

**COLUMBIA UNIVERSITY**  
IN THE CITY OF NEW YORK

DEPARTMENT OF INDUSTRIAL ENGINEERING AND OPERATIONS RESEARCH

June 21st, 2020

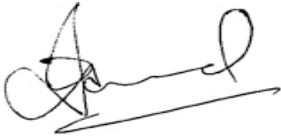
Dear APS Student Paper Competition Committee,

I am writing to certify the eligibility for the work of Min-hwan Oh submitted to the APS Student Paper Competition.

This work is titled “Sparsity-Agnostic Lasso Bandit” and is co-authored by Min-hwan Oh, Assaf Zeevi and myself. In this work, we consider a fundamental problem of sequential decision making with high-dimensional contexts.

I, along with Min-hwan, attest that all of the eligibility conditions for this competition are satisfied. In particular, I confirm that Min-hwan’s contribution comprises the majority of the paper. I look forward to hearing back from you. Thank you very much for your consideration.

Sincerely,



Garud N. Iyengar

Tang Family Professor  
Department of Industrial Engineering and Operations Research  
Data Science Institute  
Columbia University



Min-hwan Oh

Ph.D. Candidate  
IEOR Department  
Columbia University

**Title:** Sparsity-Agnostic Lasso Bandit

**Abstract:**

We consider a stochastic contextual bandit problem where the dimension  $d$  of the feature vectors is potentially large, however, only a sparse subset of features of cardinality  $s_0 \ll d$  affect the reward function. Essentially all existing algorithms for sparse bandits require a priori knowledge of the value of the sparsity index  $s_0$ . This knowledge is almost never available in practice, and misspecification of this parameter can lead to severe deterioration in the performance of existing methods. The main contribution of this paper is to propose an algorithm that does not require prior knowledge of the sparsity parameter  $s_0$  and establish tight regret bounds under mild conditions. We comprehensively evaluate our proposed algorithm numerically and show that it consistently outperforms existing methods, even when the correct sparsity parameter is revealed to these methods but is kept hidden from our algorithm.

**Key words:** Contextual Bandit, High-dimensional Statistics, Lasso

# Sparsity-Agnostic Lasso Bandit

Garud Iyengar

Columbia University, New York, NY 10027, garud@ieor.columbia.edu

Min-hwan Oh

Columbia University, New York, NY 10027, m.oh@columbia.edu

Assaf Zeevi

Columbia University, New York, NY 10027, assaf@gsb.columbia.edu

We consider a stochastic contextual bandit problem where the dimension  $d$  of the feature vectors is potentially large, however, only a sparse subset of features of cardinality  $s_0 \ll d$  affect the reward function. Essentially all existing algorithms for sparse bandits require a priori knowledge of the value of the sparsity index  $s_0$ . This knowledge is almost never available in practice, and misspecification of this parameter can lead to severe deterioration in the performance of existing methods. The main contribution of this paper is to propose an algorithm that does *not* require prior knowledge of the sparsity parameter  $s_0$  and establish tight regret bounds under mild conditions. We comprehensively evaluate our proposed algorithm numerically and show that it consistently outperforms existing methods, even when the correct sparsity parameter is revealed to these methods, but is kept hidden from our algorithm.

*Key words:* Contextual Bandit, High-dimensional Statistics, Lasso

---

## 1. Introduction

In classical multi-armed bandits (MAB), one of the arms is pulled in each round and a reward corresponding to the chosen arm is revealed to the decision-making agent. The rewards are, typically, independent and identically distributed samples from an arm-specific distribution. The goal of the agent is to devise a strategy for pulling arms that maximizes cumulative rewards, suitably balancing between exploration and exploitation. Linear contextual bandits ([Abe and Long 1999](#), [Auer 2002](#), [Chu et al. 2011](#)) and generalized linear contextual bandits ([Filippi et al. 2010](#), [Li et al. 2017](#)) are more recent important extensions of the basic MAB setting, where each arm  $a$  is associated with a known feature vector  $x_a \in \mathbb{R}^d$ , and the expected payoff of the arm is a (typically, monotone increasing) function of the inner product  $x_a^\top \beta^*$  with unknown parameter vector  $\beta^* \in \mathbb{R}^d$ . In these settings, pulling any arm provides some information about the unknown parameter vector, and hence, insight into the average reward of the other arms. These bandit algorithms are applicable in a variety of problem settings, such as online advertising, recommender systems, online retail and healthcare analytics, where the contextual information can be used for personalization.

In most application domains highlighted above, the feature space is high-dimensional ( $d \gg 1$ ), yet typically only a small subset of the features may influence the expected reward. That is, the

unknown parameter vector is *sparse* with only the elements corresponding to the relevant features being non-zero, i.e.,  $s_0 = \|\theta^*\|_0 \ll d$ . There is an emerging body of literature on contextual bandit problems with sparse linear reward functions (Abbasi-Yadkori et al. 2012, Gilton and Willett 2017, Bastani and Bayati 2020, Wang et al. 2018, Kim and Paik 2019) which propose methods to exploit the sparse structure under various conditions. However, there is a crucial drawback in almost all of these approaches: the algorithms require *prior* knowledge of the sparsity index  $s_0$ , information that is almost never available in practice. In the absence of such knowledge, the existing algorithms fail to fully leverage the sparse structure, and their performance does not guarantee the improvements in dimensionality-dependence which can be realized in the sparse problem setting. The purpose of this paper is to demonstrate that a relatively simple contextual bandit algorithm, that exploits  $\ell_1$ -regularized regression using Lasso (Tibshirani 1996) in a sparsity-agnostic manner, is provably near-optimal insofar as its regret performance (under suitable regularity). Our contributions are as follows:

- (a) We propose the first sparse bandit algorithm for a general set of arms<sup>1</sup> that does not require prior knowledge of the sparsity index  $s_0$ .
- (b) We first establish that the regret bound of our proposed algorithm is  $\mathcal{O}(s_0 \sqrt{T \log(dT)})$  for the two-armed case; the two-armed problem affords the most accessible exposition of the key analytical ideas. (Extensions to the general  $K$ -armed case are discussed later.) The regret bound scale in  $s_0$  and  $d$  matches the equivalent terms in the *offline* Lasso results (see the discussion in Section 5.1).
- (c) We comprehensively evaluate our algorithm on numerical experiments and show that it consistently outperforms existing methods, even when these methods are granted the prior knowledge of the correct sparsity index (and greatly outperforms them if this information is mis-specified).

Our algorithm does not rely on *forced sampling* which was used by almost all previous work, e.g., Bastani and Bayati (2020), Wang et al. (2018), Kim and Paik (2019), to satisfy certain regularity of the empirical Gram matrix. This requires prior knowledge of  $s_0$  because the forced sampling schemes, the key ideas going back to Goldenshluger and Zeevi (2013), need to be fine-tuned using the *correct* sparsity index. (See further discussions in Section 3.4.)

The rest of the paper is organized as follows. In Section 2, we review the related literature. In Section 3, we present the problem formulation and discuss the reason why the previously proposed methods require to know sparsity  $s_0$ . Section 4 describes our proposed algorithm, and Section 5 contains our main results. In Section 5, we present results characterizing the challenges of sparse bandit problem without requiring the sparsity information, and establish the upper bound on the cumulative regret for the two-armed sparse bandits. Section 6 contains the numerical experiments.

In Section 7, we extend our analysis and numerical evaluations to  $K$ -armed bandits. Section 8 presents discussion and future directions. The complete proofs and additional numerical results are provided in the appendix.

## 2. Related Work

Linear bandits and generalized linear bandits have been widely studied (Abe and Long 1999, Auer 2002, Dani et al. 2008, Rusmevichientong and Tsitsiklis 2010, Abbasi-Yadkori et al. 2011, Filippi et al. 2010, Chu et al. 2011, Agrawal and Goyal 2013, Li et al. 2017). However, these (generalized) linear contextual bandit strategies are not able to exploit sparse structure in the unknown parameter vector and hence may incur regret proportional to the full ambient dimension  $d$  rather than the sparse set of features of cardinality  $s_0$ . To exploit sparse structure, Abbasi-Yadkori et al. (2012) propose a framework to construct high probability confidence sets for online linear prediction and establish the  $\tilde{\mathcal{O}}(\sqrt{s_0 d T})$  regret bound, where  $\tilde{\mathcal{O}}$  hides logarithmic terms. However, their algorithm needs to know the sparsity  $s_0$ . Furthermore, their algorithm is not computationally efficient; an implementable version of their framework is not yet known (Section 23.5 in Lattimore and Szepesvári 2019). It is worth noting that the  $\sqrt{d}$  dependence in the regret bound is unavoidable unless additional assumptions are imposed; see Theorem 24.3 in Lattimore and Szepesvári (2019). Gilton and Willett (2017) adapt Thompson sampling (Thompson 1933) to sparse linear bandits; however, they also assume a priori knowledge of a small superset of the support for the parameter.

Bastani and Bayati (2020) address the contextual bandit problem with high-dimensional features. They propose a bandit algorithm which uses Lasso (Tibshirani 1996) to estimate the parameter of each arm separately. To ensure compatibility of the empirical Gram matrices, they adapt the forced-sampling technique in Goldenshluger and Zeevi (2013) which is now tuned using the (a priori known) sparsity index, and is implemented for each arm at predefined time points. They establish a  $\mathcal{O}(K^4 s_0^2 [\log d + \log T]^2)$  regret bound where  $K$  is the number of arms. Note that they invoke several additional assumptions: a margin condition that ensures that the density of the context distribution is bounded near the decision boundary, and arm-optimality which assumes a gap between the optimal and sub-optimal arms exists with some positive probability. In the same problem setting, Wang et al. (2018) propose an algorithm which uses forced-sampling along with the minimax concave penalty (MCP) estimator (Zhang 2010) and improve the regret bound to  $\mathcal{O}(K^3 s_0^2 [s_0 + \log d] \log T)$ . Note that Bastani and Bayati (2020) and Wang et al. (2018) achieve a poly-logarithmic dependence on  $T$  in regret, exploiting the arm optimality condition which assumes a gap between the optimal and sub-optimal arms exists with some probability. Since we do not assume such separability between arms, poly-logarithmic dependence on  $T$  is not attainable in our problem setting. Kim and Paik (2019) extend the Lasso bandit (Bastani and Bayati 2020) to linear

bandit settings and propose a different approach to address the non-compatibility of the empirical Gram matrices by using a doubly-robust technique (Bang and Robins 2005) that originates with the missing data / imputation literature. They achieve  $\mathcal{O}(s_0\sqrt{T}\log(dT))$  regret.

All of the aforementioned algorithms require that the learning agent know the sparsity  $s_0$  of the unknown parameter (or a non-trivial upper-bound on sparsity which is strictly less than  $d$ ).<sup>2</sup> That is, only when the algorithm knows  $s_0$ , it can guarantee the regret bounds mentioned above. Otherwise, the regret bounds would scale polynomially with  $d$  instead of  $s_0$  or potentially scale linearly with  $T$ . To our knowledge, the only work in sparse bandits which does not require this prior knowledge of sparsity is the work by Carpentier and Munos (2012) although the algorithm still requires to know the  $\ell_2$ -norm of the unknown parameter. However, their analysis uses a non-standard definition of noise and is restricted to the case where the set of arms is the  $\ell_2$  unit ball and fixed over time, a structure they exploit in a significant manner, and which limits the scope of their algorithm.

### 3. Preliminaries

#### 3.1. Notation

For a vector  $x \in \mathbb{R}^d$ , we use  $\|x\|_1$  and  $\|x\|_2$  to denote its  $\ell_1$ -norm and  $\ell_2$  norm respectively, the notation  $\|x\|_0$  is reserved for the cardinality of the set of non-zero entries of that vector. The minimum and maximum singular values of a matrix  $V$  are written as  $\lambda_{\min}(V)$  and  $\lambda_{\max}(V)$  respectively. For two symmetric matrices  $V$  and  $W$  of the same dimensions,  $V \succcurlyeq W$  means that  $V - W$  is positive semi-definite. We define  $[n]$  for a positive integer  $n$  to be a set containing positive integers up to  $n$ , i.e.,  $\{1, 2, \dots, n\}$ . For a real-valued function  $f$ , we use  $\dot{f}$  and  $\ddot{f}$  to denote its first and second derivatives.

#### 3.2. Generalized Linear Contextual Bandits

We consider the stochastic generalized linear bandit problem with  $K$  arms. Let  $T$  be the problem horizon, namely the number of rounds to be played. In round  $t \in [T]$ , the learning agent observes a context consisting of a set of  $K$  feature vectors  $\mathcal{X}_t = \{X_{t,i} \in \mathbb{R}^d \mid i \in [K]\}$ , where the tuple  $\mathcal{X}_t$  is drawn i.i.d. from an unknown joint distribution  $p_{\mathcal{X}}$ . Note that the feature vectors for different arms are allowed to be correlated. Each feature vector  $X_{t,i}$  is associated with an unknown stochastic reward  $Y_{t,i} \in \mathbb{R}$ . The agent selects one arm, denoted by  $a_t \in [K]$  and observes the reward  $Y_t := Y_{t,a_t}$  corresponding to the chosen arm's feature  $X_t := X_{t,a_t}$  as a bandit feedback. The policy consists of the sequence of actions  $\pi = \{a_t : t = 1, 2, \dots\}$  and is non-anticipating, namely each action only depends on past observations and actions.

In this work, we consider the generalized linear model (GLM) in which there is an unknown parameter  $\beta^* \in \mathbb{R}^d$  and a fixed increasing function  $\mu : \mathbb{R} \rightarrow \mathbb{R}$  (also known as *inverse link function*)

such that  $\mathbb{E}[Y|X] = \mu(X^\top \beta^*)$ , where  $X$  is the chosen arm's feature and  $Y$  is the corresponding reward. The GLM can be written as

$$Y_t = \mu(X_t^\top \beta^*) + \epsilon_t$$

where  $\{\epsilon_t, t \in [T]\}$  are independent zero-mean noise. Widely used examples are  $\mu(z) = z$ , which covers the linear model, and  $\mu(z) = 1/(1 + e^{-z})$  which is the logistic model. The parameter  $\beta^*$  and the feature vectors  $\{x_{t,i}\}$  are potentially high-dimensional (large  $d$ ) but  $\beta^*$  is *sparse*, that is, the number of non-zero elements is small with  $\|\beta^*\|_0 = s_0$  where  $s_0 \ll d$ . It is important to note that the agent *does not* know  $s_0$  or the support of  $\beta^*$ .

We assume that there is an increasing sequence of sigma fields  $\{\mathcal{F}_\tau\}$  such that  $\epsilon_t$  is  $\mathcal{F}_t$ -measurable with  $\mathbb{E}[\epsilon_t | \mathcal{F}_{t-1}] = 0$ . In our problem,  $\mathcal{F}_t$  is the sigma-field generated by random variables of chosen actions  $\{X_1, \dots, X_t\}$  and their corresponding rewards  $\{Y_1, \dots, Y_t\}$ . We assume the noise  $\epsilon_t$  is sub-Gaussian with parameter  $\sigma$ , where  $\sigma$  is a positive absolute constant, i.e.,  $\mathbb{E}[e^{\zeta \epsilon_t}] \leq e^{\zeta^2 \sigma^2 / 2}$  for all  $\zeta \in \mathbb{R}$ . In practice, when we have bounded reward  $Y_t$ , the noise  $\epsilon_t$  is also bounded and hence satisfies the sub-Gaussian assumption with appropriate  $\sigma$  value.

The agent's goal is to maximize the cumulative expected reward  $\mathbb{E}[\sum_{t=1}^T \mu(X_{t,a_t}^\top \beta^*)]$  over  $T$  rounds. Let  $a_t^* = \arg \max_{i \in [K]} \mu(X_{t,i}^\top \beta^*)$  denote the optimal arm for round  $t$ . Then, the expected cumulative *regret*  $\mathcal{R}_T$  of policy  $\pi$  is defined as

$$\mathcal{R}_T(\pi) := \sum_{t=1}^T \mathbb{E} \left[ \mu(X_{t,a_t^*}^\top \beta^*) - \mu(X_{t,a_t}^\top \beta^*) \right].$$

Hence, maximizing the expected cumulative rewards of policy  $\pi$  over  $T$  rounds is equivalent to minimizing the cumulative regret  $\mathcal{R}_T(\pi)$ .

### 3.3. Lasso for Generalized Linear Models

Suppose we have samples  $Y_1, \dots, Y_n$  and corresponding features  $X_1, \dots, X_n$ . The log-likelihood function of  $\beta$  under the canonical generalized linear model is

$$\log \mathcal{L}_n(\beta) := \sum_{j=1}^n \left[ \frac{Y_j X_j^\top \beta - m(X_j^\top \beta)}{g(\eta)} - h(Y_j, \eta) \right].$$

Here,  $\eta \in \mathbb{R}^+$  is a known scale parameter;  $m(\cdot)$ ,  $g(\cdot)$  and  $h(\cdot)$  are normalization functions where  $m(\cdot)$  is infinitely differentiable satisfying

$$\dot{m}(X^\top \beta^*) = \mathbb{E}[Y|X] = \mu(X^\top \beta^*).$$

Therefore, the Lasso (Tibshirani 1996) estimate for the generalized linear model can be defined as

$$\hat{\beta}_n \in \arg \min_{\beta} \{ \ell_n(\beta) + \lambda \|\beta\|_1 \} \quad (1)$$

where  $\ell_n(\beta) := -\frac{1}{n} \sum_{j=1}^n [Y_j X_j^\top \beta - m(X_j^\top \beta)]$  and  $\lambda$  is a penalty parameter. Lasso is known to be an efficient (offline) tool for estimating the high-dimensional linear regression parameter. The “fast convergence” property of Lasso is guaranteed when data are i.i.d. and when the observed covariates are not highly correlated. The restricted eigenvalue condition (Bickel et al. 2009, Raskutti et al. 2010), the compatibility condition (Van De Geer et al. 2009), and the restricted isometry property (Candes et al. 2007) have been used to ensure that such high correlations are avoided. In sequential learning settings, however, these conditions are often violated because the observations are adapted to the past and the feature variables of the chosen arms converge to a small region of the feature space as the learning agent updates its arm selection policy.

### 3.4. Why do existing sparse bandit algorithms require prior knowledge of the sparsity index?

The primary reason that a priori knowledge of sparsity is assumed throughout most of the literature is, roughly speaking, to ensure suitable “size” of the confidence bounds and concentration. For example, Abbasi-Yadkori et al. (2012) require the parameter  $s_0$  to explicitly construct a high probability confidence set with its radius proportional to  $s_0$  rather than  $d$ . The recently proposed bandit algorithms of Bastani and Bayati (2020), Kim and Paik (2019) and the variant with MCP estimator in Wang et al. (2018) employ a logic that is similar in spirit (though different in execution). Specifically, the compatibility condition or restricted eigenvalue condition is assumed to hold only for the theoretical Gram matrix, and the empirical Gram matrix may not satisfy such condition (the difficulty in controlling that is due to the non-i.i.d. adapted samples of the feature variables). As a remedy to this issue, Bastani and Bayati (2020) and Wang et al. (2018) utilize the forced-sampling technique of Goldenshluger and Zeevi (2013) to obtain a “sufficient” number of i.i.d. samples and use that to show that the empirical Gram matrices concentrate in the vicinity of the theoretical Gram matrix, and hence, satisfy the compatibility condition after a sufficient amount of forced-sampling. The forced-sampling duration needs to be predefined and scales at least polynomially in the sparsity  $s_0$  to ensure concentration of the Gram matrices. That is, if the algorithm does not know  $s_0$ , the forced-sampling duration will have to scale polynomially in  $d$ . Kim and Paik (2019) propose an alternative to forced sampling that builds on doubly-robust techniques used in the missing data literature; however, their algorithm involves random arm selection with a probability that is calibrated using  $s_0$ , and initial uniform sampling whose duration requires knowledge of  $s_0$  and scales polynomially with  $s_0$  in order to establish their regret bounds. The sensitivity to the sparsity index specification is also evident in cases where its value is *misspecified* which may result in severe deterioration in the performance of the algorithm (see further discussion in Section 5.1).



The key observation in our analysis is that, under some mild conditions, i.i.d. samples, which are the key output of the forced sampling scheme, are in fact not essential. We show that the empirical Gram matrix satisfies the required regularity after a sufficient number of rounds, provided the theoretical Gram matrix is also regular; the details of this analysis are in Section 5. Numerical experiments support this findings, and moreover, demonstrate that the performance of the algorithm can be superior to forced-sampling-based schemes that are tuned with foreknowledge of the parameter  $s_0$ .

## 4. Algorithm

Our proposed SPARSITY-AGNOSTIC (SA) LASSO BANDIT algorithm for high-dimensional GLM bandits is displayed in Algorithm 1. As the name suggests, our algorithm does not require prior knowledge of the sparsity  $s_0$ . It relies on Lasso for parameter estimation, and does not explicitly use exploration strategies or forced-sampling. Instead, in each round we choose an arm which maximizes the inner product of a feature vector and the Lasso estimate. After observing the reward, we update the regularization parameter  $\lambda_t$  and update the Lasso estimate  $\hat{\beta}_t$  which minimizes the penalized negative log-likelihood function defined in (1).

SA LASSO BANDIT requires only one input parameter: the subgaussian parameter  $\sigma$ , which is commonly required for almost all parametric bandit algorithms. (Note that, in comparison, Kim and Paik (2019) require three tuning parameters and Bastani and Bayati (2020) and Wang et al. (2018) require four tuning parameters, including the unknown sparsity parameter  $s_0$ ). This is a great advantage in terms of practicality since tuning parameters, while achieving low regret, are challenging to specify in online learning settings.

---

### Algorithm 1 SA LASSO BANDIT

---

- 1: **Input parameter:**  $\sigma$
  - 2: **for** all  $t = 1$  to  $T$  **do**
  - 3:   Observe  $X_{t,i}$  for all  $i \in [K]$
  - 4:   Compute  $a_t = \arg \max_{i \in [K]} X_{t,i}^\top \hat{\beta}_t$
  - 5:   Pull arm  $a_t$  and observe  $Y_t$
  - 6:   Update  $\lambda_t \leftarrow 2\sigma \sqrt{\frac{4 \log t + 2 \log d}{t}}$
  - 7:   Update  $\hat{\beta}_{t+1} \leftarrow \arg \min_{\beta} \{\ell_t(\beta) + \lambda_t \|\beta\|_1\}$
  - 8: **end for**
- 

Algorithm 1 may appear to be an *exploration-free* greedy algorithm (e.g., Bastani et al. 2017). However, this is not the case. Algorithm 1 does *not* take the *greedy* action corresponding to the

*greedy* maximum likelihood estimate based only on a few samples. Exploration in algorithms can manifest itself in many ways. For example, in upper-confidence bound (UCB) algorithms, the estimate is the parameter value in a high-probability confidence ellipsoid around the *greedy* maximum likelihood estimate that maximizes the reward. Once the UCB estimate is chosen, the action selection is greedy with respect to the parameter estimate — similarly in Thompson sampling, action selection is also greedy with respect to a sampled parameter. The UCB algorithms carefully control the size of the ellipsoid to ensure convergence. Thus, exploration is loosely equivalent to regularizing the MLE to reduce the impact of the variance in the initial measurements. The algorithm we propose computes the parameter estimate by regularizing the MLE with a sparsifying norm, and then, as in UCB, takes a greedy action with respect to this regularized parameter estimate. We adjust the penalty associated with the sparsifying norm over time at a suitable rate, and hence, ensure that our estimate is consistent as we collect more samples. (This adjustment and specification do not require knowledge of  $s_0$ .) In fact, an inadequate choice of this penalty parameter would lead to large regret, which is analogous to poor choices of confidence widths in UCB.

## 5. Regret Analysis

In this section, we establish an upper bound on the expected regret of SA LASSO BANDIT for the two-armed generalized linear bandits. We focus on the two-arm case primarily for clarity and accessibility of key analysis ideas. We later extend our analysis to the  $K$ -armed case with  $K \geq 3$  in Section 7. It is important to note that our proposed algorithm does not change with the number of arms. Now, we first provide a few definitions and assumptions used throughout the analysis. We start with an assumption standard in the (generalized) linear bandit literature.

**ASSUMPTION 1 (Feature set and link function).** *For all features  $X \in \mathcal{X}$ ,  $\|X\|_2 \leq 1$ . For the link function,  $\kappa_{\min} \leq \dot{\mu}(x^\top \beta) \leq \kappa_{\max}$  for all  $x$  and  $\beta$ .*

Clearly for the linear link function,  $\kappa_{\min} = \kappa_{\max} = 1$ . For the logistic link function,  $\kappa_{\max} = 1/4$ .

**DEFINITION 1 (ACTIVE SET AND SPARSITY INDEX).** Let  $S_0 := \{j : \beta_j^* \neq 0\}$  denote the set of indices for which  $\beta_j^*$ 's are non-zero, referred to as the active set, and put  $s_0 = |S_0|$  to be the sparsity index.

For the active set  $S_0$ , and an arbitrary vector  $\beta \in \mathbb{R}^d$ , we can define

$$\beta_{j,S_0} := \beta_j \mathbb{1}\{j \in S_0\}, \quad \beta_{j,S_0^c} := \beta_j \mathbb{1}\{j \notin S_0\}.$$

Thus,  $\beta_{S_0} = [\beta_{1,S_0}, \dots, \beta_{d,S_0}]^\top$  has zero elements outside the set  $S_0$  and the elements of  $\beta_{S_0^c}$  can only be non-zero in the complement of  $S_0$ . We define the set of vectors

$$\mathbb{C}(S_0) := \{\beta \in \mathbb{R}^d \mid \|\beta_{S_0^c}\|_1 \leq 3\|\beta_{S_0}\|_1\}. \quad (2)$$

Let  $\mathbf{X} \in \mathbb{R}^{K \times d}$  be the design matrix where each row is a feature vector for each arm. (Although we focus on  $K = 2$  case in this section, the terms and the assumptions defined here also apply to the case of  $K \geq 3$ .) Then, similar to the previous literature in sparse estimation and specifically in sparse bandits (Bastani and Bayati 2020, Wang et al. 2018, Kim and Paik 2019), we assume that the following compatibility condition is satisfied for the theoretical Gram matrix  $\Sigma := \frac{1}{K} \mathbb{E}[\mathbf{X}^\top \mathbf{X}]$ .

**ASSUMPTION 2 (Compatibility condition).** *For active set  $S_0$ , there exists compatibility constant  $\phi_0^2 > 0$  such that*

$$\phi_0^2 \|\beta_{S_0}\|_1^2 \leq s_0 \beta^\top \Sigma \beta \quad \text{for all } \beta \in \mathbb{C}(S_0).$$

We add to this the following mild assumption that is more specific to our analysis.

**ASSUMPTION 3 (Relaxed symmetry).** *For a joint distribution  $p_{\mathcal{X}}$ , there exists  $\rho_0 < \infty$  such that  $\frac{p_{\mathcal{X}}(-\mathbf{x})}{p_{\mathcal{X}}(\mathbf{x})} \leq \rho_0$  for all  $\mathbf{x}$ .*

### Discussion of the assumptions.

Assumption 1 is the standard regularity assumption used in GLM bandit literature (Filippi et al. 2010, Li et al. 2017, Kveton et al. 2019). It is important to note that unlike the existing GLM bandit algorithms which need to know  $\kappa_{\min}$  a priori, our proposed algorithm does not require either knowledge of  $\kappa_{\min}$  or  $\kappa_{\max}$ . The compatibility condition in Assumption 2 is analogous to the standard positive-definite assumption on the Gram matrix for the ordinary least squares estimator for linear models but is less restrictive. The compatibility condition allows collinearity in the covariates matrix and implies that truly active parameters are not too correlated. As mentioned above, the compatibility condition is a standard assumption in the sparse bandits literature (Bastani and Bayati 2020, Wang et al. 2018, Kim and Paik 2019). Assumption 3 states that the joint distribution  $p_{\mathcal{X}}$  can be skewed but this skewness is bounded. Obviously, if  $p_{\mathcal{X}}$  is symmetrical, we have  $\rho_0 = 1$ . Assumption 3 is satisfied for a large class of continuous and discrete distributions, e.g., elliptical distributions including Gaussian, truncated Gaussian, uniform distribution, and Rademacher distribution.

### 5.1. Regret Bound for SA Lasso Bandit

**THEOREM 1 (Regret bound for two arms).** *Suppose  $K = 2$  and Assumptions 1-3 hold. Then the expected cumulative regret of the SA LASSO BANDIT policy  $\pi$  over horizon  $T$  is upper-bounded by*

$$\mathcal{R}_T(\pi) \leq 4\kappa_{\max} + \frac{2 \log(2d^2) + 2}{C_0(\phi_0, s_0)^2} + \frac{32\kappa_{\max}\rho_1\sigma s_0 \sqrt{T \log(dT)}}{\kappa_{\min}\phi_0^2}$$

where  $C_0(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_1}\right)$ .

### Discussion of Theorem 1.

In terms of key problem primitives, Theorem 1 establishes  $\mathcal{O}(s_0\sqrt{T\log(dT)})$  regret without any prior knowledge on  $s_0$ . The bound shows that the regret of our algorithm grows at most logarithmically in feature dimension  $d$ . The key takeaway from this theorem is that SA LASSO BANDIT is sparsity-agnostic and is able to achieve *correct* dependence on parameters  $d$  and  $s_0$ . Based on the offline Lasso convergence results under the compatibility condition (e.g., Theorem 6.1 in Bühlmann and Van De Geer 2011), we believe that the dependence on  $d$  and  $s_0$  in Theorem 1 is best possible.<sup>3</sup>

The regret bound in Theorem 1 is tighter than the previously known bound in the same problem setting (Kim and Paik 2019) although direct comparison is not immediate, given the difference in assumptions involved — compared to Kim and Paik (2019), we require Assumption 3 whereas they assume the sparsity index  $s_0$  is known. Having said that, the numerical experiments in Section 6 support our theoretical claims and provide additional evidence that our proposed algorithm compares very favorably to other existing methods (which are tuned with the knowledge of the correct  $s_0$ ), and moreover, the performance is not sensitive to the assumptions that were imposed primarily for technical tractability purposes. As mentioned earlier, the previous work on sparse bandits (Bastani and Bayati 2020, Wang et al. 2018, Kim and Paik 2019) require the knowledge of sparsity. In the absence of such knowledge, if sparsity is underspecified, then these algorithms would suffer a regret linear in  $T$ . On the other hand, if the sparsity is overspecified, the regret of these algorithms scales with  $d$  instead of  $s_0$ . Our proposed algorithm does not require such prior knowledge, hence there is no risk of under or over specification, and yet our analysis provides a sharper regret guarantee. Furthermore, our result also suggests that even when the sparsity is known, random sampling to satisfy the compatibility condition, invoked by all existing sparse bandit algorithms to date, can be wasteful since said conditions may be already satisfied even in the absence of such sampling. This finding is also supported by the numerical experiments in Section 6 and Section 7.2. We provide the outline of the proof and the key lemmas in the following section.

### 5.2. Challenges and Proof Outlines

There are two essential challenges that prevent us from fully benefiting from the fast convergence property of Lasso:

- (i) The samples are not i.i.d., therefore the standard Lasso oracle inequality does not hold.
- (ii) Empirical Gram matrices do not necessarily satisfy the compatibility condition even under Assumption 2. This is because the selected feature variables for which the rewards are observed do not provide an “even” representation for the entire distribution.

To resolve (i), we provide a Lasso oracle inequality for GLM with non-i.i.d. adapted samples under the compatibility condition in Lemma 1. For (ii), we aim to provide a remedy without using the

knowledge of sparsity or without using i.i.d. samples. Hence, this poses a greater challenge. In Section 5.2.2, we address this issue by showing that the empirical Gram matrix behaves “nicely” even when we choose arms adaptively without deliberate random sampling. In particular, we show that adapted Gram matrices can be controlled by the theoretical Gram matrix and the empirical Gram matrix concentrates properly around the adapted Gram matrix as we collect more samples. Connecting this matrix concentration to the corresponding compatibility constants, we show that empirical Gram matrix satisfies the compatibility condition with high probability.

**5.2.1. Lasso Oracle Inequality for GLM with Non-i.i.d. Data.** We present an oracle inequality for the Lasso estimator for GLM under non-i.i.d. data. This is a generalization of the standard Lasso oracle inequality (Bühlmann and Van De Geer 2011) that allows adapted sequences of observations. This result may be of independent interest.

**LEMMA 1 (Oracle inequality).** *Suppose the compatibility condition holds for the empirical covariance matrix  $\hat{\Sigma}_t = \frac{1}{t} \sum_{\tau=1}^t X_\tau X_\tau^\top$  with active set  $S_0$  and compatibility constant  $\phi_t$ . For some  $\delta \in (0, 1)$ , define the regularization parameter*

$$\lambda_t := 2\sigma \sqrt{\frac{2[\log(2/\delta) + \log d]}{t}}.$$

*Then with probability at least  $1 - \delta$ , the Lasso estimate  $\hat{\beta}_t$  defined in (1) satisfies*

$$\|\hat{\beta}_t - \beta^*\|_1 \leq \frac{4s_0\lambda_t}{\kappa_{\min}\phi_t^2}.$$

Note that here we assume that the compatibility condition holds for the empirical Gram matrix  $\hat{\Sigma}_t$ . In the next section, we show that this holds with high probability. The Lasso oracle inequality holds without further assumptions on the underlying parameter  $\beta^*$  or its support. Therefore, if we show that  $\hat{\Sigma}_t$  satisfies the compatibility condition absent knowledge of  $s_0$ , then the remainder of the result does not require this knowledge as well.

**5.2.2. Compatibility Condition and Matrix Concentration.** We first define the generic compatibility constant for matrix  $M$ .

**DEFINITION 2.** The compatibility constant of  $M$  over  $S_0$  is

$$\phi^2(M, S_0) := \min_{\beta} \left\{ \frac{s_0 \beta^\top M \beta}{\|\beta_{S_0}\|_1^2} : \|\beta_{S_0^c}\|_1 \leq 3\|\beta_{S_0}\|_1 \neq 0 \right\}.$$

Hence, it suffices to show  $\phi^2(M, S_0) > 0$  in order to show that matrix  $M$  satisfies the compatibility condition. Note that this definition holds for any index set. In this section, however, we will focus on the active index set  $S_0$  of the parameter  $\beta^*$ . Also, note that the constant 3 in the inequality is for ease of exposition and may be replaced by a different value, but then one has to adjust the choice of the regularization parameter accordingly. Now, under Assumption 2, the theoretical Gram matrix  $\Sigma = \frac{1}{K} \mathbb{E}[\mathbf{X}^\top \mathbf{X}]$  satisfies the compatibility condition i.e.,  $\phi_0^2 = \phi^2(\Sigma, S_0) > 0$ .

DEFINITION 3. We define the *adapted* Gram matrix as  $\Sigma_t := \frac{1}{t} \sum_{\tau=1}^t \mathbb{E}[X_\tau X_\tau^\top | \mathcal{F}_\tau]$  and the *empirical* Gram matrix as  $\hat{\Sigma}_t := \sum_{\tau=1}^t X_\tau X_\tau^\top$ .

For each  $\mathbb{E}[X_\tau X_\tau^\top | \mathcal{F}_\tau]$  in  $\Sigma_t$ , the history  $\mathcal{F}_\tau$  affects how the feature vector  $X_\tau$  is chosen. More specifically, our algorithm uses  $\mathcal{F}_\tau$  to compute  $\hat{\beta}_\tau$  and then chooses arm  $a_\tau$  such that its feature  $x_{a_\tau}$  maximizes  $x_{a_\tau}^\top \hat{\beta}_\tau$ . Therefore, we can rewrite  $\Sigma_t$  as

$$\Sigma_t = \frac{1}{t} \sum_{\tau=1}^t \sum_{i=1}^2 \mathbb{E}_{\mathcal{X}_t} [X_{\tau,i} X_{\tau,i}^\top \mathbb{1}\{X_{\tau,i} = \arg \max_{X \in \mathcal{X}_\tau} X^\top \hat{\beta}_\tau\} | \hat{\beta}_\tau].$$

Since the compatibility condition is satisfied only for the theoretical Gram matrix  $\Sigma$  and we need to show the empirical Gram matrix  $\hat{\Sigma}_t$  satisfies the compatibility condition, the adapted Gram matrix  $\Sigma_t$  serves as a bridge between  $\Sigma$  and  $\hat{\Sigma}_t$  in our analysis. We first lower-bound the compatibility constant  $\phi^2(\Sigma_t, S_0)$  in terms of  $\phi^2(\Sigma, S_0)$  so that we can show that  $\Sigma_t$  satisfies the compatibility condition as long as  $\Sigma$  satisfies the compatibility condition. Then, we show that  $\hat{\Sigma}_t$  concentrates around  $\Sigma_t$  with high probability and that such matrix concentration guarantees the compatibility condition of  $\hat{\Sigma}_t$ .

In Lemma 2, we show that  $\Sigma_t$  can be controlled in terms of the theoretical Gram matrix  $\Sigma$ , which allows us to link the compatibility constant of  $\Sigma$  to compatibility constant of  $\Sigma_t$ . Note that Lemma 2 shows the result for any fixed vector  $\beta$ ; hence can be applied to  $\mathbb{E}[X_\tau X_\tau^\top | \mathcal{F}_\tau]$ .

LEMMA 2. For a fixed vector  $\beta \in \mathbb{R}^d$ , we have

$$\sum_{i=1}^2 \mathbb{E} \left[ X_{t,i} X_{t,i}^\top \mathbb{1}\{X_{t,i} = \arg \max_{X \in \mathcal{X}_t} X^\top \beta\} \right] \succcurlyeq \frac{\Sigma}{\rho_1}.$$

Therefore, we have  $\Sigma_t \succcurlyeq \frac{\Sigma}{\rho_1}$  which implies that  $\phi^2(\Sigma_t, S_0) \geq \frac{\phi^2(\Sigma, S_0)}{\rho_1} > 0$ , i.e.,  $\Sigma_t$  satisfies the compatibility condition. Note that both  $\Sigma$  and  $\Sigma_t$  can be singular. In Lemma 3, we show that  $\hat{\Sigma}_t$  concentrates to  $\Sigma_t$  with high probability. This result is crucial in our analysis since it allows the matrix concentration without using i.i.d. samples. The proof of Lemma 3 utilizes a new Bernstein-type inequality for adapted samples (Lemma EC.5 in the appendix) which may be of independent interest.

LEMMA 3. For  $t \geq \frac{2 \log(2d^2)}{C(\phi_0, s_0)^2}$  where  $C(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_1}\right)$ , we have

$$\mathbb{P} \left( \|\Sigma_t - \hat{\Sigma}_t\|_\infty \geq \frac{\phi_0^2}{32s_0\rho_1} \right) \leq \exp \left\{ -\frac{tC(\phi_0, s_0)^2}{2} \right\}.$$

Then, we invoke the following corollary to use the matrix concentration results to ensure the compatibility condition for  $\hat{\Sigma}_t$ .

**COROLLARY 1 (Corollary 6.8, Bühlmann and Van De Geer (2011)).** *Suppose that  $\Sigma_0$ -compatibility condition holds for the index set  $S$  with cardinality  $s = |S|$ , with compatibility constant  $\phi^2(\Sigma_0, S)$ , and that  $\|\Sigma_1 - \Sigma_0\|_\infty \leq \Delta$ , where  $32s\Delta \leq \phi^2(\Sigma_0, S)$ . Then, for the set  $S$ , the  $\Sigma_1$ -compatibility condition holds as well, with  $\phi^2(\Sigma_1, S) \geq \phi^2(\Sigma_0, S)/2$ .*

In order to satisfy the hypotheses for Lemma 3 and Corollary 1, we define the *initial* period  $t < T_0 := \frac{2\log(2d^2)}{C(\phi_0, s_0)^2}$  during which the compatibility condition for the empirical Gram matrix is not guaranteed, and the event

$$\mathcal{E}_t := \left\{ \|\Sigma_t - \hat{\Sigma}_t\|_\infty \leq \frac{\phi_0^2}{32s_0\rho_1} \right\}.$$

Then for all  $t \geq \lceil T_0 \rceil$  and  $\Sigma_t$  for which event  $\mathcal{E}_t$  holds, we have

$$\phi_t^2 := \phi^2(\hat{\Sigma}_t, S_0) \geq \frac{\phi^2(\Sigma_t, S_0)}{2} \geq \frac{\phi_0^2}{2\rho_1} > 0.$$

Hence, the compatibility condition is satisfied for the empirical Gram matrix without using sparsity information.

**5.2.3. Proof Sketch of Theorem 1** We combine the results above to analyze the regret bound of SA LASSO BANDIT shown in Theorem 1. First, we divide the planning horizon  $[T]$  into three groups:

- (a)  $(t \leq T_0)$ . Here the compatibility condition is not guaranteed to hold.
- (b)  $(t > T_0)$  such that  $\mathcal{E}_t$  holds.
- (c)  $(t > T_0)$  such that  $\mathcal{E}_t$  does not hold.

These sets are disjoint, hence we bound the regret contribution from each separately and obtain an upper bound on the overall regret. It is important to note that SA LASSO BANDIT does not rely in any way on this partitioning, and this is introduced purely for the purpose of analysis. Set (a) is an initial period in which we do not have guarantees for the compatibility condition. Therefore, we cannot apply the Lasso convergence result; hence we can incur  $\mathcal{O}(s_0^2 \log d)$  regret. Set (b) is where the compatibility condition is satisfied; hence the Lasso oracle inequality in Lemma 1 can apply. In fact, this group can be further divided to two cases: (b-1) when the high-probability Lasso result holds and (b-2) when it does not, where the regret of (b-2) can be bounded by  $\mathcal{O}(1)$ . For (b-1), using the Lasso convergence result and summing the regret over the planning horizon gives  $\mathcal{O}(s_0 \sqrt{T \log(dT)})$  regret, which is the leading factor in the regret bound of Theorem 1. Lastly, (c) contains the failure events of Lemma 3 whose regret is  $\mathcal{O}(s_0^2)$ . The complete proofs of the theorem and lemmas are presented in the Appendix.

### 5.3. Regret under the Restricted Eigenvalue Condition

In our analysis so far, we have presented the main results under the compatibility condition in order to be consistent with previous results in the sparse bandit literature. In this section, we present the regret bound for SA LASSO BANDIT under the restricted eigenvalue (RE) condition and briefly discuss its implication in terms of potentially matching lower bounds. Similar to the analysis under the compatibility condition, we assume that the RE condition is satisfied only for the theoretical Gram matrix  $\Sigma = \frac{1}{K}\mathbb{E}[\mathbf{X}^\top \mathbf{X}]$ .

**ASSUMPTION 4 (RE condition).** *For active set  $S_0$  and  $\Sigma$ , there exists restricted eigenvalue  $\phi_1 > 0$  such that  $\phi_1^2 \|\beta\|_2^2 \leq \beta^\top \Sigma \beta$  for all  $\beta \in \mathbb{C}(S_0)$  defined in (2).*

The RE condition is very similar to the compatibility condition in Assumption 2 but uses  $\ell_2$  norm. Based on this condition, we can show the following regret bound.

**THEOREM 2 (Regret bound under RE condition).** *Suppose  $K = 2$  and Assumptions 1, 3, and 4 hold. Then the expected regret of the SA LASSO BANDIT policy is  $\mathcal{O}(\sqrt{s_0 T \log(dT)})$ .*

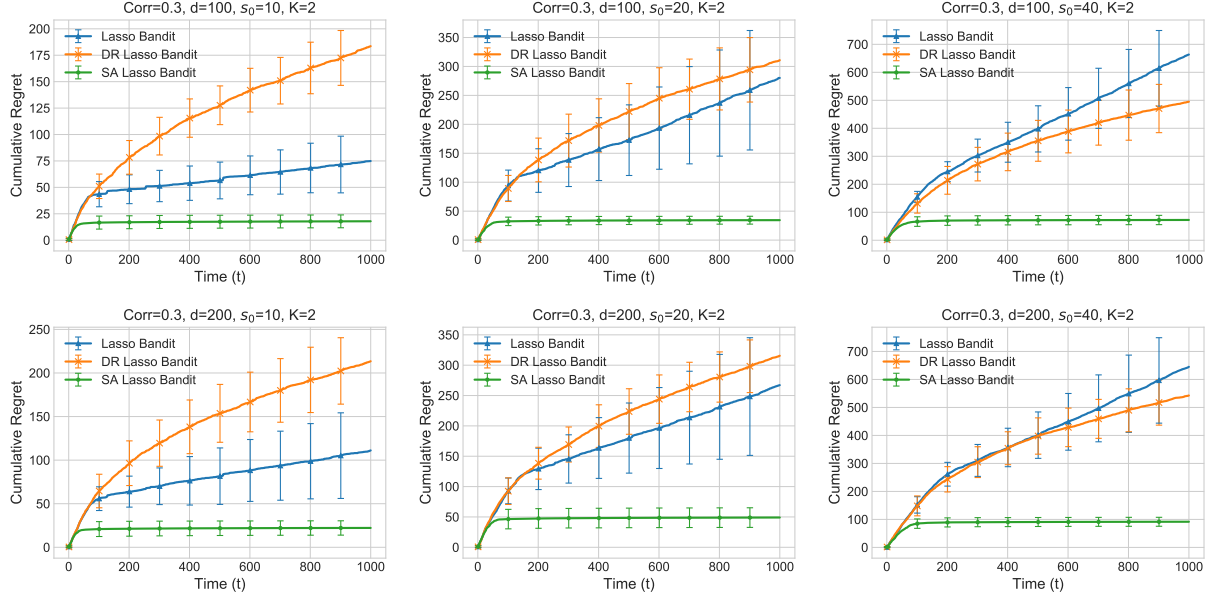
Theorem 2 establishes  $\mathcal{O}(\sqrt{s_0 T \log(dT)})$  regret without any prior knowledge on  $s_0$ . The regret upper-bound based on the RE condition still enjoys logarithmic dependence on  $d$  and furthermore sub-linear dependence on  $s_0$ . Compared to Theorem 1, the regret bound in Theorem 2 is smaller by  $\sqrt{s_0}$  factor, which is again consistent with the offline Lasso results under the RE condition (Theorem 7.19 in Wainwright 2019). The difference in the regret bounds in Theorem 1 and Theorem 2 is due to the inherent difference in the strengths of the RE condition and the compatibility condition, i.e., the RE condition is slightly stronger than the compatibility condition.

The RE condition is more directly analogous (than the compatibility condition) to the standard positive-definiteness assumption for covariance matrices in generalized linear bandits (Li et al. 2017). That is, the RE condition is equivalent to positive-definite covariance when  $s_0 = d$ , i.e., non-sparse settings. Chu et al. (2011) showed  $\Omega(\sqrt{dT})$  minimax lower bound of the expected regret for (non-sparse) linear bandits with finite arms, which is a special case of the generalized linear bandits considered in this paper. Hence, we conjecture that  $\mathcal{O}(\sqrt{s_0 T \log(dT)})$  regret is the best possible up to logarithm factor under the RE condition (and so is  $\mathcal{O}(s_0 \sqrt{T \log(dT)})$  regret under the compatibility condition). Here, while we present these conjectures, we do not claim our results as the minimax regret. In fact, we point out that the notion of minimax regret is much more delicate in sparse bandits than in non-sparse linear bandits. We discuss this briefly in Section 8.

## 6. Numerical Experiments

We conduct numerical experiments to evaluate SA LASSO BANDIT and compare with existing sparse bandit algorithms: DR LASSO BANDIT (Kim and Paik 2019) and LASSO BANDIT (Bastani

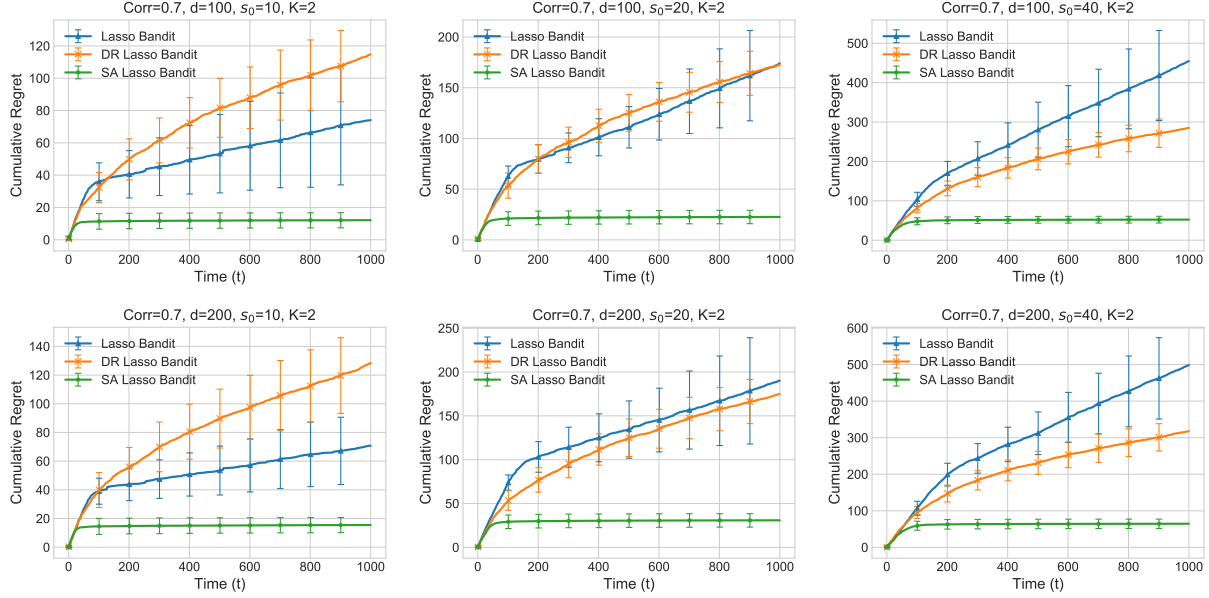




**Figure 1** The plots show the  $t$ -round regret of SA Lasso Bandit (Algorithm 1), DR Lasso Bandit (Kim and Paik 2019), and Lasso Bandit (Bastani and Bayati 2020) for  $K = 2$ ,  $d = 100$  (first row) and  $d = 200$  (second row) with varying sparsity  $s_0 = 10, 20, 40$  under a weak correlation.

and Bayati 2020) in two-armed contextual bandits. We follow the experimental setup of Kim and Paik (2019) to evaluate algorithms under different levels of correlation between arms. Although we consider  $K = 2$  case in this section, the experimental setup introduced here also applies to numerical evaluations for  $K \geq 3$  arm case in Section 7. For each dimension  $i \in [d]$ , we sample each element of the feature vectors  $[X_{t,1}^{(i)}, \dots, X_{t,K}^{(i)}]$  from multivariate Gaussian distribution  $\mathcal{N}(\vec{0}_K, V)$  where covariance matrix  $V$  is defined as  $V_{i,i} = 1$  for all diagonal elements  $i, i$  and  $V_{i,j} = \rho^2$  for all off-diagonal elements  $i \neq j \in [d]$ . Hence, for  $\rho^2 > 0$ , feature vectors for each arm are allowed to be correlated. We consider two levels of correlation with  $\rho^2 = 0.3$  (weak correlation) and  $\rho^2 = 0.7$  (strong correlation). In these two sets of experiments, we consider feature dimensions  $d = 100$  and  $d = 200$ . For comparison, we use a linear reward with the linear link function  $\mu(z) = z$  since both LASSO BANDIT and DR LASSO BANDIT are proposed in linear reward settings. We generate  $\beta^*$  with varying sparsity  $s_0 = \|\beta^*\|_0$ . For a given  $s_0$ , we generate the non-zero elements of  $\beta^*$  from a uniform distribution in  $[0, 1]$ . For each case with different experimental configurations, we conduct 20 independent runs, and report the average of the cumulative regret for each of the algorithms. The error bars represent the standard deviations.

DR LASSO BANDIT is proposed for the same problem setting as ours. Therefore, it does not require any modifications. However, the problem setting of LASSO BANDIT is different from ours: it assumes that the context variable is the same for all arms but the underlying parameter differs for



**Figure 2** The plots show the  $t$ -round regret of SA Lasso Bandit (Algorithm 1), DR Lasso Bandit (Kim and Paik 2019), and Lasso Bandit (Bastani and Bayati 2020) for  $K = 2$ ,  $d = 100$  (first row) and  $d = 200$  (second row) with varying sparsity  $s_0 = 10, 20, 40$  under a strong correlation.

each arm. As was done in the experiments of Kim and Paik (2019), LASSO BANDIT can be applied in our setting by constructing a  $Kd$ -dimensional context vector  $X_t = [X_{t,1}^\top, \dots, X_{t,K}^\top]^\top \in \mathbb{R}^{Kd}$  and  $Kd$ -dimensional parameter  $\beta_i^*$  for each arm  $i$  where  $\beta_i^* = [\beta^{*\top} \mathbb{1}(i=1), \dots, \beta^{*\top} \mathbb{1}(i=K)]^\top \in \mathbb{R}^{Kd}$ . Note that despite the concatenation, the effective dimension of the unknown parameter  $\beta_i^*$  remains the same as far as estimation is concerned. We defer the other details of the experimental setup and additional results to the appendix.

It is important to note that we report the performances of the benchmarks (DR LASSO BANDIT and LASSO BANDIT) assuming that they have access to correct sparsity  $s_0$  and this information is kept hidden from ours. Despite such advantage to the benchmarks, the experiment results shown in Figure 1 and Figure 2 demonstrate that SA LASSO BANDIT outperforms the other methods by significant margins consistently across various instances of experiments. We also verify that the performance of our proposed algorithm is the least sensitive and scales very well with changes in problem instances, which suggests that our algorithm is very effective for various high-dimensional bandit problem instances with a sparse structure. Regret scalability on sparsity  $s_0$  appears to be at most linear while dependence on feature dimension  $d$  appears to be very minimal in most of the instances, which is consistent with our theoretical findings. We also observe that a higher correlation between arms (feature vectors) improves the overall performances of the algorithms. This finding is further evidenced by the experiments for the  $K$ -armed case. We discuss this phenomenon in detail in Section 7.

## 7. Extension to $K$ Arms

Thus far, we have presented our main results in two-armed bandit settings which highlight the main challenges of sparse bandit problems without prior knowledge of sparsity. In this section, we extend our regret analysis to the case of  $K \geq 3$  arms. Also, we present additional numerical experiments for  $K$ -armed bandits.

### 7.1. Regret Analysis for $K$ Arms

Recall that SA LASSO BANDIT is valid for any number of arms; hence it does not require an algorithmic modification for the case of  $K \geq 3$  arms. The analysis of SA LASSO BANDIT for the  $K$ -armed case tackles largely the same challenges described in Section 5.2: that is, the need for a Lasso convergence result for adapted samples and ensuring the compatibility condition without knowing  $s_0$  (and without relying on i.i.d. samples). The former challenge is again taken care of by the Lasso convergence result in Lemma 1. However, the latter issue is more subtle in the  $K$ -armed case than in the two-armed case. In particular, when controlling the adapted Gram matrix  $\Sigma_t$  with the theoretical Gram matrix  $\Sigma$ , the Gram matrix for the unobserved features could be incomparable with the Gram matrix for the observed features. For this issue, we introduce an additional regularity condition, which we denote as the “balanced covariance” condition.

**ASSUMPTION 5 (Balanced covariance).** *Consider a permutation  $(i_1, \dots, i_K)$  of  $(1, \dots, K)$ . For any integer  $k \in \{2, \dots, K-1\}$  and fixed vector  $\beta$ , there exists  $C_{\mathcal{X}} < \infty$  such that*

$$\mathbb{E} [X_{i_k} X_{i_k}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_k}^\top \beta\}] \preceq C_{\mathcal{X}} \mathbb{E} [(X_{i_1} X_{i_1}^\top + X_{i_K} X_{i_K}^\top) \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}].$$

This balanced covariance condition implies that there is “sufficient randomness” in the observed features compared to non-observed features. The exact value of  $C_{\mathcal{X}}$  depends on the joint distribution of  $\mathcal{X}$  including the correlation between arms. In general, the more positive correlation, the smaller  $C_{\mathcal{X}}$  (obviously, with an extreme case of perfectly correlated arms having a constant  $C_{\mathcal{X}}$ ). In the case of independent arms, one can show that both a multivariate Gaussian distribution and a uniform distribution on a  $d$ -dimensional sphere satisfy Assumption 5 with  $C_{\mathcal{X}} = 1$ . For an arbitrary independent distribution for each arm, it holds with  $C_{\mathcal{X}} = \binom{K-1}{K_0}$  where  $K_0 = \lceil (K-1)/2 \rceil$ . It is important to note that even in this pessimistic case,  $C_{\mathcal{X}}$  does not exhibit dependence on dimensionality  $d$  or sparsity  $s_0$ . These are formalized in Proposition EC.1 in the appendix.<sup>4</sup> This balanced covariance condition is somewhat similar to “positive-definiteness” condition for observed contexts in previous bandit literature (e.g., Goldenshluger and Zeevi (2013), Bastani et al. (2017)). However, notice that we allow the covariance matrices on both sides of the inequality to be singular. Hence, the positive-definiteness condition for observed context in our setting may not hold even when the balanced covariance condition holds. While this condition admittedly originates from our proof

technique, it also provides potential insights on learnability of problem instances. That is,  $C_{\mathcal{X}}$  close to infinity implies that the distribution of feature vectors is heavily skew toward a particular direction. Hence, learning algorithms may require much more samples to learn the unknown parameter, which leads to larger regret. It is important to note that our algorithm does not require any prior information on  $C_{\mathcal{X}}$ . Under this suitable regularity, we present the regret bound for  $K$ -armed sparse bandits.

**THEOREM 3 (Regret bound for  $K$  arms).** *Suppose  $K \geq 3$  and Assumptions 1-3, and 5 hold. Then the expected cumulative regret of the SA LASSO BANDIT policy  $\pi$  over horizon  $T$  is upper-bounded by*

$$\mathcal{R}_T(\pi) \leq 4\kappa_{\max} + \frac{2\log(2d^2) + 2}{C_1(\phi_0, s_0)^2} + \frac{64\kappa_{\max}\rho_1 C_{\mathcal{X}}\sigma s_0 \sqrt{T\log(dT)}}{\kappa_{\min}\phi_0^2}$$

where  $C_1(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_1 C_{\mathcal{X}}}\right)$ .

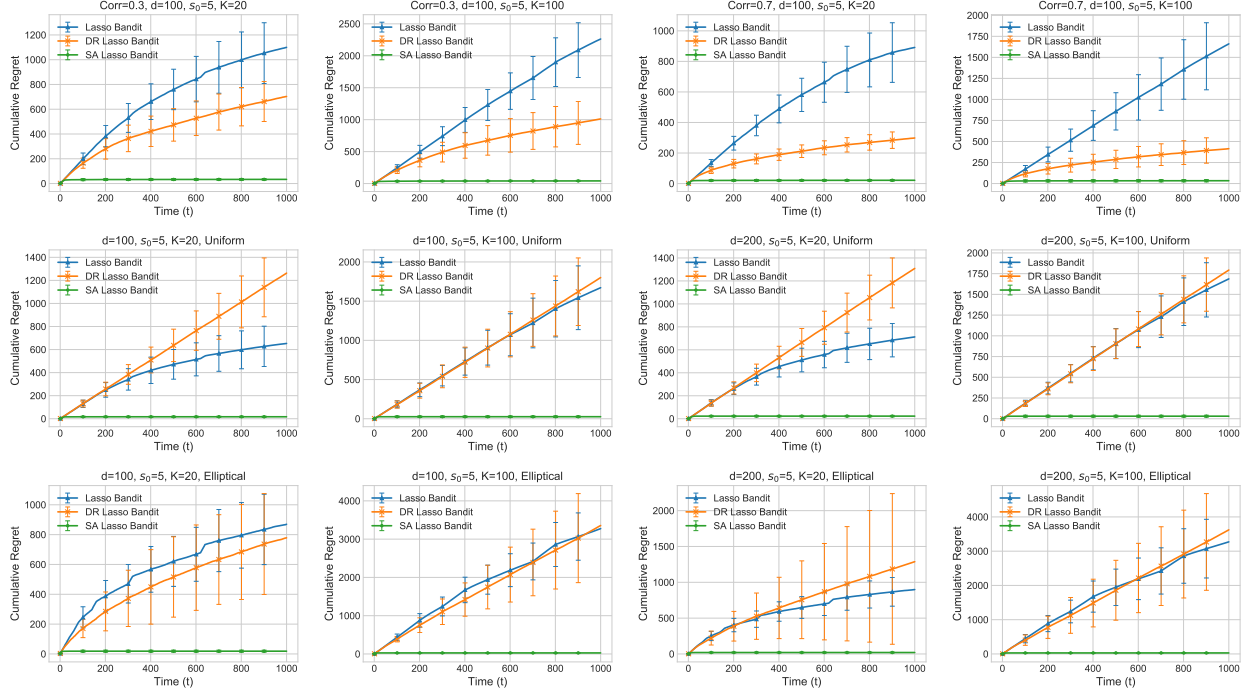
Theorem 3 establishes  $\mathcal{O}(s_0\sqrt{T\log(dT)})$  regret without prior knowledge on  $s_0$ , achieving the same rate as Theorem 1 in terms of the key problem primitives. The proof of Theorem 3 largely follows that of Theorem 1. The main difference is how we control the adapted Gram matrix  $\Sigma_t$  with the theoretical Gram matrix  $\Sigma$ . Under the balanced covariance condition, we can ensure the lower bound of the adapted Gram matrix as a function of the theoretical Gram matrix, which is analogous to the result in Lemma 2. In particular, we can show that for a fixed vector  $\beta \in \mathbb{R}^d$ ,

$$\sum_{i=1}^K \mathbb{E}_{\mathcal{X}_t} \left[ X_{t,i} X_{t,i}^\top \mathbb{1}\{X_{t,i} = \arg \max_{X \in \mathcal{X}_t} X^\top \beta\} \right] \succcurlyeq \frac{\Sigma}{2\rho_1 C_{\mathcal{X}}}.$$

The formal result is presented in Lemma EC.7 in the appendix along with its proof. Then we again invoke the matrix concentration result in Lemma 3 to connect the compatibility constant of empirical Gram matrix  $\hat{\Sigma}_t$  to that of  $\Sigma_t$ , and eventually to the theoretical Gram matrix  $\Sigma$ . Thus, we ensure the compatibility condition of  $\hat{\Sigma}_t$ . The additional cost of regret analysis in the  $K$ -armed case compared to the two-armed case is essentially dependence on  $C_{\mathcal{X}}$  for the balanced covariance condition.

## 7.2. Numerical Experiments for $K$ Arms

We now validate the performance of SA LASSO BANDIT in  $K$ -armed sparse bandit settings via additional numerical experiments and provide comparison with the existing sparse bandit algorithms. The setup of the experiments is identical to the setup described in Section 6. We perform evaluations under various instances. In particular, we focus on the performances of algorithms as the number of arms increases. Additionally, to investigate the effect of the balanced covariance



**Figure 3** Evaluations on  $K$ -armed bandits. The plots show the  $t$ -round regret of SA Lasso Bandit (Algorithm 1), DR Lasso Bandit (Kim and Paik 2019), and Lasso Bandit (Bastani and Bayati 2020) with varying number of arms  $K = 20, 100$  and different distributions. In the first row, features are drawn from a multivariate Gaussian distribution with weak and strong correlation levels. The second rows shows evaluations with features drawn from the uniform distribution on a unit sphere. In the third row, features are drawn from a non-Gaussian elliptical distribution.

condition, we evaluate algorithms on features drawn from a non-Gaussian elliptical distribution, for which we do not have a tight bound of  $C_{\mathcal{X}}$  as well as the uniform distribution.

Figure 3 shows the sample results of the numerical evaluations (averaged over 20 independent runs per problem instance), and the additional results are also presented in the appendix. The experiments results provide the convincing evidence that the performance of our proposed algorithm is superior to the existing sparse bandit methods that we compare with. Again, SA LASSO BANDIT outperforms the existing sparse bandit algorithms by significant margins. Furthermore, SA LASSO BANDIT is much more practical and easy to implement with a minimal number of a hyperparameter (only noise variance parameter is needed for our algorithm). We again observe that under strong correlation, algorithms generally perform better compared to weak correlation instances. This strong correlation would imply a smaller  $C_{\mathcal{X}}$  as briefly discussed earlier when we introduce the balanced covariance condition. Hence, the results are consistent with our theoretical findings. Note that strong correlation does not immediately imply that the performances are

generally better since it potentially decreases the value of compatibility constant. Thus, the regret would increase with an increase in correlation as far as the compatibility condition is concerned. However, as evidenced by our experiments, there appears to be an offsetting effect, which we argue can potentially be explained by the balanced covariance condition. As for features i.i.d. from the uniform distribution and non-Gaussian elliptical distributions,<sup>5</sup> while the performance of existing algorithms (e.g., DR LASSO BANDIT from [Kim and Paik \(2019\)](#)) deteriorates significantly with the change of feature distributions, SA LASSO BANDIT remains very robust, and still exhibits superior performances.

## 8. Conclusion

In this paper, we study a high-dimensional contextual bandit problem with sparse structure. In particular, we address the fundamental issue of learning algorithms having to know a priori the sparsity of the unknown parameter. We propose and analyze an algorithm that circumvent this critical issue. The proposed algorithm achieves a tight regret upper bound which depends on a logarithmic function of the feature dimension which matches the scalability of the offline Lasso convergence results. The algorithm attains this sharp result without knowing the sparsity of the unknown parameter, overcoming weaknesses of the existing algorithms. We demonstrate that our proposed algorithm significantly outperforms the benchmark, supporting the theoretical claims. As we conclude, we outline some of future directions.

### Minimax Regret in Sparse Bandits

Minimax regret in sparse bandits is more delicate to define than in (non-sparse) linear or GLM bandits. To illustrate this challenge, consider the following. If the nature is allowed to freely choose  $s_0 \in [d]$ , it can force the regret for any sparse bandit algorithm to be polynomial in  $d$  by choosing  $s_0 = d$ . On the other hand, if we limit the nature to choose  $s_0 \in [1, s_{\max}]$ , it will choose  $s_0 = s_{\max}$ , and therefore, sparse bandit algorithms can assume that the sparsity index  $s_0$  is known and set to be  $s_{\max}$  if  $s_{\max}$  is also known to algorithms. Thus, it is not clear how to define a minimax criterion in a manner that does not reveal the dominating choice for the nature, and thus forces learning algorithm to play a strategy which protects against a range of values of the sparsity index. It would be interesting to further study what are the right set of assumptions one need to impose to have a non-trivial minimax regret for sparse bandit problems.

### Reinforcement Learning with High-Dimensional Covariates

Another compelling direction is to extend our analysis and proposed approach to reinforcement learning with high-dimensional context or with high-dimensional function approximation. A main challenge in this direction appears to be the need for an algorithm to be optimistic. To our knowledge, almost all reinforcement learning algorithms with provable efficiency rely on the principle of

optimism. As we have discussed in this paper, in order to be optimistic in the tightest sense under sparse structure, the knowledge on sparsity is generally needed. Designing an algorithm that learns an efficient policy without requiring such prior knowledge on sparsity would be very interesting.

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvari. Online-to-confidence-set conversions and application to sparse stochastic bandits. In *Artificial Intelligence and Statistics*, pages 1–9, 2012.
- Naoki Abe and Philip M Long. Associative reinforcement learning using linear probabilistic concepts. In *International Conference on Machine Learning*, pages 3–11, 1999.
- Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135, 2013.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Heejung Bang and James M Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.
- Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. *Operations Research*, 68(1):276–294, 2020.
- Hamsa Bastani, Mohsen Bayati, and Khashayar Khosravi. Mostly exploration-free algorithms for contextual bandits. *arXiv preprint arXiv:1704.09011*, 2017.
- Peter J Bickel, Ya’acov Ritov, Alexandre B Tsybakov, et al. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.
- Stamatis Cambanis, Steel Huang, and Gordon Simons. On the theory of elliptically contoured distributions. *Journal of Multivariate Analysis*, 11(3):368–385, 1981.
- Emmanuel Candes, Terence Tao, et al. The dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *The annals of Statistics*, 35(6):2313–2351, 2007.
- Alexandra Carpentier and Rémi Munos. Bandit theory meets compressed sensing for high dimensional stochastic linear bandit. In *Artificial Intelligence and Statistics*, pages 190–198, 2012.
- Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214, 2011.

- Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. In *Proceedings of the 21st Annual Conference on Learning Theory*, page 355–366, 2008.
- Sarah Filippi, Olivier Cappe, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The generalized linear case. In *Advances in Neural Information Processing Systems*, pages 586–594, 2010.
- Davis Gilton and Rebecca Willett. Sparse linear contextual bandits via relevance vector machines. In *2017 International Conference on Sampling Theory and Applications (SampTA)*, pages 518–522. IEEE, 2017.
- Alexander Goldenshluger and Assaf Zeevi. A linear response bandit problem. *Stochastic Systems*, 3(1):230–261, 2013.
- Gi-Soo Kim and Myunghee Cho Paik. Doubly-robust lasso bandit. In *Advances in Neural Information Processing Systems*, pages 5869–5879, 2019.
- Branislav Kveton, Manzil Zaheer, Csaba Szepesvari, Lihong Li, Mohammad Ghavamzadeh, and Craig Boutilier. Randomized exploration in generalized linear bandits. *arXiv preprint arXiv:1906.08947*, 2019.
- Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press (preprint), 2019.
- Lihong Li, Yu Lu, and Dengyong Zhou. Provably optimal algorithms for generalized linear contextual bandits. In *International Conference on Machine Learning*, pages 2071–2080, 2017.
- Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Restricted eigenvalue properties for correlated gaussian designs. *Journal of Machine Learning Research*, 11(Aug):2241–2259, 2010.
- Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Sara A Van De Geer, Peter Bühlmann, et al. On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, 3:1360–1392, 2009.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Xue Wang, Mingcheng Wei, and Tao Yao. Minimax concave penalized multi-armed bandit model with high-dimensional covariates. In *International Conference on Machine Learning*, pages 5200–5208, 2018.
- Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894–942, 2010.



## EC.1. Proofs of Lemmas for Theorem 1

### EC.1.1. Proof of Lemma 1

The proof follows from modifying the proof of the standard Lasso oracle inequality ([Bühlmann and Van De Geer 2011](#)) using martingale theory. This is also a generalized version of Proposition 1 in [Bastani and Bayati \(2020\)](#) which shows the tail inequality for the linear model with adapted samples. Recall from (1) that the negative log-likelihood of the GLM is

$$\ell_t(\beta) = -\frac{1}{t} \sum_{\tau=1}^t [Y_\tau X_\tau^\top \beta - m(X_\tau^\top \beta)]$$

where  $m$  is a normalizing function with its gradient  $\dot{m}(X^\top \beta) = \mu(X^\top \beta)$ . Now, we denote the expectation of  $\ell_t(\beta)$  over  $Y$  by  $\bar{\ell}_t(\beta)$ :

$$\bar{\ell}_t(\beta) := \mathbb{E}_Y[\ell_t(\beta)] = -\frac{1}{t} \sum_{\tau=1}^t [\mu(X_\tau^\top \beta^*) X_\tau^\top \beta - m(X_\tau^\top \beta)].$$

Note that  $\nabla_\beta \bar{\ell}_t(\beta) = -\frac{1}{t} \sum_{\tau=1}^t [\mu(X_\tau^\top \beta^*) - \mu(X_\tau^\top \beta)] X_\tau$ . Hence, we have  $\nabla_\beta \bar{\ell}_t(\beta^*) = \vec{0}$  which implies that  $\beta^* = \arg \min_\beta \bar{\ell}_t(\beta)$  given the fact that  $m$  is convex in the GLM. Hence, for any parameter  $\beta \in \mathbb{R}^d$ , the excess risk is defined as

$$\mathcal{E}(\beta) := \bar{\ell}_t(\beta) - \bar{\ell}_t(\beta^*).$$

Note that by definition,  $\mathcal{E}(\beta) \geq 0$ , for all  $\beta \in \mathbb{R}^d$  (with  $\mathcal{E}(\beta^*) = 0$ ). The Lasso estimate  $\hat{\beta}_t$  for the GLM is given by the minimization of the penalized negative log-likelihood

$$\hat{\beta}_t := \arg \min_{\beta} \{ \ell_t(\beta) + \lambda_t \|\beta\|_1 \}$$

where  $\lambda$  is the penalty parameter whose value needs to be chosen to control the noise of the model. Now, we define the empirical process of the problem as

$$v_t(\beta) := \ell_t(\beta) - \bar{\ell}_t(\beta).$$

Note that the randomness in  $\{Y_\tau\}$  still plays a role on  $\ell_t(\beta)$  and hence on  $v_t(\beta)$ . Then by the definition of  $\hat{\beta}_t$ , we have

$$\ell_t(\hat{\beta}_t) + \lambda_t \|\hat{\beta}_t\|_1 \leq \ell_t(\beta^*) + \lambda_t \|\beta^*\|_1. \quad (\text{EC.1})$$

Adding and subtracting terms, we have

$$\ell_t(\hat{\beta}_t) - \bar{\ell}_t(\hat{\beta}_t) + \bar{\ell}_t(\hat{\beta}_t) - \bar{\ell}_t(\beta^*) + \lambda_t \|\hat{\beta}_t\|_1 \leq \ell_t(\beta^*) - \bar{\ell}_t(\beta^*) + \lambda_t \|\beta^*\|_1.$$

Rearranging terms gives the following “basic inequality” for the GLM

$$\mathcal{E}(\hat{\beta}_t) + \lambda_t \|\hat{\beta}_t\|_1 \leq -[v_t(\hat{\beta}_t) - v_t(\beta^*)] + \lambda_t \|\beta^*\|_1. \quad (\text{EC.2})$$

The basic inequality implies that to provide an upper-bound for the penalized excess risk, we need to control the behavior of the increments of the empirical process  $[v_t(\hat{\beta}_t) - v_t(\beta^*)]$ . And we will bound this in terms of  $\|\hat{\beta}_t - \beta^*\|_1$ . Essentially,  $[v_t(\hat{\beta}_t) - v_t(\beta^*)]$  is where the random noise plays a role and with large enough penalization (suitably large  $\lambda$ ) we can control such randomness in the empirical process. We define the event of the empirical process being controlled by the penalization.

$$\mathcal{T} := \{|v_t(\hat{\beta}_t) - v_t(\beta^*)| \leq \lambda_0 \|\hat{\beta}_t - \beta^*\|_1\}. \quad (\text{EC.3})$$

Lemma EC.1 ensures that we can control this empirical process with a high probability. Hence, in the rest of the proof, we restrict ourselves to the case where the empirical process behaves well, i.e., event  $\mathcal{T}$  in (EC.3) holds.

LEMMA EC.1. *Assume  $X_t$  satisfies  $\|X_t\| \leq 1$  for all  $t$ . If  $\lambda_0 = \sigma \sqrt{\frac{2[\log(2/\delta) + \log d]}{t}}$ , then with probability at least  $1 - \delta$  we have*

$$|v_t(\hat{\beta}_t) - v_t(\beta^*)| \leq \lambda_0 \|\hat{\beta}_t - \beta^*\|_1. \quad (\text{EC.4})$$

On  $\mathcal{T}$ , for  $\lambda_t \geq 2\lambda_0$ , we have

$$2\mathcal{E}(\hat{\beta}_t) + 2\lambda_t \|\hat{\beta}_t\|_1 \leq \lambda_t \|\hat{\beta}_t - \beta^*\|_1 + 2\lambda_t \|\beta^*\|_1. \quad (\text{EC.5})$$

Let  $\hat{\beta} := \hat{\beta}_t$  for brevity. Using the active set  $S_0$ , we can define the following:

$$\beta_{j,S_0} := \beta_j \mathbb{1}\{j \in S_0\} \quad \beta_{j,S_0^c} := \beta_j \mathbb{1}\{j \notin S_0\} \quad (\text{EC.6})$$

so that  $\beta_{S_0} = [\beta_{1,S_0}, \dots, \beta_{d,S_0}]^\top$  has zero elements outside the set  $S_0$  and the elements of  $\beta_{S_0^c}$  can only be non-zero in the complement of  $S_0$ . We can then lower-bound  $\|\hat{\beta}\|_1$  using the triangle inequality,

$$\begin{aligned} \|\hat{\beta}\|_1 &= \|\hat{\beta}_{S_0}\|_1 + \|\hat{\beta}_{S_0^c}\|_1 \\ &\geq \|\beta_{S_0}^*\|_1 - \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 + \|\hat{\beta}_{S_0^c}\|_1. \end{aligned}$$

Also, we can rewrite

$$\begin{aligned} \|\hat{\beta} - \beta^*\|_1 &= \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 + \|\hat{\beta}_{S_0^c} - \beta_{S_0^c}^*\|_1 \\ &= \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 + \|\hat{\beta}_{S_0^c}\|_1. \end{aligned}$$

Then we continue from (EC.5)

$$\begin{aligned} 2\mathcal{E}(\hat{\beta}) + 2\lambda_t \|\beta_{S_0}^*\|_1 - 2\lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 + 2\lambda_t \|\hat{\beta}_{S_0^c}\|_1 &\leq \lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 + \lambda_t \|\hat{\beta}_{S_0^c}\|_1 + 2\lambda_t \|\beta^*\|_1 \\ &= \lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 + \lambda_t \|\hat{\beta}_{S_0^c}\|_1 + 2\lambda_t \|\beta_{S_0}^*\|_1. \end{aligned}$$

Therefore, we have

$$\begin{aligned} 0 \leq 2\mathcal{E}(\hat{\beta}) &\leq 3\lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 - \lambda_t \|\hat{\beta}_{S_0^c} - \beta_{S_0^c}^*\|_1 \\ &= \lambda_t \left( 3\|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 - \|\hat{\beta}_{S_0^c} - \beta_{S_0^c}^*\|_1 \right) \end{aligned} \quad (\text{EC.7})$$

Then the compatibility condition can be applied to the vector  $\hat{\beta} - \beta^*$  which gives

$$\|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1^2 \leq s_0 (\hat{\beta} - \beta^*)^\top \hat{\Sigma} (\hat{\beta} - \beta^*) / \phi_t^2. \quad (\text{EC.8})$$

From (EC.7), we have

$$2\mathcal{E}(\hat{\beta}) + \lambda_t \|\hat{\beta}_{S_0^c} - \beta_{S_0^c}^*\|_1 \leq 3\lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1.$$

Therefore, we have

$$\begin{aligned} 2\mathcal{E}(\hat{\beta}) + \lambda_t \|\hat{\beta} - \beta^*\|_1 &= 2\mathcal{E}(\hat{\beta}) + \lambda_t \|\hat{\beta}_{S_0^c} - \beta_{S_0^c}^*\|_1 + \lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 \\ &\leq 3\lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 + \lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 \\ &= 4\lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 \\ &\leq 4\lambda_t \sqrt{s_0 (\hat{\beta} - \beta^*)^\top \hat{\Sigma} (\hat{\beta} - \beta^*) / \phi_t} \\ &\leq \kappa_{\min} (\hat{\beta} - \beta^*)^\top \hat{\Sigma} (\hat{\beta} - \beta^*) + \frac{4\lambda^2 s_0}{\kappa_{\min} \phi_t^2} \\ &\leq 2\mathcal{E}(\hat{\beta}) + \frac{4\lambda^2 s_0}{\kappa_{\min} \phi_t^2} \end{aligned}$$

where the second inequality is from applying the compatibility condition (EC.8) and the third inequality is by using  $4uv \leq u^2 + 4v^2$  with  $u = \sqrt{\kappa_{\min} (\hat{\beta} - \beta^*)^\top \hat{\Sigma} (\hat{\beta} - \beta^*)}$  and  $v = \frac{\lambda_t \sqrt{s_0}}{\phi_t \sqrt{\kappa_{\min}}}$ . The last inequality is from Lemma EC.2. Hence, rearranging gives

$$\|\hat{\beta} - \beta^*\|_1 \leq \frac{4s_0 \lambda_t}{\kappa_{\min} \phi_t^2}.$$

This completes the proof.

### EC.1.2. Proof of Lemma EC.1

Let  $\mathcal{F}_t$  be the sigma algebra generated by  $\{X_1, Y_1, \dots, X_t, Y_t\}$ .

By the definition of  $v_t(\beta)$  and  $\ell_t(\beta)$ , we have

$$\begin{aligned} v_t(\beta) &= \ell_t(\beta) - \ell(\beta) \\ &= -\frac{1}{t} \sum_{\tau=1}^t [Y_\tau X_\tau^\top \beta - m(X_\tau^\top \beta)] + \frac{1}{t} \sum_{\tau=1}^t [\mu(X_\tau^\top \beta^*) X_\tau^\top \beta - m(X_\tau^\top \beta)] \\ &= -\frac{1}{t} \sum_{\tau=1}^t [Y_\tau X_\tau^\top \beta - \mu(X_\tau^\top \beta^*) X_\tau^\top \beta] \\ &= -\frac{1}{t} \sum_{\tau=1}^t \epsilon_\tau X_\tau^\top \beta \end{aligned}$$

where the last equality uses the definition of  $\epsilon_\tau$ . Then, the empirical process is

$$v_t(\hat{\beta}_t) - v_n(\beta^*) = -\frac{1}{t} \sum_{\tau=1}^t \epsilon_\tau X_\tau^\top (\hat{\beta}_t - \beta^*).$$

Applying Hölder's inequality, we have

$$|v_t(\hat{\beta}_t) - v_t(\beta^*)| \leq \frac{1}{t} \left\| \sum_{\tau=1}^t \epsilon_\tau X_\tau \right\|_\infty \|\hat{\beta}_t - \beta^*\|_1.$$

Then controlling the empirical process reduces to controlling  $\frac{1}{t} \left\| \sum_{\tau=1}^t \epsilon_\tau X_\tau \right\|_\infty$ .

Using a union bound, we have

$$\begin{aligned} \mathbb{P} \left( \frac{1}{t} \left\| \sum_{\tau=1}^t \epsilon_\tau X_\tau \right\|_\infty \leq \lambda_0 \right) &= 1 - \mathbb{P} \left( \frac{1}{t} \left\| \sum_{\tau=1}^t \epsilon_\tau X_\tau \right\|_\infty > \lambda_0 \right) \\ &\geq 1 - \sum_{j=1}^d \mathbb{P} \left( \frac{1}{t} \left| \sum_{\tau=1}^t \epsilon_\tau X_\tau^{(j)} \right| > \lambda_0 \right) \end{aligned}$$

where  $X_\tau^{(j)}$  is the  $j$ -th element of  $X_\tau$ . For each  $j \in [d]$ , and  $\tau \in [t]$ , we let  $Z_\tau^{(j)} := \epsilon_\tau X_\tau^{(j)}$ . Note that each  $Z_\tau^{(j)}$  for  $j \in [d]$  is a martingale difference sequence adapted to the filtration  $\mathcal{F}_1 \subset \dots \subset \mathcal{F}_\tau$  since  $\mathbb{E}[\epsilon_\tau X_\tau^{(j)} | \mathcal{F}_\tau] = 0$  for each  $j$ . Note that each  $X_\tau^{(j)}$  is a bounded random variable with  $|X_\tau^{(j)}| \leq \|X_\tau\|_\infty \leq \|X_\tau\| \leq 1$ . Then from the fact that  $\epsilon_\tau$  is  $\sigma^2$ -sub-Gaussian, it follows that  $Z_\tau^{(j)}$  is also  $\sigma^2$ -sub-Gaussian. That is, for any  $\alpha \in \mathbb{R}$ ,

$$\begin{aligned} \mathbb{E} [\exp\{\alpha Z_\tau^{(j)}\} | \mathcal{F}_{\tau-1}] &= \mathbb{E} [\exp\{(\alpha X_\tau^{(j)}) \epsilon_\tau\} | \mathcal{F}_{\tau-1}] \\ &\leq \mathbb{E}_{X_\tau} \left[ \exp \left\{ \frac{(\alpha X_\tau^{(j)} \sigma)^2}{2} \right\} | \mathcal{F}_{\tau-1} \right] \\ &\leq \exp \left\{ \frac{\alpha^2 \sigma^2}{2} \right\}. \end{aligned}$$

Then, using Bernstein Concentration in Lemma EC.11, we have

$$\mathbb{P} \left( \left| \sum_{\tau=1}^t \epsilon_\tau X_\tau^{(j)} \right| > t\lambda_0 \right) \leq 2 \exp \left\{ -\frac{t^2 \lambda_0^2}{2t\sigma^2} \right\} \leq 2 \exp \left\{ -\frac{t\lambda_0^2}{2\sigma^2} \right\}$$

So, with  $\lambda_0 = \sigma \sqrt{\frac{2[\log(2/\delta) + \log d]}{t}}$ , we have

$$\mathbb{P} \left( \frac{1}{t} \left\| \sum_{\tau=1}^t \epsilon_\tau X_\tau \right\|_\infty \leq \lambda_0 \right) \geq 1 - 2d \exp \left\{ \log \frac{\delta}{2} - \log d \right\} = 1 - \delta.$$

LEMMA EC.2. *The excess risk is lower-bounded by*

$$\mathcal{E}(\hat{\beta}_t) \geq \frac{\kappa_{\min}}{2} (\hat{\beta}_t - \beta^*)^\top \hat{\Sigma} (\hat{\beta}_t - \beta^*).$$

By the definition of the excess risk  $\mathcal{E}(\beta)$ , we have

$$\begin{aligned}\mathcal{E}(\beta) &= \bar{\ell}_t(\beta) - \bar{\ell}_t(\beta^*) \\ &= -\frac{1}{t} \sum_{\tau=1}^t [\mu(X_\tau^\top \beta^*) X_\tau^\top \beta - m(X_\tau^\top \beta)] + \frac{1}{t} \sum_{\tau=1}^t [\mu(X_\tau^\top \beta^*) X_\tau^\top \beta^* - m(X_\tau^\top \beta^*)].\end{aligned}$$

Since  $\dot{m}(\cdot) = \mu(\cdot)$ , we have  $\nabla_\beta \bar{\ell}_t(\beta^*) = \vec{0}$ . Hence, the gradient of the excess risk  $\nabla_\beta \mathcal{E}(\beta)$  and the Hessian are given as

$$\begin{aligned}\nabla_\beta \mathcal{E}(\beta) &= -\frac{1}{t} \sum_{\tau=1}^t [\mu(X_\tau^\top \beta^*) X_\tau - \mu(X_\tau^\top \beta) X_\tau] \\ H_\mathcal{E}(\beta) &:= \nabla_\beta^2 \mathcal{E}(\beta) = \frac{1}{t} \sum_{\tau=1}^t \dot{\mu}(X_\tau^\top \beta) X_\tau X_\tau^\top.\end{aligned}$$

Using the Taylor expansion, with  $\bar{\beta} = c\beta^* + (1-c)\hat{\beta}$  for some  $c \in (0, 1)$

$$\mathcal{E}(\hat{\beta}_t) = \mathcal{E}(\beta^*) + \nabla_\beta \mathcal{E}(\beta^*)^\top (\hat{\beta}_t - \beta^*) + \frac{1}{2} (\hat{\beta}_t - \beta^*)^\top H_\mathcal{E}(\bar{\beta}) (\hat{\beta}_t - \beta^*). \quad (\text{EC.9})$$

Note that by the definition of  $\beta^*$ , we have  $\mathcal{E}(\beta^*) = 0$  and  $\nabla_\beta \mathcal{E}(\beta^*) = \nabla_\beta \ell(\beta^*) = \vec{0}$ . Hence, combining with the definition of the Hessian, we have

$$\begin{aligned}\mathcal{E}(\hat{\beta}_t) &= \frac{1}{2} (\hat{\beta}_t - \beta^*)^\top \left[ \frac{1}{t} \sum_{\tau=1}^t \dot{\mu}(X_\tau^\top \beta) X_\tau X_\tau^\top \right] (\hat{\beta}_t - \beta^*) \\ &\geq \frac{\kappa_{\min}}{2} (\hat{\beta}_t - \beta^*)^\top \hat{\Sigma} (\hat{\beta}_t - \beta^*)\end{aligned}$$

where the last inequality is from Assumption 1 and  $\hat{\Sigma} = \frac{1}{t} \sum_{\tau=1}^t X_\tau X_\tau^\top$ .

### EC.1.3. Proof of Lemma 2

Consider  $\mathcal{X} = \{X_1, X_2\}$ . Let the joint distribution of  $x_1, x_2$  as  $p_{\mathcal{X}}(x_1, x_2)$ . Then we have

$$\begin{aligned}\mathbb{E}[\mathbf{X}^\top \mathbf{X}] &= \int (x_1 x_1^\top + x_2 x_2^\top) p_{\mathcal{X}}(x_1, x_2) dx_1, x_2 \\ &= \int x_1 x_1^\top [\mathbb{1}\{(x_1 - x_2)^\top \beta \geq 0\} + \mathbb{1}\{(x_1 - x_2)^\top \beta \leq 0\}] p_{\mathcal{X}}(x_1, x_2) dx_1, x_2 \\ &\quad + \int x_2 x_2^\top [\mathbb{1}\{(x_1 - x_2)^\top \beta \geq 0\} + \mathbb{1}\{(x_1 - x_2)^\top \beta \leq 0\}] p_{\mathcal{X}}(x_1, x_2) dx_1, x_2\end{aligned}$$

Let's first look at the first integral.

$$\begin{aligned}&\int x_1 x_1^\top [\mathbb{1}\{(x_1 - x_2)^\top \beta \geq 0\} + \mathbb{1}\{(x_1 - x_2)^\top \beta \leq 0\}] p_{\mathcal{X}}(x_1, x_2) dx_1, x_2 \\ &= \int x_1 x_1^\top [\mathbb{1}\{(x_1 - x_2)^\top \beta \geq 0\} p_{\mathcal{X}}(x_1, x_2) + \mathbb{1}\{-(x_1 - x_2)^\top \beta \geq 0\} p_{\mathcal{X}}(x_1, x_2)] dx_1, x_2 \\ &\preceq \int x_1 x_1^\top \mathbb{1}\{(x_1 - x_2)^\top \beta \geq 0\} p_{\mathcal{X}}(x_1, x_2) dx_1, x_2\end{aligned}$$

$$\begin{aligned}
& + \rho_0 \int x_1 x_1^\top \mathbb{1} \{-(x_1 - x_2)^\top \beta \geq 0\} p_{\mathcal{X}}(-x_1, -x_2) dx_1, x_2 \\
& = \int x_1 x_1^\top \mathbb{1} \{(x_1 - x_2)^\top \beta \geq 0\} p_{\mathcal{X}}(x_1, x_2) dx_1, x_2 \\
& \quad + \rho_0 \int x_1 x_1^\top \mathbb{1} \{(x_1 - x_2)^\top \beta \geq 0\} p_{\mathcal{X}}(x_1, x_2) dx_1, x_2 \\
& = (1 + \rho_0) \int x_1 x_1^\top \mathbb{1} \{(x_1 - x_2)^\top \beta \geq 0\} p_{\mathcal{X}}(x_1, x_2) dx_1, x_2 \\
& = (1 + \rho_0) \mathbb{E} \left[ X_1 X_1^\top \mathbb{1} \{X_1 = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right]
\end{aligned}$$

where the inequality follows from Assumption 3. Likewise, we can show for the second integral that

$$\begin{aligned}
& \int x_2 x_2^\top [\mathbb{1} \{(x_1 - x_2)^\top \beta \geq 0\} + \mathbb{1} \{(x_1 - x_2)^\top \beta \leq 0\}] p_{\mathcal{X}}(x_1, x_2) dx_1, x_2 \\
& = (1 + \rho_0) \mathbb{E} \left[ X_2 X_2^\top \mathbb{1} \{X_2 = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right].
\end{aligned}$$

Hence,

$$\mathbb{E}[\mathbf{X}^\top \mathbf{X}] = (1 + \rho_0) \left( \mathbb{E} \left[ X_1 X_1^\top \mathbb{1} \{X_1 = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right] + \mathbb{E} \left[ X_2 X_2^\top \mathbb{1} \{X_2 = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right] \right).$$

Therefore, with the fact that  $\rho_0 \geq 1$ , we have

$$\sum_{i=1}^2 \mathbb{E} \left[ X_i X_i^\top \mathbb{1} \{X_i = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right] \succcurlyeq \frac{2}{1 + \rho_0} \cdot \frac{\mathbb{E}[\mathbf{X}^\top \mathbf{X}]}{2} \succcurlyeq \frac{\Sigma}{\rho_0}.$$

#### EC.1.4. Bernstein-type Inequality for Adapted Samples

In this section, we derive a Bernstein-type inequality for adapted samples which is shown in Lemma EC.5. We first define the following function of a random variable  $X_t$  which is used throughout this section.

DEFINITION EC.1. For all  $i, j$  with  $1 \leq i \leq j \leq d$ , we define  $\gamma_t^{ij}(X_t)$  to be a real-value function which take random variable  $X_t \in \mathbb{R}^d$  as input:

$$\gamma_t^{ij}(X_t) := X_t^{(i)} X_t^{(j)} - \mathbb{E}[X_t^{(i)} X_t^{(j)} \mid \mathcal{F}_{t-1}] \tag{EC.10}$$

where  $X_t^{(i)}$  is the  $i$ -th element of  $X_t$ .

It is easy to see that  $\mathbb{E}[\gamma_t^{ij}(X_t) \mid \mathcal{F}_{t-1}] = 0$  and  $\mathbb{E}[|\gamma_t^{ij}(X_t)|^m \mid \mathcal{F}_{t-1}] \leq 1$  for all integer  $m \geq 2$ .

LEMMA EC.3 (**Bühlmann and Van De Geer (2011), Lemma 14.1**). *Let  $Z_t \in \mathbb{R}$  be a random variable with  $\mathbb{E}[Z_t \mid \mathcal{F}_{t-1}] = 0$ . Then it holds that*

$$\log \mathbb{E} [e^{Z_t} \mid \mathcal{F}_{t-1}] \leq \mathbb{E} [e^{|Z_t|} \mid \mathcal{F}_{t-1}] - 1 - \mathbb{E} [|Z| \mid \mathcal{F}_{t-1}].$$

The proof follows directly from the proof of Lemma 14.1 in [Bühlmann and Van De Geer \(2011\)](#), applying their result to a conditional expectation. For any  $c > 0$ ,

$$\begin{aligned} \exp(Z_t - c) - 1 &\leq \frac{\exp(Z_t) - 1}{1 + c} \\ &= \frac{e^{Z_t} - 1 - Z_t + Z_t - c}{1 + c} \\ &\leq \frac{e^{|Z_t|} - 1 - |Z_t| + Z_t - c}{1 + c}. \end{aligned}$$

Let  $c = \mathbb{E}[e^{|Z_t|} | \mathcal{F}_{t-1}] - 1 - \mathbb{E}[|Z_t| | \mathcal{F}_{t-1}]$ . Hence, since  $\mathbb{E}[Z_t | \mathcal{F}_{t-1}] = 0$ ,

$$\mathbb{E}[\exp(Z_t - c) | \mathcal{F}_{t-1}] - 1 \leq \frac{\mathbb{E}[e^{|Z_t|} | \mathcal{F}_{t-1}] - 1 - \mathbb{E}[|Z_t| | \mathcal{F}_{t-1}] - c}{1 + c} = 0.$$

LEMMA EC.4. *Suppose  $\mathbb{E}[\gamma_t^{ij}(X) | \mathcal{F}_{t-1}] = 0$  and  $\mathbb{E}[|\gamma_t^{ij}(X)|^m | \mathcal{F}_{t-1}] \leq m!$  for all integer  $m \geq 2$ , all  $t \geq 1$  and all  $1 \leq i \leq j \leq d$ . Then, for  $L > 1$  we have*

$$\mathbb{E} \left[ \exp \left( \frac{1}{L} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right) \right] \leq \exp \left( \frac{\tau}{L(L-1)} \right).$$

$$\begin{aligned} \mathbb{E} \left[ \exp \left( \frac{1}{L} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right) \right] &= \mathbb{E} \left[ \mathbb{E} \left[ \exp \left( \frac{1}{L} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right) | \mathcal{F}_{\tau-1} \right] \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ \exp \left( \frac{\gamma_{\tau}^{ij}(X_{\tau})}{L} \right) | \mathcal{F}_{\tau-1} \right] \exp \left( \frac{1}{L} \sum_{t=1}^{\tau-1} \gamma_t^{ij}(X_t) \right) \right] \\ &\leq e^{\frac{1}{L(L-1)}} \mathbb{E} \left[ \exp \left( \frac{1}{L} \sum_{t=1}^{\tau-1} \gamma_t^{ij}(X_t) \right) \right] \end{aligned}$$

where the inequality is from Lemma EC.3 and noting that

$$\begin{aligned} \log \mathbb{E} \left[ \exp \left( \frac{\gamma_{\tau}^{ij}(X_{\tau})}{L} \right) | \mathcal{F}_{\tau-1} \right] &\leq \mathbb{E} \left[ e^{|\gamma_{\tau}^{ij}(X_{\tau})|/L} - 1 - \frac{|\gamma_{\tau}^{ij}(X_{\tau})|}{L} | \mathcal{F}_{\tau-1} \right] \\ &= \mathbb{E} \left[ \sum_{m=2}^{\infty} \frac{|\gamma_{\tau}^{ij}(X_{\tau})|^m}{L^m m!} | \mathcal{F}_{\tau-1} \right] \\ &= \sum_{m=2}^{\infty} \frac{\mathbb{E}[|\gamma_{\tau}^{ij}(X_{\tau})|^m | \mathcal{F}_{\tau-1}]}{L^m m!} \\ &\leq \frac{1}{L(L-1)} \end{aligned}$$

Then, repeatedly applying this to the rest of the sum  $\frac{1}{L} \sum_{t=1}^{\tau-1} \gamma_t^{ij}(X_t)$ , we have

$$\mathbb{E} \left[ \exp \left( \frac{1}{L} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right) \right] \leq \exp \left( \frac{\tau}{L(L-1)} \right)$$

LEMMA EC.5.

$$\mathbb{P} \left( \max_{1 \leq i \leq j \leq d} \left| \frac{1}{\tau} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right| \geq w + \sqrt{2w} + \sqrt{\frac{4 \log(2d^2)}{\tau} + \frac{2 \log(2d^2)}{\tau}} \right) \leq \exp \left( -\frac{\tau w}{2} \right)$$

Using the Chernoff bound and Lemma EC.4, for any  $L > 1$  we have

$$\begin{aligned} \mathbb{P} \left( \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \geq a \right) &= \mathbb{P} \left( \exp \left\{ \frac{1}{L} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right\} \geq \exp \left( \frac{a}{L} \right) \right) \\ &\leq \frac{\mathbb{E} \left[ \exp \left( \frac{1}{L} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right) \right]}{\exp \left( \frac{a}{L} \right)} \\ &\leq \exp \left( -\frac{a}{L} \right) \exp \left( \frac{\tau}{L(L-1)} \right) \\ &= \exp \left( -\frac{a}{L} + \frac{\tau}{L(L-1)} \right). \end{aligned}$$

Here,  $L = \frac{\tau+a+\sqrt{\tau^2+\tau a}}{a}$  minimizes the right hand side above for  $L > 1$ . Therefore,

$$\begin{aligned} \mathbb{P} \left( \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \geq a \right) &\leq \exp \left\{ -\frac{a^2}{\tau + a + \sqrt{\tau^2 + \tau a}} + \frac{\tau a^2}{(\tau + a + \sqrt{\tau^2 + \tau a})(\tau + \sqrt{\tau^2 + \tau a})} \right\} \\ &= \exp \left\{ -\left( \frac{\sqrt{1+a/\tau}}{1 + \sqrt{1+a/\tau}} \right) \frac{a^2}{\tau + a + \sqrt{\tau^2 + \tau a}} \right\} \\ &\leq \exp \left\{ -\frac{a^2}{2(\tau + a + \sqrt{\tau^2 + \tau a})} \right\} \\ &\leq \exp \left\{ -\frac{a^2}{2(\tau + a + \sqrt{\tau^2 + 2\tau a})} \right\}. \end{aligned}$$

Choosing  $a = \tau(w + \sqrt{2w})$  gives

$$\mathbb{P} \left( \frac{1}{\tau} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \geq w + \sqrt{2w} \right) \leq \exp \left\{ -\frac{\tau w}{2} \right\}. \quad (\text{EC.11})$$

Then for the maximal inequality, we first apply the union bound to (EC.11).

$$\begin{aligned} \mathbb{P} \left( \max_{1 \leq i \leq j \leq d} \left| \frac{1}{\tau} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right| \geq w + \sqrt{2w} \right) &\leq \sum_{1 \leq i \leq j \leq d} 2\mathbb{P} \left( \frac{1}{\tau} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \geq w + \sqrt{2w} \right) \\ &\leq 2d^2 \exp \left\{ -\frac{\tau w}{2} \right\} \\ &= \exp \left\{ -\frac{\tau w}{2} + \log(2d^2) \right\}. \end{aligned}$$

Then,

$$\mathbb{P} \left( \max_{1 \leq i \leq j \leq d} \left| \frac{1}{\tau} \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right| \geq w + \sqrt{2w} + \sqrt{\frac{4 \log(2d^2)}{\tau} + \frac{2 \log(2d^2)}{\tau}} \right)$$



$$\begin{aligned}
&\leq \mathbb{P} \left( \max_{1 \leq i \leq j \leq d} \frac{1}{\tau} \left| \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right| \geq \left( w + \frac{2\log(2d^2)}{\tau} \right) + \sqrt{2 \left( w + \frac{2\log(2d^2)}{\tau} \right)} \right) \\
&= \mathbb{P} \left( \max_{1 \leq i \leq j \leq d} \frac{1}{\tau} \left| \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right| \geq w' + \sqrt{2w'} \right) \\
&\leq \exp \left\{ -\frac{\tau w'}{2} + \log(2d^2) \right\} \\
&= \exp \left\{ -\frac{\tau w}{2} \right\}
\end{aligned}$$

where  $w' = w + \frac{2\log(2d^2)}{\tau}$ .

### EC.1.5. Proof of Lemma 3

Notice the difference between the unconditional theoretical Gram matrix  $\Sigma$  and its adapted version  $\mathbb{E}[X_t X_t^\top | \mathcal{F}_{t-1}]$  which is a conditional covariance matrix conditioned on the history  $\mathcal{F}_{t-1}$ . Recall that from Algorithm 1, in each round  $t$  we choose  $X_t$  given the history  $\mathcal{F}_{t-1}$ . More precisely, we compute  $\hat{\beta}_t$  based on  $\mathcal{F}_{t-1}$  and choose  $X_t$  which maximizes the product  $X_t^\top \hat{\beta}_t$ , i.e.,  $\arg \max_{X \in \mathcal{X}_t} X^\top \hat{\beta}_t$  where  $\mathcal{X}_t = \{X_{t,1}, X_{t,2}\}$ . Hence, we can write  $\mathbb{E}[X_t X_t^\top | \mathcal{F}_{t-1}]$  as the following:

$$\mathbb{E}[X_t X_t^\top | \mathcal{F}_{t-1}] = \sum_{i=1}^2 \mathbb{E}_{\mathcal{X}_t} \left[ X_{\tau,i} X_{\tau,i}^\top \mathbb{1} \{ X_{\tau,i} = \arg \max_{X \in \mathcal{X}_t} X^\top \hat{\beta}_t \} \mid \hat{\beta}_t \right].$$

From Lemma 2, it follows that

$$\mathbb{E}[X_t X_t^\top | \mathcal{F}_{t-1}] \succcurlyeq \frac{\Sigma}{\rho_0}.$$

Now, taking an average over  $t$  gives,

$$\Sigma_\tau = \frac{1}{\tau} \sum_{t=1}^{\tau} \mathbb{E}[X_t X_t^\top | \mathcal{F}_{t-1}] \succcurlyeq \frac{\Sigma}{\rho_0}.$$

Then, we define  $\tilde{\beta}$  corresponding to compatibility constant  $\phi^2(\Sigma_\tau, S_0)$ , that is,

$$\tilde{\beta} := \arg \min_{\beta} \left\{ \frac{\beta^\top \Sigma_\tau \beta}{\|\beta_{S_0}\|_1^2} : \|\beta_{S_0^c}\|_1 \leq 3\|\beta_{S_0}\|_1 \neq 0 \right\}.$$

Therefore, it follows that

$$\frac{\tilde{\beta}^\top \Sigma_\tau \tilde{\beta}}{\|\tilde{\beta}_{S_0}\|_1^2} \geq \frac{\tilde{\beta}^\top \Sigma \tilde{\beta}}{\rho_0 \|\tilde{\beta}_{S_0}\|_1^2} \geq \frac{\phi_0^2}{\rho_0} \quad (\text{EC.12})$$

where the second inequality is by the compatibility condition on  $\Sigma$ . Thus,  $\Sigma_\tau$  satisfies the compatibility condition with compatibility constant  $\phi_{\Sigma_\tau}^2 := \frac{\phi_0^2}{\rho_0}$ .

Now, noting that  $\|\Sigma_\tau - \hat{\Sigma}_\tau\|_\infty = \max_{1 \leq i \leq j \leq d} \frac{1}{\tau} \left| \sum_{t=1}^{\tau} \gamma_t^{ij}(X_t) \right|$  for  $\gamma_t^{ij}(\cdot)$  defined in (EC.10), we can use a Bernstein-type inequality for adapted samples in Lemma EC.5 to get

$$\mathbb{P} \left( \|\Sigma_\tau - \hat{\Sigma}_\tau\|_\infty \geq w + \sqrt{2w} + \sqrt{\frac{4\log(2d^2)}{\tau}} + \frac{2\log(2d^2)}{\tau} \right) \leq \exp \left( -\frac{\tau w}{2} \right)$$

For  $\tau \geq \frac{2\log(2d^2)}{C(\phi_0, s_0)^2}$  where  $C(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_0}\right)$ , letting  $w = C(\phi_0, s_0)^2$  gives

$$\begin{aligned} w + \sqrt{2w} + \sqrt{\frac{4\log(2d^2)}{\tau}} + \frac{2\log(2d^2)}{\tau} &\leq 2\left(C(\phi_0, s_0)^2 + \sqrt{2}C(\phi_0, s_0)\right) \\ &\leq 4C(\phi_0, s_0) \\ &\leq \frac{\phi_0^2}{32s_0\rho_0} \\ &= \frac{\phi_{\Sigma_\tau}^2}{32s_0} \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{P}\left(\|\Sigma_\tau - \hat{\Sigma}_\tau\|_\infty \geq \frac{\phi_{\Sigma_\tau}^2}{32s_0}\right) &\leq \mathbb{P}\left(\|\Sigma_\tau - \hat{\Sigma}_\tau\|_\infty \geq w + \sqrt{2w} + \sqrt{\frac{4\log(2d^2)}{\tau}} + \frac{2\log(2d^2)}{\tau}\right) \\ &\leq \exp\left(-\frac{\tau w}{2}\right) \\ &= \exp\left\{-\frac{\tau C(\phi_0, s_0)^2}{2}\right\}. \end{aligned}$$

Then we can use Corollary 1 (Bühlmann and Van De Geer (2011), Corollary 6.8) to show that the empirical Gram matrix  $\hat{\Sigma}_\tau$  satisfies the compatibility condition as long as  $\Sigma_\tau$  satisfies the compatibility condition. From (EC.12), we know  $\Sigma_\tau$  satisfies the compatibility condition with compatibility constant  $\frac{\phi_0^2}{\rho_0}$ . In particular, combining Lemma 3 and Corollary 1, it follows that given  $\|\Sigma_t - \hat{\Sigma}_t\|_\infty \leq \frac{\phi_0^2}{32s_0\rho_0}$  for  $t \geq \lceil T_0 \rceil$ , we have

$$\phi^2(\hat{\Sigma}_t, S_0) \geq \frac{\phi^2(\Sigma_t, S_0)}{2} \geq \frac{\phi_0^2}{2\rho_0} > 0.$$

That is,  $\hat{\Sigma}_\tau$  satisfies the compatibility condition with compatibility constant which is at least  $\frac{\phi_0^2}{2\rho_0} > 0$ .

## EC.2. Proof of Theorem 1

First, let  $T_0 := \frac{2\log(2d^2)}{C(\phi_0, s_0)^2}$  where  $C(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_0}\right)$ . Also, we define the high probability event  $\mathcal{E}_t = \left\{\|\Sigma_t - \hat{\Sigma}_t\|_\infty \geq \frac{\phi_0^2}{32s_0\rho_0}\right\}$ . Hence, on this event  $\mathcal{E}_t$ , if  $t \geq T_0$ , then  $\phi_t^2 \geq \frac{\phi_0^2}{2\rho_0}$ , i.e., the compatibility condition holds in round  $t$ . Slightly overloading the subscript for brevity, let  $X_t := X_{t, a_t}$  be a feature of the arm chosen in round  $t$  and  $X_{a_t^*} := X_{t, a_t^*}$  be the feature of the optimal arm in round  $t$ . Now, we look at the (non-expected) immediate regret  $\text{reg}(t)$  with  $\mathcal{R}(t) = \mathbb{E}[\text{reg}(t)]$  at round  $t$ :

$$\begin{aligned} \text{reg}(t) &\leq \mathbb{1}(t \leq T_0) + \text{reg}(t)\mathbb{1}(t > T_0, \mathcal{E}_t) + \mathbb{1}(t > T_0, \mathcal{E}_t^c) \\ &= \mathbb{1}(t \leq T_0) + \text{reg}(t)\mathbb{1}\left(\mu(X_t^\top \hat{\beta}_t) \geq \mu(X_{a_t^*}^\top \hat{\beta}_t), t > T_0, \mathcal{E}_t\right) + \mathbb{1}(t > T_0, \mathcal{E}_t^c) \end{aligned}$$

Note that

$$\begin{aligned}
\mathbb{P}\left(\mu(X_t^\top \hat{\beta}_t) \geq \mu(X_{a_t^*}^\top \hat{\beta}_t)\right) &= \mathbb{P}\left(\mu(X_t^\top \hat{\beta}_t) - \mu(X_{a_t^*}^\top \hat{\beta}_t) + \text{reg}(t) \geq \text{reg}(t)\right) \\
&= \mathbb{P}\left((\mu(X_t^\top \hat{\beta}_t) - \mu(X_t^\top \beta^*)) - (\mu(X_{a_t^*}^\top \hat{\beta}_t) - \mu(X_{a_t^*}^\top \beta^*)) \geq \text{reg}(t)\right) \\
&\leq \mathbb{P}\left(|\mu(X_t^\top \hat{\beta}_t) - \mu(X_t^\top \beta^*)| + |\mu(X_{a_t^*}^\top \hat{\beta}_t) - \mu(X_{a_t^*}^\top \beta^*)| \geq \text{reg}(t)\right) \\
&\leq \mathbb{P}\left(\kappa_{\max} \|\hat{\beta}_t - \beta^*\|_1 \|X_t\|_\infty + \kappa_{\max} \|\hat{\beta}_t - \beta^*\|_1 \|X_{a_t^*}\|_\infty \geq \text{reg}(t)\right) \\
&\leq \mathbb{P}\left(2\kappa_{\max} \|\hat{\beta}_t - \beta^*\|_1 \geq \text{reg}(t)\right).
\end{aligned}$$

For an arbitrary constant  $d_t > 0$ , we can continue with  $\mathcal{R}(t) = \mathbb{E}[\text{reg}(t)]$  for  $t > T_0$ .

$$\begin{aligned}
\mathcal{R}(t) &\leq \mathbb{E}\left[\text{reg}(t) \mathbb{1}\left(2\kappa_{\max} \|\hat{\beta}_t - \beta^*\|_1 \geq \text{reg}(t), \mathcal{E}_t\right)\right] + \mathbb{P}(\mathcal{E}_t^c) \\
&= \mathbb{E}\left[\text{reg}(t) \mathbb{1}\left(2\kappa_{\max} \|\hat{\beta}_t - \beta^*\|_1 \geq \text{reg}(t), \text{reg}(t) \leq \kappa_{\max} d_t, \mathcal{E}_t\right)\right] \\
&\quad + \mathbb{E}\left[\text{reg}(t) \mathbb{1}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq \frac{\text{reg}(t)}{\kappa_{\max}}, \frac{\text{reg}(t)}{\kappa_{\max}} > d_t, \mathcal{E}_t\right)\right] + \mathbb{P}(\mathcal{E}_t^c) \\
&\leq \kappa_{\max} d_t + \kappa_{\max} \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) + \mathbb{P}(\mathcal{E}_t^c).
\end{aligned}$$

Summing over  $T$  rounds, we have

$$\sum_{t=1}^T \mathcal{R}(t) \leq T_0 + \kappa_{\max} \sum_{t=\lceil T_0 \rceil}^T d_t + \kappa_{\max} \sum_{t=\lceil T_0 \rceil}^T \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) + \sum_{t=\lceil T_0 \rceil}^T \mathbb{P}(\mathcal{E}_t^c).$$

Let  $d_t := \frac{2s_0\lambda_t}{\kappa_{\min}\phi_t^2} = \frac{4\sigma s_0}{\kappa_{\min}\phi_t^2} \sqrt{\frac{4\log t + 2\log d}{t}}$ . From Lemma 1, we have

$$\mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) \leq \frac{2}{t^2}.$$

for all  $t$ . Therefore, it follows that

$$\sum_{t=\lceil T_0 \rceil}^T \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) \leq \sum_{t=1}^T \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) \leq \frac{\pi^2}{3}.$$

For  $t \geq T_0$ , we have  $\phi_t^2 \geq \frac{\phi_0^2}{2\rho_0}$  provided that event  $\mathcal{E}_t$  holds. Hence, we have

$$\begin{aligned}
\sum_{t=\lceil T_0 \rceil}^T d_t &= \sum_{t=\lceil T_0 \rceil}^T \frac{4\sigma s_0}{\kappa_{\min}\phi_t^2} \sqrt{\frac{4\log t + 2\log d}{t}} \\
&\leq \sum_{t=\lceil T_0 \rceil}^T \frac{8\rho_0\sigma s_0}{\kappa_{\min}\phi_0^2} \sqrt{\frac{4\log t + 2\log d}{t}} \\
&\leq \frac{8\rho_0\sigma s_0\sqrt{4\log T + 2\log d}}{\kappa_{\min}\phi_0^2} \sum_{t=\lceil T_0 \rceil}^T \frac{1}{\sqrt{t}} \\
&\leq \frac{8\rho_0\sigma s_0\sqrt{4\log T + 2\log d}}{\kappa_{\min}\phi_0^2} \sum_{t=1}^T \frac{1}{\sqrt{t}} \\
&\leq \frac{16\rho_0\sigma s_0\sqrt{4\log T + 2\log d}}{\kappa_{\min}\phi_0^2} \sqrt{T}
\end{aligned} \tag{EC.13}$$

where the last inequality is from the fact that  $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq \int_{t=0}^T \frac{1}{\sqrt{t}} = 2\sqrt{T}$ .

Finally, for  $\sum_{t=T_0}^T \mathbb{P}(\mathcal{E}_t^c)$  we have from Lemma 3:

$$\begin{aligned} \sum_{t=\lceil T_0 \rceil}^T \mathbb{P}(\mathcal{E}_t^c) &\leq \sum_{t=\lceil T_0 \rceil}^T \mathbb{P}\left(\|\Sigma_t - \hat{\Sigma}_t\|_\infty \geq \frac{\phi_0^2}{32s_0\rho_0}\right) \\ &\leq \sum_{t=\lceil T_0 \rceil}^T \exp\left\{-\frac{tC(\phi_0, s_0)^2}{2}\right\} \\ &\leq \sum_{t=1}^{\infty} \exp\left\{-\frac{tC(\phi_0, s_0)^2}{2}\right\} \\ &\leq \frac{2}{C(\phi_0, s_0)^2}. \end{aligned}$$

### EC.3. Proof of Theorem 2

The proof follows similar arguments as the proof of Theorem 1. The key difference is that the RE condition involves  $\ell_2$  norm and therefore the analysis requires the Lasso oracle inequality of GLM in  $\ell_2$  norm, which we provide as an extension of Lemma 1.

**COROLLARY EC.1.** *Assume that the RE condition holds for  $\hat{\Sigma}_t$  with active set  $S_0$  and restricted eigenvalue  $\phi_t$ . For some  $\delta \in (0, 1)$ , let the regularization parameter  $\lambda_t$  be*

$$\lambda_t := \Sigma \sqrt{\frac{2[\log(2/\delta) + \log d]}{t}}.$$

*Then with probability at least  $1 - \delta$ , we have*

$$\|\hat{\beta}_t - \beta^*\|_2 \leq \frac{3\sqrt{s_0}\lambda_t}{\kappa_{\min}\phi_t^2}.$$

Continuing from (EC.7) in Lemma 1, the RE condition can be applied to the vector  $\hat{\beta} - \beta^*$  which gives

$$\|\hat{\beta} - \beta^*\|_2^2 \leq \frac{(\hat{\beta} - \beta^*)^\top \hat{\Sigma}_t (\hat{\beta} - \beta^*)}{\phi_t^2}. \quad (\text{EC.14})$$

Again from (EC.7), we can use the margin condition in Lemma EC.2

$$\begin{aligned} 3\lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 &\geq 2\mathcal{E}(\hat{\beta}_n) \\ &\geq \kappa_{\min}(\hat{\beta} - \beta^*)^\top \hat{\Sigma}_t (\hat{\beta} - \beta^*) \\ &\geq \kappa_{\min}\phi_t^2 \|\hat{\beta} - \beta^*\|_2^2 \end{aligned}$$

where the last inequality is from (EC.14) applying the RE condition. Then, it follows that

$$\kappa_{\min}\phi_t^2 \|\hat{\beta} - \beta^*\|_2^2 \leq 3\lambda_t \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_1 \leq 3\lambda_t \sqrt{s_0} \|\hat{\beta}_{S_0} - \beta_{S_0}^*\|_2 \leq 3\lambda_t \sqrt{s_0} \|\hat{\beta} - \beta^*\|_2.$$

Hence, dividing the both sides by  $\|\hat{\beta} - \beta^*\|_2$  and rearranging gives

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{3\sqrt{s_0}\lambda_t}{\kappa_{\min}\phi_t^2}.$$

This complete the proof.

### EC.3.1. Ensuring the RE Condition for the Empirical Gram Matrix

To distinguish from the compatibility constant, we introduce the definition of a generic restricted eigenvalue as  $\phi_{\text{RE}}^2(\Sigma, S)$ .

DEFINITION EC.2. The restricted eigenvalue of  $M$  over  $S_0$  is

$$\phi_{\text{RE}}^2(M, S_0) := \min_{\beta} \left\{ \frac{\beta^\top M \beta}{\|\beta\|_2^2} : \|\beta_{S_0^c}\|_1 \leq 3\|\beta_{S_0}\|_1 \neq 0 \right\}.$$

Note that Assumption 4 only provides the RE condition for the theoretical Gram matrix  $\Sigma$ . Then, we follow the same arguments as in the analysis under the compatibility condition to show that  $\phi_{\text{RE}}^2(\Sigma_t, S_0) \geq \frac{\phi_{\text{RE}}^2(\Sigma, S_0)}{\rho_0} > 0$ , i.e.,  $\Sigma_t$  satisfies the RE condition. Then using Lemma 3, we can show that  $\hat{\Sigma}_t$  concentrates to  $\Sigma_t$  with a high probability. The following lemma (similar to Corollary 1) ensures the RE condition of  $\hat{\Sigma}_t$  conditioned on the matrix concentration of the empirical Gram matrix  $\hat{\Sigma}_t$ .

LEMMA EC.6. Suppose that the RE condition holds for  $\Sigma_0$  and the index set  $S$  with cardinality  $s = |S|$ , with restricted eigenvalue  $\phi_{\text{RE}}^2(\Sigma_0, S) > 0$ , and that  $\|\Sigma_1 - \Sigma_0\|_\infty \leq \Delta$ , where  $32s\Delta \leq \phi_{\text{RE}}^2(\Sigma_0, S)$ . Then, for the set  $S$ , the RE condition holds as well for  $\Sigma_1$ , with  $\phi_{\text{RE}}^2(\Sigma_1, S) \geq \phi_{\text{RE}}^2(\Sigma_0, S)/2$ .

The proof is an adaptation of Lemma 6.17 in [Bühlmann and Van De Geer \(2011\)](#) to the RE condition.

$$\begin{aligned} |\beta^\top \Sigma_1 \beta - \beta^\top \Sigma_0 \beta| &= |\beta^\top (\Sigma_1 - \Sigma_0) \beta| \\ &\leq \|\Sigma_1 - \Sigma_0\|_\infty \|\beta\|_1^2 \\ &\leq \Delta \|\beta\|_1^2 \end{aligned}$$

For  $\beta$  such that  $\|\beta_{S^c}\| \leq 3\|\beta_S\|$ , we have the RE condition satisfied for  $\Sigma_0$ . Hence, we have

$$\|\beta\|_1 \leq 4\|\beta_S\|_1 \leq 4\sqrt{s}\|\beta_S\|_2 \leq 4\sqrt{s}\|\beta\|_2 \leq \frac{4\sqrt{s_0\beta^\top \Sigma_0 \beta}}{\phi_{\text{RE}}(\Sigma_0, S)}.$$

Therefore, it follows that

$$|\beta^\top \Sigma_1 \beta - \beta^\top \Sigma_0 \beta| \leq \frac{16s\Delta\beta^\top \Sigma_0 \beta}{\phi_{\text{RE}}^2(\Sigma_0, S)}.$$

Since  $\beta^\top \Sigma_0 \beta > 0$ , dividing the both sides by  $\beta^\top \Sigma_0 \beta$  gives

$$\left| \frac{\beta^\top \Sigma_1 \beta}{\beta^\top \Sigma_0 \beta} - 1 \right| \leq \frac{16s\Delta}{\phi_{\text{RE}}^2(\Sigma_0, S)}$$

Now, since  $32s\Delta \leq \phi_{\text{RE}}^2(\Sigma_0, S)$ , it follows that

$$\frac{1}{2} \cdot \frac{\beta^\top \Sigma_0 \beta}{\|\beta\|_2^2} \leq \frac{\beta^\top \Sigma_1 \beta}{\|\beta\|_2^2} \leq \frac{3}{2} \cdot \frac{\beta^\top \Sigma_0 \beta}{\|\beta\|_2^2}.$$

Hence,

$$\phi_{\text{RE}}^2(\Sigma_1, S) \geq \frac{\phi_{\text{RE}}^2(\Sigma_0, S)}{2}$$

### EC.3.2. Proof of Theorem 2

The proof of Theorem 2 follows the similar arguments as the proof of Theorem 1. The only difference is that we use  $\ell_2$  error bound  $\|\hat{\beta}_t - \beta^*\|_2$  instead of  $\|\hat{\beta}_t - \beta^*\|_1$ . First, note that

$$\begin{aligned} \mathbb{P}\left(\mu(X_t^\top \hat{\beta}_t) \geq \mu(X_{a_t^*}^\top \hat{\beta}_t)\right) &\leq \mathbb{P}\left(|\mu(X_t^\top \hat{\beta}_t) - \mu(X_t^\top \beta^*)| + |\mu(X_{a_t^*}^\top \hat{\beta}_t) - \mu(X_{a_t^*}^\top \beta^*)| \geq \text{reg}(t)\right) \\ &\leq \mathbb{P}\left(\kappa_{\max} \|\hat{\beta}_t - \beta^*\|_2 \|X_t\|_2 + \kappa_{\max} \|\hat{\beta}_t - \beta^*\|_2 \|X_{a_t^*}\|_2 \geq \text{reg}(t)\right) \\ &\leq \mathbb{P}\left(2\kappa_{\max} \|\hat{\beta}_t - \beta^*\|_2 \geq \text{reg}(t)\right). \end{aligned}$$

For an arbitrary constant  $d_t > 0$ , we can continue with  $\mathbb{E}[\text{reg}(t)]$  for  $t > T_0$ .

$$\mathcal{R}(t) \leq \kappa_{\max} d_t + \kappa_{\max} \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_2 \geq d_t, \mathcal{E}_t\right) + \mathbb{P}(\mathcal{E}_t^c).$$

Hence, the cumulative regret would give

$$\sum_{t=1}^T \mathcal{R}(t) \leq T_0 + \kappa_{\max} \sum_{t=\lceil T_0 \rceil}^T d_t + \kappa_{\max} \sum_{t=\lceil T_0 \rceil}^T \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_2 \geq d_t, \mathcal{E}_t\right) + \sum_{t=\lceil T_0 \rceil}^T \mathbb{P}(\mathcal{E}_t^c).$$

Let  $d_t := \frac{3\sqrt{s_0}\lambda_t}{2\kappa_{\min}\phi_t^2} = \frac{3\sigma}{\kappa_{\min}\phi_t^2} \sqrt{\frac{s_0(4\log t + 2\log d)}{t}}$ . From Lemma 1, we have

$$\mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) \leq \frac{2}{t^2}.$$

for all  $t$ . Therefore, it follows that

$$\sum_{t=\lceil T_0 \rceil}^T \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) \leq \sum_{t=1}^T \mathbb{P}\left(2\|\hat{\beta}_t - \beta^*\|_1 \geq d_t, \mathcal{E}_t\right) \leq \frac{\pi^2}{3}.$$

For  $t \geq T_0$ , we have  $\phi_t^2 \geq \frac{\phi_1^2}{2\rho_0}$  provided that event  $\mathcal{E}_t$  holds. Hence, we have

$$\begin{aligned} \sum_{t=\lceil T_0 \rceil}^T d_t &= \sum_{t=\lceil T_0 \rceil}^T \frac{3\sigma s_0}{\kappa_{\min}\phi_t^2} \sqrt{\frac{4\log t + 2\log d}{t}} \\ &\leq \sum_{t=\lceil T_0 \rceil}^T \frac{6\rho_0\sigma s_0}{\kappa_{\min}\phi_1^2} \sqrt{\frac{4\log t + 2\log d}{t}} \\ &\leq \frac{6\rho_0\sigma s_0 \sqrt{4\log T + 2\log d}}{\kappa_{\min}\phi_1^2} \sum_{t=1}^T \frac{1}{\sqrt{t}} \\ &\leq \frac{12\rho_0\sigma s_0 \sqrt{4\log T + 2\log d}}{\kappa_{\min}\phi_1^2} \sqrt{T} \end{aligned}$$

where the last inequality is from the fact that  $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq \int_{t=0}^T \frac{1}{\sqrt{t}} = 2\sqrt{T}$ . Combining all the results with bounds on  $T_0$  and  $\sum_{t=\lceil T_0 \rceil}^T \mathbb{P}(\mathcal{E}_t^c)$  from Theorem 1, the expected regret under the RE condition is bounded by

$$\mathcal{R}_T \leq \frac{\kappa_{\max}\pi^2}{3} + \frac{2\log(2d^2) + 2}{C_1(\phi_1, s_0)^2} + \frac{24\kappa_{\max}\rho_0\sigma\sqrt{s_0 T \log(dT)}}{\kappa_{\min}\phi_1^2}$$

where  $C_1(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_1^2}{128s_0\rho_0}\right)$ .

## EC.4. Regret Analysis for $K$ -Armed Case

As discussed in Section 7, the analysis for the  $K$ -armed bandit mostly follows the proof of the two-armed bandit analysis in Section 5. Assuming the compatibility condition of the empirical Gram matrix  $\hat{\Sigma}_t$ , the Lasso oracle inequality for adapted samples in Lemma 1 can be directly applied. Hence, what we have left is ensuring the compatibility condition of  $\hat{\Sigma}_t$ . As before, for each  $\mathbb{E}[X_\tau X_\tau^\top | \mathcal{F}_\tau]$  in  $\Sigma_t$ , the history  $\mathcal{F}_\tau$  affects how feature vector  $X_\tau$  is chosen. Similar to the two-armed bandit case, we rewrite  $\Sigma_t$  as

$$\Sigma_t = \frac{1}{t} \sum_{\tau=1}^t \sum_{i=1}^K \mathbb{E}[X_{\tau,i} X_{\tau,i}^\top \mathbb{1}\{X_{\tau,i} = \arg \max_{X \in \mathcal{X}_\tau} X^\top \hat{\beta}_\tau\} | \hat{\beta}_\tau].$$

Recall that the compatibility condition is only assumed for the theoretical Gram matrix  $\Sigma$  (Assumption 2). Again, the adapted Gram matrix  $\Sigma_t$  is used to bridge  $\Sigma$  and  $\hat{\Sigma}_t$  to ensure the compatibility of  $\hat{\Sigma}_t$ . The key difference between the two-armed bandit analysis and the  $K$ -armed bandit analysis lies in how  $\Sigma_t$  is controlled by  $\Sigma$ . In particular, under the balanced covariance condition in Assumption 5, we show the following lemma which is a generalization of Lemma 2.

LEMMA EC.7. *Suppose Assumption 5 holds. For a fixed vector  $\beta \in \mathbb{R}^d$ , we have*

$$\sum_{i=1}^K \mathbb{E}[X_{t,i} X_{t,i}^\top \mathbb{1}\{X_{t,i} = \arg \max_{X \in \mathcal{X}_t} X^\top \beta\}] \succcurlyeq \frac{\Sigma}{2\rho_0 C_{\mathcal{X}}}.$$

With this result, we can lower-bound the compatibility constant  $\phi^2(\Sigma_t, S_0)$  of the adapted Gram matrix in terms of the compatibility constant  $\phi^2(\Sigma, S_0)$  for the theoretical Gram matrix. That is, we have  $\Sigma_t \succcurlyeq \frac{\Sigma}{2\rho_0 C_{\mathcal{X}}}$  which implies that  $\phi^2(\Sigma_t, S_0) \geq \frac{\phi^2(\Sigma, S_0)}{2\rho_0 C_{\mathcal{X}}} > 0$ . Hence,  $\Sigma_t$  satisfies the compatibility condition.

Then, we can show that  $\hat{\Sigma}_t$  concentrates to  $\Sigma_t$  with a high probability which directly follows from applying Lemma 2, which is formally stated as follows.

COROLLARY EC.2. *For  $t \geq \frac{2\log(2d^2)}{C(\phi_0, s_0)^2}$  where  $C(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_0 C_{\mathcal{X}}}\right)$ , we have*

$$\mathbb{P}\left(\|\Sigma_t - \hat{\Sigma}_t\|_\infty \geq \frac{\phi_0^2}{32s_0\rho_0 C_{\mathcal{X}}}\right) \leq \exp\left\{-\frac{tC(\phi_0, s_0)^2}{2}\right\}.$$

Now, we can invoke Corollary 1 to connect this matrix concentration result to guaranteeing the compatibility condition of  $\hat{\Sigma}_t$ . Therefore,  $\hat{\Sigma}_t$  satisfies the compatibility condition with compatibility constant  $\phi_t^2 = \frac{\phi_0^2}{4\rho_0 C_{\mathcal{X}}} > 0$ . The rest of the proof of Theorem 3 directly follows the proof of Theorem 1 using this compatibility constant.

**EC.4.1. Proof of Lemma EC.7**

Suppose there are  $K$  variables  $\mathcal{X} = \{X_1, \dots, X_K\}$ . Let joint distribution of  $\mathcal{X}$  as  $p_{\mathcal{X}}(x_1, \dots, x_K) = p_{\mathcal{X}}(\mathbf{x})$  where we let  $\mathbf{x} = (x_1, \dots, x_K)$ . Then the theoretical Gram matrix is defined as

$$\begin{aligned}\mathbb{E}[\mathbf{X}^\top \mathbf{X}] &= \mathbb{E} \left[ \sum_{i=1}^K X_i X_i^\top \right] \\ &= \int (x_1 x_1^\top + \dots + x_K x_K^\top) p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x}\end{aligned}$$

Let's first focus on  $\int x_1 x_1^\top p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x}$ .

$$\begin{aligned}\int x_1 x_1^\top p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} &= \int x_1 x_1^\top \mathbb{1} \left\{ x_1 = \arg \max_{x_i \in \mathcal{X}} x_i^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} + \int x_1 x_1^\top \mathbb{1} \left\{ x_1 = \arg \min_{x_i \in \mathcal{X}} x_i^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} \\ &\quad + \int x_1 x_1^\top \mathbb{1} \left\{ x_1 \neq \arg \max_{x_i \in \mathcal{X}} x_i^\top \beta, x_1 \neq \arg \min_{x_i \in \mathcal{X}} x_i^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x}.\end{aligned}$$

We define three disjoint sets of possible orderings for  $\{1, \dots, K\}$  as follows.

DEFINITION EC.3. We define the following sets of permutations of  $(1, \dots, K)$ .

$$\begin{aligned}\mathcal{I}_1^{\max} &:= \{\text{indices } (i_1, \dots, i_K) \text{ such that } i_K = 1\} \\ \mathcal{I}_1^{\min} &:= \{\text{indices } (i_1, \dots, i_K) \text{ such that } i_1 = 1\} \\ \mathcal{I}_1^{\text{mid}} &:= \{\text{indices } (i_1, \dots, i_K) \text{ such that } i_1 \neq 1 \text{ and } i_K \neq 1\}.\end{aligned}$$

Then, for  $\int x_1 x_1^\top \mathbb{1} \{x_1 = \arg \min_{x_i \in \mathcal{X}} x_i^\top \beta\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x}$ , we can write

$$\int x_1 x_1^\top \mathbb{1} \left\{ x_1 = \arg \min_{x_i \in \mathcal{X}} x_i^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} = \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}} \int x_1 x_1^\top \mathbb{1} \left\{ x_{i_1}^\top \beta \leq \dots \leq x_{i_K}^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x}$$

Then for any  $(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}$ ,

$$\begin{aligned}\int x_1 x_1^\top \mathbb{1} \left\{ x_{i_1}^\top \beta \leq \dots \leq x_{i_K}^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} &= \int x_1 x_1^\top \mathbb{1} \left\{ -x_{i_1}^\top \beta \geq \dots \geq -x_{i_K}^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} \\ &\leq \rho_0 \int x_1 x_1^\top \mathbb{1} \left\{ -x_{i_1}^\top \beta \geq \dots \geq -x_{i_K}^\top \beta \right\} p_{\mathcal{X}}(-\mathbf{x}) d\mathbf{x} \\ &= \rho_0 \int x_1 x_1^\top \mathbb{1} \left\{ x_{i_1}^\top \beta \geq \dots \geq x_{i_K}^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x}\end{aligned}$$

where the inequality is again from Assumption 3. Since the elements in  $\mathcal{I}_1^{\min}$  can be considered as reversed orderings of elements in  $\mathcal{I}_1^{\max}$  (and obviously  $|\mathcal{I}_1^{\min}| = |\mathcal{I}_1^{\max}|$ ),

$$\begin{aligned}\mathbb{E} \left[ X_1 X_1^\top \mathbb{1} \left\{ X_1 = \arg \min_{X \in \mathcal{X}} X^\top \beta \right\} \right] &= \int x_1 x_1^\top \mathbb{1} \left\{ x_1 = \arg \min_{x_i \in \mathcal{X}} x_i^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} \\ &= \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}} \int x_1 x_1^\top \mathbb{1} \left\{ x_{i_1}^\top \beta \leq \dots \leq x_{i_K}^\top \beta \right\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x}\end{aligned}$$



$$\begin{aligned}
&\preceq \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}} \rho_0 \int x_1 x_1^\top \mathbb{1}\{x_{i_1}^\top \beta \geq \dots \geq x_{i_K}^\top \beta\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} \\
&= \rho_0 \int x_1 x_1^\top \mathbb{1}\{x_1 = \arg \max_{x_i \in \mathcal{X}} x_i^\top \beta\} p_{\mathcal{X}}(\mathbf{x}) d\mathbf{x} \\
&= \rho_0 \mathbb{E} \left[ X_1 X_1^\top \mathbb{1}\{X_1 = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right]
\end{aligned}$$

From Assumption 5, we have

$$\mathbb{E} [X_1 X_1^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \preceq C_{\mathcal{X}} \mathbb{E} [(X_{i_1} X_{i_1}^\top + X_{i_K} X_{i_K}^\top) \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}].$$

Then it follows that

$$\begin{aligned}
\mathbb{E} [X_1 X_1^\top] &= \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\text{mid}}} \mathbb{E} [X_1 X_1^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}} \mathbb{E} [X_1 X_1^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&\quad + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\max}} \mathbb{E} [X_1 X_1^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&= \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\text{mid}}} \mathbb{E} [X_1 X_1^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}} \mathbb{E} [X_{i_1} X_{i_1}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&\quad + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\max}} \mathbb{E} [X_{i_K} X_{i_K}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&\preceq \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\text{mid}}} C_{\mathcal{X}} \mathbb{E} [X_{i_1} X_{i_1}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\text{mid}}} C_{\mathcal{X}} \mathbb{E} [X_{i_K} X_{i_K}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&\quad + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}} \mathbb{E} [X_{i_1} X_{i_1}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\max}} \mathbb{E} [X_{i_K} X_{i_K}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&\preceq \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\text{mid}}} C_{\mathcal{X}} \mathbb{E} [X_{i_1} X_{i_1}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\text{mid}}} C_{\mathcal{X}} \mathbb{E} [X_{i_K} X_{i_K}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&\quad + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\min}} C_{\mathcal{X}} \mathbb{E} [X_{i_1} X_{i_1}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] + \sum_{(i_1, \dots, i_K) \in \mathcal{I}_1^{\max}} C_{\mathcal{X}} \mathbb{E} [X_{i_K} X_{i_K}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \\
&= C_{\mathcal{X}} \sum_{i=1}^K \mathbb{E} \left[ X_i X_i^\top \mathbb{1}\{X_i = \arg \min_{X \in \mathcal{X}} X^\top \beta\} \right] + C_{\mathcal{X}} \sum_{i=1}^K \mathbb{E} \left[ X_i X_i^\top \mathbb{1}\{X_i = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right] \\
&\preceq C_{\mathcal{X}} (1 + \rho_0) \sum_{i=1}^K \mathbb{E} \left[ X_i X_i^\top \mathbb{1}\{X_i = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right]
\end{aligned}$$

Therefore, summing  $\mathbb{E} [X_j X_j^\top]$  over all  $j = 1, \dots, K$  gives

$$\begin{aligned}
\mathbb{E} [\mathbf{X}^\top \mathbf{X}] &= \sum_{j=1}^K \mathbb{E} [X_j X_j^\top] \\
&\preceq K C_{\mathcal{X}} (1 + \rho_0) \sum_{i=1}^K \mathbb{E} \left[ X_i X_i^\top \mathbb{1}\{X_i = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right]
\end{aligned}$$

Hence,

$$\sum_{i=1}^K \mathbb{E} \left[ X_i X_i^\top \mathbb{1} \{X_i = \arg \max_{X \in \mathcal{X}} X^\top \beta\} \right] \succcurlyeq \frac{1}{C_{\mathcal{X}}(1 + \rho_0)} \frac{\mathbb{E}[\mathbf{X}^\top \mathbf{X}]}{K} \succcurlyeq \frac{\Sigma}{2C_{\mathcal{X}}\rho_0}.$$

#### EC.4.2. Proposition EC.1

PROPOSITION EC.1. *In the case of independent arms, both a multivariate Gaussian distribution and a uniform distribution on a unit sphere satisfy Assumption 5 with  $C_{\mathcal{X}} = \mathcal{O}(1)$ . For an arbitrary distribution, it holds with  $C_{\mathcal{X}} = \binom{K-1}{K_0}$  where  $K_0 = \lceil (K-1)/2 \rceil$ .*

The proof of Proposition EC.1 involves the following few technical lemmas.

LEMMA EC.8. *Suppose each  $X_i \in \mathbb{R}^d$  is i.i.d. Gaussian with mean  $\mu$  and covariance matrix  $\Gamma$ . For any permutation  $(i_1, \dots, i_K)$  of  $(1, \dots, K)$ , any integer  $k \in \{2, \dots, K-1\}$  and fixed  $\beta$ ,*

$$\begin{aligned} \mathbb{E} \left[ X_{i_k} X_{i_k}^\top \mathbb{1} \{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\} \right] &\preceq \mathbb{E} \left[ X_{i_1} X_{i_1}^\top \mathbb{1} \{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\} \right] \\ &+ \mathbb{E} \left[ X_{i_K} X_{i_K}^\top \mathbb{1} \{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\} \right]. \end{aligned}$$

It suffices to show that for any  $y \in \mathbb{R}^d$

$$\begin{aligned} &\mathbb{E} \left[ (X_{i_k}^\top y)^2 \mathbb{1} \{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\} \right] \\ &\leq \mathbb{E} \left[ (X_{i_1}^\top y)^2 \mathbb{1} \{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\} \right] + \mathbb{E} \left[ (X_{i_K}^\top y)^2 \mathbb{1} \{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\} \right]. \end{aligned}$$

Now, we can write

$$y = \tilde{\beta}(\tilde{\beta}^\top y) + \sum_{j=1}^{d-1} g_j g_j^\top y := \tilde{\beta} w_0 + \sum_{j=1}^{d-1} g_j g_j^\top y.$$

where  $w_0 = \tilde{\beta}^\top y$  and  $\tilde{\beta} = \frac{\beta}{\|\beta\|}$  and  $[\tilde{\beta}, g_1, \dots, g_{d-1}]$  form an orthonormal basis. For  $i \in [N]$ , we can write

$$\begin{aligned} X_i^\top y &= (X_i^\top \tilde{\beta}) w_0 + X_i^\top \left( \sum_{j=1}^{d-1} g_j g_j^\top \right) y \\ &= (X_i^\top \tilde{\beta}) w_0 + \left[ \left( \sum_{j=1}^{d-1} g_j g_j^\top \right) X_i \right]^\top y. \end{aligned}$$

Then we define the following two random variables

$$U_i := X_i^\top \tilde{\beta}, \quad V_i := G X_i \tag{EC.15}$$

where  $G = \sum_{j=1}^{d-1} g_j g_j^\top$ . Then we have

$$\begin{bmatrix} U_i \\ V_i \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu^\top \tilde{\beta} \\ G\mu \end{bmatrix}, \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \right)$$

where

$$\begin{aligned} A_{11} &= \tilde{\beta}^\top \Gamma \tilde{\beta} \in \mathbb{R} \\ A_{12} &= A_{21}^\top = \tilde{\beta}^\top \Gamma G^\top \in \mathbb{R}^{1 \times d} \\ A_{22} &= G \Gamma G^\top \in \mathbb{R}^{d \times d} \end{aligned}$$

Then, we know from Lemma EC.12 that the conditional distribution  $V_i \mid U_i$  of a multivariate normal distribution is also a multivariate normal distribution. In particular,

$$V_i \mid U_i = u_i \sim \mathcal{N} \left( G\mu + A_{21}A_{11}^{-1}(u_i - \mu^\top \tilde{\beta}), B \right)$$

where  $B = A_{22} - A_{21}A_{11}^{-1}A_{12}$ . Therefore, given  $U_{i_k} = u_{i_k}$ , we can write

$$\begin{aligned} X_{i_k}^\top y &= u_{i_k} w_0 + V_{i_k}^\top y \\ &= u_{i_k} w_0 + \left( G\mu + A_{21}A_{11}^{-1}(u_{i_k} - \mu^\top \tilde{\beta}) + B^{1/2}Z \right)^\top y \end{aligned}$$

where  $Z \sim \mathcal{N}(0, I_d)$  and  $Z \perp\!\!\!\perp U_{i_k}$ . Rearranging gives

$$X_{i_k}^\top y = u_{i_k} (w_0 + A_{11}^{-1}A_{12}y) + \left( G\mu - A_{21}A_{11}^{-1}\mu^\top \tilde{\beta} \right)^\top y + Z^\top B^{1/2}y$$

Hence,  $X_{i_k}^\top y$  is a linear function of  $u_{i_k}$ . Then it follows that

$$\begin{aligned} (X_{i_k}^\top y)^2 &= \left[ u_{i_k} (w_0 + A_{11}^{-1}A_{12}y) + \left( G\mu - A_{21}A_{11}^{-1}\mu^\top \tilde{\beta} \right)^\top y + Z^\top B^{1/2}y \right]^2 \\ &\leq \max \left\{ \left[ u_{i_1} (w_0 + A_{11}^{-1}A_{12}y) + \left( G\mu - A_{21}A_{11}^{-1}\mu^\top \tilde{\beta} \right)^\top y + Z^\top B^{1/2}y \right]^2, \right. \\ &\quad \left. \left[ u_{i_K} (w_0 + A_{11}^{-1}A_{12}y) + \left( G\mu - A_{21}A_{11}^{-1}\mu^\top \tilde{\beta} \right)^\top y + Z^\top B^{1/2}y \right]^2 \right\} \\ &\leq \left[ u_{i_1} (w_0 + A_{11}^{-1}A_{12}y) + \left( G\mu - A_{21}A_{11}^{-1}\mu^\top \tilde{\beta} \right)^\top y + Z^\top B^{1/2}y \right]^2 \\ &\quad + \left[ u_{i_K} (w_0 + A_{11}^{-1}A_{12}y) + \left( G\mu - A_{21}A_{11}^{-1}\mu^\top \tilde{\beta} \right)^\top y + Z^\top B^{1/2}y \right]^2. \end{aligned}$$

Therefore, it follows that

$$\begin{aligned} &\mathbb{E} \left[ (X_{i_k}^\top y)^2 \mathbb{1} \{ X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta \} \right] \\ &\leq \mathbb{E} \left[ (X_{i_1}^\top y)^2 \mathbb{1} \{ X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta \} \right] + \mathbb{E} \left[ (X_{i_K}^\top y)^2 \mathbb{1} \{ X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta \} \right] \end{aligned}$$

Hence,

$$\mathbb{E} \left[ X_{i_k} X_{i_k}^\top \mathbb{1} \{ X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta \} \right] \preceq \mathbb{E} \left[ (X_{i_1} X_{i_1}^\top + X_{i_K} X_{i_K}^\top) \mathbb{1} \{ X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta \} \right].$$

LEMMA EC.9. Suppose  $X \in \mathbb{R}^d$  is uniformly distributed on the unit sphere  $\mathcal{S}^{d-1}$  and  $K = o(d)$ . For fixed vector  $\beta \in \mathbb{R}^d$  and a given integer  $k \in \{2, \dots, K-1\}$ ,

$$\mathbb{E} [X_{i_k} X_{i_k}^\top \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}] \preceq C_{\mathcal{X}} \mathbb{E} [(X_{i_1} X_{i_1}^\top + X_{i_K} X_{i_K}^\top) \mathbb{1}\{X_{i_1}^\top \beta < \dots < X_{i_K}^\top \beta\}].$$

where  $C_{\mathcal{X}} = \mathcal{O}(1)$ .

Here, we instead show directly

$$\mathbb{E}[XX^\top] \preceq C \left( \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \max_{X_i \in \{X_1, \dots, X_K\}} X_i^\top \tilde{\beta}\} \right] + \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \min_{X_i \in \{X_1, \dots, X_K\}} X_i^\top \tilde{\beta}\} \right] \right)$$

for some constant  $C$ . It can be shown that if  $C = \mathcal{O}(1)$ , then the claim holds with  $C_{\mathcal{X}} = \mathcal{O}(1)$ . Suppose  $X \in \mathbb{R}^d$  is uniformly distributed on the unit sphere  $\mathcal{S}^{d-1} := \{s \in \mathbb{R}^d : \|s\|_2 = 1\}$ . Then by Lemma 2 in [Cambanis et al. \(1981\)](#), we can write for each  $X_i$ ,

$$X_i \sim (B_i U_{i,1}, (1 - B_i^2)^{1/2} U_{i,2})$$

where  $B_i \sim \text{beta}(\frac{1}{2}, \frac{d-1}{2})$ ,  $U_{i,1} = \pm 1$  with probability  $\frac{1}{2}$ ,  $U_{i,2} \sim \text{unif}(\mathcal{S}^{d-2})$ .  $U_{i,1}$ ,  $U_{i,2}$  and  $B_i$  are independent of each other. Similar to the analysis of the Gaussian case, we can normalize  $\beta$  so that  $\tilde{\beta} = \frac{\beta}{\|\beta\|}$ . Without loss of generality, assume that  $\tilde{\beta} = [1, 0, \dots, 0]^\top$ . That is, only the first element is non-zero. We can do this since  $X$  is spherical and rotation invariant. Then we can write

$$\mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \max_{X_i \in \{X_1, \dots, X_K\}} X_i^\top \tilde{\beta}\} \right] = \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \max_{X_i \in \{X_1, \dots, X_K\}} X_i^{(1)}\} \right]$$

where  $X_i^{(1)}$  is the first element of  $X_i$ . Similarly,

$$\mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \min_{X_i \in \{X_1, \dots, X_K\}} X_i^\top \tilde{\beta}\} \right] = \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \min_{X_i \in \{X_1, \dots, X_K\}} X_i^{(1)}\} \right].$$

Now, from the definition of  $X$ , for  $B \sim \text{beta}(\frac{1}{2}, \frac{d-1}{2})$  we have

$$X_i X_i^\top = \begin{bmatrix} B_i^2 & B_i \sqrt{1 - B_i^2} U_{i,1} U_{i,2}^\top \\ B_i \sqrt{1 - B_i^2} U_{i,1} U_{i,2} & (1 - B_i^2) U_{i,2} U_{i,2}^\top \end{bmatrix}.$$

By the independence of  $U_1, U_2$ , and  $B$ , we have

$$\mathbb{E} [XX^\top] = \mathbb{E} \begin{bmatrix} B^2 & 0 \\ 0 & \frac{1}{d-1} (1 - B^2) I_{d-1} \end{bmatrix}.$$

By the definitions of  $B_i$  and  $U_{i,1}$ , it follows that

$$\mathbb{E} \left[ XX^\top \mathbb{1}\{B = \max_{B_i \in \{B_1, \dots, B_K\}} B_i\} \right] \preceq \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \max_{X_i \in \{X_1, \dots, X_K\}} X_i^{(1)}\} \right] + \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \min_{X_i \in \{X_1, \dots, X_K\}} X_i^{(1)}\} \right]$$

Since  $\mathbb{E}[B^2] = \frac{(\alpha+1)\alpha}{(\alpha+\beta+1)(\alpha+\beta)}$  for  $B \sim \text{beta}(\alpha, \beta)$ , we have  $\mathbb{E}[B^2] = \frac{3}{d(d+2)}$  and  $\frac{1-\mathbb{E}[B^2]}{d-1} = \frac{d+3}{d(d+2)}$  using  $\alpha = \frac{1}{2}$  and  $\beta = \frac{d-1}{2}$ . Clearly,  $\lambda_{\min}(\mathbb{E}[XX^\top]) = \frac{3}{d(d+2)}$ . Similarly, for the matrix  $\mathbb{E}[XX^\top \mathbb{1}\{B = \max_i B_i\}]$ , we have

$$\mathbb{E}[XX^\top \mathbb{1}\{B = \max_i B_i\}] = \mathbb{E} \begin{bmatrix} B^2 \mathbb{1}\{B = \max_i B_i\} & 0 \\ 0 & \frac{1}{d-1}(1-B^2) \mathbb{1}\{B = \max_i B_i\} I_{d-1} \end{bmatrix}.$$

Note that  $\mathbb{E}[B^2 \mathbb{1}\{B = \max_i B_i\}] = \sum_{j=1}^K \mathbb{E}[B_j^2 \mathbb{1}\{B_j = \max_i B_i\}] \geq \mathbb{E}[B^2]$ . Then, we need to show

$$C(1 - \mathbb{E}[B^2 \mathbb{1}\{B = \max_i B_i\}]) \geq 1 - \mathbb{E}[B^2]$$

for some  $C$ . Note that  $\mathbb{E}[B^2 \mathbb{1}\{B = \max_i B_i\}] \leq N\mathbb{E}[B^2]$ . Hence, we can show

$$C \geq \frac{1 - \mathbb{E}[B^2]}{1 - N\mathbb{E}[B^2]} = \frac{1 - \frac{3}{d(d+2)}}{1 - \frac{3K}{d(d+2)}} = \frac{d^2 + d - 3}{d^2 + d - 3K}$$

Since  $K = o(d)$ , we have  $C = \mathcal{O}(1)$ . Hence,

$$\begin{aligned} \mathbb{E}[XX^\top] &\preceq C \mathbb{E} \left[ XX^\top \mathbb{1}\{B = \max_{B_i \in \{B_1, \dots, B_K\}} B_i\} \right] \\ &\preceq C \left( \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \max_{X_i \in \{X_1, \dots, X_K\}} X_i^{(1)}\} \right] + \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \min_{X_i \in \{X_1, \dots, X_K\}} X_i^{(1)}\} \right] \right) \\ &= C \left( \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \max_{X_i \in \{X_1, \dots, X_K\}} X_i^\top \tilde{\beta}\} \right] + \mathbb{E} \left[ XX^\top \mathbb{1}\{X = \arg \min_{X_i \in \{X_1, \dots, X_K\}} X_i^\top \tilde{\beta}\} \right] \right) \end{aligned}$$

which implies  $C_X = \mathcal{O}(1)$ .

LEMMA EC.10. Consider i.i.d. arbitrary distribution  $p_X$ . Fix some vector  $\beta \in \mathbb{R}^d$ . For a given integer  $k \in \{2, \dots, K-1\}$ ,

$$\mathbb{E}[X_k X_k^\top \mathbb{1}\{X_1^\top \beta < \dots < X_k^\top \beta < \dots < X_K^\top \beta\}] \preceq C_{K,k} \mathbb{E}[(X_1 X_1^\top + X_K X_K^\top) \mathbb{1}\{X_1^\top \beta < \dots < X_K^\top \beta\}]$$

where  $C_X = \binom{K-1}{(K-1)/2}$  assuming  $K$  is odd — if  $K$  is even, we can use  $\lceil (K-1)/2 \rceil$ .

First notice that

$$\begin{aligned} &\mathbb{E}[X_k X_k^\top \mathbb{1}\{X_1^\top \beta < \dots < X_k^\top \beta < \dots < X_K^\top \beta\}] \\ &= \mathbb{E}_V[V V^\top \mathbb{E}_{X_{1:K}/X_k}[\mathbb{1}\{X_1^\top \beta < \dots < X_{k-1}^\top \beta < V^\top \beta < X_{k+1}^\top \beta < \dots < X_K^\top \beta\} | V]] \end{aligned}$$

where  $X_{1:K}/X_k$  denotes  $X_1, \dots, X_{k-1}, X_{k+1}, \dots, X_K$ . Also,

$$\begin{aligned} \mathbb{E}[X_1 X_1^\top \mathbb{1}\{X_1^\top \beta < \dots < X_K^\top \beta\}] &= \mathbb{E}_V[V V^\top \mathbb{E}_{X_{2:K}}[\mathbb{1}\{V^\top \beta < X_2^\top \beta < \dots < X_K^\top \beta\} | V]] \\ \mathbb{E}[X_K X_K^\top \mathbb{1}\{X_1^\top \beta < \dots < X_K^\top \beta\}] &= \mathbb{E}_V[V V^\top \mathbb{E}_{X_{1:K-1}}[\mathbb{1}\{X_1^\top \beta < \dots < X_{K-1}^\top \beta < V^\top \beta\} | V]] \end{aligned}$$

Let  $\psi(y) := \mathbb{P}(X^\top \beta \leq y)$  denote the CDF of  $X^\top \beta$ . Then

$$\begin{aligned} & \mathbb{P}(X_1^\top \beta < \dots < X_{k-1}^\top \beta < V^\top \beta < X_{k+1}^\top \beta < \dots < X_K^\top \beta) \\ &= \prod_{i=1}^{k-1} \mathbb{P}(X_i^\top \beta \leq V^\top \beta) \frac{1}{(k-1)!} \prod_{i=k+1}^N \mathbb{P}(X_i^\top \beta \geq V^\top \beta) \frac{1}{(K-k)!} \\ &= \frac{1}{(k-1)!(K-k)!} \psi(V^\top \beta)^{k-1} (1 - \psi(V^\top \beta))^{K-k}. \end{aligned}$$

Likewise

$$\begin{aligned} \mathbb{P}(V^\top \beta < X_2^\top \beta < \dots < X_K^\top \beta) &= \frac{1}{(K-1)!} (1 - \psi(V^\top \beta))^{K-1}, \\ \mathbb{P}(X_1^\top \beta < \dots < X_{K-1}^\top \beta < V^\top \beta) &= \frac{1}{(K-1)!} \psi(V^\top \beta)^{K-1}. \end{aligned}$$

Then, we need to show there exists  $C_{K,k}$  such that

$$\begin{aligned} & \mathbb{P}(X_1^\top \beta < \dots < X_{k-1}^\top \beta < V^\top \beta < X_{k+1}^\top \beta < \dots < X_K^\top \beta) \\ & \leq C_{K,k} [\mathbb{P}(V^\top \beta < X_2^\top \beta < \dots < X_K^\top \beta) + \mathbb{P}(X_1^\top \beta < \dots < X_{K-1}^\top \beta < V^\top \beta)] \end{aligned}$$

That is,

$$\frac{1}{(k-1)!(K-k)!} \psi(V^\top \beta)^{k-1} (1 - \psi(V^\top \beta))^{K-k} \leq \frac{C_{K,k}}{(K-1)!} [(1 - \psi(V^\top \beta))^{K-1} + \psi(V^\top \beta)^{K-1}]$$

Hence,

$$C_{K,k} \geq \binom{K-1}{k-1} \frac{\psi(V^\top \beta)^{k-1} (1 - \psi(V^\top \beta))^{K-k}}{(1 - \psi(V^\top \beta))^{K-1} + \psi(V^\top \beta)^{K-1}}$$

Since  $\psi(V^\top \beta) \in [0, 1]$ , we have

$$\frac{\psi(V^\top \beta)^{k-1} (1 - \psi(V^\top \beta))^{K-k}}{(1 - \psi(V^\top \beta))^{K-1} + \psi(V^\top \beta)^{K-1}} \leq 1$$

for all  $K$  and  $k$ . Hence, for  $C_{K,k} = \binom{K-1}{k-1}$ ,

$$\mathbb{E}[X_k X_k^\top \mathbb{1}\{X_1^\top \beta < \dots < X_k^\top \beta < \dots < X_K^\top \beta\}] \preceq C_{K,k} \mathbb{E}[(X_1 X_1^\top + X_K X_K^\top) \mathbb{1}\{X_1^\top \beta < \dots < X_K^\top \beta\}].$$

## EC.5. Other lemmas

**LEMMA EC.11 (Wainwright (2019), Theorem 2.19).** *Let  $\{Z_\tau, \mathcal{F}_\tau\}_\tau^\infty$  be a martingale difference sequence, and suppose that  $Z_\tau$  is  $\sigma$ -subgaussian in an adapted sense, i.e., for all  $\alpha \in \mathbb{R}$ ,  $\mathbb{E}[e^{\alpha Z_\tau} | \mathcal{F}_{\tau-1}] \leq e^{\alpha^2 \sigma^2 / 2}$  almost surely. Then for all  $t \geq 0$ ,  $\mathbb{P}\left[\left|\sum_{\tau=1}^K Z_\tau\right| \geq \gamma\right] \leq 2 \exp[-\gamma^2 / (2n\sigma^2)]$ .*

Note that Lemma EC.12 is a well-known result, but for the sake of completeness, we present its formal statment and proof.

LEMMA EC.12. *Let  $X \in \mathbb{R}^d$  follow a multivariate Gaussian distribution with mean  $\mu$  and covariance matrix  $\Sigma$  and consider the partition of  $X$  with*

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right).$$

*Then the conditional distribution of  $X_1$  given  $X_2$  is also a multivariate Gaussian distribution. In particular*

$$X_1 \mid X_2 = x_2 \sim \mathcal{N}(\mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}).$$

Define  $Z = X_1 + \mathbf{A}X_2$  where  $\mathbf{A} = -\Sigma_{12}\Sigma_{22}^{-1}$ . Now we can write

$$\begin{aligned} \text{cov}(Z, X_2) &= \text{cov}(X_1, X_2) + \text{cov}(\mathbf{A}X_2, X_2) \\ &= \Sigma_{12} + \mathbf{A}\text{var}(X_2) \\ &= \Sigma_{12} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{22} \\ &= 0 \end{aligned}$$

Therefore  $Z$  and  $X_2$  are not correlated and, since they are jointly normal, they are independent<sup>6</sup>.

Now, clearly we have  $\mathbb{E}(Z) = \mu_1 + \mathbf{A}\mu_2$ . Then

$$\begin{aligned} \mathbb{E}(X_1|X_2) &= \mathbb{E}(Z - \mathbf{A}X_2|X_2) \\ &= \mathbb{E}(Z|X_2) - \mathbb{E}(\mathbf{A}X_2|X_2) \\ &= \mathbb{E}(Z) - \mathbf{A}X_2 \\ &= \mu_1 + \mathbf{A}(\mu_2 - X_2) \\ &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(X_2 - \mu_2). \end{aligned}$$

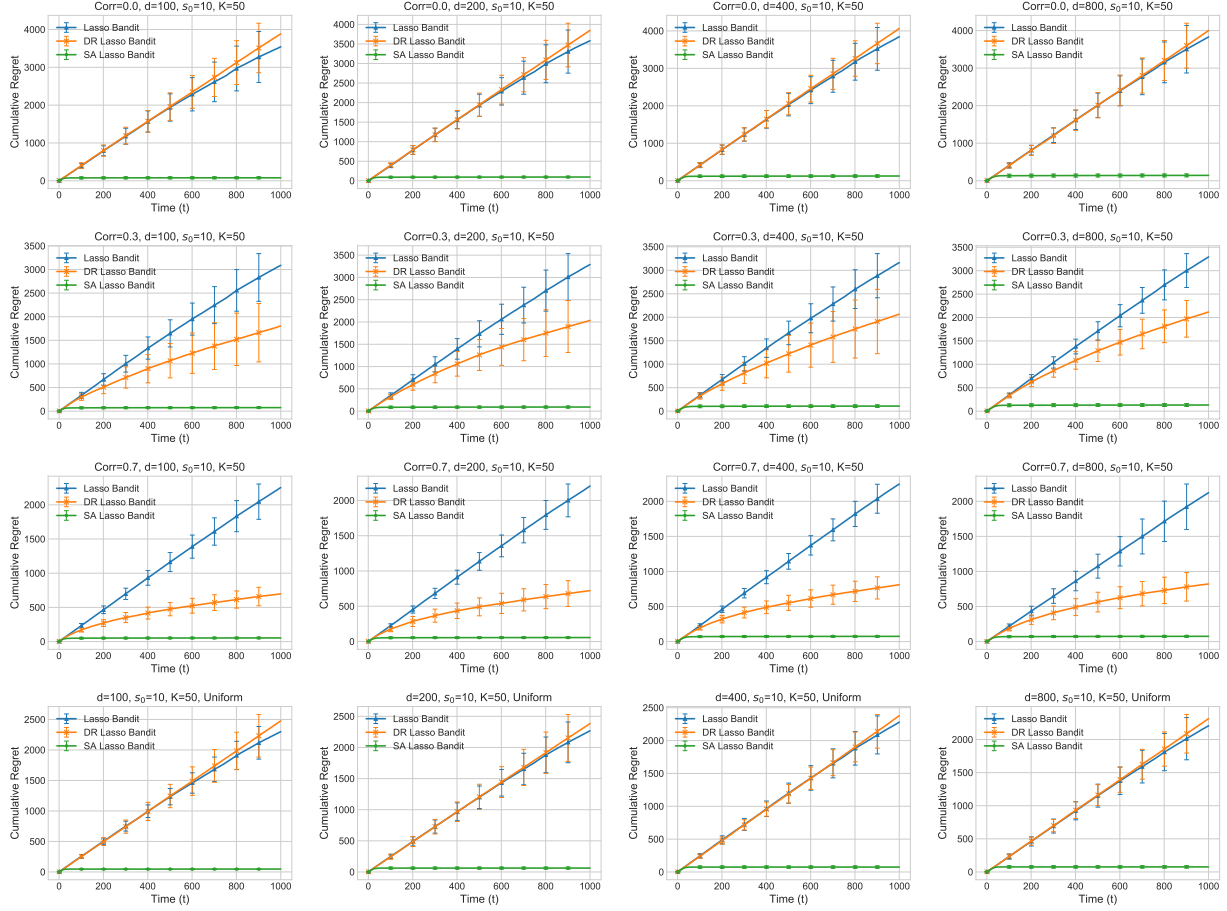
For the covariance matrix, note that

$$\begin{aligned} \text{var}(X_1|X_2) &= \text{var}(Z - \mathbf{A}X_2|X_2) \\ &= \text{var}(Z|X_2) + \text{var}(\mathbf{A}X_2|X_2) - \mathbf{A}\text{cov}(Z, -X_2) - \text{cov}(Z, -X_2)\mathbf{A}^\top \\ &= \text{var}(Z|X_2) \\ &= \text{var}(Z) \end{aligned}$$

Hence, it follows that

$$\begin{aligned} \text{var}(X_1|X_2) &= \text{var}(Z) \\ &= \text{var}(X_1 + \mathbf{A}X_2) \\ &= \text{var}(X_1) + \mathbf{A}\text{var}(X_2)\mathbf{A}^\top + \mathbf{A}\text{cov}(X_1, X_2) + \text{cov}(X_2, X_1)\mathbf{A}^\top \\ &= \Sigma_{11} + \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{22}\Sigma_{22}^{-1}\Sigma_{21} - 2\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \\ &= \Sigma_{11} + \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} - 2\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \\ &= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \end{aligned}$$

## EC.6. Additional Numerical Experiment Results



**Figure EC.1** The plots show the  $t$ -round regret of SA Lasso Bandit (Algorithm 1), DR Lasso Bandit (Kim and Paik 2019), and Lasso Bandit (Bastani and Bayati 2020) for  $K = 50$  and  $s_0 = 10$ . The first three rows are the results with features drawn from a multivariate Gaussian distribution with varying levels of correlation between arms. In the fourth row, the features are drawn from a uniform distribution on a unit sphere. For each row, we present evaluations for varying feature dimensions,  $d = 100, 200, 400, 800$ .