

---

# Best Arm Identification in Linear Bandits with Linear Dimension Dependency

---

Chao Tao<sup>1</sup> Saúl A. Blanco<sup>1</sup> Yuan Zhou<sup>1 2 3</sup>

## Abstract

We study the best arm identification problem in linear bandits, where the mean reward of each arm depends linearly on an unknown  $d$ -dimensional parameter vector  $\theta$ , and the goal is to identify the arm with the largest expected reward. We first design and analyze a novel randomized  $\theta$  estimator based on the solution to the convex relaxation of an optimal  $G$ -allocation experiment design problem. Using this estimator, we describe an algorithm whose sample complexity depends linearly on the dimension  $d$ , as well as an algorithm with sample complexity dependent on the reward gaps of the best  $d$  arms, matching the lower bound arising from the ordinary top-arm identification problem. We finally compare the empirical performance of our algorithms with other state-of-the-art algorithms in terms of both sample complexity and computational time.

## 1. Introduction

The stochastic multi-armed bandit (MAB) problems are an important framework that not only addresses the fundamental trade-off between exploration and exploitation in sequential experiments, but also crystallizes the challenge of efficiently gathering information before committing to a final decision. Many authors have investigated the *pure exploration* form of the problem with  $N$  independent arms, e.g., to identify the best arm (top arm) ((Carpentier & Locatelli, 2016; Chen et al., 2017b; Even-Dar et al., 2006; Jamieson et al., 2014; Karnin et al., 2013; Kaufmann et al., 2016)) as well as identifying the set of best (top)  $K$  arms under different metrics ((Bubeck et al., 2013; Chen et al.,

2017a; 2014; Kalyanakrishnan & Stone, 2010; Zhou et al., 2014)).

In this paper, we study the top arm identification problem in the stochastic *linear bandit* setting, introduced and studied in (Auer, 2002; Soare et al., 2014)<sup>1</sup>. In the linear bandit setting, instead of having  $N$  independent arms, each of the  $N$  arms is associated with a  $d$ -dimensional *attribute vector*, and the expected reward is the linear combination of its own attribute vector and an unknown global vector  $\theta$ . Each pull of an arm reports a stochastic reward centered at its expectation, and the goal is to pull as few times as possible to identify the arm with the highest expected reward. Due to the linear structure of the problem, pulling each arm reveals information about the global vector  $\theta$  and therefore, indirectly, about the expected reward of other arms.

Linear bandit problems find many applications in practice. For example, in online advertising, suppose the goal is to select an advertisement from a pool to maximize the likelihood for a group of unknown audience to click. The likelihood can usually be approximated by a linear function of logical combinations of various attributes associated with the audience and the ads (such as age, sex, the domain, keywords, ad genres, etc.). Now the linear top-arm identification directly addresses this problem if each ad is abstracted to be an arm with the known attribute vector. While the regret analysis of the exploration-and-exploitation linear bandit problem has been extensively studied (Abbasi-Yadkori et al., 2011; Li et al., 2010; 2017), this paper studies the pure-exploration scenario and is devoted to showing a fixed-confidence algorithm with the optimal sample complexity.

**Problem formulation.** We are given a set of  $N$  arms  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ . Each arm is associated with a  $d$ -dimensional attribute vector and we overload the notation by using  $x_i \in \mathbb{R}^d$  to denote the vector for the arm  $x_i$ . We assume  $\|x_i\|_2 \leq 1$  for all arms. We also assume w.l.o.g. that the rank of  $\mathcal{X}$  is exactly  $d$  (otherwise we can project the set of vectors to a smaller-dimensional space).

There is an unknown global vector  $\theta \in \mathbb{R}^d$  with  $\|\theta\|_2 \leq L$ . Each pull of arm  $x_i$  reports a stochastic reward  $x_i^\top \theta + \epsilon$  where  $\epsilon$  is an independent  $\kappa$ -subgaussian noise (see Sec-

---

<sup>1</sup>Department of Computer Science, Indiana University at Bloomington, Indiana, USA <sup>2</sup>Institute for Theoretical Computer Science, Shanghai University of Finance and Economics, Shanghai, China <sup>3</sup>Department of Industrial and Enterprise Systems Engineering, University of Illinois at Urbana-Champaign, Illinois, USA. Correspondence to: Yuan Zhou <yzhoucs@iu.edu>.

---

<sup>1</sup>Very recently, the problem was also studied independently in (Xu et al., 2018).

tion 2 for the necessary definitions).

Given a confidence parameter  $\delta \in (0, 1)$ , the goal is to design an algorithm with as small query complexity  $Q$  as possible that, with probability at least  $(1 - \delta)$ , sequentially makes at most  $Q$  pulls of the arms, and then identifies the arm with the highest mean reward.

**Previous results.** To make the best arm unambiguously defined, we assume  $\mathcal{X}_{[1]}^T \theta > \mathcal{X}_{[i]}^T \theta$  for  $i > 1$  where  $\mathcal{X}_{[i]}$  denotes the arm in  $\mathcal{X}$  whose mean reward (i.e.,  $\mathcal{X}_{[i]}^T \theta$ ) is the  $i$ -th largest. We let  $\Delta_i = \mathcal{X}_{[1]}^T \theta - \mathcal{X}_{[i]}^T \theta$  be the *reward gap* between arm  $\mathcal{X}_{[1]}$  and  $\mathcal{X}_{[i]}$  for  $i \geq 2$ . In (Soare et al., 2014), the authors showed a few algorithms to identify the top arm and their best query complexity is <sup>2</sup>

$$O\left(\frac{d}{\Delta_2^2}(\ln \delta^{-1} + \ln N + \ln \Delta_2^{-1}) + d^2\right). \quad (1)$$

These algorithms were combined with the Explore-Verify framework proposed in (Karnin, 2016). For very small  $\delta$  (i.e.  $\delta = (\Delta_2/N)^{\omega(1)}$ ), the sample complexity was improved (in (Karnin, 2016)) to

$$O\left(\frac{d}{\Delta_2^2}(\ln N + \ln \Delta_2^{-1}) + \rho^* \ln \delta^{-1} + d^2\right), \quad (2)$$

where  $\rho^* \geq \Delta_2^{-2}$  is an information-theoretic lower bound of the sample complexity on the input instance.

**Our contribution and techniques.** In this paper, we show linear top-arm identification algorithms with confidence  $(1 - \delta)$  and sample complexity

$$O\left(\frac{d}{\Delta_2^2}(\ln \delta^{-1} + \ln N + \ln \ln \Delta_2^{-1})\right). \quad (3)$$

The main improvement compared to (1) is the removal of the additive  $d^2$  term, so that our sample complexity is truly linear with the dimension. We achieve this goal by proposing and analyzing a novel estimator of  $\theta$  described in Section 3. Both our algorithms and (Soare et al., 2014) carefully choose the pulling strategy in order to get a better estimate for  $\theta$ .

<sup>2</sup>The termination condition of their algorithm is  $Q \leq cd(1 + \beta)(\ln Q + \ln N + \ln \delta^{-1})/\Delta_2^2$ , where  $c$  is a constant and  $\beta = d^2/Q$  results from the efficient approximation of the NP-Hard optimal  $G$ -allocation experiment design problem. Solving the inequality for  $Q$  gives the sample complexity bound in (1). The extra  $d^2$  term appears in (2) for the same reason.

We also note that there is a gap between the proofs and the algorithms in (Soare et al., 2014). More specifically, both  $G$ -Allocation and  $\mathcal{X}\mathcal{Y}$ -Allocation algorithms adopt a greedy approach to construct the set of arms to pull. However, the theoretical analysis (Lemmas 5 and 7 in Appendix C of (Soare et al., 2014)) are for a different convex-relaxation-and-rounding procedure which was proposed in (Pukelsheim, 2006).

In (Soare et al., 2014), this is done by (deterministically) approximating the optimal  $G$ -allocation experiment design (which is NP-Hard (Çivril & Magdon-Ismail, 2009)), and leads to the additional  $d^2$  term in the sample complexity. In our algorithms, however, we solve the convex relaxation of the optimal  $G$ -allocation problem, and then, based on the optimal solution, construct a novel randomized estimator of  $\theta$ . While our solution does not guarantee to solve the optimal  $G$ -allocation problem, it saves the extra  $d^2$  term.

Our algorithms also make significant improvement in terms of time complexity. The previous algorithms in (Soare et al., 2014) have to spend  $\Omega(N^2 d^2)$  time to compare the inverse of several matrices before *each* pull. In contrast, the overhead of our algorithm is to solve the convex relaxation and compute the inverse of a  $d \times d$  matrix (which costs roughly  $O(N d^2)$  time), and the remaining time cost before each pull is negligible.

We also develop fully data-dependent algorithms. For confidence parameter  $\delta$ , the sample complexity of our algorithm is

$$O\left(\sum_{i=2}^d \Delta_i^{-2}(\ln \delta^{-1} + \ln N + \ln \ln \Delta_i^{-1})\right). \quad (4)$$

We finally empirically evaluate our algorithms with extensive synthetic and real-world data-sets, and compare the sample complexity and run time with other state-of-the-art algorithms in Section 5.

## 1.1. Related Work

**Relation to optimal experiment designs.** The problem is closely related to the optimal  $G$ -allocation problem in experiment design. The greedy approach proposed in (Soare et al., 2014) was analyzed for different optimality criteria such as  $A$ -optimality ((Bian et al., 2017; Chamon & Ribeiro, 2017a)),  $E$ -optimality ((Chamon & Ribeiro, 2017a)), and  $V$ -optimality ((Chamon & Ribeiro, 2017b)). The  $A$ ,  $D$ ,  $E$ ,  $V$ , and  $G$ -allocation design problems were also studied using other approximation frameworks (e.g. (Allen-Zhu et al., 2017)). We also know that exactly solving the  $G$ -allocation design problem is NP-Hard (Çivril & Magdon-Ismail, 2009).

## Relation to the ordinary top-arm identification problem.

As pointed out in (Soare et al., 2014), the ordinary top-arm identification problem for  $d$  independent arms is a special case of the linear bandit setting when each arm  $x_i$  is associated with the  $i$ -th unit basis vector in  $\mathbb{R}^d$ . Therefore, the sample complexity lower bound (e.g. from (Chen et al., 2017b))  $\Omega(\sum_{i=2}^d \Delta_i^{-2} \ln \delta^{-1})$  directly applies to the linear bandit setting.

**Regret minimization for linear bandits.** Regret minimization, another important goal of multi-armed bandit problems, for linear (and generalized linear) bandits has been extensively studied (e.g., (Abbasi-Yadkori et al., 2011; Auer, 2002; Chu et al., 2011; Dani et al., 2008; Filippi et al., 2010; Li et al., 2017; Rusmevichientong & Tsitsiklis, 2010)). However, its pure exploration counterpart has been less investigated. While regret and sample complexity bounds are not directly comparable, the main technical difference between our algorithms and most regret minimization algorithms is that our sample strategy is computed by solving the convex relaxation of an optimal  $G$ -allocation problem beforehand, which is important for achieving linear dimensional dependence and avoiding computationally expensive matrix inverse updates after each sample.

## 2. Preliminaries

Throughout the paper, we use bold letters to denote random variables. Moreover, for any positive semi-definite matrix  $A$ ,  $\|x\|_A$  is defined to be  $\sqrt{x^T A x}$ . Furthermore, recall that if  $A$  is a  $d \times d$  positive definite matrix, then  $A^{-1}$  exists and is also positive definite. So by the Spectral Theorem, there exists an orthonormal matrix  $O$  and a diagonal matrix  $D$  such that  $A^{-1} = O D O^T$  with the elements at the diagonal of  $D$  being the eigenvalues of  $A^{-1}$ . Now we define  $A^{-1/2} = O D^{1/2} O^T$  where  $D^{1/2}$  is also a diagonal matrix with  $D_{ii}^{1/2} = \sqrt{D_{ii}}$  for all  $i \in \{1, 2, 3, \dots, d\}$ .

We now define *subgaussian random variables* and spell out some of the mathematical tools that will be used in the proofs.

**Definition 1** (Subgaussian Random Variable). *For any  $b > 0$ , a real-valued random variable  $X$  is said to be  $b$ -subgaussian if it has the property that for every  $t \in \mathbb{R}$  one has  $\mathbb{E}[e^{tX}] \leq e^{b^2 t^2 / 2}$ .*

**Proposition 2** (Special case of Lemma 2.3.18, (Stroock, 2011)). *If  $X$  is  $b$ -subgaussian, then for any  $\alpha \in \mathbb{R}$ , the random variable  $\alpha X$  is  $|\alpha|b$ -subgaussian. Moreover, if  $X_1, X_2$  are independent random variables such that  $X_i$  is  $b_i$ -subgaussian, then  $X_1 + X_2$  is  $\sqrt{b_1^2 + b_2^2}$ -subgaussian.*

**Proposition 3** (Subgaussian Tail Estimate, special case of Lemma 2.3.18, (Stroock, 2011)). *If  $X$  is  $b$ -subgaussian, then for any  $\epsilon > 0$  we have that  $\Pr[|X| \geq \epsilon] \leq 2 \exp\left(-\frac{\epsilon^2}{2b^2}\right)$ .*

We will find the following concentration inequalities useful.

**Proposition 4** (Multiplicative Chernoff Bound). *Suppose  $X_i$ 's ( $1 \leq i \leq n$ ) are independent random variables taking values in  $[0, 1]$ . Let  $X = \sum_{i=1}^n X_i$  and  $\mu = \mathbb{E}[X]$ . Then for any  $\delta \in [0, 1]$ , it holds that  $\Pr[X \geq (1 + \delta)\mu] \leq e^{-\frac{\delta^2 \mu}{3}}$ .*

**Proposition 5** (Bernstein Inequality). *Let  $X_i$  ( $1 \leq i \leq n$ ) be independent zero-mean random variables. Suppose that  $|X_i| \leq M$  almost surely, for any  $i$ . Then,*

*for all positive  $t$ , it holds that  $\Pr(\sum_{i=1}^n X_i > t) \leq \exp\left(-\frac{t^2/2}{\sum_{i=1}^n \mathbb{E}[X_i^2] + \frac{1}{3}Mt}\right)$ .*

Finally, we introduce the following theorem which will be crucially used in the design of our estimator for  $\theta$ .

**Proposition 6** (Restatement of the Equivalence-Theorem in (Kiefer & Wolfowitz, 1960)). *Given a set of  $d$ -dimensional vectors  $\mathcal{X} \subseteq \mathbb{R}^d$ , for any distribution  $\lambda$  supported on  $\mathcal{X}$  such that  $M(\lambda) = \mathbb{E}_{z \sim \lambda}[zz^T]$  is non-singular, we define  $f(x; \lambda) = x^T M(\lambda)^{-1} x$ . The following extremum problems are equivalent:*

(a) *Choosing  $\lambda$  so that  $\lambda$  maximizes  $\det M(\lambda)$ .* (5)

(b) *Choosing  $\lambda$  so that  $\lambda$  minimizes  $\max_{x \in \mathcal{X}} f(x; \lambda)$ .* (6)

*Moreover, since  $\mathbb{E}_{x \sim \lambda}[f(x; \lambda)] = d$ , it follows that  $\max_{x \in \mathcal{X}} f(x; \lambda) \geq d$ . Therefore, a sufficient condition for  $\lambda$  to satisfy (6) is*

$$\max_{x \in \mathcal{X}} f(x; \lambda) = d. \quad (7)$$

## 3. Main Theorems on the Estimator

### 3.1. The New Estimator

Let  $\lambda^*$  be the solution to the three problems defined in Proposition 6. Let  $y_1, \dots, y_n$  be  $n$  i.i.d. samples following from the distribution  $\lambda^*$ , corresponding to the  $n$  pulled arms. Let the corresponding rewards be  $r_1, \dots, r_n$  respectively. In the setting of linear bandits, we have  $r_i = y_i^T \theta + \epsilon_i$  where  $\epsilon_i$  is  $\kappa$ -subgaussian. Let  $b_i = r_i y_i$  and  $b = \sum_{i=1}^n b_i$ . Now we define our estimator  $\hat{\theta}$  to be

$$\hat{\theta} = A^{-1} b, \quad (8)$$

where  $A = nM(\lambda^*)$ .

**Remark 1.** *When  $y_1, \dots, y_n$  are the  $n$  sampled arms, the usual maximum-likelihood estimator used in literature (e.g. (Li et al., 2017; Soare et al., 2014)) is  $(\sum_{i=1}^n y_i y_i^T)^{-1} b$ . However, it is technically difficult to analyze the (spectral) concentration properties of  $(\sum_{i=1}^n y_i y_i^T)^{-1}$ , and we novelly use  $(\mathbb{E}[\sum_{i=1}^n y_i y_i^T])^{-1} = A^{-1}$  instead.*

### 3.2. The Key Lemma

For any  $\theta, x \in \mathbb{R}^d$ , and  $\ell \geq 1$ , we define a  $\theta$ -dependent error and a  $\theta$ -free error,  $\text{Err}_{\lambda^*}(x, \ell, \theta)$  and  $\text{err}_{\lambda^*}(x, \ell)$ , respectively, as follows:

$$\begin{aligned} \text{Err}_{\lambda^*}(x, \ell, \theta) &= \sqrt{\frac{2 \mathbb{E}_{z \sim \lambda^*}[(x^T M(\lambda^*)^{-1} z \cdot z^T \theta)^2]}{\ell}} \\ &\quad + \frac{2(|x^T \theta| + L \|x\|_{M(\lambda^*)^{-1}} \sqrt{d})}{3\ell} \\ &\quad + \sqrt{2\kappa} \sqrt{\frac{\|x\|_{M(\lambda^*)^{-1}}^2}{\ell} \sqrt{\frac{3d}{\ell}} + \frac{\|x\|_{M(\lambda^*)^{-1}}^2}{\ell}}, \end{aligned}$$

$$\begin{aligned} \text{and } \text{err}_{\lambda^*}(x, \ell) &= \frac{\sqrt{2}L\|x\|_{M(\lambda^*)^{-1}}}{\sqrt{\ell}} \\ &+ \frac{2L(\|x\|_2 + \|x\|_{M(\lambda^*)^{-1}}\sqrt{d})}{3\ell} \\ &+ \sqrt{2}\kappa\sqrt{\frac{\|x\|_{M(\lambda^*)^{-1}}^2}{\ell}\sqrt{\frac{3d}{\ell}} + \frac{\|x\|_{M(\lambda^*)^{-1}}^2}{\ell}}. \end{aligned}$$

We show the following key lemma.

**Lemma 7.** *When  $n \geq \ell \ln(5/\delta)$  where  $\ell \geq 3d$ , for any fixed vector (not necessarily an arm), with probability at least  $1 - \delta$ , it holds that  $|x^T(\theta - \hat{\theta})| \leq \text{Err}_{\lambda^*}(x, \ell, \theta) \leq \text{err}_{\lambda^*}(x, \ell)$ .*

*Proof.* Note that  $\mathbf{b} = \sum_{i=1}^n \mathbf{r}_i \mathbf{y}_i = \sum_{i=1}^n (\mathbf{y}_i \mathbf{y}_i^T \theta + \mathbf{y}_i \epsilon_i)$ . For any given vector  $x$ , the absolute error between the real reward and the estimated reward is

$$\begin{aligned} |x^T(\theta - \hat{\theta})| &= |x^T(\theta - A^{-1}\mathbf{b})| = \\ &\left| \sum_{i=1}^n (x^T \theta / n - x^T A^{-1} \mathbf{y}_i \mathbf{y}_i^T \theta) - \sum_{i=1}^n x^T A^{-1} \mathbf{y}_i \epsilon_i \right|. \quad (9) \end{aligned}$$

For any  $i \in \{1, 2, \dots, n\}$ , set  $\mathbf{Y}_i = x^T \theta / n - x^T A^{-1} \mathbf{y}_i \mathbf{y}_i^T \theta$  and  $\mathbf{Z}_i = x^T A^{-1} \mathbf{y}_i \epsilon_i$ , then (9) can be written as

$$\left| \sum_{i=1}^n \mathbf{Y}_i - \sum_{i=1}^n \mathbf{Z}_i \right| \leq |\mathbf{Y}| + |\mathbf{Z}|,$$

where  $\mathbf{Y} = \sum_{i=1}^n \mathbf{Y}_i$  and  $\mathbf{Z} = \sum_{i=1}^n \mathbf{Z}_i$ .

The following two claims bound  $\mathbf{Y}$  and  $\mathbf{Z}$  respectively. We provide the proofs of both claims after the proof of this lemma.

**Claim 8.** *For  $n \geq \ell \ln \delta^{-1}$ , with probability  $(1 - \delta)$ , we have  $|\mathbf{Y}| \leq \sqrt{\frac{2 \mathbb{E}_{\mathbf{z} \sim \lambda^*}[(x^T M(\lambda^*)^{-1} \mathbf{z} \cdot \mathbf{z}^T \theta)^2]}{\ell}} + \frac{2(|x^T \theta| + L\|x\|_{M(\lambda^*)^{-1}}\sqrt{d})}{3\ell}$ .*

**Claim 9.** *For  $n \geq \ell \ln(4/\delta)$  where  $\ell \geq 3d$ , with probability  $(1 - \delta)$ , we have  $|\mathbf{Z}| \leq \sqrt{2}\kappa\sqrt{\frac{\|x\|_{M(\lambda^*)^{-1}}^2}{\ell}\sqrt{\frac{3d}{\ell}} + \frac{\|x\|_{M(\lambda^*)^{-1}}^2}{\ell}}$ .*

The first inequality of Lemma 7 follows immediately by combining Claim 8 and Claim 9, and a union bound. The second inequality of Lemma 7 holds because of  $\|\mathbf{z}\|_2 \leq 1$ ,  $\|\theta\|_2 \leq L$ , and the definition of  $M(\lambda^*)$ .  $\square$

*Proof of Claim 8.* Let  $t$  denote  $\|x\|_{M(\lambda^*)^{-1}}^2$  and  $s$  denote  $\mathbb{E}_{\mathbf{z} \sim \lambda^*}[(x^T M(\lambda^*)^{-1} \mathbf{z} \cdot \mathbf{z}^T \theta)^2]$ . Recall that  $\mathbf{Y}_i = x^T \theta / n - x^T A^{-1} \mathbf{y}_i \mathbf{y}_i^T \theta$ .

It is easy to verify that  $\mathbb{E}[\mathbf{Y}_i] = 0$ . We also have

$$\begin{aligned} |\mathbf{Y}_i| &\leq |x^T \theta / n| + |x^T A^{-1} \mathbf{y}_i \mathbf{y}_i^T \theta| \\ &\leq |x^T \theta| / n + L|x^T A^{-1} \mathbf{y}_i|, \end{aligned}$$

where the last inequality is due to  $\|\theta\|_2 \leq L$  and  $\|\mathbf{y}_i\|_2 \leq 1$ . By Cauchy-Schwartz inequality, we have

$$\begin{aligned} |x^T A^{-1} \mathbf{y}_i| &= |\langle A^{-1/2} x, A^{-1/2} \mathbf{y}_i \rangle| \\ &\leq \|A^{-1/2} x\|_2 \cdot \|A^{-1/2} \mathbf{y}_i\|_2 \\ &= \sqrt{x^T A^{-1} x} \sqrt{\mathbf{y}_i^T A^{-1} \mathbf{y}_i} \leq \sqrt{td/n^2}, \quad (10) \end{aligned}$$

where the last inequality is due to the definition of  $A$  and Proposition 6. Hence it holds that

$$|\mathbf{Y}_i| \leq |x^T \theta| / n + L\sqrt{td/n^2}. \quad (11)$$

Next, we bound  $\sum_{i=1}^n \mathbb{E}[\mathbf{Y}_i^2]$

$$\begin{aligned} &= \sum_{i=1}^n \mathbb{E}[(x^T \theta)^2 / n^2 - 2x^T \theta / n \cdot x^T A^{-1} \mathbf{y}_i \mathbf{y}_i^T \theta \\ &\quad + (x^T A^{-1} \mathbf{y}_i \cdot \mathbf{y}_i^T \theta)^2] \\ &\leq (x^T \theta)^2 / n - \sum_{i=1}^n 2x^T \theta / n \cdot x^T A^{-1} (\mathbb{E} \mathbf{y}_i \mathbf{y}_i^T) \theta \\ &\quad + \sum_{i=1}^n \mathbb{E}[(x^T A^{-1} \mathbf{y}_i \cdot \mathbf{y}_i^T \theta)^2] \\ &= - (x^T \theta)^2 / n + \mathbb{E}_{\mathbf{z} \sim \lambda^*}[(x^T M(\lambda^*)^{-1} \mathbf{z} \cdot \mathbf{z}^T \theta)^2] / n \\ &\leq s/n, \quad (12) \end{aligned}$$

where the second inequality is due to  $\|\theta\|_2 \leq L$  and  $\|\mathbf{y}_i\|_2 \leq 1$ , and the third equality is due to the definition of  $A$ .

Since  $\mathbb{E}[\mathbf{Y}_i] = 0$ , the  $\mathbf{Y}_i$ 's are independent, and also we have (11) and (12), applying Proposition 5, it holds that

$$\begin{aligned} \Pr[|\mathbf{Y}| > \epsilon] &\leq \exp\left(-\frac{\epsilon^2/2}{s/n + \frac{1}{3}(|x^T \theta|/n + L\sqrt{td/n^2})\epsilon}\right) \\ &= \exp\left(-\frac{n\epsilon^2}{2s + \frac{2}{3}(|x^T \theta| + L\sqrt{td})\epsilon}\right). \end{aligned}$$

By letting  $\epsilon = \sqrt{\frac{2s}{\ell}} + \frac{2(|x^T \theta| + L\sqrt{td})}{3\ell}$ , it holds that

$$\Pr\left[|\mathbf{Y}| > \sqrt{\frac{2s}{\ell}} + \frac{2(|x^T \theta| + L\sqrt{td})}{3\ell}\right] \leq \delta$$

for  $n \geq \ell \ln \delta^{-1}$ , which concludes the proof.  $\square$

*Proof of Claim 9.* Let  $t$  denote  $\|x\|_{M(\lambda^*)}^2$ . Recall that  $\mathbf{Z}_i = x^T A^{-1} \mathbf{y}_i \epsilon_i$ .

By Proposition 2,  $\mathbf{Z}_i$  can be seen as a  $(\kappa \sigma_i)$ -subgaussian random variable where  $\sigma_i = \sqrt{x^T A^{-1} \mathbf{y}_i \mathbf{y}_i^T A^{-1} x}$ . Again by Proposition 2,  $\mathbf{Z}$  is a  $(\kappa \sigma)$ -subgaussian random variable where  $\sigma = \sqrt{\sum_{i=1}^n \sigma_i^2}$ .

For a given  $\epsilon$  and a realization of  $\sigma^2$ , by Proposition 3 we know  $\Pr[|\mathbf{Z}| > \epsilon] \leq 2 \exp(-\epsilon^2/(2\kappa^2 \sigma^2))$ . When  $\sigma^2 \leq \frac{\epsilon^2}{2\kappa^2 \ln(4/\delta)}$ , we have  $\Pr[|\mathbf{Z}| > \epsilon] \leq \delta/2$ , which means

$$\Pr \left[ |\mathbf{Z}| > \epsilon \mid \sigma^2 \leq \frac{\epsilon^2}{2\kappa^2 \ln(4/\delta)} \right] \leq \delta/2. \quad (13)$$

Next, we give a bound on  $\Pr \left[ \sigma^2 \leq \frac{\epsilon^2}{2\kappa^2 \ln(4/\delta)} \right]$ . Note that

$$\begin{aligned} \mathbb{E}[\sigma_i^2] &= \mathbb{E}[x^T A^{-1} \mathbf{y}_i \mathbf{y}_i^T A^{-1} x] \\ &= x^T A^{-1} \mathbb{E}[\mathbf{y}_i \mathbf{y}_i^T] A^{-1} x = t/n^2 \end{aligned} \quad (14)$$

and by (10) we have

$$\sigma_i^2 = |x^T A^{-1} \mathbf{y}_i|^2 \leq td/n^2. \quad (15)$$

Let  $\lambda_i = \sigma_i^2 n^2/(td)$ . By (14) we have  $\mathbb{E}[\lambda_i] = 1/d$  and by (15) we have  $\lambda_i \in [0, 1]$ . Now applying Proposition 4, for  $\tau \in [t/n, 2t/n]$ , we have

$$\begin{aligned} \Pr[\sigma^2 > \tau] &= \Pr \left[ \sum_{i=1}^n \lambda_i > \tau n^2/(td) \right] \\ &\leq \exp \left( -\frac{(\tau n/t - 1)^2 \cdot n/d}{3} \right). \end{aligned}$$

Hence, by letting  $\tau = \frac{t}{n} \sqrt{\frac{3d \ln(2/\delta)}{n}} + \frac{t}{n}$ , it holds that

$$\Pr \left[ \sigma^2 \leq \frac{t}{n} \sqrt{\frac{3d \ln(2/\delta)}{n}} + \frac{t}{n} \right] \geq 1 - \delta/2. \quad (16)$$

Therefore, combining (13) and (16), we have  $\Pr[|\mathbf{Z}| \leq \epsilon]$

$$\begin{aligned} &\geq \Pr \left[ |\mathbf{Z}| \leq \epsilon \mid \sigma^2 \leq \frac{\epsilon^2}{2\kappa^2 \ln(4/\delta)} \right] \\ &= \Pr \left[ |\mathbf{Z}| \leq \epsilon \mid \sigma^2 \leq \frac{\epsilon^2}{2\kappa^2 \ln(4/\delta)} \right] \Pr \left[ \sigma^2 \leq \frac{\epsilon^2}{2\kappa^2 \ln(4/\delta)} \right] \\ &\geq (1 - \delta/2)(1 - \delta/2) \geq 1 - \delta, \end{aligned}$$

whenever  $\frac{\epsilon^2}{2\kappa^2 \ln(4/\delta)} \geq \frac{t}{n} \sqrt{\frac{3d \ln(2/\delta)}{n}} + \frac{t}{n}$ .

By letting  $\epsilon = \sqrt{2\kappa} \sqrt{\frac{t}{\ell} \sqrt{\frac{3d}{\ell}} + \frac{t}{\ell}}$ , we get

$\Pr \left[ |\mathbf{Z}| \leq \sqrt{2\kappa} \sqrt{\frac{t}{\ell} \sqrt{\frac{3d}{\ell}} + \frac{t}{\ell}} \right] \geq 1 - \delta$  for  $n \geq \ell \ln(4/\delta)$ , which concludes the proof of this claim.  $\square$

### 3.3. Error Bounds for Every Arm

We first simplify  $\text{err}_{\lambda^*}(x, \ell)$  and get the following lemma.

**Lemma 10.** Let  $c_0 = \max\{4L^2, 3\kappa^2, 3\}$ , when  $n \geq c_0 \ell \ln(5/\delta)$  where  $\ell \geq d$ , with probability  $(1 - \delta)$ , we have  $|x^T(\theta - \hat{\theta})| \leq \sqrt{\frac{2\|x\|_2^2 + 2\sqrt{d}\|x\|_2\|x\|_{M(\lambda^*)}^{-1} + (4+2\sqrt{d/\ell})\|x\|_{M(\lambda^*)}^2}{\ell}}$ .

Please refer to Appendix A for the proof of Lemma 10.

If we further make assumptions on  $\|x\|_{M(\lambda^*)}^{-1}$ , we have

**Lemma 11.** When  $n \geq c_0 \ell \ln(5/\delta)$  where  $\ell \geq d$ , and  $\|x\|_{M(\lambda^*)}^{-1} \leq c_1 \sqrt{d}$  where  $\|x\|_2 \leq c_1$ , we have  $\Pr \left[ |x^T(\theta - \hat{\theta})| \leq c_1 \sqrt{\frac{2+6d+2d\sqrt{d/\ell}}{\ell}} \right] \geq 1 - \delta$ .

Note that  $\|x\|_{M(\lambda^*)}^{-1} \leq \sqrt{d}$  and  $\|x\|_2 \leq 1$  for every  $x \in \mathcal{X}$  (by Proposition 6). Applying Lemma 11 for every arm in  $\mathcal{X}$  and via a union bound, we have

**Theorem 12.** When  $n \geq c_0 \cdot \frac{2+(6+\epsilon)d}{\epsilon^2} \ln(5N/\delta)$  where  $\epsilon \leq 3$ , we have  $\Pr \left[ |x^T \theta - x^T \hat{\theta}| \leq \epsilon, \forall x \in \mathcal{X} \right] \geq 1 - \delta$ .

Please refer to Appendix B for the proof of Theorem 12.

Next, we apply Lemma 7 for every  $y \in \mathcal{Y}$  where  $\mathcal{Y}$  is the set  $\{x - x' \mid x, x' \in \mathcal{X}\}$ . Via a union bound, we have

**Theorem 13.** When  $n \geq \ell \ln(5N^2/(2\delta))$  where  $\ell \geq 3d$ , with probability  $(1 - \delta)$ , we have  $|y^T \theta - y^T \hat{\theta}| \leq \text{Err}_{\lambda^*}(y, \ell, \theta) \leq \text{err}_{\lambda^*}(y, \ell), \forall y \in \mathcal{Y}$ .

For every  $y = x - x' \in \mathcal{Y}$  (where  $x, x' \in \mathcal{X}$ ), by Proposition 6, we have  $\|y\|_{M(\lambda^*)}^{-1} = \|x - x'\|_{M(\lambda^*)}^{-1} \leq \|x\|_{M(\lambda^*)}^{-1} + \|x'\|_{M(\lambda^*)}^{-1} \leq 2\sqrt{d}$ . Note that  $\|y\|_2 \leq 2$ . Similar to Theorem 12, we apply Lemma 11 for every  $y \in \mathcal{Y}$  and via a union bound, we have

**Theorem 14.** When  $n \geq 4c_0 \frac{2+(6+\epsilon)d}{\epsilon^2} \ln(5N^2/(2\delta))$  where  $\epsilon \leq 7$ , with probability  $(1 - \delta)$ , we have  $|y^T \theta - y^T \hat{\theta}| \leq \text{err}_{\lambda^*}(y, \ell) \leq \epsilon, \forall y \in \mathcal{Y}$ , where  $\ell = 4c_0 \frac{2+(6+\epsilon)d}{\epsilon^2}$ .

### 3.4. Bounds for a Significant Number of Arms

We also prove the following theorem and corollary which work better when we have a significant number of (or infinitely many) arms.

**Theorem 15.** When  $n \geq c_0 \ell \ln(5d/\delta)$  where  $\ell \geq d$ , with probability  $1 - \delta$ , we have  $|x^T(\theta - \hat{\theta})| \leq \sqrt{\frac{(2+d)d + (5d+2d\sqrt{d/\ell})\|x\|_{M(\lambda^*)}^2}{\ell}}, \forall x \in \mathbb{R}^d : \|x\|_2 \leq 1$ .

Please refer to Appendix C for the proof of Theorem 15.

**Corollary 16.** When  $n \geq c_0 \frac{(4d+(6+\epsilon)d^2)}{\epsilon^2} \ln(5d/\delta)$  where  $\epsilon \leq 3\sqrt{d}$ , we have  $\Pr \left[ |x^T(\theta - \hat{\theta})| \leq \epsilon, \forall x \in \mathcal{X} \right] \geq 1 - \delta$ .

*Proof.* Note that  $\|x\|_{M(\lambda^*)^{-1}}^2 \leq d$  for all  $x \in \mathcal{X}$ . We prove the corollary by applying Theorem 15 with  $\ell = (4d + (6 + \epsilon)d^2)/\epsilon^2$ .  $\square$

## 4. Algorithms

Given a set  $S$  of arms, we use  $\lambda_S^*$  to represent the solution to the three problems defined in Proposition 6, which can be computed by solving the optimization problem stated in (5).

### 4.1. The Estimator for $\theta$

First we define a procedure,  $\text{VECTOREST}(S, n)$ , that aims to estimate the underlying unknown  $\theta$  given a set  $S$  of arms and  $n$  samples.  $\text{VECTOREST}(S, n)$  utilizes the estimator in (8), and its described in Algorithm 1.

---

#### Algorithm 1 $\text{VECTOREST}(S, n)$

---

- 1: **Input:** A set  $S$  of arms, and  $n$  samples.
  - 2: Let  $y_1, \dots, y_n$  be the  $n$  samples acquired from  $S$  according to the distribution  $\lambda_S^*$ .
  - 3: Pull arms  $y_1, \dots, y_n$ , and suppose their corresponding rewards are  $r_1, \dots, r_n$  respectively.
  - 4: Compute  $A \leftarrow n \cdot \sum_{x \in S} \lambda_S^*(x) x x^T$  and  $b \leftarrow \sum_{i=1}^n r_i y_i$ .
  - 5:  $\hat{\theta} \leftarrow A^{-1} b$ .
  - 6: **Output:** The estimate  $\hat{\theta}$ .
- 

### 4.2. $\mathcal{X}$ -dependent Algorithm

We now describe an iterative algorithm that eliminates  $(|S| - p)$  suboptimal arms from  $S$  with probability at least  $(1 - \delta)$ . We call this algorithm  $\text{ELIMTIL}_p$  and present the details in Algorithm 2. In essence, during the  $r$ -th iteration,  $\text{ELIMTIL}_p$  uses  $\text{VECTOREST}$  to get an  $\frac{\epsilon_r}{2}$ -close estimate of  $x^T \theta$  for all  $x \in S_r$  with probability at least  $1 - \delta_r$ , and then discards all arms whose estimated mean rewards are  $\epsilon_r$  worse than that of the highest estimated mean.  $\text{ELIMTIL}_p$  continues until at most  $p$  arms remain. By letting  $\epsilon_r$  and  $\delta_r$  decrease exponentially, we are able to keep the best arm in the set of output arms with probability at least  $(1 - \delta)$ .

---

#### Algorithm 2 $\text{ELIMTIL}_p(S, \delta)$

---

- 1: **Input:** Arms set  $S$  and confidence level  $\delta$ .
  - 2: Initialize  $S_1 \leftarrow S, r \leftarrow 1$ .
  - 3: **while**  $|S_r| > p$  **do**
  - 4:   Set  $\epsilon_r \leftarrow 1/2^r, \delta_r \leftarrow 6/\pi^2 \cdot \delta/r^2$ .
  - 5:    $\hat{\theta}_r \leftarrow \text{VECTOREST}\left(S, c_0 \frac{2+(6+\epsilon_r/2)d}{(\epsilon_r/2)^2} \ln \frac{5|S|}{\delta_r}\right)$ .
  - 6:   Select arm  $a_r \leftarrow \arg\max_{x \in S_r} x^T \hat{\theta}_r$ .
  - 7:    $S_{r+1} \leftarrow S_r \setminus \{x \in S_r | x^T \hat{\theta}_r < x_{a_r}^T \hat{\theta}_r - \epsilon_r\}$ .
  - 8:    $r \leftarrow r + 1$ .
  - 9: **Output:** Set  $S_r$ .
- 

**Lemma 17.** *With probability at least  $(1 - \delta)$ ,  $\text{ELIMTIL}_p(S, \delta)$  satisfies the following properties: (i) the procedure outputs a set of at most  $p$  arms with the best arm included; and (ii) the sample complexity is  $O\left(\frac{c_0 d}{\Delta_{p+1}^2} (\ln \delta^{-1} + \ln |S| + \ln \ln \Delta_{p+1}^{-1})\right)$ .*

Please refer to Appendix D for the proof of Lemma 17.

As a consequence of Lemma 17, we get the following corollary by setting  $S = \mathcal{X}$ .

**Corollary 18.** *With probability at least  $(1 - \delta)$ ,  $\text{ELIMTIL}_1(\mathcal{X}, \delta)$  outputs the best arm, and the sample complexity is  $O\left(\frac{c_0 d}{\Delta_2^2} (\ln \delta^{-1} + \ln |\mathcal{X}| + \ln \ln \Delta_2^{-1})\right)$ .*

If we use the bound in Theorem 15 and Corollary 16 instead, we have the following statement.

**Corollary 19.** *There exists an algorithm that, with probability at least  $(1 - \delta)$ , outputs the best arm, and uses at most  $O\left(\frac{c_0 d^2}{\Delta_2^2} (\ln \delta^{-1} + \ln d + \ln \ln \Delta_2^{-1})\right)$  samples.*

### 4.3. $\mathcal{Y}$ -dependent Algorithm

We now present the algorithm  $\mathcal{Y}$ - $\text{ELIMTIL}_p$ , which has the same asymptotic sample complexity guarantee as  $\text{ELIMTIL}_p$ , but performs more efficiently in practice. This improvement is achieved thanks to the following two ideas.

First, when comparing the mean rewards of two arms (namely  $x$  and  $x'$ ), instead of checking the confidence intervals of the two arms, we turn to the confidence interval of  $y^T \theta$  (where  $y = x - x'$ ) and check the corresponding confidence interval.

Second, by Lemma 7, at the  $r$ -th iteration, the confidence interval of  $y^T \theta$  has length  $\text{Err}_{\lambda_T^*}(y, \ell_r, \theta) \leq \text{err}_{\lambda_T^*}(y, \ell_r)$  where the latter quantity is an upper estimate of the first one and does not need the knowledge about  $\theta$  in advance. At the  $r$ -th iteration, we use  $\hat{\theta}_r$  as an estimate of  $\theta$  and define

$$\begin{aligned} \widehat{\text{Err}}_{\lambda_T^*}(y, \ell_r) = & \sqrt{\frac{2 \mathbb{E}_{z \sim \lambda_T^*} [(y^T M(\lambda_T^*)^{-1} z \cdot (|z^T \hat{\theta}_r| + \epsilon_r/2))^2]}{\ell}} \\ & + \frac{2(|y^T \hat{\theta}_r| + \epsilon_r/2) + L \|y\|_{M(\lambda_T^*)^{-1}} \sqrt{d}}{3\ell} \\ & + \sqrt{2\kappa} \sqrt{\frac{\|y\|_{M(\lambda_T^*)^{-1}}^2 \sqrt{3d/\ell}}{\ell} + \frac{\|y\|_{M(\lambda_T^*)^{-1}}^2}{\ell}}. \end{aligned}$$

We will be able to show that  $\text{Err}_{\lambda_T^*}(y, \ell_r, \theta) \leq \widehat{\text{Err}}_{\lambda_T^*}(y, \ell_r)$  holds with high probability and therefore we also use  $\widehat{\text{Err}}_{\lambda_T^*}(y, \ell_r)$  as an upper estimate of the length of the confidence interval which in many cases proves to be tighter than  $\text{err}_{\lambda_T^*}(y, \ell_r)$ .

We present the details of  $\mathcal{Y}$ -ELIMTIL<sub>p</sub> in Algorithm 3. Note that our theoretical analysis does not require  $\hat{\theta}_r$  for different values of  $r$  to be independent. Therefore, the samples may be reused across different invocations of VECTOREST.

---

**Algorithm 3**  $\mathcal{Y}$ -ELIMTIL<sub>p</sub>( $S, T, \delta$ )
 

---

- 1: **Input:** Set  $S$  of active arms (the ones not yet eliminated), set  $T \supseteq S$  of arms those can be pulled, and confidence level  $\delta$ .
  - 2: Initialize  $S_1 \leftarrow S, r \leftarrow 1$ .
  - 3: **while**  $|S_r| > p$  **do**
  - 4:   Set  $\delta_r \leftarrow 6/\pi^2 \cdot \delta/r^2, \ell_r \leftarrow 4c_0(2 + (6 + 4/1.1^r)d)(1.1^r/4)^2$ .
  - 5:    $\hat{\theta}_r \leftarrow \text{VECTOREST}(T, \ell_r \ln(5|S|^2/(2\delta_r)))$ .
  - 6:   Select arm  $a_r \leftarrow \arg\max_{x \in S_r} x^T \hat{\theta}_r$ .
  - 7:    $S_{r+1} \leftarrow S_r \setminus \{x \in S_r | (x_{a_r} - x)^T \hat{\theta}_r > \min\{\text{err}_{\lambda_T^*}(x_{a_r} - x, \ell_r), \widehat{\text{Err}}_{\lambda_T^*}(x_{a_r} - x, \ell_r)\}\}$ .
  - 8:    $r \leftarrow r + 1$ .
  - 9: **Output:** Set  $S_r$ .
- 

**Lemma 20.** *With probability at least  $(1 - \delta)$ ,  $\mathcal{Y}$ -ELIMTIL<sub>p</sub>( $S, T, \delta$ ) satisfies the following properties: (i) the procedure outputs a set of at most  $p$  arms in  $S$  with the best arm included; and (ii) the sample complexity is  $O\left(\frac{c_0 d}{\Delta_{p+1}^2} (\ln \delta^{-1} + \ln |S| + \ln \ln \Delta_{p+1}^{-1})\right)$ .*

Please refer to Appendix E for the proof of Lemma 20.

Lemma 20 gives the following corollary by setting  $S$  and  $T$  to  $\mathcal{X}$ .

**Corollary 21.** *With probability at least  $(1 - \delta)$ ,  $\mathcal{Y}$ -ELIMTIL<sub>1</sub>( $\mathcal{X}, \mathcal{X}, \delta$ ) outputs the best arm, and the sample complexity is  $O\left(\frac{c_0 d}{\Delta_2^2} (\ln \delta^{-1} + \ln |\mathcal{X}| + \ln \ln \Delta_2^{-1})\right)$ .*

#### 4.4. Fully Data-dependent Algorithm

Using the  $\mathcal{Y}$ -ELIMTIL algorithm as a subroutine, we introduce an algorithm that outputs the best arm with probability at least  $(1 - \delta)$  and whose sample complexity depends on  $\Delta_1, \dots, \Delta_d$ . To achieve this goal, the algorithm runs in rounds. During each round  $r$ , it invokes  $\mathcal{Y}$ -ELIMTIL to identify the top- $(d/2^r)$  arms and discards the remaining arms. We call this algorithm ALBA and include the details in Algorithm 4.

---

**Algorithm 4** Adaptive Linear Best Arm, ALBA( $\mathcal{X}, \delta$ )
 

---

- 1: **Input:** Arms set  $\mathcal{X}$  and confidence level  $\delta$ .
  - 2: Initialize  $S_0 \leftarrow \mathcal{X}$ .
  - 3: **for**  $r \leftarrow 0$  **to**  $\lfloor \log_2 d \rfloor$  **do**
  - 4:   Set  $\delta_r \leftarrow 6/\pi^2 \cdot \delta/((r+1)^2)$ .
  - 5:    $S_{r+1} \leftarrow \mathcal{Y}\text{-ELIMTIL}_{\lfloor d/2^r \rfloor}(S_r, \mathcal{X} \cap \text{span}(S_r), \delta_r)$ .<sup>3</sup>
  - 6:    $r \leftarrow r + 1$ .
  - 7: **Output:** Best arm in  $\mathcal{X}$ .
- 

**Theorem 22.** *With probability at least  $(1 - \delta)$ , the following are true: (i) ALBA( $\mathcal{X}, \delta$ ) outputs the best arm; and (ii) the sample complexity is  $O\left(\sum_{i=2}^d \frac{c_0}{\Delta_i^2} (\ln \delta^{-1} + \ln |\mathcal{X}| + \ln \ln \Delta_i^{-1})\right)$ .*

Please refer to Appendix F for the proof of Theorem 22.

## 5. Experiments

We test our algorithms  $\mathcal{Y}$ -ELIMTIL<sub>1</sub>( $\mathcal{X}, \mathcal{X}, \delta$ ) and ALBA( $\mathcal{X}, \delta$ ), and compare them with the state-of-the-art algorithms  $\mathcal{X}\mathcal{Y}$ -Allocation (Figure 2 of (Soare et al., 2014)) and  $\mathcal{X}\mathcal{Y}$ -Adaptive (Figure 3 of (Soare et al., 2014)). We do not include the Explore-Verify algorithm in (Karnin, 2016) as its main contribution is considered to be theoretical.

In the implementation of our algorithms, we use the *entropic mirror descent* method introduced in (Beck & Teboulle, 2003) to compute the optimal solution  $\lambda_{\mathcal{X}}^*$  defined in Proposition 6(a) (see Appendix G for details). In the  $\mathcal{X}\mathcal{Y}$ -Adaptive algorithm, we set  $\alpha = 1/10$  following the choice made in (Soare et al., 2014). All algorithms are implemented in Python 3, and are tested without parallelization.

We test the algorithms using both synthetic (similar to that in (Soare et al., 2014), and random data) and real-world data. We fix the confidence parameter  $\delta = 0.05$ , and report the total number of samples and time used by each algorithm and their empirical failure probabilities (i.e. to fail to identify the best arm). For each setting, the reported numbers are averaged over 100 runs.

### 5.1. Synthetic Data Set 1

In this experiment, we consider a set  $\mathcal{X} \subset \mathbb{R}^d$  of arms, with  $|\mathcal{X}| = d + 1$ .  $\mathcal{X}$  includes the canonical basis  $\{e_1, \dots, e_d\}$  of  $\mathbb{R}^d$ , and an additional arm  $x_{d+1} = [\cos(\omega), \sin(\omega), 0, \dots, 0]^T$ . We choose  $\theta = [2, 0, \dots, 0]^T$ , and fix  $\omega = 0.1$ , so that  $x_1$  is the best arm and  $\Delta_2 = (x_1 - x_{d+1})^T \theta$  is the minimum reward gap. Also, we set the noise  $\epsilon \sim \mathcal{N}(0, 1)$  independently for each pull. We test the algorithms for  $d = 2, \dots, 10$ , and report the number of samples used in Figure 1(a). In Table 1, we summarize the runtime needed by the algorithms, from which we can see the improvement made by our algorithms. The empirical failure probability for each of the algorithms is 0.

### 5.2. Synthetic Data Set 2: Random Vectors

In this data set, the feature vectors of the  $|\mathcal{X}| = 100$  arms are independently uniformly random sampled from  $\mathbb{S}^{d-1}$ , the unit sphere centered at the origin. We pick the two closest arms  $x, y \in \mathcal{X}$  and set  $\theta = x + \alpha \cdot (y - x)$  with

<sup>3</sup>Changing this step to  $S_{r+1} \leftarrow \text{ELIMTIL}_{\lfloor d/2^r \rfloor}(S_r, \delta_r)$  does not make a difference in the theoretical guarantee of the algorithm (i.e., Theorem 22). However,  $\mathcal{Y}$ -ELIMTIL leads to better empirical performance.

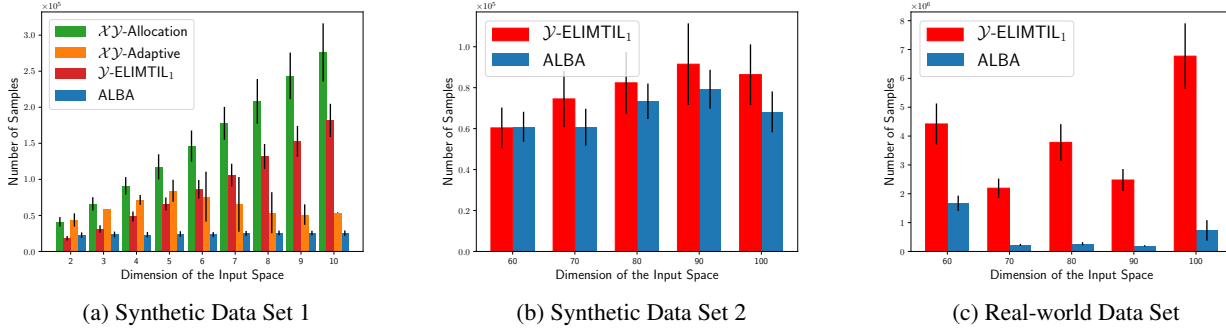


Figure 1: Number of samples needed to identify the best arm for each data set. Half of each vertical black line represents the standard deviation. Notice that the sample complexity may not be monotonically increasing on  $d$ . For (c), the  $\Delta_2$  corresponding to  $d = 60$  is much smaller than the other cases (see Table 3), and so it may need more samples.

$\alpha = 0.01$  so that  $x$  becomes the best arm. We also set the noise  $\epsilon \sim \mathcal{N}(0, 1)$  independently for each pull. We report the samples used by our algorithms in Figure 1(b) and more information (e.g. runtime and  $\Delta_2$  of each data point) in Table 2. The empirical failure probabilities for both algorithms are 0. No results are reported for  $\mathcal{XY}$ -Allocation and  $\mathcal{XY}$ -Adaptive as both failed to terminate within one hour.

### 5.3. Real-world Data Set

We now work with candidate arms generated from the advertisement placement dataset provided in (Lefortier et al., 2016). Given dimensional parameter  $d$ , we select  $d$  features with the highest variance and then randomly pick  $|\mathcal{X}| = 200$  identical vectors from the data set and normalize every vector to unit vectors. We also pick the two closest arms  $x, y \in \mathcal{X}$  and set  $\theta = x + \alpha \cdot (y - x)$  with  $\alpha = 0.01$  so that  $x$  becomes the best arm. We also set the noise  $\epsilon \sim \mathcal{N}(0, 1)$  independently for each pull. We now report the samples used by our algorithms in Figure 1(c) and more information (e.g. runtime and  $\Delta_2$  of each data point) in Table 3. The empirical failure probabilities for both algorithms are 0. Again, no results are shown for  $\mathcal{XY}$ -Allocation and  $\mathcal{XY}$ -Adaptive as they failed to terminate within one hour.

## 6. Conclusion

Via exploiting the global linear structure of the problem, we have shown that the sample complexity of identifying the best arm in linear bandit only depends on the reward gaps of the top  $d$  arms (where  $d$  is the dimension and up to a logarithmic factor). The experimental results also demonstrate the substantial improvement made by our algorithms in terms of both sample complexity and computational time.

However, it remains open to design algorithms with *instance-wise optimal* sample complexity. In (Soare et al., 2014), the authors proposed  $\rho^* = \rho^*(Y^*) = \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{y \in Y^*} \frac{y^T E_{z \sim \lambda} [zz^T]^{-1} y}{(y^T \theta)^2}$  (where  $Y^* = \{x -$

$\mathcal{X}_{[1]} \mid x \neq \mathcal{X}_{[1]}\}$ ) as an information-theoretic lower bound on every input instance. It is a very interesting question to explore whether this lower bound is achievable by an algorithm, or some stronger lower bound exists.

Table 1: The average runtime (in seconds) needed for Synthetic Data Set 1.

$d$	$\mathcal{XY}$ -ALLOC.	$\mathcal{XY}$ -ADAPT.	$\mathcal{Y}$ -ELIMTIL <sub>1</sub>	ALBA
2	11.3	6.2	0.2	0.3
3	35.7	12.2	0.2	0.3
4	56.3	19.6	0.2	0.3
5	122.8	26.4	0.2	0.3
6	229.0	22.7	0.2	0.3
7	699.0	18.9	0.2	0.3
8	1015.3	27.5	0.2	0.3
9	1391.7	37.2	0.2	0.3
10	1869.0	44.4	0.2	0.3

Table 2: Experimental results for Synthetic Data Set 2.

$d$	$\Delta_2$	#samples used		Runtime (secs)	
		$\frac{d \cdot \Delta_2^{-2}}{d \cdot \Delta_2^{-2}}$			
		$\mathcal{Y}$ -ELIMTIL <sub>1</sub>	ALBA	$\mathcal{Y}$ -ELIMTIL <sub>1</sub>	ALBA
60	0.52	268.56	270.99	5.2	5.6
70	0.55	318.94	259.73	3.8	4.6
80	0.57	340.17	302.95	4.6	5.3
90	0.59	353.94	306.51	4.1	4.6
100	0.67	383.91	302.99	2.4	3.7

Table 3: Experimental results for the Real-world Data Set.

$d$	$\Delta_2$	#samples used		Runtime (secs)	
		$\frac{d \cdot \Delta_2^{-2}}{d \cdot \Delta_2^{-2}}$			
	$(\times 10^{-3})$	$\mathcal{Y}$ -ELIMTIL <sub>1</sub>	ALBA	$\mathcal{Y}$ -ELIMTIL <sub>1</sub>	ALBA
60	3.386	0.84	0.32	18.6	22.7
70	9.850	3.04	0.32	43.8	46.8
80	7.567	2.71	0.19	20.0	22.9
90	9.655	2.56	0.20	39.4	42.1
100	9.850	6.57	0.71	20.0	21.0



## References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Allen-Zhu, Z., Li, Y., Singh, A., and Wang, Y. Near-optimal discrete optimization for experimental design: A regret minimization approach. *arXiv preprint arXiv:1711.05174*, 2017.
- Auer, P. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Beck, A. and Teboulle, M. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- Bian, A. A., Buhmann, J. M., Krause, A., and Tschitschek, S. Guarantees for greedy maximization of non-submodular functions with applications. In *Proceedings of International Conference of Machine Learning (ICML)*, 2017.
- Bubeck, S., Wang, T., and Viswanathan, N. Multiple identifications in multi-armed bandits. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 258–265, 2013.
- Carpentier, A. and Locatelli, A. Tight (lower) bounds for the fixed budget best arm identification bandit problem. In *Conference on Learning Theory*, pp. 590–604, 2016.
- Chamon, L. and Ribeiro, A. Approximate supermodularity bounds for experimental design. In *Advances in Neural Information Processing Systems*, pp. 5407–5416, 2017a.
- Chamon, L. F. and Ribeiro, A. Greedy sampling of graph signals. *arXiv preprint arXiv:1704.01223*, 2017b.
- Chen, J., Chen, X., Zhang, Q., and Zhou, Y. Adaptive multiple-arm identification. In *Proceedings of International Conference on Machine Learning (ICML)*, 2017a.
- Chen, L., Li, J., and Qiao, M. Towards instance optimal bounds for best arm identification. In *Proceedings of the Conference on Learning Theory (COLT)*, 2017b.
- Chen, S., Lin, T., King, I., Lyu, M. R., and Chen, W. Combinatorial pure exploration of multi-armed bandits. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, 2014.
- Chu, W., Li, L., Reyzin, L., and Schapire, R. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 208–214, 2011.
- Civril, A. and Magdon-Ismail, M. On selecting a maximum volume sub-matrix of a matrix and related problems. *Theoretical Computer Science*, 410(47-49):4801–4811, 2009.
- Dani, V., Hayes, T. P., and Kakade, S. M. Stochastic linear optimization under bandit feedback. In *Proceedings of the Conference on Learning Theory (COLT)*, 2008.
- Even-Dar, E., Mannor, S., and Mansour, Y. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7:1079–1105, 2006.
- Filippi, S., Cappe, O., Garivier, A., and Szepesvári, C. Parametric bandits: The generalized linear case. In *Advances in Neural Information Processing Systems*, pp. 586–594, 2010.
- Jamieson, K., Malloy, M., Nowak, R., and Bubeck, S. lil’ UCB : An optimal exploration algorithm for multi-armed bandits. In *Proceedings of the Conference on Learning Theory (COLT)*, 2014.
- Kalyanakrishnan, S. and Stone, P. Efficient selection of multiple bandit arms: Theory and practice. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pp. 511–518, 2010.
- Karnin, Z., Koren, T., and Somekh, O. Almost optimal exploration in multi-armed bandits. In *Proceedings of International Conference on Machine Learning (ICML)*, 2013.
- Karnin, Z. S. Verification based solution for structured MAB problems. In *Advances in Neural Information Processing Systems*, pp. 145–153, 2016.
- Kaufmann, E., Cappé, O., and Garivier, A. On the complexity of best arm identification in multi-armed bandit models. *Journal of Machine Learning Research*, 17:1–42, 2016.
- Kiefer, J. and Wolfowitz, J. The equivalence of two extremum problems. *Canadian Journal of Mathematics*, 12(5):363–365, 1960.
- Lefortier, D., Swaminathan, A., Gu, X., Joachims, T., and de Rijke, M. Large-scale validation of counterfactual learning methods: A test-bed. In *NIPS Workshop on Inference and Learning of Hypothetical and Counterfactual Interventions in Complex Systems*, 2016.
- Li, L., Chu, W., Langford, J., and Schapire, R. E. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th International Conference on World Wide Web (WWW)*, pp. 661–670, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-799-8. doi: 10.1145/1772690.1772758.

- Li, L., Lu, Y., and Zhou, D. Provable optimal algorithms for generalized linear contextual bandits. In *Proceedings of International Conference of Machine Learning (ICML)*, 2017.
- Pukelsheim, F. *Optimal design of experiments*. SIAM, 2006.
- Rusmevichientong, P. and Tsitsiklis, J. N. Linearly parameterized bandits. *Mathematics of Operations Research*, 35 (2):395–411, 2010.
- Soare, M., Lazaric, A., and Munos, R. Best-arm identification in linear bandits. In *Advances in Neural Information Processing Systems*, pp. 828–836, 2014.
- Stroock, D. W. *Probability theory: an analytic view*. Cambridge University Press, Cambridge, second edition, 2011. ISBN 978-0-521-13250-3.
- Xu, L., Honda, J., and Sugiyama, M. A fully adaptive algorithm for pure exploration in linear bandits. In Storkey, A. and Perez-Cruz, F. (eds.), *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 843–851, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR.
- Zhou, Y., Chen, X., and Li, J. Optimal PAC multiple arm identification with applications to crowdsourcing. In *Proceedings of International Conference on Machine Learning (ICML)*, pp. 217–225, 2014.