Min-hwan Oh¹ Garud Iyengar² Assaf Zeevi²

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 index s_0 and establish tight regret bounds on its performance under mild conditions. We also comprehensively evaluate our proposed algorithm numerically and show that it consistently outperforms existing methods, even when the correct sparsity index is revealed to them but is kept hidden from our algorithm.

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 & 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 $x_a^T \beta^*$ for a fixed and unknown parameter vector $\beta^* \in \mathbb{R}^d$. Unlike the traditional MAB problem, here pulling any one arm provides some information about the unknown parameter vector, and hence, insight into the average reward of all the other arms. These contextual bandit algorithms are applicable in a variety of problem settings, such as recommender systems, assortment selection in online retail, and healthcare analytics (Li et al., 2010; Oh & Iyengar, 2019; Tewari & Murphy, 2017), where the contextual information can be used for personalization and generalization.

In most application domains highlighted above, the feature space is high-dimensional $(d \gg 1)$, yet typically only a small subset of the features influence the expected reward. That is, the unknown parameter vector is sparse with only elements corresponding to the relevant features being non-zero, i.e., the sparsity index $s_0 = \|\beta^*\|_0 \ll d$, where the zero norm $||x||_0$ counts non-zero entries in the vector x. There is an emerging body of literature on contextual bandit problems with sparse linear reward functions (Abbasi-Yadkori et al., 2012; Gilton & Willett, 2017; Bastani & Bayati, 2020; Wang et al., 2018; Kim & Paik, 2019) which propose methods to exploit the sparse structure under various conditions. However, there is a crucial shortcoming 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 (and can lead to extremely poor performance if s_0 is underspecified). The purpose of this paper is to demonstrate that a relatively simple contextual bandit algorithm that exploits ℓ_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 general sparse bandit algorithm that does not require prior knowledge of the sparsity index s_0 .
- (b) We establish that the regret bound of our proposed algorithm is $\mathcal{O}(s_0\sqrt{T\log(dT)})$ for the two-armed case, which affords the most accessible exposition of the key analytical ideas. (Extensions to the general *K*-armed

^{*}Equal contribution ¹Seoul National University, Seoul, South Korea ²Columbia University, New York, NY, USA. Correspondence to: Min-hwan Oh <minoh@snu.ac.kr>.

Proceedings of the 38th International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

case are discussed later.) The regret bound scale in s_0 and d matches the equivalent terms in the *offline* Lasso results (see the discussions in Section 4.1).

(c) We comprehensively evaluate our algorithm on numerical experiments and show that it consistently outperforms existing methods, even when these methods are granted prior knowledge of the correct sparsity index (and can greatly outperform them if this information is misspecified).

The salient feature of our algorithm is that it does not rely on *forced sampling* which was used by almost all previous work, e.g., Bastani & Bayati (2020); Wang et al. (2018); Kim & Paik (2019), to satisfy certain regularity of the empirical Gram matrix. Forced sampling requires prior knowledge of s_0 because such schemes, the key ideas of which go back to Goldenshluger & Zeevi (2013), need to be fine-tuned using the *correct* sparsity index. (See further discussions in Section 2.4.)

2. Preliminaries

2.1. Notation

For a vector $x \in \mathbb{R}^d$, we use $||x||_1$ and $||x||_2$ to denote its ℓ_1 -norm and ℓ_2 norm respectively, the notation $||x||_0$ is reserved for the cardinality of the set of non-zero entries of that vector. We define [n] for a positive integer n to be a set containing positive integers up to n, i.e., $\{1, 2, ..., n\}$. For a real-valued function f, we use \dot{f} and \ddot{f} to denote its first and second derivatives.

2.2. 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 each 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. over $t \in [T]$ from an unknown joint distribution with probability density $p_{\mathcal{X}}$ with respect to the Lebesgue measure. 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 then 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, ... and is non-anticipating, namely each action only depends on past observations and actions.

In this work, we assume that the reward $Y_{t,i}$ of arm *i* is given by a generalized linear model (GLM), i.e.

$$Y_{t,i} = \mu(X_{t,i}^{\top}\beta^*) + \epsilon_{t,i}$$

where $\mu : \mathbb{R} \to \mathbb{R}$ (also known as *inverse link function*) is a *known* increasing function, $\beta^* \in \mathbb{R}^d$ is an *unknown* parameter, and each $\epsilon_{t,i}$ is an independent zero-mean noise. Therefore, $\mathbb{E}[Y_{t,i}|X_{t,i} = x] = \mu(x^{\top}\beta^*)$ for all $i \in [K]$ and $t \in [T]$. Widely used examples for μ are $\mu(z) = z$ which corresponds to the linear model, and $\mu(z) = 1/(1 + e^{-z})$ which corresponds to the logistic model. The parameter β^* and the feature vectors $\{X_{t,i}\}$ are potentially highdimensional, i.e., $d \gg 1$, but β^* is *sparse*, that is, the number of non-zero elements in β^* , $s_0 = \|\beta^*\|_0 \ll d$. It is important to note that the agent *does not* know s_0 or the support of the unknown parameter β^* .

We assume that there is an increasing sequence of sigma fields $\{\mathcal{F}_t\}$ such that each $\epsilon_{t,i}$ is \mathcal{F}_t -measurable with $\mathbb{E}[\epsilon_{t,i}|\mathcal{F}_{t-1}] = 0$. In our problem, \mathcal{F}_t is the sigma-field generated by random variables of chosen actions $\{a_1, ..., a_t\}$, their features $\{X_{1,a_1}, ..., X_{t,a_t}\}$, and the corresponding rewards $\{Y_{1,a_1}, ..., Y_{t,a_t}\}$. We assume the noise $\epsilon_{t,i}$ for all $i \in [K]$ is sub-Gaussian with parameter σ , where σ is a positive absolute constant, i.e., $\mathbb{E}[e^{\alpha \epsilon_{t,i}}] \leq e^{\alpha^2 \sigma^2/2}$ for all $\alpha \in \mathbb{R}$. In practice, for bounded reward $Y_{t,i}$, the noise $\epsilon_{t,i}$ is also bounded and hence satisfies the sub-Gaussian assumption with an appropriate σ 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 each round t. Then, the expected cumulative *regret* of policy $\pi = \{a_1, ..., a_T\}$ is defined as

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

Hence, maximizing the expected cumulative rewards of policy π over T rounds is equivalent to minimizing the cumulative regret $\mathcal{R}^{\pi}(T)$. Note that all the expectations and probabilities throughout the paper are with respect to feature vectors and noise unless explicitly stated otherwise.

2.3. Lasso for Generalized Linear Models

For given samples $Y_1, ..., Y_n$ and corresponding features $X_1, ..., X_n$, the Lasso (Tibshirani, 1996) estimate for the generalized linear model can be defined as

$$\hat{\beta}_n \in \operatorname*{argmin}_{\beta} \left\{ \ell_n(\beta) + \lambda \|\beta\|_1 \right\} \tag{1}$$

where $\ell_n(\beta) := -\frac{1}{n} \sum_{j=1}^n [Y_j X_j^\top \beta - m(X_j^\top \beta)]$, $m(\cdot)$ is infinitely differentiable with $\dot{m}(X^\top \beta^*) = \mathbb{E}[Y|X] = \mu(X^\top \beta^*)$, and λ 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 & Bühlmann, 2009), and the restricted isometry property (Candes & Tao, 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.

2.4. 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 & Bayati, 2020; Kim & 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 & Bayati, 2020) and (Wang et al., 2018) utilize the forced-sampling technique of (Goldenshluger & 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 vicinty 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 & 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 4. 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 .

3. Proposed Algorithm

Our proposed SPARSITY-AGNOSTIC (SA) LASSO BANDIT algorithm for high-dimensional GLM bandits is summarized in Algorithm 1. As the name suggests, our algorithm does not require prior knowledge of the sparsity index 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 λ_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 λ_0 . We show in Section 4 that $\lambda_0 = 2\sigma x_{\max}$ where x_{\max} is a bound on the ℓ_2 -norm of the feature vectors $X_{t,i}$. Thus, λ_0 does *not* depend on the sparsity index s_0 or the underlying parameter β^* . (Note that, in comparison, Kim & Paik (2019) require three tuning parameters, and Bastani & Bayati (2020) and Wang et al. (2018) require four tuning parameters, most of which are functions of the unknown sparsity index s_0 .) It is worth noting that tuning parameters, while helping to achieve low regret, are challenging to specify in online learning settings. In contrast, our proposed algorithm is practical and easy to implement.

Algorithm 1 SA LASSO BANDIT	
1: Input parameter: λ_0	
2: for all $t = 1$ to T do	
3: Observe $X_{t,i}$ for all $i \in [K]$	
4: Compute $a_t = \operatorname{argmax}_{i \in [K]} X_{t,i}^{\top} \hat{\beta}_t$	
5: Pull arm a_t and observe Y_t	
6: Update $\lambda_t \leftarrow \lambda_0 \sqrt{\frac{4\log t + 2\log d}{t}}$	
7: Update $\hat{\beta}_{t+1} \leftarrow \operatorname{argmin}_{\beta} \{ \ell_t(\beta) + \lambda_t \ \beta \ _1 \}$	
8: end for	

Discussion of the algorithm. Algorithm 1 may appear to be an *exploration-free* greedy algorithm (see e.g., Bastani et al. 2020), but this is not the case. To better understand this we will compare the steps in Algorithm 1 to upperconfidence bound (UCB) algorithms. A UCB algorithm constructs a high-probability confidence ellipsoid around a *greedy* maximum likelihood estimate and chooses a parameter value within the ellipse that maximizes the reward. Once the UCB estimate is chosen, the action selection is greedy with respect to the parameter estimate.¹ The UCB algorithms regularize parameter estimates by carefully controlling the size of the confidence ellipsoid to ensure convergence, thus, exploration is loosely equivalent to regularizing the parameter estimate. The algorithm we propose also 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 parameter associated with the sparsifying norm over time at carefully specified rate in order to ensure that our estimate is consistent as we collect more samples. (This adjustment and specification do not require knowledge of sparsity s_0 .) Incorrect choice for the penalty parameter would lead to large regret, which is analogous to poor choice of confidence widths in UCB.

4. 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 $K \ge 3$ in Section 5. It is important to note that our proposed algorithm does not change with the number of arms. We start with an assumption standard in the (generalized) linear bandit literature.

Assumption 1 (Feature set and parameter). There exists a positive constant x_{\max} such that $||x||_2 \le x_{\max}$ for all $x \in \mathcal{X}_t$ and all t, and a positive constant b such that $||\beta^*||_2 \le b$.

Assumption 2 (Link function). There exist $\kappa_0 > 0$ and $\kappa_1 < \infty$ such that the derivative $\dot{\mu}(\cdot)$ of the link function satisfies $\kappa_0 \leq \dot{\mu}(x^{\top}\beta) \leq \kappa_1$ for all x and β .

Clearly for the linear link function, $\kappa_0 = \kappa_1 = 1$. For the logistic link function, we have $\kappa_1 = 1/4$.

Definition 1 (Active set and sparsity index). The active set $S_0 := \{j : \beta_j^* \neq 0\}$ is the set of indices j for which β_j^* is non-zero, and the sparsity index $s_0 = |S_0|$ denotes the cardinality of the active set S_0 .

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, $\beta_{S_0} = [\beta_{1,S_0}, ..., \beta_{d,S_0}]^\top$ has zero elements outside the set S_0 and the components of $\beta_{S_0^c}$ can only be nonzero in the complement of S_0 . Let $\mathbb{C}(S_0)$ denote the set of vectors

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

¹Likewise, in Thompson sampling (Thompson, 1933), the agent chooses the greedy action for the sampled parameter.

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

Assumption 3 (Compatibility condition). For active set S_0 , there exists compatibility constant $\phi_0^2 > 0$ such that

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

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

Assumption 4 (Relaxed symmetry). For a joint distribution $p_{\mathcal{X}}$, there exists $\nu < \infty$ such that $\frac{p_{\mathcal{X}}(-\mathbf{x})}{p_{\mathcal{X}}(\mathbf{x})} \leq \nu$ for all \mathbf{x} .

Discussion of the assumptions. Assumptions 1 and 2 are the standard regularity assumptions used in the GLM bandit literature (Filippi et al., 2010; Li et al., 2017; Kveton et al., 2020). It is important to note that unlike the existing GLM bandit algorithms which explicitly use the value of κ_0 , our proposed algorithm does not use κ_0 or κ_1 — this information is only needed to establish the regret bound. The compatibility condition in Assumption 3 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 ensures that truly active components of the parameter vector are not "too correlated." As mentioned above, the compatibility condition is a standard assumption in the sparse bandit literature (Bastani & Bayati, 2020; Wang et al., 2018; Kim & Paik, 2019). Assumption 4 states that the joint distribution p_{χ} can be skewed but this skewness is bounded. Obviously, if $p_{\mathcal{X}}$ is symmetrical, we have $\nu = 1$. Assumption 4 is satisfied for a large class of continuous and discrete distributions, e.g., elliptical distributions including Gaussian and truncated Gaussian distributions, multi-dimensional uniform distribution, and Rademacher distribution. Note that in the non-sparse low dimensional setting (i.e., d = s), the relaxed symmetry in Assumption 4 together with the positive definiteness of the theoretical Gram matrix is equivalent to the covariate diversity condition introduced in (Bastani et al., 2020). However, in the sparse high-dimensional setting considered here, the relaxed symmetry does not imply diversity in all covariates. Consequently, the greedy parameter estimation approach proposed by (Bastani et al., 2020) is not guaranteed to achieve a sublinear regret. As in the case of κ_0 and κ_1 in Assumption 2, the parameter ν is only needed to establish the regret bound, our proposed algorithm does not require knowledge of ν .

4.1. Regret Bound for SA LASSO BANDIT

Theorem 1 (Regret bound for two arms). Suppose K = 2and Assumptions 2-4 hold. Then the expected regret of the SA LASSO BANDIT policy (π) over horizon T is upperbounded by

$$\mathcal{R}^{\pi}(T) \leq 4\kappa_{\max} + \frac{2\log(2d^2) + 2}{C_0(\phi_0, s_0)^2} + \frac{32\kappa_{\max}\rho_0\sigma s_0\sqrt{T\log(dT)}}{\kappa_{\min}\phi_0^2}$$

where $C_0(\phi_0, s_0) = \min\left(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_0}\right)$.

Discussion of Theorem 1. In terms of key problem primitives, Theorem 1 establishes a regret bound of $\mathcal{O}(s_0\sqrt{T\log(dT)})$ without any prior knowledge on s_0 . The bound shows that the regret of SA LASSO BANDIT 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 & Van De Geer, 2011)), we believe that the dependence on d and s_0 in Theorem 1 is best possible.²

The regret bound in Theorem 1 is tighter than the previously known bound in the same problem setting (Kim & Paik, 2019) although direct comparison is not immediate, given the difference in assumptions involved -- compared to (Kim & Paik, 2019), we require Assumption 4 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 several of our assumptions that were imposed primarily for technical tractability purposes. As mentioned earlier, the previous work on sparse bandits (Bastani & Bayati, 2020; Wang et al., 2018; Kim & 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 additional experiments in the appendix. We provide the outline of the proof and the key lemmas in the following section.

4.2. 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 induced by our bandit policy 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 3. 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 the 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 4.2.2, 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 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 the empirical Gram matrix satisfies the compatibility condition with high probability.

4.2.1. 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 & 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 ϕ_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 proba-

²Since the horizon T does not exist in offline Lasso results, it is not straightforward to see whether \sqrt{T} dependence can be improved comparing only with the offline Lasso results. Clearly, without an additional assumption on the separability of the arms, we know that poly-logarithmic scalability in T is not feasible. We briefly discuss our conjecture in comparison with the lower bound result in the non-sparse linear bandits in Secton B.1 in the appendix where we discuss the regret bound under the RE condition.

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

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

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 β^* 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.

4.2.2. Compatibility Condition and Matrix Concentration.

For matrix M, we define $\phi^2(M, S_0) := \min_{\beta} \{s_0\beta^{\top}M\beta/\|\beta_{S_0}\|_1^2 : \|\beta_{S_0}\|_1 \leq 3\|\beta_{S_0}\|_1 \neq 0\}$ as the (generic) compatibility constant. Hence, it suffices to show $\phi^2(M, S_0) > 0$ in order to show that matrix M satisfies the compatibility condition. Now, under Assumption 3, 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 2. 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 Σ_t , the history \mathcal{F}_{τ} affects how the feature vector X_{τ} is chosen. More specifically, our algorithm uses \mathcal{F}_{τ} to compute $\hat{\beta}_{\tau}$ and then chooses arm a_{τ} such that its (realized) feature $x_{a_{\tau}}$ maximizes $x_{a_{\tau}}^{\top}\hat{\beta}_{\tau}$. Therefore, we can rewrite Σ_t as

$$\Sigma_t = \frac{1}{t} \sum_{\tau=1}^t \sum_{i=1}^2 \mathbb{E}_{\mathcal{X}_{\tau}} \left[X_{\tau,i} X_{\tau,i}^\top \mathbb{1} \{ X_{\tau,i} = \underset{X \in \mathcal{X}_{\tau}}{\operatorname{argmax}} X^\top \hat{\beta}_{\tau} \} | \hat{\beta}_{\tau} \right]$$

Since the compatibility condition is satisfied only for the theoretical Gram matrix Σ and we need to show the empirical Gram matrix $\hat{\Sigma}_t$ satisfies the compatibility condition, the adapted Gram matrix Σ_t serves as a bridge between Σ 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 Σ_t satisfies the compatibility condition as long as Σ satisfies the compatibility condition. Then, we show that $\hat{\Sigma}_t$ concentrates around Σ_t with high probability and that such matrix concentration guarantees the compatibility condition of $\hat{\Sigma}_t$.

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

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

$$\sum_{i=1}^2 \mathbb{E} \Big[X_{t,i} X_{t,i}^\top \mathbb{1} \{ X_{t,i} = \operatorname{argmax}_{X \in \mathcal{X}_t} X^\top \beta \} \Big] \succcurlyeq \frac{\Sigma}{\rho_0}.$$

Therefore, we have $\Sigma_t \succeq \frac{\Sigma}{\rho_0}$ which implies that $\phi^2(\Sigma_t, S_0) \ge \frac{\phi^2(\Sigma, S_0)}{\rho_0} > 0$, i.e., Σ_t satisfies the compatibility condition. Note that both Σ and Σ_t can be singular. In Lemma 3, we show that $\hat{\Sigma}_t$ concentrates to Σ_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 8 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 $C(\phi_0, s_0) = \min(\frac{1}{2}, \frac{\phi_0^2}{128s_0\rho_0})$, we have

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

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 & Van De Geer, 2011)). Suppose that Σ_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 Σ_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 $t < T_0 := 2\log(2d^2)/C(\phi_0, s_0)^2$ during which the compatibility condition for the empirical Gram matrix is not guaranteed, and the event $\mathcal{E}_t := \{ \|\Sigma_t - \hat{\Sigma}_t\|_{\infty} \le \phi_0^2/(32s_0\rho_0) \}$. Then for all $t \ge [T_0]$ and Σ_t for which event \mathcal{E}_t holds, we have

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

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

4.2.3. 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 time 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 BAN-DIT Algorithm does not rely in any way on this partitioning - it is introduced purely for the purpose of analysis. Set (a) is the initial period over 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 time 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 proofs of the lemmas are deferred to the appendix.

5. Extension to K Arms

Recall that SA LASSO BANDIT is valid for any number of arms; hence, no modifications are required to extend the algorithm to $K \geq 3$ arms. The analysis of SA LASSO BANDIT for the K-armed case tackles largely the same challenges described in Section 4.2: 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 Σ_t with the theoretical Gram matrix Σ , the Gram matrix for the unobserved feature vectors could be incomparable with the Gram matrix for the observed feature vectors. For this issue, we introduce an additional regularity condition, which we denote as the "balanced covariance" condition.

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

$$\begin{split} & \mathbb{E}\left[X_{i_k}X_{i_k}^{\top}\mathbbm{1}\{X_{i_1}^{\top}\beta < \ldots < X_{i_K}^{\top}\beta\}\right] \\ & \preccurlyeq C_{\mathcal{X}}\mathbb{E}\left[(X_{i_1}X_{i_1}^{\top} + X_{i_K}X_{i_K}^{\top})\mathbbm{1}\{X_{i_1}^{\top}\beta < \ldots < X_{i_K}^{\top}\beta\}\right] \end{split}$$

In Algorithm 1 we observe only the reward corresponding to arm i_1 , and Assumption 4 implies that we have some control on the arm i_K . This balanced covariance condition implies that there is "sufficient randomness" in the observed features compared to non-observed features. The exact value of C_X depends on the joint distribution of \mathcal{X} including the correlation between arms. In general, the more positive the corre-

lation, the smaller $C_{\mathcal{X}}$ (obviously, with an extreme case of perfectly correlated arms having a constant $C_{\mathcal{X}}$ independent of any problem parameters). When the arms are independent and identically distributed, Assumption 5 holds with $C_{\mathcal{X}} = \mathcal{O}(1)$ for both the multivariate Gaussian distribution and a uniform distribution on a sphere, and for an arbitrary independent distribution for each arm, Assumption 5 holds for $C_{\mathcal{X}} = \binom{K-1}{K_0}$ 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 the sparsity index s_0 . These are formalized in Proposition 1 in the appendix.³ This balanced covariance condition is somewhat similar to "positive-definiteness" condition for observed contexts in the bandit literature (e.g., Goldenshluger & Zeevi (2013); Bastani et al. (2020)). 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 skewed toward a particular direction. Hence, learning algorithms may require many more samples to learn the unknown parameter, leading to larger regret. It is important to note that our algorithm does not require any prior information on $C_{\mathcal{X}}$. The regret bound for the K-armed sparse bandits under Assumption 5 is as follows.

Theorem 2 (Regret bound for *K* arms). Suppose $K \ge 3$ and Assumptions 1-4, and 5 hold. Let $\lambda_0 = 2\sigma x_{\text{max}}$. Then the expected cumulative regret of the SA LASSO BANDIT policy π over horizon $T \ge 1$ is upper-bounded by

$$\begin{aligned} \mathcal{R}^{\pi}(T) &\leq 4\kappa_{1} + \frac{4\kappa_{1}x_{\max}b(\log(2d^{2})+1)}{C_{1}(s_{0})^{2}} \\ &+ \frac{64\kappa_{1}\nu C_{\mathcal{X}}\sigma x_{\max}s_{0}\sqrt{T\log(dT)}}{\kappa_{0}\phi_{0}^{2}} \end{aligned}$$
where $C_{1}(s_{0}) = \min\left(\frac{1}{2}, \frac{\phi_{0}^{2}}{256s_{0}\nu C_{\mathcal{X}}x_{\max}^{2}}\right).$

Theorem 2 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. Since both multivariate Gaussian distributions and uniform distributions satisfy Assumption 4 with $\nu = 1$ and Assumption 5 with $C_{\mathcal{X}} = \mathcal{O}(1)$, the regret bound in Theorem 2 still holds

³While it is not our primary goal to derive general tight bounds on $C_{\mathcal{X}}$, we acknowledge that the bound on $C_{\mathcal{X}}$ for an arbitrary distribution for independent arms is very loose, and is the result of conservative analysis driven by lack of information on $p_{\mathcal{X}}$. Numerical evaluation on distributions other than Gaussian and uniform distributions, detailed in Section 5, buttress this point and indicate that the dependence on K is no greater than linear.

Sparsity-Agnostic Lasso Bandit



Figure 1. The evaluations of SA LASSO BANDIT (Algorithm 1), DR LASSO BANDIT (Kim & Paik, 2019), and LASSO BANDIT (Bastani & Bayati, 2020). The first row shows results for features drawn from a multivariate Gaussian distribution with varying correlation between arms. The second and third rows show results for uniform and non-Gaussian elliptical distributions respectively. The results provide clear evidence that SA LASSO BANDIT outperforms the benchmarks across various experiments.

under Assumptions 1-3 for these distributions. Therefore, to our knowledge, this is the first sparsity-agnostic regret bound for a general *K*-armed high-dimensional contextual bandit algorithm even for the Gaussian distribution or uniform distribution.

The proof of Theorem 2 largely follows that of Theorem 1. The main difference is how we control the adapted Gram matrix Σ_t with the theoretical Gram matrix Σ . 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 show that for a fixed $\beta \in \mathbb{R}^d$,

$$\sum_{i=1}^{K} \mathbb{E}_{\mathcal{X}_{t}} \Big[X_{t,i} X_{t,i}^{\top} \mathbb{1} \{ X_{t,i} = \underset{X \in \mathcal{X}_{t}}{\operatorname{argmax}} X^{\top} \beta \} \Big] \succcurlyeq (2\nu C_{\mathcal{X}})^{-1} \Sigma$$

v

The formal result is presented in Lemma 10 in the appendix along with its proof. Next, 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 Σ_t , and eventually to the theoretical Gram matrix Σ . Thus, we ensure the compatibility condition of $\hat{\Sigma}_t$. The additional regret in the *K*-armed case as compared to the two-armed case is essentially a scaling by C_X to ensure the balanced covariance condition.

6. Experiments

We conduct numerical experiments to evaluate SA LASSO BANDIT and compare with existing sparse bandit algorithms: DR LASSO BANDIT (Kim & Paik, 2019) and LASSO BANDIT (Bastani & Bayati, 2020). For each case with different experimental configurations, we conduct 20 independent runs. For performance evaluations, we report the average of the cumulative regret for each of the algorithms. The error bars represent the standard deviations. Each row of the plots show experiments using different distributions for feature vectors. Additional results are presented in the appendix. SA LASSO BANDIT exhibits superior performances across different distributions as well as other problem parameters.

The results provide convincing evidence that the performance of our proposed algorithm is superior to the existing sparse bandit methods that we compare with. SA LASSO BANDIT outperforms the existing sparse bandit algorithms by significant margins, even though the correct sparsity index s_0 is revealed to these algorithms and kept hidden from SA LASSO BANDIT. Furthermore, SA LASSO BANDIT is much more practical and simple to implement with a minimal number of a hyperparameter.

References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Abbasi-Yadkori, Y., Pal, D., and Szepesvari, C. Onlineto-confidence-set conversions and application to sparse stochastic bandits. In *Artificial Intelligence and Statistics*, pp. 1–9, 2012.
- Abe, N. and Long, P. M. Associative reinforcement learning using linear probabilistic concepts. In *International Conference on Machine Learning*, pp. 3–11, 1999.
- Agrawal, S. and Goyal, N. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pp. 127–135, 2013.
- Auer, P. Using confidence bounds for exploitationexploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Bang, H. and Robins, J. M. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61 (4):962–973, 2005.
- Bastani, H. and Bayati, M. Online decision making with high-dimensional covariates. *Operations Research*, 68 (1):276–294, 2020.
- Bastani, H., Bayati, M., and Khosravi, K. Mostly exploration-free algorithms for contextual bandits. *Management Science*, 2020.
- Bickel, P. J., Ritov, Y., and Tsybakov, A. B. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- Bühlmann, P. and Van De Geer, S. Statistics for highdimensional data: methods, theory and applications. Springer Science & Business Media, 2011.
- Cambanis, S., Huang, S., and Simons, G. On the theory of elliptically contoured distributions. *Journal of Multivariate Analysis*, 11(3):368–385, 1981.
- Candes, E. and Tao, T. The dantzig selector: Statistical estimation when p is much larger than n. *The annals of Statistics*, 35(6):2313–2351, 2007.
- Carpentier, A. and Munos, R. Bandit theory meets compressed sensing for high dimensional stochastic linear bandit. In *Artificial Intelligence and Statistics*, pp. 190– 198, 2012.

- Chu, W., Li, L., Reyzin, L., and Schapire, R. Contextual bandits with linear payoff functions. In *Proceedings* of the Fourteenth International Conference on Artificial Intelligence and Statistics, pp. 208–214, 2011.
- Dani, V., Hayes, T. P., and Kakade, S. M. Stochastic linear optimization under bandit feedback. In *Proceedings* of the 21st Annual Conference on Learning Theory, pp. 355–366, 2008.
- Filippi, S., Cappe, O., Garivier, A., and Szepesvári, C. Parametric bandits: The generalized linear case. In Advances in Neural Information Processing Systems, pp. 586–594, 2010.
- Gilton, D. and Willett, R. Sparse linear contextual bandits via relevance vector machines. In 2017 International Conference on Sampling Theory and Applications (SampTA), pp. 518–522. IEEE, 2017.
- Goldenshluger, A. and Zeevi, A. A linear response bandit problem. *Stochastic Systems*, 3(1):230–261, 2013.
- Kim, G.-S. and Paik, M. C. Doubly-robust lasso bandit. In Advances in Neural Information Processing Systems, pp. 5869–5879, 2019.
- Kveton, B., Zaheer, M., Szepesvari, C., Li, L., Ghavamzadeh, M., and Boutilier, C. Randomized exploration in generalized linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pp. 2066–2076, 2020.
- Lattimore, T. and Szepesvári, C. Bandit Algorithms. Cambridge University Press (preprint), 2019.
- Li, L., Chu, W., Langford, J., and Schapire, R. E. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pp. 661–670. ACM, 2010.
- Li, L., Lu, Y., and Zhou, D. Provably optimal algorithms for generalized linear contextual bandits. In *International Conference on Machine Learning*, pp. 2071–2080, 2017.
- Oh, M.-h. and Iyengar, G. Thompson sampling for multinomial logit contextual bandits. In Advances in Neural Information Processing Systems, pp. 3151–3161, 2019.
- Raskutti, G., Wainwright, M. J., and Yu, B. Restricted eigenvalue properties for correlated gaussian designs. *Journal of Machine Learning Research*, 11(Aug):2241–2259, 2010.
- Rusmevichientong, P. and Tsitsiklis, J. N. Linearly parameterized bandits. *Mathematics of Operations Research*, 35 (2):395–411, 2010.

- Tewari, A. and Murphy, S. A. From ads to interventions: Contextual bandits in mobile health. In *Mobile Health*, pp. 495–517. Springer, 2017.
- Thompson, W. R. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B* (*Methodological*), 58(1):267–288, 1996.
- Van De Geer, S. A. and Bühlmann, P. On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, 3:1360–1392, 2009.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Wang, X., Wei, M., and Yao, T. Minimax concave penalized multi-armed bandit model with high-dimensional covariates. In *International Conference on Machine Learning*, pp. 5200–5208, 2018.
- Zhang, C.-H. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2): 894–942, 2010.