Similarly under  $\theta_-$ , the optimal value is  $(4\varepsilon/\rho)(-p_{\min})$  and  $\varepsilon$ -optimality of  $\hat{x}$  implies

$$-\frac{4\varepsilon}{\rho} \langle \hat{x}, u \rangle = \langle \hat{x}, \theta_- \rangle \geq \frac{4\varepsilon}{\rho} (-p_{\min}) - \varepsilon \quad \Rightarrow \quad \langle \hat{x}, u \rangle \leq p_{\min} + \frac{\rho}{4} \leq \langle b, u \rangle + \frac{\rho}{4}.$$

Let

$$s := \frac{\langle a, u \rangle + \langle b, u \rangle}{2}.$$

Since  $\langle a, u \rangle - \langle b, u \rangle = \rho$ , the two implications based on  $\varepsilon$ -optimality of  $\hat{x}$  become:

$$\text{Under } \theta_+ : \langle \hat{x}, u \rangle \geq s + \frac{\rho}{4}, \quad \text{Under } \theta_- : \langle \hat{x}, u \rangle \leq s - \frac{\rho}{4},$$

each with probability at least  $1 - \delta$ . Therefore the decision rule

$$\hat{H} = \begin{cases} \theta_+, & \text{if } \langle \hat{x}, u \rangle \geq s, \\ \theta_-, & \text{otherwise,} \end{cases}$$

distinguishes  $\theta_+$  from  $\theta_-$  with error at most  $\delta$  under each hypothesis.

Let  $\mathbb{P}_+$  and  $\mathbb{P}_-$  be the probability laws under  $\theta_+$  and  $\theta_-$ . By Bretagnolle–Huber inequality, we get

$$\mathbb{P}_+(\hat{H} = \theta_-) + \mathbb{P}_-(\hat{H} = \theta_+) \geq \frac{1}{2} \exp(-D_{\text{KL}}(\mathbb{P}_+ \parallel \mathbb{P}_-)).$$

Since each error is at most  $\delta$ , the left-hand side is at most  $2\delta$ , hence

$$D_{\text{KL}}(\mathbb{P}_+ \parallel \mathbb{P}_-) \geq \log\left(\frac{1}{4\delta}\right).$$

For Gaussian noise and adaptive actions  $x_t \in \mathcal{X} \subseteq \mathbb{B}_2^d$ ,

$$\begin{aligned} D_{\text{KL}}(\mathbb{P}_+ \parallel \mathbb{P}_-) &= \frac{1}{2} \sum_{t=1}^n \mathbb{E}_+[\langle x_t, \theta_+ - \theta_- \rangle^2] = \frac{1}{2} \sum_{t=1}^n \mathbb{E}_+ \left[ \left\langle x_t, \frac{8\varepsilon}{\rho} u \right\rangle^2 \right] \\ &\leq \frac{1}{2} \sum_{t=1}^n \left( \frac{8\varepsilon}{\rho} \right)^2 = \frac{32n\varepsilon^2}{\rho^2}, \end{aligned}$$

using  $|\langle x_t, u \rangle| \leq \|x_t\|_2 \|u\|_2 \leq 1$ . Combining the last two inequalities, we have

$$n \geq \frac{\rho^2}{32\varepsilon^2} \log\left(\frac{1}{4\delta}\right).$$

Finally,  $\rho^2 \geq 2 - \sqrt{2}$  by (6) yields

$$n \geq \frac{2 - \sqrt{2}}{32\varepsilon^2} \log\left(\frac{1}{4\delta}\right). \quad \blacksquare$$

Define the alternatives

$$\theta^{(j)} := 3\varepsilon v_j, \quad j = 1, \dots, m,$$

and consider the hypothesis test

$$H_0 : \theta = 0, \quad H_1 : \theta = \theta^{(J)} \text{ where } J \sim \mathcal{D}.$$

Let  $\mathbb{P}_0$  be the law under  $H_0$ ,  $\mathbb{P}_1$  the law under  $H_1$ , and  $\mathbb{P}_{\theta^{(j)}}$  the law under  $\theta = \theta^{(j)}$ . For Gaussian noise, convexity of KL and the chain rule yield, for any (possibly adaptive) choice of actions  $x_t \in \mathcal{X}$ ,

$$D_{\text{KL}}(\mathbb{P}_0 \parallel \mathbb{P}_1) \leq \sum_{j=1}^m \frac{c_j}{d} D_{\text{KL}}(\mathbb{P}_0 \parallel \mathbb{P}_{\theta^{(j)}}) = \frac{1}{2} \sum_{t=1}^N \mathbb{E} \left[ \langle x_t, \theta^{(J)} \rangle^2 \right],$$

where  $N$  is the total number of samples used by the algorithm solving the hypothesis testing problem, and the expectation is under  $\mathbb{P}_0$  and the distribution  $\mathcal{D}$ . Using (5) and  $\theta^{(j)} = 3\varepsilon v_j$ ,

$$\begin{aligned} \mathbb{E}[\langle x_t, \theta^{(j)} \rangle^2] &= 9\varepsilon^2 \mathbb{E}[\langle x_t, v_j \rangle^2] = 9\varepsilon^2 x_t^\top \mathbb{E}[v_j v_j^\top] x_t \\ &= 9\varepsilon^2 x_t^\top \left(\frac{1}{d} I_d\right) x_t = \frac{9\varepsilon^2}{d} \|x_t\|_2^2 \leq \frac{9\varepsilon^2}{d}. \end{aligned}$$

Therefore, we have

$$D_{\text{KL}}(\mathbb{P}_0 \|\mathbb{P}_1) \leq \frac{9N\varepsilon^2}{2d}. \quad (7)$$

By Bretagnolle–Huber inequality, we have

$$\mathbb{P}_0(\widehat{H} = H_1) + \mathbb{P}_1(\widehat{H} = H_0) \geq \frac{1}{2} \exp(-D_{\text{KL}}(\mathbb{P}_0 \|\mathbb{P}_1)) \geq \frac{1}{2} \exp\left(-\frac{9N\varepsilon^2}{2d}\right).$$

In particular, if  $N \leq \frac{2d}{9\varepsilon^2} \log\left(\frac{1}{10\delta}\right)$ , then the sum of the two errors exceeds  $4\delta$ , so at least one error exceeds  $2\delta$ . Hence any algorithm with error at most  $2\delta$  under both  $H_0$  and  $H_1$  requires

$$N = \Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right). \quad (8)$$

Consider an algorithm Alg that satisfies the conditions of our theorem statement. We now solve the above hypothesis test using Alg.

Run Alg for  $n$  rounds to obtain  $\widehat{x} \in \mathcal{X}$ . Then take

$$n_{\text{est}} := \frac{8 \log(2/\delta)}{\varepsilon^2}$$

additional samples using the constant action  $x_t = \widehat{x}$  and let  $\widehat{v}$  be the empirical mean of these additional observations. Output  $\widehat{H} = H_1$  if  $\widehat{v} > \varepsilon$  and  $\widehat{H} = H_0$  otherwise.

Gaussian concentration implies

$$|\widehat{v} - \langle \widehat{x}, \theta \rangle| \leq \varepsilon/2 \quad \text{with probability at least } 1 - \delta.$$

Under  $H_0$ ,  $\langle \widehat{x}, \theta \rangle = 0$ , so  $\widehat{v} \leq \varepsilon/2$  with probability at least  $1 - \delta$ , hence  $\widehat{H} = H_0$  with probability at least  $1 - \delta$ .

Under  $H_1$ , let us condition on  $J = j$ . Since  $v_j \in \mathcal{X} \cap \mathbb{S}^{d-1}$ ,

$$\max_{x \in \mathcal{X}} \langle x, \theta^{(j)} \rangle \geq \langle v_j, 3\varepsilon v_j \rangle = 3\varepsilon.$$

Therefore Alg outputs  $\widehat{x}$  such that with probability at least  $1 - \delta$ ,

$$\langle \widehat{x}, \theta^{(j)} \rangle \geq 3\varepsilon - \varepsilon = 2\varepsilon.$$

Therefore, with probability at least  $1 - 2\delta$ , we have  $\widehat{v} \geq 3\varepsilon/2 > \varepsilon$  and therefore  $\widehat{H} = H_1$ .

Thus we solve the hypothesis testing problem with probability at least  $1 - 2\delta$  under both  $H_0$  and  $H_1$  using

$$N = n + n_{\text{est}}$$

samples.

Finally, Lemma 20 implies  $n \geq c \log(2/\delta)/\varepsilon^2$  for an absolute constant  $c > 0$ . Hence  $n_{\text{est}} \leq (8/c)n$  and so  $N \leq (1 + 8/c)n$ . Combining with (8) yields

$$n = \Omega\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right).$$

■

### B.3. Gaussian Width Lower Bound

Fix  $\varepsilon \in (0, 1]$ ,  $\delta \in (0, 1)$  and an  $(\varepsilon, \delta)$ -PAC algorithm. Let the algorithm run for  $T := \frac{H_1 \log(1/\delta) + H_2}{\varepsilon^2}$  rounds and output  $\hat{x} \in \mathcal{X}$ . Define the simple regret for any  $\theta$  as

$$R_T(\theta) := \max_{x \in \mathcal{X}} \langle x, \theta \rangle - \langle \hat{x}, \theta \rangle \geq 0.$$

Also define the quantity

$$Z(\theta) := \max_{x, y \in \mathcal{X}} \langle x - y, \theta \rangle,$$

Since  $\hat{x} \in \mathcal{X}$ , we have  $R_T(\theta) \leq Z(\theta)$  for every  $\theta$ .

As the algorithm is  $(\varepsilon, \delta)$ -PAC, we have  $\mathbb{P}(R_T(\theta) > \varepsilon) \leq \delta$ . Hence, taking expectation over  $\theta \sim \mathcal{N}(0, \Sigma)$  where  $\Sigma = \tau^2 A^{-1}$ , we have

$$\mathbb{E}_{\theta \sim \mathcal{N}(0, \Sigma)}[R_T(\theta)] \leq \varepsilon + \delta \cdot \mathbb{E}_{\theta \sim \mathcal{N}(0, \Sigma)}[Z(\theta)] = \varepsilon + 2\tau \cdot \delta \cdot w(\mathcal{X}; A)$$

Recall that for any  $T$  and any algorithm, we have:

$$\mathbb{E}_{\theta}[R_T(\theta)] \geq \frac{\tau(1 - \tau)}{1 + \tau^2} w(\mathcal{X}; A)$$

Hence, using the above two inequalities and the facts that  $w(\mathcal{X}; A) \geq w(\mathcal{X})/\sqrt{T}$  and  $H_1 \leq H_2$ , we have the following by setting  $\delta = 0.1$

$$\frac{0.12 \cdot w(\mathcal{X})}{\sqrt{T}} \leq \varepsilon \quad \Rightarrow \quad H_2 \geq \frac{0.0144}{\log(10) + 1} \cdot w(\mathcal{X})^2$$

### B.4. Singular Case of the Gaussian Width Lower Bound

In this section, we prove the following result.

**Theorem 21** *Assume  $\mathcal{X} \subset \mathbb{R}^d$  is finite with  $\text{span}(\mathcal{X}) = \mathbb{R}^d$ . Consider  $\delta \in (0, 1/2)$  and a fixed design  $x_1, \dots, x_T$  such that  $\sum_{t=1}^T x_t x_t^\top$  is singular. Then for every (possibly randomized) non-adaptive procedure  $\mathcal{A}$  that outputs  $\hat{x} \in \mathcal{X}$  based on this fixed design, there exists a  $\theta \in \mathbb{R}^d$  such that*

$$\mathbb{P}(\max_{x \in \mathcal{X}} \langle x - \hat{x}, \theta \rangle > \varepsilon) > \delta.$$

**Proof** Since  $A := \sum_{t=1}^T x_t x_t^\top$  is singular, pick  $v \neq 0$  such that  $Av = 0$ . Then we have:

$$0 = v^\top Av = \sum_{t=1}^T (v^\top x_t)^2.$$

Therefore, for all  $t \in [T]$ , we have:

$$\langle x_t, v \rangle = 0$$

Consider two reward vectors

$$\theta_{+1} = \alpha v, \quad \theta_{-1} = -\alpha v,$$

with  $\alpha > 0$ . Then for every  $t$  and  $s \in \{-1, +1\}$ , we have

$$y_t = \langle x_t, \theta_s \rangle + \eta_t = s \cdot \alpha \langle x_t, v \rangle + \eta_t = \eta_t$$

Hence the distribution of the algorithm's output  $\hat{x}$  is identical under  $\theta_{+1}$  and  $\theta_{-1}$ .

Define

$$a := \max_{x \in \mathcal{X}} \langle x, v \rangle, \quad b := \min_{x \in \mathcal{X}} \langle x, v \rangle.$$

We claim  $a > b$ . For the sake of contradiction, assume  $a = b$ . In this case  $\langle x, v \rangle$  is constant over  $\mathcal{X}$ . However,  $\langle x_t, v \rangle$  and  $x_t \in \mathcal{X}$  which implies that  $\langle x, v \rangle = 0$  for all  $x \in \mathcal{X}$ . This implies  $v \perp \text{span}(\mathcal{X})$ , contradicting  $\text{span}(\mathcal{X}) = \mathbb{R}^d$ . Thus

$$\Delta(v) := a - b > 0.$$

Let us define regret as  $r(\hat{x}, \theta) = \max_{x \in \mathcal{X}} \langle x - \hat{x}, \theta \rangle$ . We now compute regret under each  $\theta_s$  for  $s \in \{-1, +1\}$ .

- For  $\theta_{+1} = \alpha v$ ,

$$r(\hat{x}, \theta_{+1}) = \max_{x \in \mathcal{X}} \langle x - \hat{x}, \alpha v \rangle = \alpha (a - \langle \hat{x}, v \rangle),$$

- For  $\theta_{-1} = -\alpha v$ ,

$$r(\hat{x}, \theta_{-1}) = \max_{x \in \mathcal{X}} \langle x - \hat{x}, -\alpha v \rangle = \alpha (\langle \hat{x}, v \rangle - b),$$

Adding the two bounds gives

$$r(\hat{x}, \theta_{+1}) + r(\hat{x}, \theta_{-1}) = \alpha(a - b) = \alpha \Delta(v).$$

Therefore, if  $\alpha = \frac{4\varepsilon}{\Delta(v)}$ , we have

$$\max \left\{ r(\hat{x}, \theta_{+1}), r(\hat{x}, \theta_{-1}) \right\} \geq \frac{\alpha \Delta(v)}{2} = 2\varepsilon.$$

This implies that  $\max_{s \in \{-1, +1\}} \mathbb{P}(\max_{x \in \mathcal{X}} \langle x - \hat{x}, \theta_s \rangle > \varepsilon) \geq \frac{1}{2} > \delta$ . ■

### Appendix C. Adaptive Version of Our Non-Adaptive Algorithm.

In this section, we make our non-adaptive algorithm from Section 2.1 adaptive. Let  $\mathcal{R} = (R_1, R_2, \dots, R_d)$  be a partition of  $\mathcal{X}$  into  $d$  regions. Let us define

$$w(\lambda, \mathcal{X}, \mathcal{R}) = \max_{i \in [d]} \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in R_i} \langle x, A(\lambda; \mathcal{X})^{-1/2} \eta \rangle \right]$$

Let  $i_* \in [d]$  be the index such that  $x_* \in R_{i_*}$ .

Let  $\lambda_1$  be a distribution over  $\mathcal{X}$  that minimizes the expression  $\mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} [\max_{x \in \mathcal{X}} \langle x, A(\lambda; \mathcal{X})^{-1/2} \eta \rangle]$ . Similarly let  $\lambda_2$  be a distribution over  $\mathcal{X}$  that minimizes the expression  $\max_{x \in \mathcal{X}} \|x\|_{A(\lambda; \mathcal{X})^{-1}}^2$ . Assume that  $\lambda_1$  and  $\lambda_2$  are finite supported (such minimizers always exist). Let  $\lambda_0$  be a distribution over  $\mathcal{X}$  such that we sample from  $\lambda_1$  with probability 1/2 and we sample from  $\lambda_2$  with probability 1/2.

As  $\frac{1}{2}A(\lambda_2; \mathcal{X}) \prec A(\lambda_0; \mathcal{X})$ , for all  $x \in \mathcal{X}$ , we have  $x^\top A(\lambda_0; \mathcal{X})^{-1}x \leq 2x^\top A(\lambda_2; \mathcal{X})^{-1}x$ . As  $\lambda_2$  is a G-optimal design, we have  $\|x\|_{A(\lambda_0; \mathcal{X})^{-1}}^2 \leq 2d$ .

As  $\frac{1}{2}A(\lambda_1; \mathcal{X}) \prec A(\lambda_0; \mathcal{X})$ , due to Sudakov-Fernique inequality, we have

$$w(\lambda_0, \mathcal{X}, \mathcal{R}) \leq \sqrt{2} \cdot w(\lambda_1, \mathcal{X}, \mathcal{R}) \leq \sqrt{2} \cdot w(\mathcal{X}).$$

Now consider a fixed design  $x_1, x_2, \dots, x_T \in \mathcal{X}$  such that  $\tau(A_T, \mathcal{R}) \leq 2\tau(A(\lambda_0; \mathcal{X}), \mathcal{R})$  where  $A_T := \frac{1}{T} \sum_{i=1}^T x_i x_i^\top$  and  $\tau(A, \mathcal{R}) := \max_{i \in [d]} \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} [\max_{x \in R_i} x^\top A^{-1/2} \eta]^2 + 2 \max_{x \in \mathcal{X}} \|x\|_{A^{-1}}^2 \log(4/\delta)$ . Such a fixed design exists for any  $T \geq 180d$  due to Lemma 17. Due to the calculations above, we have  $\tau(A_T) \leq 4w(\lambda_1, \mathcal{X}, \mathcal{R})^2 + 8d \log(4/\delta)$ .

Let  $y_t = \langle x_t, \theta \rangle + \eta_t$  denote the noisy rewards for our fixed design where  $\eta_t \sim \mathcal{N}(0, 1)$ . Let  $\hat{\theta} = (\sum_{t=1}^T x_t x_t^\top)^{-1} \sum_{t=1}^T x_t y_t$ . Now we observe that  $\hat{\theta}$  is distributionally equivalent to  $\theta + (\sum_{t=1}^T x_t x_t^\top)^{-1/2} \eta$  where  $\eta \sim \mathcal{N}(0, I_d)$ . This enables us to apply Borell-TIS inequality on the gaussian process  $V_x = x^\top (\hat{\theta} - \theta)$  over  $R_{i_*}$  and obtain the following with probability  $1 - \delta/2$ :

$$\left| \max_{x \in R_{i_*}} \langle x, \hat{\theta} - \theta \rangle \right| \leq \frac{1}{\sqrt{T}} \cdot \mathbb{E}_{\eta \sim \mathcal{N}(0, I_d)} \left[ \max_{x \in R_{i_*}} x^\top A_T^{-1/2} \eta \right] + \frac{1}{\sqrt{T}} \cdot \sqrt{2 \max_{x \in \mathcal{X}} \|x\|_{A_T^{-1}}^2 \log(4/\delta)} \leq \sqrt{\frac{2\tau(A_T)}{T}}$$

Let us compute  $x^{(1)}, x^{(2)}, \dots, x^{(d)}$  such that  $x^{(i)} \in \arg \max_{x \in R_{i_*}} \langle x, \hat{\theta} \rangle$ . If we choose  $T = 1440(\frac{d}{\varepsilon^2} \log(4/\delta) + w(\lambda_1, \mathcal{X}, \mathcal{R})^2/\varepsilon^2)$ , then  $|\max_{x \in R_{i_*}} \langle x, \hat{\theta} - \theta \rangle| \leq \varepsilon/4$  with probability  $1 - \delta/2$ . This implies that  $\langle x^{(i_*)}, \theta \rangle \geq \langle x_*, \theta \rangle - \varepsilon/2$  with probability  $1 - \delta/2$ .

If we run the Median Elimination algorithm from [Even-Dar et al. \(2002\)](#) by treating each  $x^{(i)}$  as an arm of a stochastic MAB instance with mean  $\langle x^{(i)}, \theta \rangle$ , we get an index  $\hat{i}$  after additional  $\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples such that with probability  $1 - \delta/2$ , we have  $\max_{i \in [d]} \langle x^{(i)} - x^{(\hat{i})}, \theta \rangle \leq \varepsilon/2$ . We output  $x^{(\hat{i})}$  as the candidate best arm for  $\mathcal{X}$  and due to union bound, we have  $\langle x_* - x^{(\hat{i})}, \theta \rangle \leq \varepsilon$ .

In certain cases, if the partition is chosen appropriately, our adaptive algorithm can improve the sample complexity by logarithmic factors. For instance, consider a finite set  $\mathcal{X} \subset \mathbb{R}^d$ . For simplicity of presentation, let us assume that  $|\mathcal{X}|$  is multiple of  $d$ . Let  $\mathcal{R} = (R_1, \dots, R_d)$  be a partition of  $\mathcal{X}$  such that for all  $i \in [d]$ , we have  $|R_i| = |\mathcal{X}|/d$ . Using [Proposition 25](#), we can show that  $w(\lambda_1, \mathcal{X}, \mathcal{R}) \leq O(\sqrt{d \log(|\mathcal{X}|/d)})$ . This implies that our adaptive algorithm has an improved sample complexity of  $O\left(\frac{d \log(|\mathcal{X}|/d)}{\varepsilon^2}\right)$ .

## Appendix D. $\ell_2$ Norm Estimation Algorithm

Let  $\mathcal{X}$  denote the unit  $\ell_2$  ball in  $\mathbb{R}^d$ , and let  $\theta \in \mathbb{R}^d$  be the reward vector. Our goal is to estimate  $r := \|\theta\|_2$  using only samples obtained by querying points in  $\mathcal{X}$ . In [Appendix D.1](#), we analyze [Algorithm 1](#). In [Appendix D.2](#), we study how to estimate  $\|\theta\|_2$  up to a constant multiplicative factor. In [Appendix D.3](#), we consider the regime where  $\|\theta\|_2 \geq \sqrt{d}$ . Finally, in [Appendix D.4](#), we analyze a meta-algorithm that handles all regimes and provides a universal guarantee.

### D.1. Estimating with the Help of a Multiplicative Estimate $r_0$

Begin by defining  $r := \|\theta\|_2$ . We now assume that we are given a value  $r_0$  such that  $\varepsilon < r_0 < 2\sqrt{d}$  and  $\frac{r}{2} < r_0 \leq 2r$ . In the next section, we show how to satisfy these assumptions. Now we analyse [Algorithm 1](#).

In this section we work with sub-exponential random variables. A mean zero  $U$  is sub-exponential with parameters  $(V, b)$  if

$$\mathbb{E}[e^{\lambda U}] \leq \exp\left(\frac{\lambda^2 V}{2}\right) \quad \text{for } |\lambda| \leq \frac{1}{b}$$

If  $X \sim \mathcal{N}(\mu, \sigma^2)$  and  $Y := X^2 - \mathbb{E}[X^2]$ , then  $Y$  is a sub-exponential with parameters  $(V, b) = (8(\sigma^4 + \mu^2 \sigma^2), 4\sigma^2)$  due to [Lemma 9](#).

If  $U_1, \dots, U_K$  are independent mean-zero, sub-exponential random variables with parameters  $(V_1, b), \dots, (V_K, b)$  respectively, then we get the following due to [Lemma 10](#):

$$\Pr\left(\left|\frac{1}{K} \sum_{s=1}^K U_s\right| > t\right) \leq 2 \exp\left(-cK \min\left\{\frac{t^2}{\bar{V}}, \frac{t}{b}\right\}\right)$$

where  $\bar{V} = \frac{1}{K} \sum_{k=1}^K V_k$  and  $c > 0$  is some absolute constant.

We now begin with some high-level analysis. Recall  $x^{(k)} = \frac{(\varepsilon_1, \dots, \varepsilon_d)}{\sqrt{d}}$  where  $\varepsilon_i \stackrel{iid}{\sim} \{\pm 1\}$  and  $\mu_k = \langle x^{(k)}, \theta \rangle$ . Then we have the following:

$$\mathbb{E}[\mu_k^2] = \theta^\top \mathbb{E} \left[ x^{(k)} (x^{(k)})^\top \right] \theta = \theta^\top \left( \frac{1}{d} I_d \right) \theta = \frac{r^2}{d}$$

Now conditioning on  $x^{(k)}$  (hence  $\mu_k$ ), we have

$$\bar{y}_k \sim \mathcal{N} \left( \mu_k, \frac{1}{s} \right) \Rightarrow \mathbb{E}[\bar{y}_k^2 \mid \mu_k] = \mu_k^2 + \frac{1}{s} \Rightarrow \mathbb{E}[Z_k \mid \mu_k] = d\mu_k^2$$

We will now aim to bound

$$\bar{Z} - r^2 = \underbrace{\left( \frac{1}{K} \sum_{k=1}^K d\mu_k^2 - r^2 \right)}_{\text{Term 1}} + \underbrace{\left( \frac{1}{K} \sum_{k=1}^K (Z_k - d\mu_k^2) \right)}_{\text{Term 2}}$$

**Bounding Term 1:** Note that if  $r = 0$ , then Term I =  $d\mu_k^2 - r^2 = 0$ . Hence, the non-trivial case is when  $r > 0$  which we focus on below. Let  $X_k := d\mu_k^2 - r^2$ . We now show that  $X_k$  is sub-exponential with parameters

$$(V, b) = (Cr^4, Cr^2)$$

for some absolute constant  $C$  with the help of Hanson–Wright inequality (see Lemma 11).

Recall that

$$\mu_k = \left\langle \frac{1}{\sqrt{d}} \varepsilon^{(k)}, \theta \right\rangle, \quad \varepsilon_i^{(k)} \stackrel{iid}{\sim} \{\pm 1\}$$

Now observe that:

$$d\mu_k^2 = \left( \varepsilon^{(k)} \right)^\top A \varepsilon^{(k)}, \quad \text{where } A := \theta \theta^\top.$$

As  $\mathbb{E}[d\mu_k^2] = r^2$ , we apply the Hanson–Wright inequality (we apply Lemma 11 with  $K = 1$ , as Rademacher variables have unit  $\psi_2$ -norm) to get:

$$\Pr(|X_k| > t) \leq 2 \exp \left( -c \min \left( \frac{t^2}{\|A\|_F^2}, \frac{t}{\|A\|} \right) \right)$$

As  $A = \theta \theta^\top$ , we have

$$\|A\|_F^2 = \sum_i \theta_i^4 + 2 \sum_{i < j} \theta_i^2 \theta_j^2 = \left( \sum_i \theta_i^2 \right)^2 = \|\theta\|_2^4 = r^4$$

Note that  $A$  is a rank-1 matrix, and

$$A\theta = \theta \theta^\top \theta = \|\theta\|_2^2 \theta$$

Hence,  $\theta$  is an eigenvector of  $A$  with eigenvalue  $\|\theta\|_2^2$ .

As the operator norm is equal to the largest eigenvalue, we have

$$\|A\| = \|\theta\|_2^2 = r^2$$

Hence, we have:

$$\Pr(|X_k| > t) \leq 2 \exp \left( -c \min \left( \frac{t^2}{r^4}, \frac{t}{r^2} \right) \right)$$

Using the above inequality and applying Lemma 22, we get that  $X_k$  is sub-exponential with parameters  $(Cr^4, Cr^2)$  for some constant  $C$ .

Hence, we have:

$$\Pr \left( \left| \frac{1}{K} \sum_{k=1}^K X_k \right| > t \right) \leq 2 \exp \left( -\frac{c}{C} \cdot K \min \left( \frac{t^2}{r^4}, \frac{t}{r^2} \right) \right) \quad (9)$$

Choosing  $t = \frac{r\varepsilon}{4}$ , and  $K = c_1 r_0^2 \varepsilon^{-2} \log \frac{4}{\delta}$  for some large constant  $c_1$ , we get:

$$\left| \frac{1}{K} \sum_{k=1}^K d\mu_k^2 - r^2 \right| \leq \frac{r\varepsilon}{4} \quad \text{with probability} \geq 1 - \frac{\delta}{4}$$

**Bounding Term 2:** Let  $W_k := Z_k - d\mu_k^2 = d(\bar{y}_k^2 - \frac{1}{s} - \mu_k^2)$ . Conditioning on  $\mu_k$ , we have:

$$\sqrt{d} \cdot \bar{y}_k \mid \mu_k \sim \mathcal{N}(\sqrt{d} \cdot \mu_k, \sigma^2), \quad \text{with } \sigma^2 = \frac{d}{s}$$

Hence  $W_k \mid \mu_k$  is sub-exponential with parameters:

$$V(\mu_k) = 8 \left( \frac{d^2}{s^2} + \mu_k^2 \cdot \frac{d^2}{s} \right), \quad b = \frac{4d}{s}$$

Conditioning on  $\mu_{1:K} := \mu_1, \dots, \mu_K$ , we get:

$$\Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \mid \mu_{1:K} \right) \leq 2 \exp \left( -cK \min \left( \frac{t^2}{\bar{V}}, \frac{t}{b} \right) \right)$$

where  $\bar{V} := \frac{1}{K} \sum_{k=1}^K V(\mu_k) = 8 \left( \frac{d^2}{s^2} + \frac{d^2}{s} \cdot \frac{1}{K} \sum_{k=1}^K \mu_k^2 \right)$ . We now aim to upper bound  $\bar{V}$ .

Recall that  $X_k = d\mu_k^2 - r^2$  is sub-exponential with parameters  $(Cr^4, Cr^2)$ .

Hence, plugging  $t = r^2/2$  into Eq. (9) yields the following, for a sufficiently large constant  $c_1$ :

$$\Pr \left( \left| \frac{1}{K} \sum_{k=1}^K d\mu_k^2 - r^2 \right| > \frac{r^2}{2} \right) \leq \frac{\delta}{8}$$

Hence, with probability at least  $1 - \frac{\delta}{8}$ , we have  $\bar{V} \leq \tau := 8 \left( \frac{d^2}{s^2} + \frac{3dr^2}{2s} \right)$ .

$$\begin{aligned} \Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \right) &= \Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t, \bar{V} > \tau \right) + \Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t, \bar{V} \leq \tau \right) \\ &\leq \Pr(\bar{V} > \tau) + \Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \mid \bar{V} \leq \tau \right) \cdot \Pr(\bar{V} \leq \tau) \\ &\leq \frac{\delta}{8} + \underbrace{\Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \mid \bar{V} \leq \tau \right)}_{\text{Term III}} \end{aligned}$$

Now we have:

$$\text{Term III} = \Pr \left( \left| \frac{1}{K} \sum_{k=1}^K W_k \right| > t \mid \bar{V} \leq \tau \right) \leq 2 \exp \left( -cK \min \left( \frac{t^2}{\tau}, \frac{t}{b} \right) \right) \quad (10)$$

As  $s = \frac{c_0 d}{r_0^2}$ , we have the following for some absolute constant  $C'$ :

$$b = \left( \frac{4d}{s} \right) = \left( \frac{4}{c_0} \cdot r_0^2 \right) \quad \text{and} \quad \tau = 8 \left( \frac{r_0^4}{c_0^2} + \frac{3r^2 r_0^2}{2c_0} \right) \leq C' r_0^4$$

Also, with

$$K = c_1 r_0^2 \varepsilon^{-2} \log \frac{4}{\delta}, \quad t = \frac{r\varepsilon}{4} \quad \text{and} \quad \frac{r}{2} \leq r_0 \leq 2r$$

we get that the Term III is at most  $\frac{\delta}{8}$  after plugging the above parameters into Eq. (10).

Hence, with probability at least  $1 - \delta/2$ , we have  $|\bar{Z} - r^2| \leq \frac{r\varepsilon}{2}$ .

Let us now assume that  $|\bar{Z} - r^2| \leq \frac{r\varepsilon}{2}$ . As  $\bar{Z} \geq 0$  we have  $\hat{r} = \sqrt{\bar{Z}}$ , which implies the following:

$$|\hat{r} - r| = \left| \sqrt{\bar{Z}} - r \right| = |\bar{Z} - r^2| \Big/ \left( \sqrt{\bar{Z}} + r \right) \leq \frac{\frac{r\varepsilon}{2}}{r} = \frac{\varepsilon}{2} < \varepsilon$$

**D.2. Estimating  $r$  up to a Constant Factor**


---

**Algorithm 2:** Adaptive multi-scale test for  $r = \|\theta\|_2$ 


---

**Input :**  $\varepsilon \in (0, 1]$ ,  $\delta \in (0, 1/3)$ , dimension  $d$ , large absolute constants  $c_0, c_1 > 0$ **Output:** An estimate  $r_0$ 

```

1 Function Statistic ( $t_j, \delta_j$ )
2    $s_j \leftarrow c_0 \frac{d}{t_j^2}$ 
3    $K_j \leftarrow c_1 \log\left(\frac{1}{\delta_j}\right)$ 
4   for  $k \leftarrow 1$  to  $K_j$  do
5     Draw a Rademacher unit vector  $x^{(k)} = \frac{(\varepsilon_1, \dots, \varepsilon_d)}{\sqrt{d}}$  where  $\varepsilon_i \stackrel{iid}{\sim} \{\pm 1\}$ 
6     for  $\ell \leftarrow 1$  to  $s_j$  do
7       Observe  $y_{k,\ell} = \langle x^{(k)}, \theta \rangle + \eta_{k,\ell}$  where  $\eta_{k,\ell} \sim \mathcal{N}(0, 1)$  i.i.d.
8     end
9      $\bar{y}_k \leftarrow \frac{1}{s_j} \sum_{\ell=1}^{s_j} y_{k,\ell}$ 
10     $Z_k \leftarrow d \left( \bar{y}_k^2 - \frac{1}{s} \right)$ 
11  end
12   $U(t_j) \leftarrow \frac{1}{K_j} \sum_{k=1}^{K_j} Z_k$ 
13  return  $U(t_j)$ 
14 end

15 Function Test ( $t_j, \delta_j$ )
16   $U(t_j) \leftarrow$  Statistic ( $t_j, \delta_j$ )
17  if  $U(t_j) \geq \frac{3}{2} t_j^2$  then
18    return  $H_1$  //  $H_1$  is the hypothesis that  $r \geq 2t_j$ 
19  else
20    return  $H_0$  //  $H_0$  is the hypothesis that  $r \leq t_j$ 
21  end
22 end

23  $j \leftarrow 0$ 
24 while  $2^j \cdot \varepsilon < 2\sqrt{d}$  do
25    $t_j \leftarrow 2^j \cdot \varepsilon$ 
26    $\delta_j \leftarrow \frac{\delta}{2^{j+2}}$ 
27   outcome  $\leftarrow$  Test ( $t_j, \delta_j$ )
28   if outcome =  $H_0$  then
29     break
30   end
31    $j \leftarrow j + 1$ 
32 end
33  $r_0 \leftarrow t_j$ 
34 return  $r_0$ 

```

---

We aim to output  $r_0$  such that, with high probability, it satisfies meaningful properties that depend on the value of  $r$ . We discuss these properties in detail toward the end of this section.

Now we analyse the algorithm. Fix  $j \geq 0$ . Recall that under  $\text{Test}(t_j, \delta_j)$ , we return the hypothesis  $H_1$  ( $r \geq 2t_j$ ) if  $U(t_j) \geq \frac{3}{2}t_j^2$  otherwise we return the hypothesis  $H_0$  ( $r \leq t_j$ ).

We now compute the probability of error under the hypothesis  $H_0$  ( $r \leq t_j$ ). Error under the hypothesis  $H_0$  ( $r \leq t_j$ ) means:

$$U(t_j) \geq \frac{3}{2}t_j^2 \quad \text{whereas} \quad \mathbb{E}[U(t_j)] = r^2 \leq t_j^2$$

This requires an upward deviation of at least:

$$U(t_j) - \mathbb{E}[U(t_j)] \geq \Delta_0 := \frac{3}{2}t_j^2 - r^2 \geq \frac{1}{2}t_j^2$$

Recall:

$$U(t_j) - r^2 = \left( \frac{1}{K_j} \sum_k d\mu_k^2 - r^2 \right) + \frac{1}{K_j} \sum_k (Z_k - d\mu_k^2)$$

Let:

$$a_0 := \frac{1}{8}t_j^2, \quad \nu_0 := \frac{1}{8}t_j^2$$

Then the error under  $H_0$  satisfies:

$$\{\text{error under } H_0\} \subseteq \left\{ \left| \frac{1}{K_j} \sum_k d\mu_k^2 - r^2 \right| > a_0 \right\} \cup \left\{ \left| \frac{1}{K_j} \sum_k (Z_k - d\mu_k^2) \right| > \nu_0 \right\}$$

Recall that  $X_k := d\mu_k^2 - r^2$  is sub-exponential with parameters

$$(V, b) = (Cr^4, Cr^2)$$

for some absolute constant  $C$ . Now observe that

Hence, we have the following for some absolute constant  $c$ :

$$\Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} X_k \right| > a_0 \right) \leq 2 \exp \left( -cK_j \min \left( \frac{a_0^2}{r^4}, \frac{a_0}{r^2} \right) \right)$$

Note that:

$$\min \left\{ \frac{a_0^2}{r^4}, \frac{a_0}{r^2} \right\} = \min \left\{ \frac{t_j^4}{64r^4}, \frac{t_j^2}{8r^2} \right\} \geq \frac{1}{64}$$

As  $K_j = c_1 \log \left( \frac{1}{\delta_j} \right)$ , then for a large universal constant  $c_1$ , we have the following:

$$\left| \frac{1}{K_j} \sum_{k=1}^{K_j} d\mu_k^2 - r^2 \right| \leq a_0 \quad \text{with probability} \geq 1 - \frac{\delta_j}{4}$$

Let  $W_k := Z_k - d\mu_k^2$

Conditioning on  $\mu_k$ , we have:

$$\sqrt{d} \cdot \bar{y}_k \mid \mu_k \sim \mathcal{N}(\sqrt{d} \cdot \mu_k, \sigma^2), \quad \text{with } \sigma^2 = \frac{d}{s_j}$$

Recall  $W_k \mid \mu_k$  is sub-exponential with parameters:

$$V(\mu_k) = 8 \left( \frac{d^2}{s_j^2} + \mu_k^2 \cdot \frac{d^2}{s_j} \right), \quad b = \frac{4d}{s_j}$$

and therefore we had the following for some absolute constant  $c$ :

$$\Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} W_k \right| > \nu_0 \right) \leq \Pr(\bar{V} > \tau) + \underbrace{2 \exp \left( -cK_j \min \left( \frac{\nu_0^2}{\tau}, \frac{\nu_0}{b} \right) \right)}_{\text{Term I}}$$

where

$$\bar{V} := \frac{1}{K_j} \sum_{k=1}^{K_j} V(\mu_k), \quad \nu_0 = \frac{t_j^2}{8}, \quad b = \frac{4t_j^2}{c_0}, \quad \tau := 8 \left( \frac{d^2}{s_j^2} + \frac{3dr^2}{2s_j} \right) \leq C' t_j^4$$

for some absolute constant  $C'$ .

Recall that  $X_k = d\mu_k^2 - r^2$  is sub-exponential with parameters  $(Cr^4, Cr^2)$ . As  $K_j = c_1 \log \left( \frac{1}{\delta_j} \right)$ , we have for some absolute constant  $c$  and for a large universal constant  $c_1$ :

$$\Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} d\mu_k^2 - r^2 \right| > \frac{r^2}{2} \right) \leq 2 \exp(-cK_j) \leq \frac{\delta_j}{4}$$

Hence, we have  $\Pr(\bar{V} > \tau) \leq \frac{\delta_j}{4}$ .

As  $\frac{\nu_0}{b} = \frac{c_0}{32}, \frac{\nu_0^2}{\tau} \geq \frac{1}{64C'}$ , and  $K_j = c_1 \log \left( \frac{1}{\delta_j} \right)$ , then for a large universal constant  $c_1$ , we have the following:

$$\Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} W_k \right| > \nu_0 \right) \leq \Pr(\bar{V} > \tau) + \text{Term I} \leq \frac{\delta_j}{4} + \frac{\delta_j}{4} = \frac{\delta_j}{2}.$$

Hence, we have the following:

$$\Pr(\text{error under } H_0) \leq \Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} X_k \right| > a_0 \right) + \Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} W_k \right| > \nu_0 \right) \leq \frac{3}{4} \delta_j$$

We next compute the probability of error under the hypothesis  $H_1$  ( $r \geq 2t_j$ ). Error under the hypothesis  $H_1$  ( $r \geq 2t_j$ ) means:

$$U(t_j) \leq \frac{3}{2} t_j^2 \quad \text{whereas} \quad \mathbb{E}[U(t_j)] = r^2 \geq 4t_j^2$$

This requires a downward deviation of at least:

$$\mathbb{E}[U(t_j)] - U(t_j) \geq \Delta_1 := r^2 - \frac{3}{2} t_j^2 \geq \frac{5}{8} r^2$$

Recall:

$$U(t_j) - r^2 = \left( \frac{1}{K_j} \sum_k d\mu_k^2 - r^2 \right) + \frac{1}{K_j} \sum_k (Z_k - d\mu_k^2)$$

Let:

$$a_1 := \frac{1}{8}r^2, \quad \nu_1 := \frac{1}{8}r^2$$

Then the error under  $H_1$  satisfies:

$$\{\text{error under } H_1\} \subseteq \left\{ \left| \frac{1}{K_j} \sum_k d\mu_k^2 - r^2 \right| > a_1 \right\} \cup \left\{ \left| \frac{1}{K_j} \sum_k (Z_k - d\mu_k^2) \right| > \nu_1 \right\}$$

Recall that  $X_k := d\mu_k^2 - r^2$  is sub-exponential with parameters

$$(V, b) = (Cr^4, Cr^2)$$

for some absolute constant  $C$ . Now observe that

Hence, we have the following for some absolute constant  $c$ :

$$\Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} X_k \right| > a_1 \right) \leq 2 \exp \left( -cK_j \min \left( \frac{a_1^2}{r^4}, \frac{a_1}{r^2} \right) \right)$$

Note that:

$$\min \left\{ \frac{a_1^2}{r^4}, \frac{a_1}{r^2} \right\} = \min \left\{ \frac{r^4}{64r^4}, \frac{r^2}{8r^2} \right\} = \frac{1}{64}$$

As  $K_j = c_1 \log \left( \frac{1}{\delta_j} \right)$ , then for a large universal constant  $c_1$ , we have the following:

$$\left| \frac{1}{K_j} \sum_{k=1}^{K_j} d\mu_k^2 - r^2 \right| \leq a_1 \quad \text{with probability} \geq 1 - \frac{\delta_j}{4}$$

Let  $W_k := Z_k - d\mu_k^2$

Conditioning on  $\mu_k$ , we have:

$$\sqrt{d} \cdot \bar{y}_k \mid \mu_k \sim \mathcal{N}(\sqrt{d} \cdot \mu_k, \sigma^2), \quad \text{with } \sigma^2 = \frac{d}{s_j}$$

Recall  $W_k \mid \mu_k$  is sub-exponential with parameters:

$$V(\mu_k) = 8 \left( \frac{d^2}{s_j^2} + \mu_k^2 \cdot \frac{d^2}{s_j} \right), \quad b = \frac{4d}{s_j}$$

and therefore we had the following for some absolute constant  $c$ :

$$\Pr \left( \left| \frac{1}{K_j} \sum_{k=1}^{K_j} W_k \right| > \nu_1 \right) \leq \Pr(\bar{V} > \tau) + \underbrace{2 \exp \left( -cK_j \min \left( \frac{\nu_1^2}{\tau}, \frac{\nu_1}{b} \right) \right)}_{\text{Term II}}$$

where

$$\bar{V} := \frac{1}{K_j} \sum_{k=1}^{K_j} V(\mu_k), \quad \nu_1 = \frac{r^2}{8}, \quad b = \frac{4t_j^2}{c_0} \leq \frac{r^2}{c_0}, \quad \tau := 8 \left( \frac{d^2}{s_j^2} + \frac{3dr^2}{2s_j} \right) \leq C'r^4$$

for some absolute constant  $C'$ .

Recall that  $X_k = d\mu_k^2 - r^2$  is sub-exponential with parameters  $(Cr^4, Cr^2)$ . As  $K_j = c_1 \log\left(\frac{1}{\delta_j}\right)$ , we have the following for some absolute constant  $c$  and for a large universal constant  $c_1$ :

$$\Pr\left(\left|\frac{1}{K_j} \sum_{k=1}^{K_j} d\mu_k^2 - r^2\right| > \frac{r^2}{2}\right) \leq 2 \exp(-cK_j) \leq \frac{\delta_j}{4}$$

Hence, we have  $\Pr(\bar{V} > \tau) \leq \frac{\delta_j}{4}$ .

As  $\frac{\nu_1}{b} \geq \frac{c_0}{8}$ ,  $\frac{\nu_1^2}{\tau} \geq \frac{1}{64C'}$ , and  $K_j = c_1 \log\left(\frac{1}{\delta_j}\right)$ , then for a large universal constant  $c_1$ , we have the following:

$$\Pr\left(\left|\frac{1}{K_j} \sum_{k=1}^{K_j} W_k\right| > \nu_1\right) \leq \Pr(\bar{V} > \tau) + \text{Term II} \leq \frac{\delta_j}{4} + \frac{\delta_j}{4} = \frac{\delta_j}{2}.$$

Hence, we have:

$$\Pr(\text{error under } H_1) \leq \frac{\delta_j}{2} + \frac{\delta_j}{4} = \frac{3}{4}\delta_j$$

We now establish the guarantees of our algorithm. We divide it into 4 cases.

**Case 1:**  $r \leq \varepsilon$ : In this case, for  $j = 0$ , we have  $r \leq t_j$  which implies that the hypothesis  $H_0$  holds. Which implies with probability at least  $1 - \frac{3\delta}{16}$ , we return  $r_0 = \varepsilon$ .

**Case 2:**  $\varepsilon < r < \sqrt{d}$ : In this case, for all  $j$  such that  $2t_j \leq r$ ,  $H_1$  always holds. Consider the index  $j_*$  such that  $t_{j_*} \leq r < 2t_{j_*}$ . If we return  $r_0 = t_{j_*}$ , then  $r/2 < r_0 \leq r$ . If we instead proceed the index  $j_* + 1$ , then  $H_0$  holds and if we terminate and return  $r_0 = t_{j_*+1} = 2t_{j_*}$ , then  $r < r_0 \leq 2r$ . Hence, with probability at least  $1 - \sum_{j=1}^{\infty} 3\delta_j/4 = 1 - \frac{3\delta}{8}$ , we return  $r/2 < r_0 \leq 2r$ .

**Case 3:**  $\sqrt{d} \leq r < 4\sqrt{d}$ : In this case, for all  $j$  such that  $2t_j \leq r$ ,  $H_1$  always holds. If there exists an index  $j_*$  such that  $t_{j_*} < 2\sqrt{d}$  and  $t_{j_*} \leq r < 2t_{j_*}$  and if we return  $r_0 = t_{j_*}$ , then  $r/2 < r_0 \leq r$ . If we instead proceed to the index  $j_* + 1$ , then  $H_0$  holds and if we terminate and return  $r_0 = t_{j_*+1} = 2t_{j_*}$ , then  $r < r_0 \leq 2r$ . Hence, in this scenario, with probability at least  $1 - \sum_{j=1}^{\infty} 3\delta_j/4 = 1 - \frac{3\delta}{8}$ , we return  $r/2 < r_0 \leq 2r$ .

On the other hand, if there is no index  $j_*$  such that  $t_{j_*} < 2\sqrt{d}$  and  $t_{j_*} \leq r < 2t_{j_*}$ , then  $H_1$  always holds for all  $j$  such that  $2^j \varepsilon < 2\sqrt{d}$ . This implies that with probability at least  $1 - \sum_{j=1}^{\infty} 3\delta_j/4 = 1 - \frac{3\delta}{8}$ , we return  $r_0 \geq 2\sqrt{d}$ .

**Case 4:**  $r \geq 4\sqrt{d}$ : In this case, for all  $j$  such that  $2^j \varepsilon < 2\sqrt{d}$ ,  $H_1$  always holds. This implies that with probability at least  $1 - \sum_{j=1}^{\infty} 3\delta_j/4 = 1 - \frac{3\delta}{8}$ , we return  $r_0 \geq 2\sqrt{d}$ .

Now we establish the sample complexity of our algorithm. The sample complexity is upper bounded as

$$\sum_{j=0}^{\lceil \log_2(2\sqrt{d}/\varepsilon) \rceil} s_j t_j \leq \mathcal{O}\left(\sum_{j=0}^{\infty} (d/\varepsilon)^2 2^{-j} \log(2^{j+2}/\delta)\right) = \mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right).$$

### D.3. Estimation in the Large-Norm Regime

Recall that observe

$$y_t = \langle x_t, \theta \rangle + \eta_t, \quad t = 1, \dots, n,$$

where  $\theta \in \mathbb{R}^d$  and the noise variables are i.i.d.  $\eta_t \sim \mathcal{N}(0, 1)$ . Let  $r := \|\theta\|_2$ . In this section we assume the ‘‘large norm’’ condition

$$r^2 \geq d. \tag{11}$$

Our goal is to construct an estimator  $\hat{r}$  such that, for suitable absolute constant  $C > 0$ ,

$$n \geq C \frac{d}{\varepsilon^2} \log \frac{1}{\delta} \quad \text{and} \quad \Pr(|\hat{r} - r| > \varepsilon) \leq \delta$$

for all  $\theta$  satisfying (11). We first describe our algorithm below.

---

**Algorithm 3:** Algorithm for Large-Norm Regime

---

**Input** : dimension  $d$ , sample size  $n$

```

1 for  $t \leftarrow 1$  to  $n$  do
2   | Draw a Rademacher unit vector  $x_t = \frac{(\varepsilon_1, \dots, \varepsilon_d)}{\sqrt{d}}$  where  $\varepsilon_i \stackrel{iid}{\sim} \{\pm 1\}$ .
3   | Observe  $y_t$  from the model  $y_t = \langle x_t, \theta \rangle + \eta_t$  with  $\eta_t \sim \mathcal{N}(0, 1)$  i.i.d.
4 end
5 Form  $X \in \mathbb{R}^{n \times d}$  whose  $t$ -th row is  $x_t^\top$ , and  $y \leftarrow (y_1, \dots, y_n)^\top$ 
6 if  $X^\top X$  is invertible then
7   |  $\hat{\theta} \leftarrow (X^\top X)^{-1} X^\top y$ 
8   |  $\Sigma \leftarrow (X^\top X)^{-1}$ 
9 end
10 else
11   |  $\hat{r} \leftarrow \sqrt{d}$ 
12   | return  $\hat{r}$ 
13 end
14  $\hat{r} \leftarrow \sqrt{\max\{\|\hat{\theta}\|_2^2 - \text{tr}(\Sigma), 0\}}$ 
15 return  $\hat{r}$ 
    
```

---

We now start analysing our algorithm. Let us consider the case when  $X^\top X$  is invertible. We later show that this holds with high probability. Recall that the least-squares estimator is defined as:

$$\hat{\theta} := (X^\top X)^{-1} X^\top y.$$

Define

$$\Sigma := (X^\top X)^{-1}, \quad \Delta := \hat{\theta} - \theta = \Sigma X^\top \eta.$$

Conditioning on  $X$  and using that  $\eta \sim \mathcal{N}(0, I_n)$ , we have

$$\Delta \mid X \sim \mathcal{N}(0, \Sigma). \tag{12}$$

Conditioning on  $X$  and using (12), we get

$$\mathbb{E}[\Delta \mid X] = 0, \quad \mathbb{E}[\|\Delta\|_2^2 \mid X] = \text{tr}(\Sigma).$$

We wish to estimate  $r^2 = \|\theta\|_2^2$ . We first have the following:

$$\|\hat{\theta}\|_2^2 = \|\theta + \Delta\|_2^2 = \|\theta\|_2^2 + 2\theta^\top \Delta + \|\Delta\|_2^2 = r^2 + 2\theta^\top \Delta + \|\Delta\|_2^2.$$

Define the debiased estimator

$$\hat{R} := \|\hat{\theta}\|_2^2 - \text{tr}(\Sigma). \tag{13}$$

Then, we have

$$\mathbb{E}[\hat{R} \mid X] = r^2, \quad \mathbb{E}[\hat{R}] = r^2.$$

Let

$$Z := \widehat{R} - r^2 = 2\theta^\top \Delta + (\|\Delta\|_2^2 - \text{tr}(\Sigma)). \quad (14)$$

We will now derive a tail bound for  $Z$  conditional on  $X$ . Let  $g \sim \mathcal{N}(0, I_d)$ , and observe that  $\Delta$  is distributionally equivalent to  $\Sigma^{1/2}g$ . Then (14) becomes distributionally equivalent to

$$\tilde{Z} = 2\theta^\top \Sigma^{1/2}g + (g^\top \Sigma g - \text{tr}(\Sigma)).$$

Define

$$L := 2\theta^\top \Sigma^{1/2}g, \quad Q := g^\top \Sigma g - \text{tr}(\Sigma),$$

so that

$$\tilde{Z} = L + Q.$$

Conditioning on  $X$ , the random variable  $L$  is Gaussian with mean zero and variance

$$\text{Var}(L | X) = 4\|\Sigma^{1/2}\theta\|_2^2 = 4\theta^\top \Sigma \theta.$$

Therefore for any  $u > 0$ ,

$$\Pr(|L| \geq u | X) \leq 2 \exp\left(-\frac{u^2}{8\theta^\top \Sigma \theta}\right). \quad (15)$$

We now bound the term  $Q$ . We use the Hanson–Wright inequality in the Gaussian case: if  $g \sim \mathcal{N}(0, I_d)$  and  $A \in \mathbb{R}^{d \times d}$  is symmetric, then there exists an absolute constant  $c_0 > 0$  such that for all  $u > 0$ ,

$$\Pr(|g^\top A g - \text{tr}(A)| \geq u) \leq 2 \exp\left(-c_0 \min\left\{\frac{u^2}{\|A\|_F^2}, \frac{u}{\|A\|_{\text{op}}}\right\}\right). \quad (16)$$

Applying (16) with  $A = \Sigma$ , and using  $\|\Sigma\|_F^2 = \text{tr}(\Sigma^2)$ , we obtain

$$\Pr(|Q| \geq u | X) \leq 2 \exp\left(-c_0 \min\left\{\frac{u^2}{\text{tr}(\Sigma^2)}, \frac{u}{\|\Sigma\|_{\text{op}}}\right\}\right). \quad (17)$$

For any  $t > 0$ ,

$$\{|\tilde{Z}| \geq t\} = \{L + Q \geq t\} \subseteq \left\{|L| \geq \frac{t}{2}\right\} \cup \left\{|Q| \geq \frac{t}{2}\right\},$$

hence, by the union bound and the distributional equivalence between  $\tilde{Z}$  and  $Z$ , we get

$$\Pr(|Z| \geq t | X) \leq \Pr\left(|L| \geq \frac{t}{2} | X\right) + \Pr\left(|Q| \geq \frac{t}{2} | X\right). \quad (18)$$

Using (15) with  $u = t/2$ , we get

$$\Pr\left(|L| \geq \frac{t}{2} | X\right) \leq 2 \exp\left(-\frac{(t/2)^2}{8\theta^\top \Sigma \theta}\right) = 2 \exp\left(-\frac{t^2}{32\theta^\top \Sigma \theta}\right).$$

Using (17) with  $u = t/2$ , we get

$$\Pr\left(|Q| \geq \frac{t}{2} | X\right) \leq 2 \exp\left(-c_0 \min\left\{\frac{t^2}{4\text{tr}(\Sigma^2)}, \frac{t}{2\|\Sigma\|_{\text{op}}}\right\}\right).$$

Thus (18) implies

$$\Pr(|Z| \geq t | X) \leq 2 \exp\left(-\frac{t^2}{32\theta^\top \Sigma \theta}\right) + 2 \exp\left(-c_0 \min\left\{\frac{t^2}{4\text{tr}(\Sigma^2)}, \frac{t}{2\|\Sigma\|_{\text{op}}}\right\}\right). \quad (19)$$

Define

$$A_1 := \frac{t^2}{32\theta^\top \Sigma \theta}, \quad A_2 := c_0 \frac{t^2}{4 \operatorname{tr}(\Sigma^2)}, \quad A_3 := c_0 \frac{t}{2\|\Sigma\|_{\text{op}}}.$$

Then (19) can be written as

$$\Pr(|Z| \geq t \mid X) \leq 2e^{-A_1} + 2e^{-\min\{A_2, A_3\}} \leq 4e^{-\min\{A_1, A_2, A_3\}}.$$

For any  $a, b > 0$  we have

$$\min\left\{\frac{t^2}{a}, \frac{t^2}{b}\right\} = \frac{t^2}{\max\{a, b\}} \geq \frac{t^2}{a+b}.$$

Applying this with  $a = 32\theta^\top \Sigma \theta$  and  $b = \frac{4}{c_0} \operatorname{tr}(\Sigma^2)$ , and absorbing constants, we obtain for some absolute  $c > 0$ :

$$\Pr(|Z| \geq t \mid X) \leq 4 \exp\left(-c \min\left\{\frac{t^2}{4\theta^\top \Sigma \theta + 2 \operatorname{tr}(\Sigma^2)}, \frac{t}{\|\Sigma\|_{\text{op}}}\right\}\right). \quad (20)$$

Standard results on sub-Gaussian random matrices from the Section 4.7 of [Vershynin \(2018\)](#) gives the following. Since  $\mathbb{E}[x_t x_t^\top] = I_d/d$  and the rows are i.i.d. sub-Gaussian, there exists an absolute constant  $C_0 > 1$  such that if

$$n \geq C_0(d + \log \frac{2}{\delta}), \quad (21)$$

then with probability at least  $1 - \delta/2$ ,

$$\frac{1}{2d}I_d \preceq \frac{X^\top X}{n} \preceq \frac{2}{d}I_d. \quad (22)$$

Multiplying (22) by  $n$  and inverting the inequalities yields

$$\frac{d}{2n}I_d \preceq \Sigma \preceq \frac{2d}{n}I_d. \quad (23)$$

From (23) we obtain, on this ‘‘good event’’ for  $X$ ,

$$\|\Sigma\|_{\text{op}} = \lambda_{\max}(\Sigma) \leq \frac{2d}{n}, \quad (24)$$

$$\theta^\top \Sigma \theta \leq \|\Sigma\|_{\text{op}} \|\theta\|_2^2 \leq \frac{2d}{n} r^2, \quad (25)$$

$$\operatorname{tr}(\Sigma^2) = \sum_{i=1}^d \lambda_i^2 \leq d \lambda_{\max}^2 \leq d \left(\frac{2d}{n}\right)^2 = \frac{4d^3}{n^2}. \quad (26)$$

We now use (25) and (26) in (20). On the good event for  $X$ ,

$$\begin{aligned} 4\theta^\top \Sigma \theta + 2 \operatorname{tr}(\Sigma^2) &\leq 4 \cdot \frac{2d}{n} r^2 + 2 \cdot \frac{4d^3}{n^2} \\ &= \frac{8d}{n} r^2 + \frac{8d^3}{n^2}. \end{aligned}$$

Using  $r^2 \geq d$  and  $n \geq d$ , we have

$$\frac{d^3}{n^2} \leq \frac{d^2}{n} \leq \frac{d}{n} r^2,$$

hence

$$\frac{8d^3}{n^2} \leq 8 \frac{d}{n} r^2.$$

Therefore

$$4\theta^\top \Sigma \theta + 2 \operatorname{tr}(\Sigma^2) \leq \frac{8d}{n} r^2 + 8 \frac{d}{n} r^2 = 16 \frac{d}{n} r^2. \quad (27)$$

Let us choose

$$t = \varepsilon r.$$

On the good event,

$$\min \left\{ \frac{t^2}{4\theta^\top \Sigma \theta + 2 \operatorname{tr}(\Sigma^2)}, \frac{t}{\|\Sigma\|_{\text{op}}} \right\} \geq \frac{\varepsilon^2 r^2}{16 \frac{d}{n} r^2} = \frac{n\varepsilon^2}{16d}.$$

Using (20) we obtain the following for some absolute constant  $c' > 0$ :

$$\Pr(|Z| \geq \varepsilon r \mid X) \leq 4 \exp\left(-c' \frac{n\varepsilon^2}{d}\right) \quad \text{on the good event for } X. \quad (28)$$

If we choose

$$n \geq C_1 \frac{d}{\varepsilon^2} \log \frac{4}{\delta}$$

for a sufficiently large absolute constant  $C_1$ , then the right-hand side of (28) is at most  $\delta/2$ . Combining with (21) and the fact that the good event for  $X$  has probability at least  $1 - \delta/2$ , we conclude that

$$\Pr(|\widehat{R} - r^2| \geq \varepsilon r) \leq \delta. \quad (29)$$

Recall that our estimator for  $r$  is

$$\widehat{r} := \sqrt{\max\{\widehat{R}, 0\}}. \quad (30)$$

On the event

$$|\widehat{R} - r^2| \leq \varepsilon r$$

and assuming  $\varepsilon \leq r/2$  (which trivially holds for  $d \geq 4$ ), we have

$$\widehat{R} \geq r^2 - \varepsilon r \geq r^2 - \frac{r^2}{2} = \frac{r^2}{2} > 0,$$

so  $\widehat{r} = \sqrt{\widehat{R}}$ , and therefore

$$|\widehat{r} - r| = |\sqrt{\widehat{R}} - \sqrt{r^2}| = \frac{|\widehat{R} - r^2|}{\sqrt{\widehat{R}} + r} \leq \frac{|\widehat{R} - r^2|}{r} \leq \varepsilon.$$

Hence

$$\{|\widehat{R} - r^2| \leq \varepsilon r\} \subseteq \{|\widehat{r} - r| \leq \varepsilon\}.$$

Combining with (29), we obtain

$$\Pr(|\widehat{r} - r| > \varepsilon) \leq \Pr(|\widehat{R} - r^2| > \varepsilon r) \leq \delta.$$

Finally, we needed

$$n \geq C_0(d + \log \frac{2}{\delta}) \quad \text{and} \quad n \geq C_1 \frac{d}{\varepsilon^2} \log \frac{4}{\delta},$$

so for some absolute constant  $C_2 > 0$ ,

$$n \geq C_2 \frac{d}{\varepsilon^2} \log \frac{1}{\delta} \implies \Pr(|\widehat{r} - \|\theta\|_2| > \varepsilon) \leq \delta$$

for all  $\theta$  satisfying  $r^2 \geq d$ .

#### D.4. Meta Algorithm for $\ell_2$ -Norm Estimation

In section, we describe a meta algorithm that uses Algorithms 1, 2, and 3.

---

**Algorithm 4:** Meta-algorithm to estimate  $\|\theta\|_2^2$

---

**Input** : dimension  $d$

```

1 Run the Algorithm 2 and receive the estimate  $r_0$ 
2 if  $r_0 = \varepsilon$  then
3   |  $\hat{r} \leftarrow r_0$ 
4   | return  $\hat{r}$ 
5 end
6 if  $r_0 \geq 2\sqrt{d}$  then
7   | Run the Algorithm 3 and receive the estimate  $\hat{r}$ 
8   | return  $\hat{r}$ 
9 end
10 if  $\varepsilon < r_0 < 2\sqrt{d}$  then
11   | Run the Algorithm 1 and receive the estimate  $\hat{r}$ 
12   | return  $\hat{r}$ 
13 end
    
```

---

Now we claim that Algorithm 4 takes  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples and returns an estimate  $\hat{r}$  such that with probability at least  $1 - \delta$  we have  $|r - \hat{r}| \leq \varepsilon$ , where  $r := \|\theta\|_2$ .

First, consider the case where  $r \leq \varepsilon$ . We showed in Appendix D.2 that with probability at least  $1 - \delta$ , we have  $r_0 = \varepsilon$ . Hence, probability at least  $1 - \delta$  we have  $|r - \hat{r}| \leq \varepsilon$ .

Next, consider the case where  $\varepsilon < r < \sqrt{d}$ . We showed in Appendix D.2 that with probability at least  $1 - \delta/2$ , we have  $r/2 < r_0 \leq 2r$  and  $r_0 < 2\sqrt{d}$ . Conditioned on this good event, if  $r_0 = \varepsilon$ , we have  $\varepsilon < r < 2\varepsilon$  and therefore  $|\hat{r} - r| \leq \varepsilon$ . On the other hand, conditioned on this good event, if  $r_0 > \varepsilon$  Algorithm 1 is run and it returns  $\hat{r}$  such that probability at least  $1 - \delta/2$  we have  $|r - \hat{r}| \leq \varepsilon$ . Due to union bound, we return  $\hat{r}$  such that with probability at least  $1 - \delta$  we have  $|r - \hat{r}| \leq \varepsilon$ .

Next, consider the case where  $\sqrt{d} \leq r < 4\sqrt{d}$ . Due to the analysis in Appendix D.2, with probability at least  $1 - \delta/2$  we either run Algorithm 1 with an estimate  $r_0$  satisfying  $r/2 < r_0 \leq 2r$  or run the Algorithm 3. In either case, the algorithm being run returns  $\hat{r}$  such that probability at least  $1 - \delta/2$  we have  $|r - \hat{r}| \leq \varepsilon$ . Due to union bound, we return  $\hat{r}$  such that with probability at least  $1 - \delta$  we have  $|r - \hat{r}| \leq \varepsilon$ .

Finally, consider the case where  $r \geq 4\sqrt{d}$ . Due to the analysis in Appendix D.2, with probability at least  $1 - \delta/2$  we have  $r_0 \geq 2\sqrt{d}$ . Hence, conditioned on this good event, Algorithm 3 is run and it returns  $\hat{r}$  such that probability at least  $1 - \delta/2$  we have  $|r - \hat{r}| \leq \varepsilon$ . Due to union bound, we return  $\hat{r}$  such that with probability at least  $1 - \delta$  we have  $|r - \hat{r}| \leq \varepsilon$ .

The sample complexity guarantee follows as each algorithm run in the meta-algorithms takes  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2}\right)$  samples.

#### Appendix E. Tail to Sub-Exponential Technical Lemma

**Lemma 22** *Assume  $U$  is mean-zero and for all  $t \geq 0$  and some absolute constant  $c > 0$ , we have:*

$$\Pr(|U| > t) \leq 2 \exp\left(-c \min\left(\frac{t^2}{V}, \frac{t}{b}\right)\right).$$

*Then there exists an absolute constant  $C > 0$  such that  $U$  is sub-exponential with parameters*

$$(C(V + b^2), Cb),$$

i.e.

$$\mathbb{E}[e^{\lambda U}] \leq \exp\left(\frac{\lambda^2 C(V + b^2)}{2}\right) \quad \text{for all } |\lambda| \leq \frac{1}{Cb}.$$

**Proof** First, observe that for all  $t \geq 0$ , we have:

$$\Pr(|U| > t) \leq 2e^{-ct^2/V} + 2e^{-ct/b}. \quad (31)$$

For any integer  $k \geq 2$ ,

$$\mathbb{E}|U|^k = k \int_0^\infty t^{k-1} \Pr(|U| > t) dt,$$

so by (31), we have:

$$\mathbb{E}|U|^k \leq 2k \int_0^\infty t^{k-1} e^{-ct^2/V} dt + 2k \int_0^\infty t^{k-1} e^{-ct/b} dt.$$

Using the definition of gamma function, we have the following:

$$\int_0^\infty t^{k-1} e^{-ct^2/V} dt = \frac{1}{2} \left(\frac{V}{c}\right)^{k/2} \Gamma\left(\frac{k}{2}\right),$$

$$\int_0^\infty t^{k-1} e^{-ct/b} dt = \left(\frac{b}{c}\right)^k \Gamma(k).$$

Hence

$$\mathbb{E}|U|^k \leq k \left(\frac{V}{c}\right)^{k/2} \Gamma\left(\frac{k}{2}\right) + 2k \left(\frac{b}{c}\right)^k \Gamma(k). \quad (32)$$

Using the properties of gamma function from [Gubner \(2021\)](#), we have the following for all  $x \geq 1$ :

$$\Gamma(x) \leq 3x^{x-\frac{1}{2}}e^{-x}. \quad (33)$$

Apply (33) to  $x = k/2 \geq 1$  and  $x = k \geq 2$ :

$$\Gamma(k/2) \leq 3 \left(\frac{k}{2}\right)^{\frac{k}{2}-\frac{1}{2}} e^{-k/2}, \quad \Gamma(k) \leq 3k^{k-\frac{1}{2}}e^{-k}.$$

Plugging into (32) and taking  $k$ -th roots gives:

$$(\mathbb{E}|U|^k)^{1/k} \leq \left[ k \left(\frac{V}{c}\right)^{k/2} \Gamma(k/2) \right]^{1/k} + \left[ 2k \left(\frac{b}{c}\right)^k \Gamma(k) \right]^{1/k} \leq 2\sqrt{\frac{Vk}{c}} + \frac{2bk}{c} \quad (34)$$

From (34) we also get, using  $(x + y)^k \leq 2^{k-1}(x^k + y^k)$ ,

$$\mathbb{E}|U|^k \leq 2^{k-1} \left( (2\sqrt{Vk/c})^k + (2bk/c)^k \right) \quad (35)$$

Since  $\mathbb{E}U = 0$ ,

$$\mathbb{E}e^{\lambda U} = 1 + \sum_{k \geq 2} \frac{\lambda^k \mathbb{E}[U^k]}{k!} \leq 1 + \sum_{k \geq 2} \frac{|\lambda|^k \mathbb{E}|U|^k}{k!}.$$

Using (35), we get

$$\mathbb{E}e^{\lambda U} \leq 1 + \frac{1}{2} \sum_{k \geq 2} \frac{(4|\lambda|)^k \left( (V/c)^{k/2} k^{k/2} + (b/c)^k k^k \right)}{k!} = 1 + S_1 + S_2,$$

where

$$S_1 := \frac{1}{2} \sum_{k \geq 2} \frac{(4|\lambda|\sqrt{V/c})^k k^{k/2}}{k!}, \quad S_2 := \frac{1}{2} \sum_{k \geq 2} \frac{(4|\lambda|b/c)^k k^k}{k!}.$$

First, we bound the term  $1 + S_2$ . Use  $k! \geq (k/e)^k$ , so  $k^k/k! \leq e^k$ . Then

$$S_2 \leq \frac{1}{2} \sum_{k \geq 2} (4e|\lambda|b/c)^k.$$

If

$$|\lambda| \leq \frac{c}{8eb}, \tag{36}$$

then  $r := 4e|\lambda|b/c \leq 1/2$ , so

$$S_2 \leq \frac{1}{2} \sum_{k \geq 2} r^k = \frac{1}{2} \cdot \frac{r^2}{1-r} \leq r^2 \leq 16e^2 \frac{\lambda^2 b^2}{c^2}.$$

Hence, using  $1 + x \leq e^x$ , we get:

$$1 + S_2 \leq \exp\left(16e^2 \frac{\lambda^2 b^2}{c^2}\right). \tag{37}$$

Next, we bound the term  $1 + S_1$ . Let

$$a := 4|\lambda|\sqrt{V/c}.$$

For  $k \geq 2$ , letting  $m = \lfloor k/2 \rfloor$ , we have:

$$k! \geq m! \left(\frac{k}{2}\right)^{k-m}.$$

This implies

$$\frac{k^{k/2}}{k!} \leq \frac{k^{k/2}}{m!(k/2)^{k-m}} \leq \frac{2^{\lfloor k/2 \rfloor}}{m!} \leq \frac{2^{k/2+1}}{m!}. \tag{38}$$

Now we consider both case on  $k$ , one where  $k$  is odd and one where  $k$  is even:

- $k = 2m$  with  $m \geq 1$ : by (38), we have

$$\frac{a^{2m}(2m)^m}{(2m)!} \leq \frac{a^{2m}2^m}{m!}.$$

- $k = 2m + 1$  with  $m \geq 1$ : by (38), we have

$$\frac{a^{2m+1}(2m+1)^{m+1/2}}{(2m+1)!} \leq \frac{a^{2m+1}2^{m+1}}{m!}.$$

Therefore

$$S_1 \leq \frac{1}{2} \sum_{m \geq 1} \frac{(2a^2)^m}{m!} + a \sum_{m \geq 1} \frac{(2a^2)^m}{m!} = \left(a + \frac{1}{2}\right)(e^{2a^2} - 1). \quad (39)$$

Using the inequalities  $(x - \frac{1}{2})^2 + \frac{1}{4} \geq 0$  and  $e^x \geq 1 + x$ , we have the following:

$$a + \frac{1}{2} \leq 1 + a^2 \leq e^{a^2}. \quad (40)$$

Combining (39) and (40), we get:

$$S_1 \leq e^{a^2}(e^{2a^2} - 1) = e^{3a^2} - e^{a^2} \leq e^{3a^2} - 1,$$

so

$$1 + S_1 \leq e^{3a^2} = \exp\left(48 \frac{\lambda^2 V}{c}\right). \quad (41)$$

Since  $S_1, S_2 \geq 0$ ,

$$1 + S_1 + S_2 \leq (1 + S_1)(1 + S_2).$$

Using (37) and (41), whenever (36) holds, we have

$$\mathbb{E}e^{\lambda U} \leq \exp\left(48 \frac{\lambda^2 V}{c}\right) \exp\left(16e^2 \frac{\lambda^2 b^2}{c^2}\right) = \exp\left(\frac{\lambda^2}{2} \left(\frac{96}{c} V + \frac{32e^2}{c^2} b^2\right)\right).$$

Define

$$b' := \frac{8e}{c} b, \quad V' := \frac{96}{c} V + \frac{32e^2}{c^2} b^2.$$

Then (6) is exactly  $|\lambda| \leq 1/b'$ , and we have shown

$$\mathbb{E}[e^{\lambda U}] \leq \exp\left(\frac{\lambda^2 V'}{2}\right) \quad \text{for all } |\lambda| \leq \frac{1}{b'}.$$

So  $U$  is sub-exponential with parameters  $(V', b')$ .

Finally, if one desires a single constant multiplying both  $(V + b^2)$  and  $b$ , one can take

$$C := \max\left\{\frac{96}{c}, \frac{32e^2}{c^2}, \frac{8e}{c}\right\},$$

so that  $V' \leq C(V + b^2)$  and  $b' \leq Cb$ . This gives exactly the stated form:

$$\mathbb{E}[e^{\lambda U}] \leq \exp\left(\frac{\lambda^2 C(V + b^2)}{2}\right) \quad \text{for } |\lambda| \leq \frac{1}{Cb}.$$

■

## Appendix F. Gaussian Width Properties

For a design distribution  $\lambda$  over the set  $\mathcal{X}$ , let

$$\Sigma(\lambda) := \mathbb{E}_{x \sim \lambda}[xx^\top], \quad f(\lambda) := \mathbb{E}_{g \sim \mathcal{N}(0, I_d)} \left[ \sup_{x \in \mathcal{X}} \langle x, \Sigma(\lambda)^{-1/2} g \rangle \right].$$

Equivalently, we have

$$f(\lambda) = \mathbb{E}_{g \sim \mathcal{N}(0, I_d)} \left[ \sup_{x \in \mathcal{X}} \langle \Sigma(\lambda)^{-1/2} x, g \rangle \right].$$

Recall that  $w(\mathcal{X}) := \inf_{\lambda \in \Delta_{\mathcal{X}}} f(\lambda)$ . Now we prove the following propositions.

**Proposition 23**  $w(\mathcal{X}) \leq d$ .

**Proof**

Let  $\lambda_G$  be a G-optimal design, that is a distribution that minimizes  $\sup_{x \in \mathcal{X}} x^\top \Sigma(\lambda)^{-1} x$ .

$$\lambda_G \in \arg \min_{\lambda} \sup_{x \in \mathcal{X}} x^\top \Sigma(\lambda)^{-1} x.$$

For compact  $\mathcal{X}$  with  $\text{span}(\mathcal{X}) = \mathbb{R}^d$ , we have the following:

$$\sup_{x \in \mathcal{X}} x^\top \Sigma(\lambda_G)^{-1} x \leq d.$$

Equivalently,

$$\sup_{x \in \mathcal{X}} \|\Sigma(\lambda_G)^{-1/2} x\|_2 \leq \sqrt{d}.$$

For any fixed  $g$ , we have the following:

$$\sup_{x \in \mathcal{X}} \langle x, \Sigma(\lambda_G)^{-1/2} g \rangle = \sup_{x \in \mathcal{X}} \langle \Sigma(\lambda_G)^{-1/2} x, g \rangle \leq \left( \sup_{x \in \mathcal{X}} \|\Sigma(\lambda_G)^{-1/2} x\|_2 \right) \|g\|_2 \leq \sqrt{d} \|g\|_2.$$

Take expectation:

$$f(\lambda_G) \leq \sqrt{d} \mathbb{E} \|g\|_2 \leq \sqrt{d} \sqrt{\mathbb{E} \|g\|_2^2} = d.$$

Therefore, we have the following:

$$w(\mathcal{X}) = \inf_{\lambda \in \Delta_{\mathcal{X}}} f(\lambda) \leq f(\lambda_G) \leq d.$$

■

**Proposition 24** For a finite set  $\mathcal{X}$  with  $|\mathcal{X}| = m$ , we have  $w(\mathcal{X}) \leq O(\sqrt{d \log m})$ .

**Proof** Define  $u_x := \Sigma(\lambda_G)^{-1/2} x$ . Then  $\|u_x\|_2^2 \leq d$  for all  $x \in \mathcal{X}$ .

If  $g \sim \mathcal{N}(0, I_d)$ , then for each  $x \in \mathcal{X}$ ,  $Z_x := \langle u_x, g \rangle$  is Gaussian with variance  $\|u_x\|_2^2 \leq d$ , so for  $t \geq 0$ ,

$$\Pr(Z_x \geq t) \leq \exp\left(-\frac{t^2}{2d}\right).$$

Now by applying the union bound, we get the following:

$$\Pr\left(\max_{x \in \mathcal{X}} Z_x \geq t\right) \leq m \exp\left(-\frac{t^2}{2d}\right).$$

Using the fact  $\mathbb{E}[\max_{x \in \mathcal{X}} Z_x] = \int_0^\infty \Pr(\max_{x \in \mathcal{X}} Z_x \geq t) dt$  and choosing  $t_0 = \sqrt{2d \log m}$ , we get the following:

$$\mathbb{E}[\max_{x \in \mathcal{X}} Z_x] \leq t_0 + \int_{t_0}^\infty m e^{-t^2/(2d)} dt \leq t_0 + m \cdot \frac{d}{t_0} e^{-t_0^2/(2d)} = t_0 + \frac{d}{t_0} \leq O(\sqrt{d \log m}).$$

Therefore, we have the following:

$$w(\mathcal{X}) = \inf_{\lambda \in \Delta_{\mathcal{X}}} f(\lambda) \leq f(\lambda_G) \leq O(\sqrt{d \log |\mathcal{X}|}).$$

■

Next, recall the definition of  $\lambda_1$ ,  $\mathcal{R}$ ,  $w(\lambda_1, \mathcal{R}, \mathcal{X})$  for the finite set  $\mathcal{X}$  in Appendix C.

**Proposition 25** For a finite set  $\mathcal{X}$  with  $|\mathcal{X}| = m \geq d$ , we have  $w(\lambda_1, \mathcal{R}, \mathcal{X}) \leq O(\sqrt{d \log(m/d)})$ .

**Proof** For simplicity of presentation, assume that  $m$  is a multiple of  $d$ . Fix an index  $i \in [d]$ . Define  $u_x := \Sigma(\lambda_G)^{-1/2} x$ . Then  $\|u_x\|_2^2 \leq d$  for all  $x \in \mathcal{X}$ .

If  $g \sim \mathcal{N}(0, I_d)$ , then for each  $x \in \mathcal{X}$ ,  $Z_x := \langle u_x, g \rangle$  is Gaussian with variance  $\|u_x\|_2^2 \leq d$ , so for  $t \geq 0$ ,

$$\Pr(Z_x \geq t) \leq \exp\left(-\frac{t^2}{2d}\right).$$

Now by applying the union bound, we get the following:

$$\Pr\left(\max_{x \in R_i} Z_x \geq t\right) \leq (m/d) \exp\left(-\frac{t^2}{2d}\right).$$

Using the fact  $\mathbb{E}[\max_{x \in R_i} Z_x] = \int_0^\infty \Pr(\max_{x \in R_i} Z_x \geq t) dt$  and choosing  $t_0 = \sqrt{2d \log(m/d)}$ , we get the following:

$$\mathbb{E}[\max_{x \in R_i} Z_x] \leq t_0 + \int_{t_0}^\infty (m/d) e^{-t^2/(2d)} dt \leq t_0 + (m/d) \cdot \frac{d}{t_0} e^{-t_0^2/(2d)} = t_0 + \frac{d}{t_0} \leq O(\sqrt{d \log(m/d)}).$$

Therefore, we have the following:

$$w(\lambda_1, \mathcal{R}, \mathcal{X}) \leq \max_{i \in [d]} \mathbb{E}[\max_{x \in R_i} Z_x] \leq O(\sqrt{d \log(m/d)}).$$

■

**Proposition 26**  $w(\mathcal{X}) \geq \Omega(\sqrt{d \log d})$ .

**Proof** Fix any  $\lambda$  with  $\Sigma := \Sigma(\lambda) \succ 0$ . Define the set

$$T := \{\Sigma^{-1/2} x : x \in \mathcal{X}\} \subset \mathbb{R}^d,$$

and define the Gaussian process  $Z_t := \langle g, t \rangle$  for  $g \sim \mathcal{N}(0, I_d)$  over  $T$ . Then we have:

$$f(\lambda) = \mathbb{E}\left[\sup_{t \in T} Z_t\right].$$

Let  $U := \Sigma^{-1/2}X$  where  $X \sim \lambda$ . Since  $\mathbb{E}[XX^\top] = \Sigma$ ,

$$\mathbb{E}[UU^\top] = \Sigma^{-1/2} \mathbb{E}[XX^\top] \Sigma^{-1/2} = I_d.$$

Let  $m := \lfloor d/2 \rfloor$  and  $r := \sqrt{d/2}$ . We now claim that there exists  $t_1, \dots, t_m \in T$  such that  $\|t_i - t_j\|_2 \geq r$  for all  $i \neq j$ .

We construct them greedily. First pick  $t_1$  arbitrarily. Now assume that  $t_1, \dots, t_k$  have been chosen for some  $1 \leq k < m$ . Let  $V_k := \text{span}\{t_1, \dots, t_k\}$ , let  $P_k$  be the orthogonal projection onto  $V_k$ , and let  $Q_k := I - P_k$  be the orthogonal projection onto  $V_k^\perp$ . Then  $Q_k$  is symmetric and idempotent ( $Q_k^2 = Q_k$ ). Moreover, for an orthogonal projection, the trace equals the dimension of the subspace it projects onto. Note that  $Q_k$  projects onto  $V_k^\perp$ , which has dimension  $d - \dim(V_k)$ . Therefore

$$\text{tr}(Q_k) = \dim(V_k^\perp) = d - \dim(V_k) \geq d - k,$$

since  $\dim(V_k) \leq k$ .

Since  $Q_k$  is an orthogonal projection,

$$\|Q_k U\|_2^2 = U^\top Q_k^\top Q_k U = U^\top Q_k^2 U = U^\top Q_k U.$$

Using  $\text{tr}(ab) = ab$  for scalars and cyclicity of trace, we have the following:

$$\mathbb{E}\|Q_k U\|_2^2 = \mathbb{E}[U^\top Q_k U] = \mathbb{E}[\text{tr}(Q_k U U^\top)] = \text{tr}(Q_k \mathbb{E}[U U^\top]) = \text{tr}(Q_k I_d) = \text{tr}(Q_k) \geq d - k.$$

Therefore there exists a point  $t_{k+1} \in T$  with

$$\|Q_k t_{k+1}\|_2^2 \geq d - k.$$

For any  $i \leq k$ , we have  $t_i \in V_k$ , which implies  $Q_k t_i = 0$ . As orthogonal projections cannot increase the norm, we have the following:

$$\|t_{k+1} - t_i\|_2 \geq \|Q_k(t_{k+1} - t_i)\|_2 = \|Q_k t_{k+1}\|_2 \geq \sqrt{d - k} \geq \sqrt{d/2} = r.$$

This completes the proof of our claim.

Let  $S := \{t_1, \dots, t_m\}$ . Recall that  $\|t_i - t_j\|_2 \geq r$  for all  $i \neq j$ . Now define  $Z_i := \langle g, t_i \rangle$  where  $g \sim \mathcal{N}(0, 1)$ . Then  $(Z_i)_{i=1}^m$  is a centered Gaussian vector and for  $i \neq j$ , we have:

$$\mathbb{E}[(Z_i - Z_j)^2] = \|t_i - t_j\|_2^2 \geq r^2.$$

Let  $\xi_1, \dots, \xi_m$  be i.i.d.  $\mathcal{N}(0, 1)$  and set  $Y_i := (r/\sqrt{2})\xi_i$ . Then for  $i \neq j$ , we have:

$$\mathbb{E}[(Y_i - Y_j)^2] = \frac{r^2}{2} \mathbb{E}[(\xi_i - \xi_j)^2] = \frac{r^2}{2} \cdot 2 = r^2.$$

Due to Sudakov–Fernique inequality, we have the following:

$$\mathbb{E} \left[ \sup_{t \in T} Z_t \right] \geq \mathbb{E} \left[ \max_{1 \leq i \leq m} Z_i \right] \geq \mathbb{E} \left[ \max_{1 \leq i \leq m} Y_i \right] = \frac{r}{\sqrt{2}} \mathbb{E} \left[ \max_{1 \leq i \leq m} \xi_i \right] \geq \Omega(\sqrt{d \log d}).$$

This implies that  $w(\mathcal{X}) \geq \Omega(\sqrt{d \log d})$ . ■

**Proposition 27** *When  $\mathcal{X}$  is  $\{-1, +1\}^d$ ,  $\{0, +1\}^d$  or a unit  $\ell_2$ -ball, we have  $w(\mathcal{X}) \geq \Omega(d)$ . For  $m \leq d/21$ , we have  $w(\mathcal{X}) \geq \Omega(\sqrt{md})$  when  $\mathcal{X}$  is an  $m$ -set.*

**Proof** Recall that in Section 2.1, we described an algorithm with an upper bound of  $\mathcal{O}\left(\frac{d \log(1/\delta)}{\varepsilon^2} + \frac{w(\mathcal{X})^2}{\varepsilon^2}\right)$ .

When  $\mathcal{X}$  is  $\{-1, +1\}^d$ ,  $\{0, +1\}^d$  or a unit  $\ell_2$ -ball, our adaptive lower bound results in Table 1 imply that  $w(\mathcal{X}) \geq \Omega(d)$ . Similarly when  $\mathcal{X}$  is an  $m$ -set, our adaptive lower bound results in Table 1 imply that  $w(\mathcal{X}) \geq \Omega(\sqrt{md})$  for  $m \leq d/21$ . ■

## Appendix G. Gaussian Width Calculations for Multi-Task MAB

Recall that  $d = \sum_{j=1}^m d_j$ . Consider the blocks  $B_j := \{d_{1:j-1} + 1, \dots, d_{1:j}\}$ . Recall that the action set is

$$\mathcal{X} := \left\{ x \in \{0, 1\}^d : \forall j \in [m], \sum_{i \in B_j} x_i = 1 \right\}.$$

Let  $r := d - m + 1$ . Since  $\mathcal{X}$  is  $r$ -dimensional, it is contained in an  $r$ -dimensional linear subspace  $\text{span}(\mathcal{X}) \subseteq \mathbb{R}^d$ . Let  $U \in \mathbb{R}^{d \times r}$  be a matrix whose columns form an orthonormal basis of  $\text{span}(\mathcal{X})$ , that is  $U^\top U = I_r$  and  $\text{range}(U) = \text{span}(\mathcal{X})$ . Define the coordinate representation of  $\mathcal{X}$  in  $\mathbb{R}^r$  by

$$\mathcal{X}_r := \{U^\top x : x \in \mathcal{X}\} \subseteq \mathbb{R}^r.$$

Then the maps

$$x \mapsto U^\top x \quad \text{and} \quad z \mapsto Uz$$

are inverse to each other on  $\mathcal{X}$  and  $\mathcal{X}_r$ , respectively. In particular, for every  $x \in \mathcal{X}$ ,

$$x = U(U^\top x),$$

so  $U^\top x$  is the unique coordinate vector of  $x$  w.r.t the orthonormal basis  $\{U_1, \dots, U_r\}$  where  $U_j$  is the  $j$ -th column of  $U$ .

We now construct one such matrix  $U$ . Define the vector  $\mu \in \mathbb{R}^d$  by setting, for each  $i \in \{1, \dots, d\}$ ,

$$\mu_i := \frac{1}{d_j} \quad \text{where } j \text{ is the unique index such that } i \in B_j.$$

Then, we have:

$$S := \|\mu\|_2^2 = \sum_{i=1}^d \mu_i^2 = \sum_{j=1}^m d_j \left(\frac{1}{d_j}\right)^2 = \sum_{j=1}^m \frac{1}{d_j},$$

and we define the unit vector:

$$u_0 := \frac{\mu}{\|\mu\|_2} = \frac{\mu}{\sqrt{S}}.$$

We now define a matrix  $H_n \in \mathbb{R}^{n \times (n-1)}$ . For  $n \geq 2$ , define  $(H_n)_{\ell, k}$  for  $\ell \in [n]$ ,  $k \in [n-1]$  by

$$(H_n)_{\ell, k} = \begin{cases} \frac{1}{\sqrt{k(k+1)}} & 1 \leq \ell \leq k, \\ -\frac{k}{\sqrt{k(k+1)}} & \ell = k+1, \\ 0 & \ell \geq k+2. \end{cases}$$

Then, we have:

$$H_n^\top H_n = I_{n-1}, \quad H_n^\top \mathbf{1}_n = \mathbf{0}_{n-1}, \quad H_n H_n^\top = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top. \quad (42)$$

Let  $s_j := d_{1:j-1}$ . We define the entries of the matrix  $Q_j \in \mathbb{R}^{d \times (d_j-1)}$  as:

$$(Q_j)_{i, k} = \begin{cases} (H_{d_j})_{i-s_j, k} & i \in B_j, \\ 0 & i \notin B_j, \end{cases} \quad k \in [d_j - 1]. \quad (43)$$

Now we define:

$$U := [u_0 \ Q_1 \ \cdots \ Q_m] \in \mathbb{R}^{d \times r}, \quad r := 1 + \sum_{j=1}^m (d_j - 1) = d - m + 1.$$

As  $Q_j, Q_\ell$  have disjoint supports for  $j \neq \ell$  and due to Eq. (42), we get:

$$Q_j^\top Q_j = H_{d_j}^\top H_{d_j} = I_{d_j-1}, \quad Q_j^\top Q_\ell = 0 \ (j \neq \ell).$$

For any matrix  $A$ , let  $A_{i,:}$  denote the  $i$ -th row of the matrix  $A$ . Since  $u_0$  is constant on block  $B_j$ , we have:

$$u_0^\top Q_j = \sum_{i \in B_j} (u_0)_i (Q_j)_{i,:} = \frac{1/d_j}{\sqrt{S}} \sum_{r=1}^{d_j} (H_{d_j})_{r,:} = \frac{1/d_j}{\sqrt{S}} \mathbf{1}_{d_j}^\top H_{d_j} = \mathbf{0}_{n-1}.$$

Hence

$$U^\top U = I_r.$$

Next, we prove that  $\text{range}(U) = \text{span}(\mathcal{X})$ . Towards that, define

$$R_j := \left\{ v \in \mathbb{R}^d : \text{supp}(v) \subseteq B_j, \mathbf{1}_{B_j}^\top v = 0 \right\}.$$

where  $\mathbf{1}_{B_j}$  is the indicator vector of the block  $B_j$ . From (42)–(43), we get:

$$\text{range}(Q_j) = R_j.$$

Also  $\mu \in \text{span}(\mathcal{X})$  and  $\forall j, \forall k \in [d_j], e_{s_j+k} - e_{s_j+1} \in \text{span}(\mathcal{X})$ , so  $R_j \subseteq \text{span}(\mathcal{X})$  (since  $\{e_{s_j+k} - e_{s_j+1}\}_{k=2}^{d_j}$  spans  $R_j$ ). Thus

$$\text{span}(\mu) \oplus \left( \bigoplus_{j=1}^m R_j \right) \subseteq \text{span}(\mathcal{X}).$$

where  $\oplus$  means direct sum of vector spaces.

Conversely, for any  $x \in \mathcal{X}$ , we have:

$$\forall j \in [m] : \mathbf{1}_{B_j}^\top (x - \mu) = 1 - 1 = 0 \quad \Rightarrow \quad x - \mu \in \bigoplus_{j=1}^m R_j, \quad \Rightarrow \quad x \in \text{span}(\mu) \oplus \left( \bigoplus_{j=1}^m R_j \right).$$

Hence

$$\text{span}(\mathcal{X}) = \text{span}(\mu) \oplus \left( \bigoplus_{j=1}^m R_j \right) = \text{range}(U), \text{ and indeed } \dim(\text{span}(\mathcal{X})) = r = d - m + 1.$$

Now the gaussian width term  $w(\mathcal{X})$  is defined using the coordinate representation of  $\mathcal{X}$  in  $\mathbb{R}^r$  as follows:

$$w(\mathcal{X}) := \inf_{\lambda \in \Delta_{\mathcal{X}_r}} \mathbb{E}_{g_r \sim \mathcal{N}(0, I_r)} \left[ \sup_{x \in \mathcal{X}} \langle U^\top x, A(\lambda; \mathcal{X}_r)^{-1/2} g_r \rangle \right]$$

Fix a probability distribution  $\lambda$  on  $\mathcal{X}_r$  such that, if  $Z \sim \lambda$  and we map to  $\mathcal{X}$  via

$$X := UZ \in \mathcal{X} \subseteq \mathbb{R}^d,$$

then the following holds: for each block  $B_k$ , exactly one coordinate index  $j_k \in B_k$  is selected uniformly at random, and the selections  $\{j_k\}_{k=1}^m$  are independent across blocks. In other words, viewed in  $\mathbb{R}^d$ ,  $X$

is distributed as a vector obtained by choosing one coordinate uniformly within each block, independently across blocks. Now define:

$$\Sigma_r := \mathbb{E}_{x \sim \lambda}[xx^\top],$$

where  $\Sigma_r \in \mathbb{R}^{r \times r}$ . Now the Gaussian width of  $\mathcal{X}$  is upper bounded as follows:

$$w(\mathcal{X}) \leq w(\mathcal{X}_r; \Sigma_r) := \mathbb{E}_{g_r \sim \mathcal{N}(0, I_r)} \left[ \sup_{x \in \mathcal{X}} \langle U^\top x, \Sigma_r^{-1/2} g_r \rangle \right]$$

Observe that  $\Sigma_r := U^\top \Sigma U$  where  $\Sigma := \mathbb{E}_{x \sim \lambda}[Ux(Ux)^\top]$ . Due to the definition of  $\lambda$ , for  $i \in B_j, k \in B_\ell$ , we have:

$$\Sigma_{ik} = \begin{cases} \frac{1}{d_j} \mathbf{1}\{i = k\} & j = \ell, \\ \frac{1}{d_j d_\ell} & j \neq \ell. \end{cases}$$

Let  $v \in R_j$  ( $\mathbf{1}_{B_j}^\top v = 0$ , support in  $B_j$ ). For  $i \in B_j$ ,

$$(\Sigma v)_i = \sum_{k \in B_j} \frac{1}{d_j} \mathbf{1}\{i = k\} v_k + \sum_{\ell \neq j} \sum_{k \in B_\ell} \frac{1}{d_j d_\ell} v_k = \frac{1}{d_j} v_i + 0,$$

so  $\Sigma v = \frac{1}{d_j} v$ . Therefore

$$Q_j^\top \Sigma Q_j = \frac{1}{d_j} Q_j^\top Q_j = \frac{1}{d_j} I_{d_j-1}, \quad Q_j^\top \Sigma Q_\ell = 0 \quad (j \neq \ell).$$

For  $i \in B_j$ ,

$$\begin{aligned} (\Sigma \mu)_i &= \sum_{k \in B_j} \frac{1}{d_j} \mathbf{1}\{i = k\} \frac{1}{d_j} + \sum_{\ell \neq j} \sum_{k \in B_\ell} \frac{1}{d_j d_\ell} \frac{1}{d_\ell} \\ &= \frac{1}{d_j^2} + \sum_{\ell \neq j} \frac{1}{d_j d_\ell} = \frac{1}{d_j} \sum_{\ell=1}^m \frac{1}{d_\ell} = S \mu_i. \end{aligned}$$

Thus  $\Sigma \mu = S \mu$  and, since  $\|u_0\| = 1$ ,

$$u_0^\top \Sigma u_0 = S, \quad u_0^\top \Sigma Q_j = 0.$$

Hence, in the explicit basis  $U$ ,

$$\Sigma_r := U^\top \Sigma U = \text{diag} \left( S, \frac{1}{d_1} I_{d_1-1}, \dots, \frac{1}{d_m} I_{d_m-1} \right), \quad \Sigma_r^{-1/2} = \text{diag} \left( S^{-1/2}, \sqrt{d_1} I_{d_1-1}, \dots, \sqrt{d_m} I_{d_m-1} \right). \quad (44)$$

#### D) Gaussian Width in $r = d - m + 1$ Dimensions and Evaluation

Let  $g_r \sim \mathcal{N}(0, I_r)$ , write  $g_r = (g_0, g^{(1)}, \dots, g^{(m)})$  with  $g^{(j)} \sim \mathcal{N}(0, I_{d_j-1})$  independent. Recall that

$$w(\mathcal{X}) \leq \mathbb{E}_{g_r \sim \mathcal{N}(0, I_r)} \left[ \sup_{x \in \mathcal{X}} \langle U^\top x, \Sigma_r^{-1/2} g_r \rangle \right].$$

Also, we have

$$\langle u_0, x \rangle = \frac{\langle \mu, x \rangle}{\sqrt{S}} = \frac{\sum_{j=1}^m \frac{1}{d_j}}{\sqrt{S}} = \sqrt{S}, \quad \forall x \in \mathcal{X},$$

so, from (44), we have

$$\begin{aligned} \langle U^\top x, \Sigma_r^{-1/2} g_r \rangle &= S^{-1/2} g_0 \langle u_0, x \rangle + \sum_{j=1}^m \sqrt{d_j} \langle g^{(j)}, Q_j^\top x \rangle \\ &= g_0 + \sum_{j=1}^m \sqrt{d_j} \langle g^{(j)}, Q_j^\top x \rangle. \end{aligned}$$

Therefore

$$w(\mathcal{X}) \leq \mathbb{E}[g_0] + \sum_{j=1}^m \sqrt{d_j} \mathbb{E} \left[ \max_{k \in [d_j]} \langle g^{(j)}, H_{d_j}^\top e_k \rangle \right] = \sum_{j=1}^m \sqrt{d_j} \mathbb{E} \left[ \max_{k \in [d_j]} \langle g^{(j)}, H_{d_j}^\top e_k \rangle \right]. \quad (45)$$

Let  $h^{(j)} := H_{d_j} g^{(j)} \in \mathbb{R}^{d_j}$ . Then we have

$$\langle g^{(j)}, H_{d_j}^\top e_k \rangle = e_k^\top H_{d_j} g^{(j)} = (h^{(j)})_k,$$

and by (42),

$$h^{(j)} \sim \mathcal{N}\left(0, H_{d_j} H_{d_j}^\top\right) = \mathcal{N}\left(0, I_{d_j} - \frac{1}{d_j} \mathbf{1}\mathbf{1}^\top\right).$$

If  $Z^{(j)} \sim \mathcal{N}(0, I_{d_j})$  and  $\bar{Z}^{(j)} := \frac{1}{d_j} \mathbf{1}^\top Z^{(j)}$ , then

$$\left(I_{d_j} - \frac{1}{d_j} \mathbf{1}\mathbf{1}^\top\right) Z^{(j)} = Z^{(j)} - \bar{Z}^{(j)} \mathbf{1} \sim \mathcal{N}\left(0, I_{d_j} - \frac{1}{d_j} \mathbf{1}\mathbf{1}^\top\right),$$

so

$$\max_{k \leq d_j} (h^{(j)})_k \stackrel{d}{=} \max_{k \leq d_j} ((Z^{(j)})_k - \bar{Z}^{(j)}) = \max_{k \leq d_j} (Z^{(j)})_k - \bar{Z}^{(j)}.$$

Taking expectations and using  $\mathbb{E}[\bar{Z}^{(j)}] = 0$ ,

$$\mathbb{E} \left[ \max_{k \leq d_j} (h^{(j)})_k \right] = \mathbb{E} \left[ \max_{k \leq d_j} Z_k \right] \leq \mathcal{O}(\sqrt{\log d_j}), \quad Z_k \stackrel{iid}{\sim} \mathcal{N}(0, 1).$$

Thus, due to Eq. (45) and the calculations above, we have

$$w(\mathcal{X}) \leq \mathcal{O} \left( \sum_{j=1}^m \sqrt{d_j \log d_j} \right).$$

## Appendix H. Lower Bounds for Structured Sets

First, in section H.1, we present the lower bound for the multi-task MAB problem. Next, in section H.1, we present the lower bound for hypercubes. Finally, in section H.3, we present the lower bound for  $m$ -sets.

### H.1. Lower Bound for Multi-task MAB

Recall that in the Multi-task MAB, the set of arms  $\mathcal{X}$  is defined as follows:

$$\mathcal{X} = \left\{ \mathbf{x} \in \{0, 1\}^d : \forall j \in [m] \sum_{i=d_{1:j-1}+1}^{d_{1:j}} \mathbf{x}[i] = 1 \right\}$$

where  $d_i \geq 2$ ,  $d = \sum_{i=1}^m d_i$ ,  $d_{1:j} = \sum_{i=1}^j d_i$  and  $d_{1:0} = 0$  for all  $j \in [m]$ .

We first focus on the case when  $d_j \geq 100$  for all  $j \in [m]$ .

Let us fix one regret minimizing algorithm, say  $\mathcal{A}$  and assume that  $\mathcal{A}$  is deterministic given the reward realizations. Let us assume that  $\mathcal{A}$  terminates after taking  $T := \frac{(\sum_{i=1}^m \sqrt{d_i})^2}{20000\varepsilon^2}$  samples (if the algorithm terminates before this, we can take additional dummy samples). Let  $\hat{x} \in \mathcal{X}$  be the arm recommended by  $\mathcal{A}$  after sampling the arms  $\mathbf{x}_1, \dots, \mathbf{x}_T$  and terminating.

Later in this section, we construct a set of input instances  $\mathcal{I}$ . Each instance  $I \in \mathcal{I}$  is associated with a reward vector  $\theta_I \in \mathbb{R}^d$  such that the following hold:

- If  $\mathcal{A}$  samples an arm  $x \in \mathcal{X}$ , it observes  $\langle x, \theta_I \rangle + \eta$  where  $\eta \sim \mathcal{N}(0, 1)$ .
- For all  $x \in \mathcal{X}$ , we have  $\langle x, \theta_I \rangle \geq 0$ .
- We have  $\max_{x \in \mathcal{X}} \langle x, \theta_I \rangle = 10\varepsilon$ .

Suppose we show that  $\mathbb{E}_{I \sim \text{Unif}(\mathcal{I})}[\langle \hat{x}, \theta_I \rangle] \leq 7\varepsilon$ , where  $\mathbb{E}_I[\cdot]$  is the expectation under the instance  $I$ . Then by Yao's lemma, we have for any randomized algorithm, the recommended arm  $\hat{x}$  satisfies  $\min_{I \in \mathcal{I}} \mathbb{E}_I[\langle \hat{x}, \theta_I \rangle] \leq 7\varepsilon$ . Then due to markov's inequality, we have that with probability at least 0.1, there exists an instance on which the randomized algorithm does not recommend an  $\varepsilon$ -best arm.

Now we focus on showing that  $\mathbb{E}_{I \sim \text{Unif}(\mathcal{I})}[\langle \hat{x}, \theta_I \rangle] \leq 7\varepsilon$ . We begin with construction instances  $I$  with their corresponding reward vector  $\theta_I$ . These instances include both the set of all input instances and a set of alternate instances that will be used to argue the performance of the algorithm  $\mathcal{A}$ .

Let  $\tilde{\mathcal{X}} = \left\{ \mathbf{x} \in \{0, 1\}^d : \forall j \in [m] \sum_{i=d_{1:j-1}+1}^{d_{1:j}} \mathbf{x}[i] \leq 1 \right\}$ . First we describe an instance  $I_{\tilde{\mathbf{x}}}$ , where  $\tilde{\mathbf{x}} \in \tilde{\mathcal{X}}$ . In this instance, we define the associated reward vector  $\theta_{I_{\tilde{\mathbf{x}}}}$  as follows. For all  $j \in [m]$  and  $i \in \{d_{1:j-1}+1, \dots, d_{1:j}\}$ , we have  $\theta_{I_{\tilde{\mathbf{x}}}}[i] = 10\varepsilon_j \cdot \tilde{\mathbf{x}}[i]$  where  $\varepsilon_j = \frac{\varepsilon \cdot \sqrt{d_j}}{\sum_{s \in [m]} \sqrt{d_s}}$ . Observe that for all  $x \in \mathcal{X}$ , we have  $\langle x, \theta_{I_{\tilde{\mathbf{x}}}} \rangle \geq 0$ . Also observe that  $\max_{x \in \mathcal{X}} \langle x, \theta_{I_{\tilde{\mathbf{x}}}} \rangle = 10\varepsilon$  if  $\tilde{\mathbf{x}} \in \tilde{\mathcal{X}}$ .

Let  $r_{\tilde{\mathbf{x}}}^{(j)} := \mathbb{E}_{I_{\tilde{\mathbf{x}}}}[\sum_{i=d_{1:j-1}+1}^{d_{1:j}} \tilde{\mathbf{x}}[i] \cdot \theta_{I_{\tilde{\mathbf{x}}}}[i]]$ . Observe that  $\mathbb{E}_{I_{\tilde{\mathbf{x}}}}[\langle \tilde{\mathbf{x}}, \theta_{I_{\tilde{\mathbf{x}}}} \rangle] = \sum_{j=1}^m r_{\tilde{\mathbf{x}}}^{(j)}$ .

Let  $\mathcal{I} = \bigcup_{\mathbf{x} \in \mathcal{X}} I_{\mathbf{x}}$ . Fix an index  $j \in [m]$ . We now show that  $\mathbb{E}_{I_{\mathbf{x}' \sim \text{Unif}(\mathcal{I})}}[\sum_{i=d_{1:j-1}+1}^{d_{1:j}} \hat{x}[i] \cdot \theta_{I_{\mathbf{x}'}}[i]] \leq 7\varepsilon_j$ . Let

$$\mathcal{X}^{(j)} := \left\{ \mathbf{x} \in \{0, 1\}^d : \forall i \in [m] \setminus \{j\} \sum_{s=d_{1:i-1}+1}^{d_{1:i}} \mathbf{x}[s] = 1, \sum_{s=d_{1:j-1}+1}^{d_{1:j}} \mathbf{x}[s] = 0 \right\}.$$

For any  $\mathbf{x} \in \mathcal{X}^{(j)}$ , let  $\mathbf{x}^{(i)}$  be the vector in  $\mathcal{X}$  such that  $\mathbf{x}^{(i)}[s] = \mathbf{x}[s]$  for all  $s \notin \{d_{1:j-1}+1, \dots, d_{1:j}\}$  and  $\mathbf{x}^{(i)}[d_{1:j-1}+i] = 1$ .

Let us fix  $\mathbf{x} \in \mathcal{X}^{(j)}$ . Now we claim that there is a set  $\mathcal{S}_{\mathbf{x}} \subseteq [d_j]$  with at least  $d_j/3$  indices such that for each  $i \in \mathcal{S}_{\mathbf{x}}$ , we have  $r_{\mathbf{x}^{(i)}}^{(j)} \leq \varepsilon_j$ .

Before we prove our claim, we first show that if our claim holds true, then we have that

$$\mathbb{E}_{I_{\mathbf{x}' \sim \text{Unif}(\mathcal{I})}}[\sum_{i=d_{1:j-1}+1}^{d_{1:j}} \hat{x}[i] \cdot \theta_{I_{\mathbf{x}'}}[i]] \leq 7\varepsilon_j.$$

Now we have the following:

$$\begin{aligned} \mathbb{E}_{I_{\mathbf{x}' \sim \text{Unif}(\mathcal{I})}}[\sum_{i=d_{1:j-1}+1}^{d_{1:j}} \hat{x}[i] \cdot \theta_{I_{\mathbf{x}'}}[i]] &= \frac{1}{\prod_{s=1}^m d_s} \sum_{\mathbf{x} \in \mathcal{X}^{(j)}} \sum_{i=1}^{d_j} r_{\mathbf{x}^{(i)}}^{(j)} \\ &= \frac{1}{\prod_{s=1}^m d_s} \sum_{\mathbf{x} \in \mathcal{X}^{(j)}} \left( \sum_{i \in \mathcal{S}_{\mathbf{x}}} r_{\mathbf{x}^{(i)}}^{(j)} + \sum_{i \in [d_j] \setminus \mathcal{S}_{\mathbf{x}}} r_{\mathbf{x}^{(i)}}^{(j)} \right) \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{1}{\prod_{s=1}^m d_s} \sum_{\mathbf{x} \in \mathcal{X}^{(j)}} \left( \frac{d_j}{3} \cdot \varepsilon_j + \frac{2d_j}{3} \cdot 10\varepsilon_j \right) \\
 &\leq \frac{1}{\prod_{s=1}^m d_s} \sum_{\mathbf{x} \in \mathcal{X}^{(j)}} 7d_j \varepsilon_j \\
 &= \frac{7\varepsilon_j}{\prod_{s=1}^m d_s} \cdot \prod_{s \neq j} d_s \cdot d_j \\
 &= 7\varepsilon_j
 \end{aligned}$$

Now we prove our claim. We use the adaptive KL chain rule (Lemma 12) in our analysis.

Let  $\mathbb{P}_I$  denote the probability law instances under an instance  $I$ . For an instance  $I_{\mathbf{x}^{(i)}}$ , let  $f_i(\ell_1, \dots, \ell_T)$  denote the joint PDF for the tuple of reward values observed by  $\mathcal{A}$  in each round under the probability law  $\mathbb{P}_{I_{\mathbf{x}^{(i)}}}$ . Observe that our sample space is  $\Omega = \mathbb{R}^T$ . This is a valid sample space as  $\mathcal{A}$  is deterministic and the probability density function of the reward values in round  $t$  only depends on the reward values it observed in the previous rounds. Similarly for the alternate instance  $I_{\mathbf{x}}$ , let  $f_0(\ell_1, \dots, \ell_T)$  denote the joint PDF for the tuple of loss values observed by  $\mathcal{A}$  in each round under the probability law  $\mathbb{P}_{I_{\mathbf{x}}}$ .

First observe that the instances  $I_{\mathbf{x}^{(i)}}$  and  $I_{\mathbf{x}}$  only differ at index  $d_{1:j-1} + i$ . For each  $\omega \in \Omega$ , let  $\mathbf{x}_{1,\omega}, \mathbf{x}_{2,\omega}, \dots, \mathbf{x}_{T,\omega}$  be the sequence of arms chosen by  $\mathcal{A}$  on  $\omega$ . Conditioning on a set of outcomes  $X_1 = \omega_1, X_2 = \omega_2, \dots, X_{t-1} = \omega_{t-1}$ , we have  $X_t \sim \mathcal{N}(\mu_i, 1)$  for the instance  $I_{\mathbf{x}^{(i)}}$  and  $X_t \sim \mathcal{N}(\mu_0, 1)$  for the instance  $I_{\mathbf{x}}$  where  $\mu_i - \mu_0 = 10\varepsilon_j \cdot \mathbf{x}_{t,\omega}[d_{1:j-1} + i]$ . Let  $T_i = \sum_{t=1}^T \mathbf{x}_{t,\omega}[d_{1:j-1} + i]$ . For each  $\omega \in \Omega$ , let  $T_{i,\omega} = \sum_{t=1}^T \mathbf{x}_{t,\omega}[d_{1:j-1} + i]$ . Note that  $T_i$  is a random variable and  $T_{i,\omega}$  is a fixed value. Let  $X_{-t}$  denote  $X_1, \dots, X_{t-1}$  and  $\omega_{-t}$  denote  $\omega_1, \dots, \omega_{t-1}$ . Now we have the following:

$$\begin{aligned}
 KL(f_0, f_i) &= \int_{\omega \in \Omega} f_0(\omega) \left( KL(f_0(X_1), f_i(X_1)) + \sum_{t=2}^T KL(f_0(X_t | X_{-t} = \omega_{-t}), f_i(X_t | X_{-t} = \omega_{-t})) \right) d\omega \\
 &= 50\varepsilon_j^2 \int_{\omega \in \Omega} f_0(\omega) \sum_{t=1}^T \mathbf{x}_{t,\omega}[d_{1:j-1} + i] d\omega \\
 &= 50\varepsilon_j^2 \int_{\omega \in \Omega} f_0(\omega) T_{i,\omega} d\omega \\
 &= 50\varepsilon_j^2 \cdot \mathbb{E}_{I_{\mathbf{x}}}[T_i]
 \end{aligned}$$

Now observe that  $\sum_{i=1}^{d_j} \mathbb{E}_{I_{\mathbf{x}}}[T_i] = T$ . Hence, there exists a set  $\mathcal{S}_x^{(1)} \subseteq [d_j]$  with at least  $2d_j/3$  indices such that for each  $i \in \mathcal{S}_x^{(1)}$ , we have  $\mathbb{E}_{I_{\mathbf{x}}}[T_i] \leq \frac{3T}{d_j}$ . Similarly, there exists a set  $\mathcal{S}_x^{(2)} \subseteq [d_j]$  with at least  $2d_j/3$  indices such that for each  $i \in \mathcal{S}_x^{(2)}$ , we have  $\mathbb{P}_{I_{\mathbf{x}}}[\hat{x}[d_{1:j-1} + i] = 1] \leq \frac{3}{d_j}$ . Now we define  $\mathcal{S}_x := \mathcal{S}_x^{(1)} \cap \mathcal{S}_x^{(2)}$ . Observe that  $|\mathcal{S}_x| \geq d_j/3$ . Next for each  $i \in \mathcal{S}_x$ , we have  $\mathbb{E}_{I_{\mathbf{x}}}[T_i] \leq \frac{3T}{d_j}$  and  $\mathbb{P}_{I_{\mathbf{x}}}[\hat{x}[d_{1:j-1} + i] = 1] \leq \frac{3}{d_j}$ . Now for each  $i \in \mathcal{S}_x$ , we have  $KL(f_0, f_i) \leq \frac{150\varepsilon_j^2 T}{d_j} = \frac{3}{400}$ .

Fix  $i \in \mathcal{S}_x$ . Let  $A_i$  be the event that  $\hat{x}[d_{1:j-1} + i] = 1$ . Now due to Pinsker's inequality we have the following:

$$\begin{aligned}
 \mathbb{P}_{I_{\mathbf{x}^{(i)}}}(A_i) &\leq \mathbb{P}_{I_{\mathbf{x}}}(A_i) + \sqrt{\frac{KL(f_0, f_i)}{2}} \\
 &\leq \frac{3}{d_j} + \sqrt{\frac{3}{800}}
 \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{3}{100} + \sqrt{\frac{3}{800}} && \text{(as } d_j \geq 100) \\
 &< \frac{1}{10}
 \end{aligned}$$

Hence, we have  $r_{\mathbf{x}^{(i)}}^{(j)} < \frac{1}{10} \cdot 10\varepsilon_j = \varepsilon_j$ .

Now, we conclude the proof by showing that  $\mathbb{E}_{I \sim \text{Unif}(\mathcal{I})}[\langle \hat{x}, \theta_I \rangle] \leq 7\varepsilon$ :

$$\begin{aligned}
 \mathbb{E}_{I \sim \text{Unif}(\mathcal{I})}[\langle \hat{x}, \theta_I \rangle] &= \sum_{j=1}^m \mathbb{E}_{I_{\mathbf{x}'} \sim \text{Unif}(\mathcal{I})} \left[ \sum_{i=d_{1:j-1}+1}^{d_{1:j}} \hat{x}[i] \cdot \theta_{I_{\mathbf{x}'}}[i] \right] \\
 &= 7 \sum_{j=1}^m \varepsilon_j \\
 &= 7\varepsilon \cdot \frac{\sum_{j=1}^m \sqrt{d_j}}{\sum_{s=1}^m \sqrt{d_s}} && \text{(as } \varepsilon_j = \frac{\varepsilon \cdot \sqrt{d_j}}{\sum_{s=1}^m \sqrt{d_s}}) \\
 &= 7\varepsilon
 \end{aligned}$$

### H.1.1. $d_j < 100$ CASE

For the simplicity of presentation let us focus on the case  $d_j = 2$ . This analysis can be easily extended to the other cases easily. For the case  $d_j = 2$ , the construction of the hard instances remain exactly the same as that of the  $d_j \geq 100$  case. The analysis for this case only differs at the end. Consider an index  $i \in \{1, 2\}$ . Observe that  $KL(f_0, f_i) \leq 50\varepsilon_j^2 T = \frac{1}{200}$ . Recall that  $A_i$  is the event that  $\hat{x}[d_{1:j-1} + i] = 1$ . Now due to Pinsker's inequality we have the following:

$$\mathbb{P}_{I_{\mathbf{x}^{(i)}}}(A_i) \leq \mathbb{P}_{I_{\mathbf{x}}}(A_i) + \sqrt{\frac{KL(f_0, f_i)}{2}} \leq \mathbb{P}_{I_{\mathbf{x}}}(A_i) + \frac{1}{20}$$

Recall that  $r_{\hat{x}}^{(j)} := \mathbb{E}_{I_{\hat{x}}}[\sum_{i=d_{1:j-1}+1}^{d_{1:j}} \hat{x}[i] \cdot \theta_{I_{\hat{x}}}[i]]$ . Now we have the following:

$$\begin{aligned}
 \mathbb{E}_{I_{\mathbf{x}'} \sim \text{Unif}(\mathcal{I})} \left[ \sum_{i=d_{1:j-1}+1}^{d_{1:j}} \hat{x}[i] \cdot \theta_{I_{\mathbf{x}'}}[i] \right] &= \frac{1}{\prod_{s=1}^m d_s} \sum_{\mathbf{x} \in \mathcal{X}^{(j)}} \sum_{i=1}^2 r_{\mathbf{x}^{(i)}}^{(j)} \\
 &\leq \frac{10\varepsilon_j}{\prod_{s=1}^m d_s} \sum_{\mathbf{x} \in \mathcal{X}^{(j)}} (\mathbb{P}_{I_{\mathbf{x}}}(A_1) + \mathbb{P}_{I_{\mathbf{x}}}(A_2) + 0.1) \\
 &= \frac{10\varepsilon_j}{\prod_{s=1}^m d_s} \cdot \frac{1.1}{2} \cdot \prod_{s \neq j} d_s \cdot d_j \\
 &= 5.5\varepsilon_j
 \end{aligned}$$

## H.2. Hypercubes

For the  $\{0, 1\}^d$  hypercube, we consider the family of hard instances  $\{I_{\mathbf{x}}\}_{\mathbf{x} \in \{0,1\}^d}$  with corresponding reward vectors  $\{\theta_{\mathbf{x}}\}_{\mathbf{x} \in \{0,1\}^d}$ . For each  $\mathbf{x} \in \{0, 1\}^d$ , define  $\theta_{\mathbf{x}} \in \mathbb{R}^d$  coordinate-wise by  $(\theta_{\mathbf{x}})_i = \mathbb{1}\{\mathbf{x}_i = 1\} \cdot \frac{10\varepsilon}{d} - \mathbb{1}\{\mathbf{x}_i = 0\} \cdot \frac{10\varepsilon}{d}$ . An analysis analogous to our  $d_j = 2$  argument for multi-task MAB yields a lower bound of  $\Omega(d^2/\varepsilon^2)$ . The only minor difference is that, in the multi-task setting, we reasoned about the expected value of  $\langle x, \theta \rangle$ , whereas here one must instead work with the expected gap  $\langle x, \theta \rangle - \min_{x \in \{0,1\}^d} \langle x, \theta \rangle$ .

For the  $\{-1, +1\}^d$  hypercube, we similarly consider the family of hard instances  $\{I_{\mathbf{x}}\}_{\mathbf{x} \in \{-1, +1\}^d}$  with corresponding reward vectors  $\{\theta_{\mathbf{x}}\}_{\mathbf{x} \in \{-1, +1\}^d}$ . For each  $\mathbf{x} \in \{-1, +1\}^d$ , define  $\theta_{\mathbf{x}} \in \mathbb{R}^d$  coordinate-wise by  $(\theta_{\mathbf{x}})_i = \mathbf{x}_i \cdot \frac{5\varepsilon}{d}$ . An analysis analogous to our  $d_j = 2$  argument for multi-task MAB again yields a lower bound of  $\Omega(d^2/\varepsilon^2)$ . As above, the only difference is that the argument proceeds via the expected gap  $\langle x, \theta \rangle - \min_{x \in \{-1, +1\}^d} \langle x, \theta \rangle$  rather than expectation of  $\langle x, \theta \rangle$  itself.

### H.3. $m$ -Sets Lower Bound

**Theorem 28** *Let us denote the  $m$ -sets by*

$$\mathcal{X} = \{x \in \{0, 1\}^d : \|x\|_0 = m\}.$$

*In each round  $t = 1, \dots, T$ , the learner chooses  $x_t \in \mathcal{X}$  and observes*

$$y_t = \langle x_t, \theta \rangle + \eta_t, \quad \eta_t \sim \mathcal{N}(0, 1) \text{ i.i.d.}$$

*Suppose  $d - m + 1 \geq 20m$ . Then there exists a universal constant  $c > 0$  such that for any (possibly adaptive and randomized) algorithm, if*

$$T \leq c \frac{m(d - m + 1)}{\varepsilon^2},$$

*then there exists an instance  $\theta$  for which the algorithm fails to output an  $\varepsilon$ -optimal arm with probability at least  $2/9$ .*

**Proof** We first describe the family of hard instances. For each  $S \subset [d]$  with  $|S| = m$ , define

$$\theta_i^{(S)} = \begin{cases} \Delta & i \in S, \\ 0 & i \notin S, \end{cases} \quad \text{where} \quad \Delta := \frac{10\varepsilon}{m}.$$

Let  $S$  be uniform over all  $m$ -sized subsets of  $[d]$  and the environment parameter be  $\theta = \theta^{(S)}$ . For any output  $\hat{x} \in \mathcal{X}$ , write  $\hat{S} := \text{supp}(\hat{x})$  so  $|\hat{S}| = m$ . Then

$$\langle \hat{x}, \theta^{(S)} \rangle = \Delta |\hat{S} \cap S|.$$

The optimal value is  $\max_{x \in \mathcal{X}} \langle x, \theta^{(S)} \rangle = \Delta |S| = m\Delta = 10\varepsilon$ . Thus on any instance,  $\varepsilon$ -success implies  $\langle \hat{x}, \theta^{(S)} \rangle \geq 9\varepsilon$ .

By Yao's lemma it suffices to fix an arbitrary deterministic algorithm and analyze its error under the random draw of  $S$ .

Given  $S$ , let  $I$  be uniform on  $S$  (an auxiliary random variable), and define  $B := S \setminus \{I\}$  so  $|B| = m - 1$ . Conditioning on  $(S, \hat{S})$ ,

$$\mathbb{P}(I \in \hat{S} \mid S, \hat{S}) = \frac{|\hat{S} \cap S|}{|S|} = \frac{|\hat{S} \cap S|}{m},$$

so by the tower rule

$$\mathbb{E}[|\hat{S} \cap S|] = m \mathbb{P}(I \in \hat{S}). \tag{46}$$

Therefore,

$$\mathbb{E}[\langle \hat{x}, \theta^{(S)} \rangle] = \Delta \mathbb{E}[|\hat{S} \cap S|] = m\Delta \mathbb{P}(I \in \hat{S}) = 10\varepsilon \cdot \mathbb{P}(I \in \hat{S}). \tag{47}$$

So it is enough to show  $\mathbb{P}(I \in \hat{S}) \leq 0.7$ , which would imply  $\mathbb{E}[\langle \hat{x}, \theta^{(S)} \rangle] \leq 7\varepsilon$ .

We now prove the following technical lemma.

**Lemma 29 (Equivalent sampling)** *The joint law of  $(B, I, S)$  defined above is the same as: (i) draw  $B$  uniformly among all  $(m - 1)$ -subsets of  $[d]$ , (ii) draw  $I$  uniformly from  $[d] \setminus B$ , (iii) set  $S = B \cup \{I\}$ . In particular, conditional on  $B$ ,  $I \sim \text{Unif}([d] \setminus B)$ .*

**Proof** Fix  $b \subset [d]$  with  $|b| = m - 1$  and  $i \notin b$ . Under the original procedure,

$$\mathbb{P}(B = b, I = i) = \mathbb{P}(S = b \cup \{i\}) \cdot \mathbb{P}(I = i \mid S = b \cup \{i\}) = \frac{1}{\binom{d}{m}} \cdot \frac{1}{m}.$$

Under the alternative procedure,

$$\mathbb{P}(B = b, I = i) = \mathbb{P}(B = b) \cdot \mathbb{P}(I = i \mid B = b) = \frac{1}{\binom{d}{m-1}} \cdot \frac{1}{d - m + 1} = \frac{1}{\binom{d}{m}} \cdot \frac{1}{m}.$$

Hence, both sampling procedures have the same joint law. ■

Fix any  $B$  and define the alternative instance  $\theta^{(B)}$  by  $\theta_j^{(B)} = \Delta$  if  $j \in B$  and 0 otherwise. For each  $i \notin B$ , define

$$N_i := \sum_{t=1}^T x_t(i) \quad \text{and} \quad A_i := \{i \in \hat{S}\}.$$

Because every played action has exactly  $m$  ones,

$$\sum_{i=1}^d N_i = mT \quad \Rightarrow \quad \sum_{i \notin B} \mathbb{E}_{\theta^{(B)}}[N_i] \leq mT.$$

Because  $|\hat{S}| = m$  always,

$$\sum_{i \notin B} \mathbb{P}_{\theta^{(B)}}(A_i) \leq m.$$

Let  $n := d - m + 1 = |[d] \setminus B|$ . Due to the above two inequalities, there exists a set  $G_B \subseteq [d] \setminus B$  with  $|G_B| \geq n/2$  such that for all  $i \in G_B$ ,

$$\mathbb{E}_{\theta^{(B)}}[N_i] \leq \frac{4mT}{n}, \quad \mathbb{P}_{\theta^{(B)}}(A_i) \leq \frac{4m}{n}. \quad (48)$$

Fix  $i \in G_B$ . Compare  $\theta^{(B)}$  with  $\theta^{(B \cup \{i\})}$ . They differ only in coordinate  $i$  by  $\Delta$ .

Let  $P_B$  and  $P_{B \cup \{i\}}$  be the probability law under  $\theta^{(B)}$  and  $\theta^{(B \cup \{i\})}$ , respectively. For Gaussian noise with variance 1, the one-step KL between  $\mathcal{N}(\mu, 1)$  and  $\mathcal{N}(\mu + \Delta, 1)$  equals  $\Delta^2/2$ . Using the adaptive KL chain rule (Lemma 12) in the same way as we did the multi-task MAB, we get

$$\text{KL}(P_B \| P_{B \cup \{i\}}) = \frac{\Delta^2}{2} \mathbb{E}_{\theta^{(B)}}[N_i] \leq \frac{\Delta^2}{2} \cdot \frac{4mT}{n} = 2\Delta^2 \frac{mT}{n}, \quad (49)$$

where we used (48). Pinsker's inequality gives, for the event  $A_i$ ,

$$\mathbb{P}_{\theta^{(B \cup \{i\})}}(A_i) \leq \mathbb{P}_{\theta^{(B)}}(A_i) + \sqrt{\frac{1}{2} \text{KL}(P_B \| P_{B \cup \{i\}})}.$$

Combining with (48) and (49),

$$\mathbb{P}_{\theta^{(B \cup \{i\})}}(i \in \hat{S}) \leq \frac{4m}{n} + \Delta \sqrt{\frac{mT}{n}}. \quad (50)$$

Assume  $n \geq 20m$  so that  $4m/n \leq 0.2$ . Also assume  $T \leq \frac{1}{2500} \cdot \frac{mn}{\varepsilon^2}$ . Since  $\Delta = 10\varepsilon/m$ , we have

$$\Delta \sqrt{\frac{mT}{n}} = \frac{10\varepsilon}{m} \sqrt{\frac{mT}{n}} \leq \frac{10\varepsilon}{m} \sqrt{\frac{m}{n} \cdot \frac{1}{2500} \frac{mn}{\varepsilon^2}} = 0.2.$$

Plugging into (50) yields, for every  $i \in G_B$ ,

$$\mathbb{P}_{\theta(B \cup \{i\})}(i \in \hat{S}) \leq 0.4. \quad (51)$$

Due to Lemma 29, conditional on  $B$ ,  $I$  is uniform on  $[d] \setminus B$ , and if  $I = i$  then the instance is  $\theta^{(B \cup \{i\})}$ . Therefore, we have:

$$\mathbb{P}(I \in \hat{S} \mid B) = \frac{1}{n} \sum_{i \in [d] \setminus B} \mathbb{P}_{\theta(B \cup \{i\})}(i \in \hat{S}). \quad (52)$$

Split the sum into  $G_B$  and its complement. As  $|G_B| \geq n/2$ , applying (51) for each  $i \in G_B$ , and using the trivial upper bound of 1 for  $i \notin G_B$ , we get:

$$\mathbb{P}(I \in \hat{S} \mid B) \leq \frac{1}{n} \left( \frac{n}{2} \cdot 0.4 + \frac{n}{2} \cdot 1 \right) = 0.7.$$

Averaging over  $B$  gives

$$\mathbb{P}(I \in \hat{S}) \leq 0.7. \quad (53)$$

From (47) and (53), we get:

$$\mathbb{E}[\langle \hat{x}, \theta^{(S)} \rangle] \leq 10\varepsilon \cdot 0.7 = 7\varepsilon.$$

On the other hand, success implies  $\langle \hat{x}, \theta^{(S)} \rangle \geq 9\varepsilon$ , and the value is always nonnegative, hence

$$\mathbb{E}[\langle \hat{x}, \theta^{(S)} \rangle] \geq 9\varepsilon \cdot \mathbb{P}(\text{success}).$$

Therefore  $\mathbb{P}(\text{success}) \leq 7/9$ , so  $\mathbb{P}(\text{fail}) \geq 2/9$  under uniform distribution over the hard instances. By Yao's principle, there exists a fixed instance  $\theta$  on which the algorithm fails with probability at least  $2/9$ . This proves the theorem with a suitable universal constant  $c$  (for instance  $c = 1/2500$ ).  $\blacksquare$

## Appendix I. Unit Ball Lower Bound

Fix  $\varepsilon \in (0, 1]$ ,  $\delta \in (0, 1)$  and an  $(\varepsilon, \delta)$ -PAC algorithm. Let the algorithm run for  $T := \frac{H_1 \log(1/\delta) + H_2}{\varepsilon^2}$  rounds and output  $\hat{x} \in \mathcal{X} := \mathbb{B}_d$ . Define the simple regret for any  $\theta$  as

$$R_T(\theta) := \max_{x \in \mathcal{X}} \langle x, \theta \rangle - \langle \hat{x}, \theta \rangle \geq 0.$$

Observe that  $R_T(\theta) \leq 2\|\theta\|$ .

As the algorithm is  $(\varepsilon, \delta)$ -PAC, we have  $\mathbb{P}(R_T(\theta) > \varepsilon) \leq \delta$ . Hence, taking expectation, we have

$$\mathbb{E}[R_T(\theta)] \leq \varepsilon + 2\|\theta\| \cdot \delta \quad (54)$$

Consider a small enough  $\varepsilon$  such that  $T \geq d^2$ . Now due to Chen et al. (2024), there exists a  $\theta \in \mathbb{R}^d$  such that  $\|\theta\|_2 \leq c_0 d / \sqrt{T}$  and  $\mathbb{E}[R_T(\theta)] \geq c_1 d / \sqrt{T}$  for some absolute constants  $c_0, c_1 \in (0, 1)$ . Let us consider now consider such a  $\theta$ . If we choose  $\delta = c_1/4$ , we have  $\varepsilon \geq \frac{c_1 d}{2\sqrt{T}}$  due to (54). As  $H_1 \leq H_2$  and  $T = \frac{H_1 \log(1/\delta) + H_2}{\varepsilon^2}$ , we have  $H_2 \geq c_2 d^2$  for some absolute constant  $c_2 > 0$ .

Let us also consider the case where stopping time  $\tau$  is randomized as it will be used in the next section. Due to [Chen et al. \(2024\)](#), there exists a  $\theta \in \mathbb{R}^d$  such that  $\|\theta\|_2 \leq c_0 d / \sqrt{T}$  and  $\mathbb{E}[R_T(\theta)] \geq c_1 d / \sqrt{T}$  for some absolute constants  $c_0, c_1 \in (0, 1)$ . Consider one such  $\theta$ . Let assume that  $\mathbb{E}[\tau] = T_0 := \frac{H_1 \log(1/\delta) + H_2}{\varepsilon^2}$  with  $0 < H_1 \leq H_2$ . Let us consider  $T = \frac{1}{\delta} \cdot T_0$ . If the algorithm terminates before  $T$  rounds, we pad the remaining rounds using dummy samples. Now we have the following due to Markov's inequality:

$$\mathbb{E}[R_T(\theta)] \leq \varepsilon + 2\|\theta\|(\mathbb{P}(\tau > T) + \mathbb{P}(R_T(\theta) > \varepsilon, \tau \leq T)) \leq \varepsilon + 4\delta \cdot \|\theta\|$$

Hence, for small enough absolute constant  $\delta$ , we have  $H_2 \geq c_3 \cdot d^2$  for an absolute constant  $c_3$ .

## Appendix J. Polynomial Separation Instance's Non-Adaptive Lower Bound

Recall  $\mathcal{X} = \bigcup_{i=1}^k \mathcal{X}_i \subseteq \mathbb{R}^{kd}$ , where for each  $i \in [k]$  we define the  $i$ -th block of coordinates as  $B_i := \{(i-1)d + 1, \dots, id\}$  and

$$\mathcal{X}_i := \left\{ x \in \mathbb{R}^{kd} : \text{supp}(x) \subseteq B_i, \|x\|_2 \leq 1 \right\}.$$

Consider any (possibly randomized) non-adaptive  $(\varepsilon, \delta)$ -PAC algorithm that uses  $T = \frac{H_1 \log(1/\delta) + H_2}{\varepsilon^2}$  samples with  $0 < H_1 \leq H_2$ , and let  $\tau_i := \sum_{t=1}^T \mathbb{1}\{X_t \in \mathcal{X}_i\}$  denote the (random) number of samples allocated to block  $i$ , so that  $\sum_{i=1}^k \tau_i = T$  and hence there exists  $i^* \in [k]$  with  $T_{i^*} := \mathbb{E}[\tau_{i^*}] \leq T/k$ . Fix such an  $i^*$  and define a hard instance  $\theta \in \mathbb{R}^{kd}$  supported only on block  $B_{i^*}$  as follows: set  $\theta_j = 0$  for all  $j \notin B_{i^*}$ , and on the coordinates in  $B_{i^*}$  set  $(\theta_j)_{j \in B_{i^*}} = \vartheta$  where  $\vartheta \in \mathbb{R}^d$  is chosen according to the unit-ball lower bound from the previous section ([Chen et al. \(2024\)](#) together with the padding and Markov's inequality argument), so that for some absolute constants  $c_0, c_1 \in (0, 1)$  one has  $\|\vartheta\|_2 \leq c_0 \cdot \sqrt{\delta} \cdot d / \sqrt{T_{i^*}}$  and any algorithm that receives  $T_{i^*}/\delta$  informative samples on this block satisfies

$$c_1 \cdot \sqrt{\delta} \cdot d / \sqrt{T_{i^*}} \leq \mathbb{E}[R_{T_{i^*}/\delta}(\vartheta)] \leq \varepsilon + 4\delta \cdot \|\vartheta\|$$

For our chosen  $\theta$ , any action  $x \in \mathcal{X}_j$  with  $j \neq i^*$  satisfies  $\langle x, \theta \rangle = 0$ , so samples taken outside block  $i^*$  are uninformative; moreover, if the algorithm outputs  $\hat{x} \in \mathcal{X}_j$  with  $j \neq i^*$ , then  $\langle \hat{x}, \theta \rangle = 0$ , which is the same value as outputting the zero vector on block  $i^*$ . Therefore the  $(\varepsilon, \delta)$ -PAC requirement for sufficiently small absolute constant  $\delta$  forces  $H_2 \geq c_3 \cdot kd^2$  for an absolute constant  $c_3 > 0$ .

$H_1 \geq \Omega(kd)$  follows directly from [Theorem 3](#).