
Doubly High-Dimensional Contextual Bandits: An Interpretable Model with Applications to Assortment/Pricing

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider contextual bandits that are doubly high-dimensional: both covariates
2 and actions take values in high-dimensional spaces. We propose a simple model
3 that captures the interactions between covariates and actions via a (near) low-rank
4 representation matrix. The resulting class of models is reasonably expressive while
5 remaining interpretable, and includes various structured linear bandit models as
6 particular cases. We propose a computationally tractable procedure that combines
7 an exploration/exploitation protocol with an efficient low-rank matrix estimator,
8 and we prove bounds on its regret. Simulation results show that this method
9 has lower regret than state-of-the-art methods applied to various standard bandit
10 models. We also apply our method to a real-world online retail data set involving
11 assortment and pricing; in contrast to most existing methods, our method allows
12 the assortment-pricing problem to be solved simultaneously. We demonstrate the
13 effectiveness of this joint approach for revenue maximization.

14 1 Introduction

15 The bandit problem, dating back to the seminal work of Robbins [43], is a canonical problem in
16 sequential decision-making. At each round, a decision-maker chooses an action (arm), and then
17 observes a reward. The goal is to act strategically so as to determine a near-optimal policy without
18 incurring large regret. There is now a very well-developed literature on the bandit problem, and its
19 extension to the contextual bandits; see the book by Lattimore and Szepesvári [32] and references
20 therein for more background. Contextual bandit models and algorithms play a central role in online
21 decision-making, with e-commerce and health care being two fruitful domains.

22 The classical bandit formulation involves a finite action space, but many modern applications lead to
23 settings where actions take values in continuous space, also known as the *continuum-armed bandit*
24 (e.g., Agrawal [2], Kleinberg [26]). For instance, in e-commerce, an online retailer seeks to decide
25 upon product assortment and pricing so as to maximize long-term profits [9, 45, 15, 23]. In mobile
26 health, the personal device provides exercise and dietary suggestions to improve physical and mental
27 health (see Debon et al. [14] and references therein). In both of these cases, the action vector \mathbf{a}
28 takes values in some subset of \mathbb{R}^{d_a} , where the action dimension d_a can be quite large, which we
29 refer to as the high-dimensional setting. In the contextual bandit problem, decision-makers observe
30 an additional vector $\mathbf{x} \in \mathbb{R}^{d_x}$ of features or covariates, also known as the context. In applications,
31 the covariate dimension d_x may also be high-dimensional. The reward mean is modeled as some
32 unknown function of the covariate-action pair (\mathbf{x}, \mathbf{a}) .

33 As the action-covariate dimensions d_a and d_x grow, traditional bandit algorithms suffer from the
34 curse of dimension; indeed, without some kind of structure, there are “no-free-lunch” theorems
35 showing that it is prohibitively costly, both in terms of samples and computation, to learn an optimal
36 policy [32]. This fact motivates various models that encode some form of low-dimensional structure
37 in the reward function. To date, researchers have pursued structure in the covariates and actions
38 in isolation, including sparsity for high-dimensional contextual bandit problems (e.g., Bastani and
39 Bayati [6]), or subspace structure for continuum-action bandits (e.g., Tyagi et al. [50]). This line of
40 work leaves open the following question:

41 *Are there useful models and efficient learning procedures for contextual bandits that are high-*
42 *dimensional in both actions and covariates?*

43 In this paper, we tackle this challenge by proposing a new model that captures interactions between
44 actions and covariates via an (approximately) low-rank matrix representation. Within this model
45 class, we also propose a new algorithm (Hi-CCAB) that combines low-rank estimation with an
46 exploration/exploitation strategy, and prove non-asymptotic bound on its expected regret.

47 Our reward model takes the following form: given a covariate vector $\mathbf{x} \in \mathbb{R}^{d_x}$ and an action vector
48 $\mathbf{a} \in \mathbb{R}^{d_a}$, we observe a noisy reward Y with conditional mean

$$\mathbb{E}[Y \mid \mathbf{x}, \mathbf{a}] = \mathbf{a}^T \Theta \mathbf{x},$$

49 where $\Theta \in \mathbb{R}^{d_a \times d_x}$ is an unknown representation matrix. As we discuss in the sequel, in many
50 applications, it is natural to assume that this representation matrix is relatively low-rank—say with
51 rank $r \ll \min\{d_a, d_x\}$. Given this structure, our proposed Hi-CCAB algorithm interleaves estimation
52 steps, in which the low-rank representation matrix is estimated based on data observed thus far, and
53 exploration/exploitation steps. In the proposal given here, we analyze a standard estimator based
54 on the nuclear norm relaxation of rank (cf. Chapter 10 in Wainwright [51] for details). While
55 the estimator itself is not novel, our analysis of it does require new ingredients since we apply it
56 to data adaptively collected under a bandit protocol. We further demonstrate the benefits of our
57 methodologies in e-commerce with real sales data where the online retailer needs to decide on the
58 product assortment and pricing jointly. The generality of our model makes it possible to learn policy
59 on product assortment and pricing at the same time, while previous literature mostly studies the
60 assortment and pricing problem separately.

61 **Contributions.** Let us summarize some of our main contributions:

- 62 1. We propose a new model for high-dimensional contextual bandits, in which both the covariates
63 and actions can be high-dimensional and continuous. The crux of our model is a low-rank
64 representation matrix that represents the interaction between action-covariate pairs via its left
65 and right singular vectors. This model, while quite simple, unifies a number of structured bandit
66 models analyzed in past work.
- 67 2. As we argue, an advantage of this low-rank model is its combination of prediction power with
68 a high degree of interpretability. Performing a singular value decomposition (SVD) on the
69 representation matrix yields the latent structure, with the left (respectively right) singular vectors
70 corresponding to the action (respectively covariate) space structure. In this way, our model
71 implicitly performs a form of dimension reduction in how the actions and covariates interact to
72 determine the reward function. On the other hand, given the covariate, our model is able to predict
73 the reward of an unseen arm. Both interpretability and predictive power can be tremendously
74 useful for decision-makers.
- 75 3. We propose an efficient algorithm for online learning in the active setting, referred to as the
76 **High-dimensional Contextual and High-dimensional Continuum Armed Bandit** (Hi-CCAB) by
77 adopting the low-rank matrix estimator. We further provide a non-asymptotic upper bound on the
78 expected regret of Hi-CCAB.
- 79 4. The generality of our model allows for a wide range of applications. Specifically, we apply
80 Hi-CCAB to the joint assortment and pricing problem. We show that our model reveals insights
81 for product designs, assortment, and pricing and that the assortment-pricing policy based on
82 Hi-CCAB yields sales four times as high as the original strategy.

83 **Connections to past work.** The literature on high-dimensional bandit problems has been growing,
84 and exploits a relatively mature body of statistical tools for high-dimensional problems (e.g., see the
85 book [51] and references therein). There is a line of work on contextual bandits with high-dimensional
86 covariates, including the LASSO bandit problem [1, 24, 6, 17, 39, 52], in which the mean reward is
87 assumed to be a linear function of a sparse unknown parameter vector. As we describe in the sequel,
88 these high-dimensional bandit models are special cases of the high-dimensional low-rank model in
89 this paper. Other work uses non-parametric methods, such as boosting, random forests, or neural
90 networks, for estimating the reward function [16, 55, 5, 13, 53]. Such approaches are quite different
91 in flavor, and we compare to one such method in our experimental results.

92 There are various other models and problems that have connections to but differ from the set-up in
93 this paper. For example, one line of research focuses on representation learning in linear bandits,
94 specifically for low-rank bandit models and multi-task learning where several bandits are played

concurrently. The arms for each task are embedded in the same space and share a common low-dimensional representation [30, 31, 54, 18, 35, 22]. However, this line of research does not consider contextual information, and often imposes case-specific assumptions on the action space. Among such papers, Kang et al. [22] study a trace inner product bandit with a matrix of known (low) rank r , in which the action is matrix-valued. In the supplement, they sketch some potential extensions to contextual problems, but their theory and numerical validation only applies to non-contextual bandits. Our algorithm and theory, in contrast, are designed for contextual problems, and we do not need to know the rank r of the target matrix. Our reward model is connected to but different from other papers that propose bilinear-type reward models (e.g., Jun et al. [19], Kim and Vojnovic [25], Rizk et al. [42]) in which *both* arguments of the bilinear function are part of the action. Such models can be understood as a structured linear bandit of a particular type, and unlike our models, do not capture the interaction between the covariate and action at each time step.

For continuum-action bandits, there exists a thread of literature that assumes the mean reward function is smooth and continuous on the action space in some sense, e.g., the function lies in the Lipschitz or Hölder space [2, 26, 27]. Approaches taken in these works include discretizing the action space, or using non-parametric regression to estimate the reward function, which is quite different from our model. Other work on contextual bandits with continuous actions and covariates assumes that the reward function is continuous with Lipschitz type conditions over the action-covariate space and it is restricted to relatively low-dimension [34, 46, 29]. There are a few recent papers on high-dimensional models for contextual bandits [49], but using rather different techniques with relatively strict assumptions and lacking the interpretability of our model.

Methods for low-rank matrix prediction and estimation have been studied extensively in both statistics and machine learning (e.g., [47, 41, 8, 37, 7]). Our Hi-CCAB algorithm uses least-squares with nuclear norm regularization, which is a well-known approach; however, our analysis requires some novel results so as to deal with the adaptive nature of bandit data collection.

Finally, in the field of operations research, assortment and pricing are key decisions to be made by any firm; accordingly, there is a substantial body of past work on dynamic assortment and dynamic pricing. Much of the work on assortment is based on the multinomial logit (MNL) choice model [9, 28, 45]; more recent work adopts multi-arm bandit techniques to the MNL model [12, 3, 21, 11]. For dynamic pricing, the problem usually comes with demand learning. In presence of covariates, the demand can be modeled as a parametric function [40, 4] or a nonparametric function [10] which adopt the continuum-action bandit techniques in Slivkins [46]. However, there are relatively few papers on the joint assortment-pricing problem. Recently, Miao and Chao [36] provides a solution using the MNL choice model with a finite number of actions, while our model involves infinitely many actions. In addition, their model assumes independence of the products and can only handle a small number of products. Their model can neither incorporate contextual information nor predicts new products.

Roadmap. The rest of the paper is organized as follows. Section 2 describes the problem formulation and introduces our model with two concrete examples in assortment-pricing and health care. Section 3 presents our Hi-CCAB algorithm and its convergence result. Finally, Section 4 shows the empirical results on simulated data and a case study on real sales data from one of the largest online retailers. The proof of our theorem and additional empirical results are provided in the Appendix.

Notation. We use bold lowercase for vectors and bold uppercase for matrices. We denote the ℓ_2 norm of vector \mathbf{a} by $\|\mathbf{a}\|$. For matrix \mathbf{A} , let $\|\mathbf{A}\|_F := \sum_{ij} a_{ij}^2$ be its Frobenius norm, $\|\mathbf{A}\|_2$ be its ℓ_2 spectrum norm, i.e., $\|\mathbf{A}\|_2 := \sup_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2$, and $\|\mathbf{A}\|_* := \sum_{k=1}^d s_k$ be its nuclear norm, where d is the rank and s_k 's are the singular values of \mathbf{A} . We use $\langle \mathbf{a}, \mathbf{b} \rangle := \mathbf{a}^\top \mathbf{b}$ to denote the inner product between two vectors and $\langle \mathbf{A}, \mathbf{B} \rangle := \text{trace}(\mathbf{A}^\top \mathbf{B})$ between two matrices.

2 Problem formulation

In this section, we first introduce our doubly high-dimensional contextual bandit model. We compare and contrast it with traditional bandit models, and provide intuition as to why it is suitable for two different applications (joint assortment and pricing, and mobile healthcare apps). Finally, we discuss how various structured bandit models can be seen as special cases of our proposal.

Problem setup. Our goal is to learn the action with the highest expected reward based on T samples in a sequential model for data collection. At each time t , we are allowed to choose an action \mathbf{a}_t based on the data seen to date. This chosen action applies to a batch of objects of size $L \geq 1$. At each round, we observe a collection of contexts of covariate vectors $\{\mathbf{x}_{t,j}\}_{j=1}^L$, one associated with each

of the L objects, and each lying in \mathbb{R}^{d_x} . At time t , we are allowed to make a decision based on all observations prior to time t along with the covariates at time t , we decide on an action \mathbf{a}_t that takes values in a constraint set $\mathcal{A} \subseteq \mathbb{R}^{d_a}$. After taking action \mathbf{a}_t at time t , we observe a batch of rewards

$$y_{t,j} = \mathbf{a}_t^\top \boldsymbol{\Theta} \mathbf{x}_{t,j} + \varepsilon_{t,j}, \quad \text{for } j = 1, 2, \dots, L, \quad (1)$$

where $\boldsymbol{\Theta} \in \mathbb{R}^{d_a \times d_x}$ is an unknown low-rank matrix and $\varepsilon_{t,j}$ is independent noise with $\mathbb{E}[\varepsilon_{t,j}] = 0$ and $\text{Var}[\varepsilon_{t,j}] \leq \sigma^2$.

Suppose that we run T rounds of this scheme, using some protocol π for choosing actions; doing so yields a sequence $\{\mathbf{a}_{t,\pi}\}_{t=1}^T$ of T actions. We measure the quality of this sequence—and hence the protocol π —via its *expected average regret*¹

$$\mathcal{R}^\pi(T) = \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T \sup_{\mathbf{a} \in \mathcal{A}_t} \left(\sum_{j=1}^L \mathbf{a}^\top \boldsymbol{\Theta} \mathbf{x}_{t,j} - \mathbf{a}_{t,\pi}^\top \boldsymbol{\Theta} \mathbf{x}_{t,j} \right) \right]. \quad (2)$$

In this definition, the expectation is taken with respect to $(\mathbf{x}_{t,j}, \varepsilon_{t,j})$ since $\mathbf{a}_{t,\pi}$ depends on both. Our goal is to design protocols π with low regret.

At the core is the model $\mathbb{E}[Y \mid \mathbf{x}, \mathbf{a}] = \mathbf{a}^\top \boldsymbol{\Theta} \mathbf{x}$ of the mean reward function. It is worth noting that this form of reward function generalizes various known models, including classical K -arm bandits, K -arm bandits with context, and continuum-armed bandits [43, 6, 27]. We discuss these connections in more detail at the end of this section.

The other key ingredient in our model is the low-rank representation matrix $\boldsymbol{\Theta}$. It encapsulates the effect of both the arm and covariates on the reward and exploits the low-dimensional structure in the high-dimensional actions and covariates. To understand this condition, consider a matrix $\boldsymbol{\Theta}$ that is of rank $r \ll \min\{d_a, d_x\}$. It then has a singular value decomposition of the form $\boldsymbol{\Theta} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$, where $\mathbf{S} = \text{diag}\{s_1, \dots, s_r\}$ is a diagonal matrix with the ordered singular values $s_1 \geq s_2 \geq \dots \geq s_r > 0$, and both $\mathbf{U} \in \mathbb{R}^{d_a \times r}$ and $\mathbf{V} \in \mathbb{R}^{d_x \times r}$ are matrices with orthonormal columns, corresponding to the left $\{\mathbf{u}_j\}_{j=1}^r$ and right singular vectors $\{\mathbf{v}_j\}_{j=1}^r$, respectively. With this notation, the reward function (1) can be written as

$$\mathbb{E}[Y \mid \mathbf{a}, \mathbf{x}] = \mathbf{a}^\top \boldsymbol{\Theta} \mathbf{x} = \sum_{i=1}^r s_i \langle \mathbf{a}, \mathbf{u}_i \rangle \cdot \langle \mathbf{v}_i, \mathbf{x} \rangle. \quad (3)$$

In other words, the mean reward is the summation of inner products between the action projected on the left singular vector and the covariates projected on the right singular vector, weighted by the singular values. By assuming $\boldsymbol{\Theta}$ to be low-rank, the mean reward is assumed to be governed by only a few linear combinations of the arm attributes and covariates. Hence, our model automatically explores the low-dimensional structure of the arm and contextual vectors in terms of its effect on the reward, from which we can draw interpretation and insights from the effective subspaces of both the arm and covariates.

To explain why the low-rank assumption is well-motivated in practice, let's discuss some concrete cases.

Example 1 (Assortment and Pricing). The assortment problem, which arises in retail and e-commerce, is to decide which combination of products to present at each given time while satisfying capacity constraints [28]. The closely related pricing problem is to decide the prices of the products. A given firm wants to solve these problems so as to maximize a certain objective (e.g., revenue or profit).

In these problems, the action vector consists of both prices and attributes associated with products. Attributes differ in a case-by-case manner: for clothes, the pattern, color and size are standard attributes, whereas for electronics, the technical specifications provide attributes. In our case study, as described in more detail in Section 4, we focus on boxes of instant noodles. Box j has a price p_j , and can contain packages that correspond to one of m total flavors; this can be encoded by an attribute vector $\mathbf{f}_j = (f_{j,1}, \dots, f_{j,m})$ where $f_{j,\ell}$ is the number of packages of flavor ℓ in box j . Given a total of $K \geq 1$ slots in which to present products, a given store needs to decide which products to present, along with their corresponding prices. To formalize this set-up, the action vector takes the form $\mathbf{a} = (\mathbf{f}_1, p_1, \mathbf{f}_2, p_2, \dots, \mathbf{f}_K, p_K, 1)$, and so is high-dimensional. At the same time, we observe

¹Here we have defined the average regret (via our rescaling by T), but bounds on this quantity can immediately be translated to the cumulative regret as needed.

194 covariates associated with the observations; they represent information such as geographic location,
195 seasonal information at the aggregated level, or demographic information at the user level.

196 The demand and sales of products with similar attributes react similarly to the same market conditions.
197 It is often the case that there exist latent factors of the products that govern the demand and sales.
198 Therefore, it is reasonable to parameterize the reward function (1) rather than ignoring the similarity
199 between products as in the literature [36, 21, 11]. Our model can further suggest new products rather
200 than only the products that have already been provided. ♣

201 *Example 2 (Health-care).* Bandits are used for health monitoring, in which the behavior and vital
202 measurements of users are tracked, and the system provides suggested actions (e.g., take a walk). In
203 this application, the action vector (\mathbf{a}) is high-dimensional and can include continuous values (e.g.,
204 sleeping time, length and kind of exercise, usage of social media, diet choices including energy,
205 water, protein, minerals, and nutrition intakes), and the health outcome depends on not only our
206 suggestions, but also the user's characteristics (e.g., age, gender, weight, height, basic health status,
207 and the tendency of following suggestions) as contextual variables (\mathbf{x}). Clearly, both the arm and
208 the contextual variable vectors can take continuous values and are possibly high-dimensional. The
209 classical bandit models do not fit the situation. Similar actions usually share similar effects on
210 health and the user's characteristics can typically be captured by a few latent factors. Therefore, it is
211 reasonable to assume a bilinear model with a low-rank Θ . ♣

212 To close this section, let us summarize formally some classical bandit models that are special cases of
213 our reward model (1).

- 214 1. Multi-arm bandit: For i -th arm, $\mathbf{a} = (0, 0, \dots, 1, \dots, 0)$, where 1 is in i -th element. Suppose \mathbf{x}
215 has its first element being constant. Then $\Theta_{i,1} = \mu_i$, where μ_i is the mean reward of the i -th arm,
216 and $\Theta_{i,j} = 0$ if $j \neq 1$. Clearly, the matrix Θ has rank one.
- 217 2. Multi-arm high-dimensional contextual bandit: for i -th arm, $\mathbf{a} = (0, 0, \dots, 1, \dots, 0)$, where 1
218 is in i -th element. \mathbf{x} is the contextual vector. Then $\Theta = (\beta_1, \beta_2, \dots, \beta_m)^\top$, where β_i is the
219 parameter vector corresponding to i -th arm [6].
- 220 3. Continuum-action bandits (without context): Suppose a is the arm in the original continuum arm
221 bandit, and the mean reward function is $\mu(a)$. Since all continuous functions on a bounded interval
222 can be approximated by polynomial functions to arbitrary precision, it is reasonable to assume
223 $\mu(a)$ be a polynomial of order n , which is not known precisely but known to be smaller than a
224 fixed N . Let $\mathbf{a} = (1, a, a^2, a^3, \dots, a^n, \dots, a^N)$, and suppose the first element of \mathbf{x} is constant 1,
225 then $\Theta_{i,j} = \frac{1}{i!} f^{(i)}(a)$ for $j = 1$ and $\Theta_{i,j} = 0$ for $j \neq 1$. This matrix Θ also has rank one.

226 3 Hi-CCAB algorithm and theoretical results

227 In this section, we present our learning algorithm with a regret upper bound. Specifically, we detail
228 the Hi-CCAB algorithm in Section 3.1 and establish an upper bound for its convergence rate of the
229 expected regret in Section 3.2.

230 3.1 Description of the learning algorithm

231 Our policy consists of two phases for each period $t \in [T]$: the first phase learns a low-rank represen-
232 tation and the second phase determines the assortment and the selling prices. In the first phase, our
233 policy estimates $\hat{\Theta}_t$ by a penalized least-square estimator using $(\mathbf{a}_i, \mathbf{x}_{i,j}, y_{i,j})$ for $i = 1, \dots, t$ and
234 $j = 1, \dots, L$. Based on $\hat{\Theta}_t$, we look for the optimal assortment and pricing within the action space
235 \mathcal{A}_t . Algorithm 1 describes the detailed procedure of our policy.

236 **Low-rank representation learning.** As mentioned in Section 2, both the arm and contextual vectors
237 $\mathbf{a} \in \mathbb{R}^{d_a}$ and $\mathbf{x} \in \mathbb{R}^{d_x}$ are high-dimensional, thus $\Theta \in \mathbb{R}^{d_a \times d_x}$ is also high-dimensional. Fortunately,
238 there often exists structure in both the arm and covariate space as explained in Section 1. To leverage
239 the underlying structure, we impose a low-rank assumption on Θ , which automatically explores the
240 effect of the low-rank structure and the relationships between the action and the contextual arms.

241 In order to estimate the low-rank representation of Θ at time t , one could in principle solve the rank-
242 regularized least-squares problem: $\arg \min_{\Theta} \left\{ \sum_{i=1}^t \sum_{j=1}^L (\mathbf{a}_i^\top \Theta \mathbf{x}_{i,j} - y_{i,j})^2 + \lambda_t \cdot \text{rank}(\Theta) \right\}$,
243 where $\text{rank}(\Theta)$ is the rank function, and $\lambda_t > 0$ is a regularization parameter. Rank penalization
244 leads to a non-convex problem with associated computational challenges, so that it is standard to
245 replace it with the nuclear norm so as to obtain a convex problem. Doing so in our context yields the
246 nuclear-norm regularized estimator

Algorithm 1: The Hi-CCAB Algorithm.

Result: Actions $\mathbf{a}_{t_1+1}, \dots, \mathbf{a}_T$.

Input: The number of steps for initialization t_1 , set of possible actions \mathcal{A}_{t_1} , action vectors based on domain knowledge $\{\mathbf{a}_i\}_{i=1}^{t_1}$, covariate vectors $\{\mathbf{x}_{i,j}\}_{i=1}^{t_1}$, rewards $y_{i,j}$ for $j = 1, \dots, L$, and exploration parameter h .

Initialization: $\lambda_0 \leftarrow \left\| \frac{1}{2t_1L} \sum_{i=1}^{t_1} \sum_{j=1}^L |\mathbf{a}_i^\top \hat{\Theta}_{t_1} \mathbf{x}_{i,j} - y_{i,j}| \mathbf{x}_{i,j} \mathbf{a}_i^\top \right\|_2$, $t \leftarrow t_1$.

while $t < T$ **do**

$\lambda_t \leftarrow \lambda_0 / \sqrt{t}$;

Low-rank representation learning:

$\hat{\Theta}_t \leftarrow \arg \min_{\Theta} \left\{ \frac{1}{tL} \sum_{i=1}^t \sum_{j=1}^L (\mathbf{a}_i^\top \Theta \mathbf{x}_{i,j} - y_{i,j})^2 + \lambda_t \|\Theta\|_* \right\}$;

Policy learning:

$\hat{\mathbf{a}}_{t+1} \leftarrow \arg \max_{\mathbf{a} \in \mathcal{A}_t} \sum_{j=1}^L \mathbf{a}^\top \hat{\Theta}_t \mathbf{x}_{t+1,j}$;

if $t \notin \{\lfloor w^{\frac{3}{2}} \rfloor : w \in \mathbb{Z}_+\}$ **then** *Exploitation:* $\mathbf{a}_{t+1} \leftarrow \hat{\mathbf{a}}_{t+1}$;

else *Exploration:* $\mathbf{a}_{t+1} \leftarrow \hat{\mathbf{a}}_{t+1} + \delta_{t+1}$ where $\delta_{t+1} \sim N(\mathbf{0}_{d_a}, h\mathbf{I}_{d_a})$, update \mathcal{A}_{t+1} ;

Apply action \mathbf{a}_{t+1} and observe reward $y_{t+1,j}$ for $j = 1, \dots, L$;

$t \leftarrow t + 1$;

end while

$$\hat{\Theta}_t := \arg \min_{\Theta} \left\{ \sum_{i=1}^t \sum_{j=1}^L (\mathbf{a}_i^\top \Theta \mathbf{x}_{i,j} - y_{i,j})^2 + \lambda_t \cdot \|\Theta\|_* \right\}. \quad (4)$$

247 The penalization parameter λ_t is updated in each iteration according to the schedule $\lambda_t = \lambda_0 / \sqrt{t}$,
 248 where λ_0 is the initialized penalization parameter, which can be chosen by cross-validation or guided
 249 by $\left\| \frac{1}{2t_1L} \sum_{i=1}^{t_1} \sum_{j=1}^L |\mathbf{a}_i^\top \hat{\Theta}_{t_1} \mathbf{x}_{i,j} - y_{i,j}| \mathbf{x}_{i,j} \mathbf{a}_i^\top \right\|_2$.

250 **Policy learning.** Once we estimated the low-rank representation of Θ , we can proceed to the action
 251 step. The goal of the action step is to *exploit* the knowledge we have learned, i.e., $\hat{\Theta}_t$, so as to decide
 252 on the next action \mathbf{a}_{t+1} that maximizes the reward, and at the same time to *explore* actions that better
 253 inform the true Θ , which in turn will help make better decisions to achieve higher long-term rewards.
 254 Specifically, given $\hat{\Theta}_t$ and the covariate $\mathbf{x}_{t+1,j}$ for $j = 1, \dots, L$, we look for an action $\hat{\mathbf{a}}_{t+1}$ in the
 255 action space \mathcal{A}_t that maximizes the total rewards across L objects:

$$\hat{\mathbf{a}}_{t+1} := \arg \max_{\mathbf{a} \in \mathcal{A}_t} \left\{ \sum_{j=1}^L \mathbf{a}^\top \hat{\Theta}_t \mathbf{x}_{t+1,j} \right\}. \quad (5)$$

256 At a subset of times, we further perturb $\hat{\mathbf{a}}_{t+1}$ for the purpose of exploration by adding random noise
 257 to each coordinate as follows: $\mathbf{a}_{t+1} = \hat{\mathbf{a}}_{t+1} + \delta_{t+1}$ where $\delta_{t+1} \sim N(\mathbf{0}_{d_a}, h\mathbf{I}_{d_a})$ and h is a tuning
 258 parameter. In our current algorithm, we perform this perturbation at times $t \in \{\lfloor w^{\frac{3}{2}} \rfloor : w \in \mathbb{Z}_+\}$.

259 The intuition for this particular choice ($\lfloor w^{\frac{3}{2}} \rfloor$) is to explore more in the initial stage and exploit less
 260 in the later stage of the algorithm. To be specific, there are approximately $T^{\frac{2}{3}}$ steps for exploration
 261 before time T . The density of exploration at a small time frame around T is $T^{-\frac{1}{3}}$, which goes to zero
 262 as $T \rightarrow \infty$. Note that the exponent need not be $\frac{3}{2}$, but can be any number strictly larger than 1; this
 263 choice will affect the decay rate of the regret, as will be discussed later in Remark 3.4.

264 The perturbation form in exploration can be changed as well. For each exploration step, one can
 265 let $\delta_{t+1} \sim N(\mathbf{0}_{d_a}, \text{diag}(\hat{\tau}_t))$ where each element of $\hat{\tau}_t$ is the coordinate-wise standard error of the
 266 previous actions $\{\mathbf{a}_i\}_{i=1}^t$. In this way, we avoid tuning parameter h while taking the right scale.
 267 Finally, we update the action space \mathcal{A}_{t+1} according to \mathbf{a}_{t+1} . For example, if the action space
 268 $\mathcal{A}_t \in \mathbb{R}^{d_a}$ can be defined by an upper limit $\bar{\mathbf{a}}_t$ and a lower limit $\underline{\mathbf{a}}_t$, then we simply expand the action
 269 space by pushing the boundary of each coordinate to $\mathbf{a}_{t+1,j}$ if $\mathbf{a}_{t+1,j} \notin [\underline{\mathbf{a}}_{t,j}, \bar{\mathbf{a}}_{t,j}]$ for $j = 1, \dots, d_a$.

270 **Remark 3.1** (Adaptivity and robustness). It should be noted that our algorithm is adaptive (w.r.t.
 271 T), and also robust, especially compared with explore-then-commit type algorithms. It does not
 272 require specifying the T -dependent tuning parameters (e.g., exploration period T_1), and it updates

the representation matrix across the entire time horizon, making it more suitable for online learning. Moreover, it does not require knowing or pre-specifying the target rank r .

Remark 3.2 (Interpretability). To take advantage of the interpretability of our model, we can further explore the structure of $\hat{\Theta}_t$. Specifically, we can apply singular value decomposition (SVD) on $\hat{\Theta}_t$ to explore the underlying latent structure of arms and covariates through the left and right singular vectors. One can further rotate the singular vectors using techniques in factor analysis such as Varimax [20, 44] so as to obtain a sparse/simplified loading structure for easier interpretation.

3.2 Theoretical Results

In this section, we state a theorem that provides a non-asymptotic instance-dependent bound on the expected regret associated with Algorithm 1. It shows that in the worst-case and for any dimensions, the expected regret decays to zero as $T^{-1/6} \log T$.

Our analysis applies to an instantiation of Algorithm 1 in which the exploratory actions are chosen as

$$\mathbf{a}_t = \hat{\mathbf{a}}_t + \delta_t \quad \text{for each } t \in \{\lfloor w^{\frac{3}{2}} \rfloor : w \in \mathbb{Z}_+\},$$

where $\delta_t \sim N(\mathbf{0}_{d_a}, h\mathbf{I}_{d_a})$, $\mathcal{A}_t = \{\mathbf{a} \in \mathbb{R}^{d_a} : \|\mathbf{a}\| \leq 1\}$. Finally, our statement involves a burn-in period $B_{\text{init}} = C_{h,L,\lambda_0}(d_x + d_a)^6 (\log(d_x + d_a))^3$.

Theorem 3.3. Suppose that Θ has rank r , we observe covariates $\mathbf{x}_{t,j} \stackrel{i.i.d}{\sim} N(\mathbf{0}_{d_x}, \mathbf{I}_{d_x})$, and the reward errors $\varepsilon_{t,j} \stackrel{i.i.d}{\sim} N(0, \sigma^2)$ in equation (1). Then there are universal constants $\{c_j\}_{j=1}^4$ such that for all $T \geq B_{\text{init}}$, the expected regret is bounded as

$$\begin{aligned} \mathcal{R}^\pi(T) \leq & \frac{c_1}{T} \sqrt{Ld_x} \|\Theta^*\|_2 B_{\text{init}} + \frac{c_2 \log T}{T} \sqrt{Ld_x} \|\Theta^*\|_2 \\ & + \frac{c_3}{T^{1/6}} \frac{\lambda_0 \sqrt{2rd_x L}}{h^2} + \frac{c_4 \log T}{T^{1/6}} \frac{(d_x + 2 \log L) \sigma}{h^2}. \quad (6) \end{aligned}$$

Remark 3.4 (Convergence rate). An intuitive understanding of Theorem 3.3 is that the expected regret converges to zero at least as quickly as $\frac{\log T}{T^{1/6}}$ as T tends to infinity. This rate is relatively slow, and we suspect that our regret analysis can be tightened. The convergence rate depends on the frequency of the exploration which depends on the exponent $\frac{3}{2}$ in the exploration set, $\{\lfloor w^{\frac{3}{2}} \rfloor : w \in \mathbb{Z}_+\}$. The exponent can be tuned (i.e., replaced by any number larger than 1).

Remark 3.5 (“Burn-in” term). The first term in the bound (6) is a “burn-in” term, where the algorithm is gaining knowledge of Θ from scratch. We do not impose any assumptions on these starting steps so that we have a relatively conservative “burn-in” term. In practice, we can leverage historical data to obtain an initial estimation of Θ so that the “burn-in” term can be much smaller.

The order of the “burn-in” term depends on the exponent of the w in the exploration set — the more exploration there is, the smaller the “burn-in” term. The exponent can be chosen depending on how ample the historical data is.

Remark 3.6 (Constant C_{h,L,λ_0} of B_{init}). While constant C_{h,L,λ_0} depends on h, L, λ_0 , the primary dependency is actually on h and L . The order of λ_0 in terms of dimensions and noise level is $\sigma \sqrt{d_x}$. We do not assume the order of λ_0 or bound it with a high probability bound in order to show its role in time-averaged expected cumulative regret. If we utilize the order $\sigma \sqrt{d_x}$, then C_{h,L,λ_0} can be replaced by a constant depending on h and L only.

Remark 3.7 (Dependence on dimensions d_a, d_x and rank r). When T is small, the “burn-in” term (the first term) dominates. It depends on T and the dimensions but not the rank. As T grows, the last two terms dominate. Recall from Remark 3.6 that λ_0 is of order $\sigma \sqrt{d_x}$, so the third term depends on T, d_x and r but not d_a ; it has the order $\Omega(d_x \sqrt{r} T^{-\frac{1}{6}})$. The last term is of the order $\Omega(d_x T^{-\frac{1}{6}} \log T)$. In terms of T , these last two terms are of the same order up to $\log T$.

Remark 3.8 (Adaptivity). Algorithm 1 is adaptive as it does not require knowing T a priori (except the ending point) as mentioned in Remark 3.1; moreover, the non-asymptotic bounds hold for all T . This adaptivity is of both theoretical interest and practical importance. Adaptivity overcomes the limitations of the traditional bandit framework, which possibly favors good performance at a specific T at the expense of other values; this leads to algorithms involving T -dependent tuning parameters. In practice, it is preferable to have algorithms that do not require such tuning yet consistently perform well across all T .

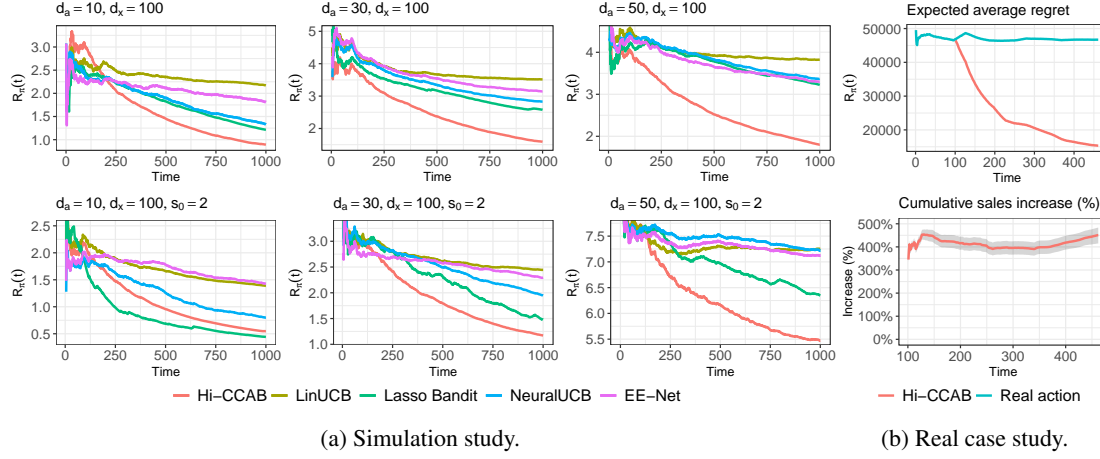


Figure 1. (a) The expected average regret under the non-sparse (first row) and sparse (second row) settings; (b) Performance of Hi-CCAB compared with real actions. The band covers the 5-th and 95-th quantiles.

319 *Remark 3.9* (Assumptions). To convey the main idea in a simple way, we have chosen to enforce
 320 relatively stringent assumptions. Neither the normality assumption nor the shape of the constraint set
 321 is essential.

322 **Proof sketch** Due to space constraints, we limit ourselves to an outline of the proof of Theorem
 323 3.3. There are two major steps: (1) bounding the estimation error for the low-rank representation
 324 matrix estimator; (2) bounding the expected regret. See Appendix A for the full proof.

325 (1) *High-probability bound on estimating Θ* : Introduce the shorthand $\Delta_t = \hat{\Theta}_t - \Theta$. We show that
 326 for a large t ,

$$P\left(\|\Delta_t\|_F \geq \frac{9t^{\frac{1}{3}}(2 + \sqrt{t})(d_x + 2 \log tL)\sigma}{th^2} + 6\lambda_0 \frac{\sqrt{2r}}{h^2 t^{\frac{1}{6}}}\right) \leq \frac{3}{t} + \frac{2}{t^2} + \frac{3}{Lt} + \frac{2}{L^3 t^3} + \frac{1}{t^2 L}.$$

327 The technical challenges in establishing this result lie in the fact that the actions taken are
 328 based on past data, and also affect future data, so that the dataset is very different from i.i.d.
 329 Thus, classical results on matrix completion, valid for i.i.d. or weakly dependent data, are
 330 not applicable here. Thus, our analysis requires some careful use of conditional expectations,
 331 martingales, and empirical process so as to separate out different sources of randomness (i.e.,
 332 $\delta_1, \dots, \delta_t, x_1, \dots, x_t$) to derive the bounds. Lemma A.1 establishes a form of restricted-
 333 strong-convexity (RSC) [38] for the objective function. Lemma A.2 establishes a Lipschitz bound
 334 for a portion of the objective function. Further analysis of the nuclear-norm-penalized sum of
 335 squares with the two lemmas and low-rank properties gives the tail bound of the estimation error.

336 (2) *Bounding the expected regret*: At each round t , we define the event $\mathcal{E}_t = \{\|\Delta_t\|_F \leq$
 337 $\frac{9t^{\frac{1}{3}}(2 + \sqrt{t})(d_x + 2 \log tL)\sigma}{th^2} + 6\lambda_0 \frac{\sqrt{2r}}{h^2 t^{\frac{1}{6}}}\}$. From step (1), we know for large t , $\mathbb{P}(\mathcal{E}_t^c) \leq \frac{3}{t} +$
 338 $\frac{2}{t^2} + \frac{3}{Lt} + \frac{2}{L^3 t^3} + \frac{1}{t^2 L^2}$. Consider the expectation of the regret on \mathcal{E}_t and \mathcal{E}_t^c separately and both
 339 terms vanish with t at the polynomial rate.

340 4 Experimental evaluations

341 In this section, we report experimental evaluations of our model and algorithm applied to both
 342 synthetic and some large-scale real-world datasets. First, we conduct simulation studies to compare
 343 the proposed Hi-CCAB with LinUCB [33], Lasso Bandit [6], NeuralUCB [55] and EE-Net [5]; we
 344 then study the joint assortment-pricing problem on the e-commerce platform for one of the largest
 345 instant noodle producers in China. Details on the tuning parameters of each algorithm and additional
 346 results comparing the low-rank bandit [22] and the case study are provided in Appendices B–C.

347 **Simulation study** We consider the multi-armed linear bandit setup, a special case of our model
 348 with $\Theta = (\beta_1, \beta_2, \dots, \beta_m)^\top$ so that each row of Θ is the parameter of each arm for the multi-arm
 349 contextual bandit. Specifically, we set the number of arms $d_a = \{10, 30, 50\}$ and the dimension of
 350 covariates $d_x = 100$. For Θ , we consider a non-sparse and sparse case. For the non-sparse case, we
 351 generate $\Theta = UDV^\top$ where $U \in \mathbb{R}^{d_a \times r}$, $V \in \mathbb{R}^{d_x \times r}$ ($r = 5$), and D is a diagonal matrix with
 352 $(1, .9, .9, .8, .5)$ as the diagonal entries. All entries of U and V are first generated from i.i.d. $N(0, 1)$,

and then applied Gram–Schmidt to make each column orthogonal. The matrix \mathbf{U} is scaled to have length $\sqrt{d_a}$ so that the rewards are comparable across different d_a ’s. For the sparse case, each row of Θ is set as zero except for $s_0 = 2$ randomly selected elements that are drawn from $N(0, 1)$. We generate the covariate $\mathbf{x} \stackrel{i.i.d}{\sim} N(\mathbf{0}, \mathbf{I}_{d_x})$ and the rewards from (1) with $\sigma = 0.1$.

Figure 1a shows the regret (averaged over 50 simulations). For the non-sparse case, Hi-CCAB converges faster than all other methods. The advantage of Hi-CCAB is more pronounced when the dimension of arms becomes larger. For the sparse case, which is not to the advantage of Hi-CCAB, when the dimension of arms is relatively small ($d_a = 10$), Lasso Bandit converges faster but the gap between Hi-CCAB and Lasso Bandit is small. As the number of arms increases, Hi-CCAB outperforms all other methods.

Assortment-pricing case study. We also applied our methods to a large-scale dataset of Chinese noodle company (as previously described in the introduction). The original data includes daily sales of 176 products across 369 cities from March 1st, 2021 to May 31st, 2022 ($T = 456$ days), aggregated by 31 provinces. Products have single or assorted flavors (13 possibilities) with varying counts. Their assortment, prices, and promotions changed daily, being uniform across locations. The maximum number of products on the homepage is $K = 30$. Considering one combination as an arm, the total possible combinations are $\binom{176}{30}$, resulting in an extremely high-dimensional arm space, which is unsuitable for most multi-arm bandit algorithms.

To apply Hi-CCAB, we specify the arms \mathbf{a}_t and the covariate vectors $\{\mathbf{x}_{t,j}\}_{j=1}^{L=31}$ at given time t following the setup in Example 1. The action is represented as

$$\mathbf{a} = (\mathbf{f}_1, \mathbf{f}_1^2, p_1, p_1^2, S_1, S_1^2, \dots, \mathbf{f}_K, \mathbf{f}_K^2, p_K, p_K^2, S_K, S_K^2, 1) \in \mathbb{R}^{2(m+2)K+1}$$

where $d_a = 2(m+2)K+1 = 901$. Here $\mathbf{f}_k = (f_{k,1}, \dots, f_{k,m})$ is a vector of non-negative integers to denote the counts of $m = 13$ flavors, p_k is the price, S_k is the indicator of promotion of product k , and \mathbf{f}_k^2 is the element-wise quadratics. The covariate $\mathbf{x}_{t,j} \in \mathbb{R}^{50}$ for location j includes dummy variables of 31 provinces, the year 2021/2022, 12 months, weekdays, and an indicator of the annual sales event on Jun 18 and Nov 11. More details are deferred to Appendix C.

To run simulations using the dataset, we first create a pseudo-truth model by estimating Θ and σ using all data of 456 days, considering them as the pseudo-ground truth. We check our model assumption (1) against our data before preceding the formal analysis and further examine the structure of the representation matrix Θ in Appendix C. We evaluate the performance of Hi-CCAB in terms of the cumulative regret (2) and the percentage gain of the cumulative sales by comparing with the original actions, since no existing bandit algorithm is applicable to this problem.

Figure 1b shows the time-averaged cumulative regret (averaged over 100 simulations) and the percentage gain in cumulative sales compared to the real sales. The expected average regret of Hi-CCAB converges to zero while that of original actions remains flat. In terms of percentage gain in cumulative sales, Hi-CCAB boosts cumulative sales by more than 4 times.

5 Discussion

The growing demand for online decision-making has led to increased interest in the bandit problem among theoreticians and practitioners. Despite the richness of the bandit literature, to date, there has been relatively little work on contextual bandits in which both the covariate and action spaces are high-dimensional. In this paper, we have argued that many applications of bandit have this “doubly” high-dimensional nature, and we have provided a structured matrix model for capturing interactions between covariates and actions. This model is reasonably general, including a number of structured bandit models as special cases, but also interpretable. We propose an efficient algorithm Hi-CCAB that interleaves steps of low-rank matrix estimation with exploration/exploitation, and we proved a non-asymptotic upper bound on its expected regret. The generality and flexibility of our model enable its application to the joint assortment-pricing problem, each of which has been studied extensively in operations research but not as a joint optimization problem. In a real case study with the largest instant noodle producers in China, our method can boost sales by a factor of four times, while also provide insights into the underlying structure of the effect on the reward of the arms and covariates such as purchasing behaviors. Future work could improve regret analysis by tightening the regret bound and establishing a potentially matching lower bound.

References

- [1] Abbasi-Yadkori, Y., Pal, D., and Szepesvari, C. (2012). Online-to-confidence-set conversions and application to sparse stochastic bandits. In *Artificial Intelligence and Statistics*, pages 1–9. PMLR.
- [2] Agrawal, R. (1995). The continuum-armed bandit problem. *SIAM journal on control and optimization*, 33(6):1926–1951.
- [3] Agrawal, S., Avadhanula, V., Goyal, V., and Zeevi, A. (2019). Mnl-bandit: A dynamic learning approach to assortment selection. *Operations Research*, 67(5):1453–1485.
- [4] Ban, G.-Y. and Keskin, N. B. (2021). Personalized dynamic pricing with machine learning: High-dimensional features and heterogeneous elasticity. *Management Science*, 67(9):5549–5568.
- [5] Ban, Y., Yan, Y., Banerjee, A., and He, J. (2022). Ee-net: Exploitation-exploration neural networks in contextual bandits.
- [6] Bastani, H. and Bayati, M. (2020). Online decision making with high-dimensional covariates. *Operations Research*, 68(1):276–294.
- [7] Cai, T. T. and Zhang, A. (2018). Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *The Annals of Statistics*, 46(1):60–89.
- [8] Candes, E. J. and Plan, Y. (2010). Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936.
- [9] Caro, F. and Gallien, J. (2007). Dynamic assortment with demand learning for seasonal consumer goods. *Management science*, 53(2):276–292.
- [10] Chen, N. and Gallego, G. (2021). Nonparametric pricing analytics with customer covariates. *Operations Research*, 69(3):974–984.
- [11] Chen, X., Shi, C., Wang, Y., and Zhou, Y. (2021). Dynamic assortment planning under nested logit models. *Production and Operations Management*, 30(1):85–102.
- [12] Chen, X. and Wang, Y. (2017). A note on a tight lower bound for mnl-bandit assortment selection models. *arXiv preprint arXiv:1709.06109*.
- [13] Chen, Y., Xie, M., Liu, J., and Zhao, K. (2022). Interconnected neural linear contextual bandits with ucb exploration. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 169–181. Springer.
- [14] Debon, R., Coleone, J. D., Bellei, E. A., and De Marchi, A. C. B. (2019). Mobile health applications for chronic diseases: A systematic review of features for lifestyle improvement. *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, 13(4):2507–2512.
- [15] den Boer, A. V. and Zwart, B. (2014). Simultaneously learning and optimizing using controlled variance pricing. *Management science*, 60(3):770–783.
- [16] Féraud, R., Allesiardo, R., Urvoy, T., and Clérot, F. (2016). Random forest for the contextual bandit problem. In *Artificial intelligence and statistics*, pages 93–101. PMLR.
- [17] Hao, B., Lattimore, T., and Wang, M. (2020). High-dimensional sparse linear bandits. *Advances in Neural Information Processing Systems*, 33:10753–10763.
- [18] Hu, J., Chen, X., Jin, C., Li, L., and Wang, L. (2021). Near-optimal representation learning for linear bandits and linear rl. In *International Conference on Machine Learning*, pages 4349–4358. PMLR.
- [19] Jun, K.-S., Willett, R., Wright, S., and Nowak, R. (2019). Bilinear bandits with low-rank structure. In *International Conference on Machine Learning*, pages 3163–3172. PMLR.
- [20] Kaiser, H. F. (1958). The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23(3):187–200.

- [21] Kallus, N. and Udell, M. (2020). Dynamic assortment personalization in high dimensions. *Operations Research*, 68(4):1020–1037.
- [22] Kang, Y., Hsieh, C.-J., and Lee, T. C. M. (2022). Efficient frameworks for generalized low-rank matrix bandit problems. *Advances in Neural Information Processing Systems*, 35:19971–19983.
- [23] Keskin, N. B. and Zeevi, A. (2014). Dynamic pricing with an unknown demand model: Asymptotically optimal semi-myopic policies. *Operations research*, 62(5):1142–1167.
- [24] Kim, G.-S. and Paik, M. C. (2019). Doubly-robust lasso bandit. *Advances in Neural Information Processing Systems*, 32.
- [25] Kim, J.-h. and Vojnovic, M. (2021). Scheduling servers with stochastic bilinear rewards. *arXiv preprint arXiv:2112.06362*.
- [26] Kleinberg, R. (2004). Nearly tight bounds for the continuum-armed bandit problem. In *Proceedings of the 17th International Conference on Neural Information Processing Systems*, NIPS’04, pages 697–704, Cambridge, MA, USA. MIT Press.
- [27] Kleinberg, R., Slivkins, A., and Upfal, E. (2019). Bandits and experts in metric spaces. *J. ACM*, 66(4).
- [28] Kök, A. G., Fisher, M. L., and Vaidyanathan, R. (2008). Assortment planning: Review of literature and industry practice. *Retail supply chain management*, 122(1):99–153.
- [29] Krishnamurthy, A., Langford, J., Slivkins, A., and Zhang, C. (2020). Contextual bandits with continuous actions: Smoothing, zooming, and adapting. *The Journal of Machine Learning Research*, 21(1):5402–5446.
- [30] Kveton, B., Szepesvári, C., Rao, A., Wen, Z., Abbasi-Yadkori, Y., and Muthukrishnan, S. (2017). Stochastic low-rank bandits. *arXiv preprint arXiv:1712.04644*.
- [31] Lale, S., Azzizadenesheli, K., Anandkumar, A., and Hassibi, B. (2019). Stochastic linear bandits with hidden low rank structure. *arXiv preprint arXiv:1901.09490*.
- [32] Lattimore, T. and Szepesvári, C. (2020). *Bandit algorithms*. Cambridge University Press.
- [33] Li, L., Chu, W., Langford, J., and Schapire, R. E. (2010). A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670.
- [34] Lu, T., Pál, D., and Pál, M. (2010). Contextual multi-armed bandits. In *Proceedings of the Thirteenth international conference on Artificial Intelligence and Statistics*, pages 485–492. JMLR Workshop and Conference Proceedings.
- [35] Lu, Y., Meisami, A., and Tewari, A. (2021). Low-rank generalized linear bandit problems. In *International Conference on Artificial Intelligence and Statistics*, pages 460–468. PMLR.
- [36] Miao, S. and Chao, X. (2021). Dynamic joint assortment and pricing optimization with demand learning. *Manufacturing & Service Operations Management*, 23(2):525–545.
- [37] Negahban, S. and Wainwright, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, 39(2):1069–1097.
- [38] Negahban, S. and Wainwright, M. J. (2012). Restricted strong convexity and (weighted) matrix completion: Optimal bounds with noise. *Journal of Machine Learning Research*, 13:1665–1697.
- [39] Papini, M., Tirinzoni, A., Restelli, M., Lazaric, A., and Pirotta, M. (2021). Leveraging good representations in linear contextual bandits. In *International Conference on Machine Learning*, pages 8371–8380. PMLR.
- [40] Qiang, S. and Bayati, M. (2016). Dynamic pricing with demand covariates. *arXiv preprint arXiv:1604.07463*.

- 492 [41] Recht, B., Fazel, M., and Parrilo, P. A. (2010). Guaranteed minimum-rank solutions of linear
493 matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501.
- 494 [42] Rizk, G., Thomas, A., Colin, I., Laraki, R., and Chevaleyre, Y. (2021). Best arm identification
495 in graphical bilinear bandits. In *International Conference on Machine Learning*, pages 9010–9019.
496 PMLR.
- 497 [43] Robbins, H. (1952). Some aspects of the sequential design of experiments. *Bulletin of the*
498 *American Mathematical Society*, 58(5):527–535.
- 499 [44] Rohe, K. and Zeng, M. (2020). Vintage factor analysis with varimax performs statistical
500 inference. *arXiv preprint arXiv:2004.05387*.
- 501 [45] Sauré, D. and Zeevi, A. (2013). Optimal dynamic assortment planning with demand learning.
502 *Manufacturing & Service Operations Management*, 15(3):387–404.
- 503 [46] Slivkins, A. (2011). Contextual bandits with similarity information. In *Proceedings of the*
504 *24th annual Conference On Learning Theory*, pages 679–702. JMLR Workshop and Conference
505 Proceedings.
- 506 [47] Srebro, N., Alon, N., and Jaakkola, T. S. (2005). Generalization error bounds for collaborative
507 prediction with low-rank matrices. In *Neural Information Processing Systems (NIPS)*, Vancouver,
508 Canada.
- 509 [48] Tropp, J. A. (2012). User-friendly tail bounds for sums of random matrices. *Foundations of*
510 *computational mathematics*, 12(4):389–434.
- 511 [49] Turğay, E., Bulucu, C., and Tekin, C. (2020). Exploiting relevance for online decision-making
512 in high-dimensions. *IEEE Transactions on Signal Processing*, 69:1438–1451.
- 513 [50] Tyagi, H., Stich, S. U., and Gärtner, B. (2016). On two continuum armed bandit problems in
514 high dimensions. *Theory of Computing Systems*, 58(1):191–222.
- 515 [51] Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48.
516 Cambridge University Press.
- 517 [52] Xu, K. and Bastani, H. (2021). Learning across bandits in high dimension via robust statistics.
518 *arXiv preprint arXiv:2112.14233*.
- 519 [53] Xu, P., Wen, Z., Zhao, H., and Gu, Q. (2022). Neural contextual bandits with deep representation
520 and shallow exploration.
- 521 [54] Yang, J., Hu, W., Lee, J. D., and Du, S. S. (2020). Impact of representation learning in linear
522 bandits. *arXiv preprint arXiv:2010.06531*.
- 523 [55] Zhou, D., Li, L., and Gu, Q. (2020). Neural contextual bandits with ucb-based exploration. In
524 *International Conference on Machine Learning*, pages 11492–11502. PMLR.

525 A Proof of Theorem 3.3

526 In this proof, we denote the true parameter as Θ^* .

527 Let $\mathcal{L}_T(\Theta) := \frac{1}{2LT} \sum_{t=1}^T \sum_{l=1}^L (\mathbf{a}_{t,l}^\top \Theta \mathbf{x}_{t,l} - y_{t,l})^2$. Then we have the following lemmas that we
528 will prove later.

529 **Lemma A.1.** Suppose all the assumptions in Theorem 3.3 holds. Denote $\mathcal{E}_t(\Delta) = \mathcal{L}_t(\Theta^* + \Delta) -$
530 $\mathcal{L}_t(\Theta^*) - \langle \nabla \mathcal{L}_t(\Theta^*), \Delta \rangle$. Then for any $t \geq 2$, with probability at least $1 - \frac{1}{Lt} - \frac{2}{t} - \frac{1}{t^2}$,

$$\mathcal{E}_t(\Delta) \geq \frac{\lfloor t^{\frac{2}{3}} \rfloor}{2t} h^2 \|\Delta\|_F^2 - 14t^{-\frac{2}{3}}(h + h^2)(2d_x + 2d_a + 6 \log t + 6 \log L)^2 \log t \|\Delta\|_2^2. \quad (7)$$

531 **Lemma A.2.** Suppose all the assumptions in Theorem 3.3 holds. With probability at least $1 - \frac{1}{t^2 L} -$
532 $\frac{2}{L^3 t^3} - \frac{2}{Lt} - \frac{1}{t} - \frac{1}{t^2}$, the following holds for all $\Delta \in \mathbb{R}^{d_a \times d_x}$

$$\begin{aligned} |\langle \nabla \mathcal{L}_t(\Theta^*), \Delta \rangle| \leq & \|\Delta\|_F \frac{\sigma(6 + 3\sqrt{t})(d_x + 2 \log tL)}{t} + \\ & \left(2h\sigma t^{-2/3} \log t \sqrt{\frac{\max\{d_a, d_x\} \log(d_a + d_x)}{L}} + \right. \\ & \left. \frac{8h\sigma}{t} \sqrt{\log(tL)} \sqrt{(d_x + 3 \log(Lt))(d_a + 3 \log t)} (\log(d_x + d_a) + 2 \log t) \right) \|\Delta\|_*. \end{aligned} \quad (8)$$

533 Recall the definition of $\hat{\Theta}_t$, we know that

$$\mathcal{L}_t(\hat{\Theta}_t) + \lambda_t \|\hat{\Theta}_t\|_* \leq \mathcal{L}_t(\Theta^*) + \lambda_t \|\Theta^*\|_*. \quad (9)$$

534 Denote $\Delta_t = \hat{\Theta}_t - \Theta^*$ and for notation simplicity we will drop the subscript t for Δ_t in the
535 following when there is no confusion. Equation (9) then implies that

$$\mathcal{E}_t(\Delta) \leq -\langle \nabla \mathcal{L}_t(\Theta^*), \Delta \rangle + \lambda_t (\|\Theta^*\|_* - \|\Theta^* + \Delta\|_*). \quad (10)$$

536 Suppose the singular value decomposition of Θ^* is $\Theta^* = U S V^\top$, where S is an $r \times r$ diagonal
537 matrix. Let U_\perp be an $d_a \times (d_a - r)$ matrix satisfying $(U, U_\perp)(U, U_\perp)^\top = I_{d_a}$. We define V_\perp
538 similarly.

539 Denote $\Delta_\perp = U_\perp^\top \Delta V_\perp$. Then $\|\Theta^* + \Delta\|_* \geq \|\Theta^* + \Delta_\perp\|_* - \|\Delta - \Delta_\perp\|_* = \|\Theta^*\|_* + \|\Delta_\perp\|_* -$
540 $\|\Delta - \Delta_\perp\|_* \geq \|\Theta^*\|_* + \|\Delta_\perp\|_* - \sqrt{2r} \|\Delta - \Delta_\perp\|_F$.

541 Going back to inequality (10), and combing with Lemma A.1 and Lemma A.2, we have, with
542 probability at least $1 - \frac{3}{t} - \frac{2}{t^2} - \frac{3}{Lt} - \frac{2}{L^3 t^3} - \frac{1}{t^2 L}$, the following holds

$$\begin{aligned} & \left(\frac{\lfloor t^{\frac{2}{3}} \rfloor}{2t} h^2 - 14t^{-\frac{2}{3}}(h + h^2)(2d_x + 2d_a + 6 \log t + 6 \log L)^2 \log t \right) \|\Delta\|_F^2 \\ & \leq \|\Delta\|_F \frac{\sigma(6 + 3\sqrt{t})(d_x + 2 \log tL)}{t} + \\ & \quad \left(\frac{8h\sigma}{t} \sqrt{\log(tL)} \sqrt{(d_x + 3 \log(Lt))(d_a + 3 \log t)} (\log(d_x + d_a) + 2 \log t) + \right. \\ & \quad \left. 2h\sigma t^{-2/3} \log t \sqrt{\frac{\max\{d_a, d_x\} \log(d_a + d_x)}{L}} \right) (\|\Delta - \Delta_\perp\|_* + \|\Delta_\perp\|_*) \\ & \quad + \lambda_0 \frac{\sqrt{t}}{t} \sqrt{2r} \|\Delta\|_F - \lambda_0 \frac{\sqrt{t}}{t} \|\Delta_\perp\|_*. \end{aligned} \quad (11)$$

543 Note that $\|\Delta - \Delta_\perp\|_* \leq \sqrt{2r} \|\Delta - \Delta_\perp\|_F$, divide both side with $\|\Delta\|_F$ and multiply both sides
544 with $3t^{\frac{1}{3}}/h^2$. Suppose B_{init} is the smallest integer such that for all $t \geq B_{\text{init}}$, the following inequalities

545 hold

$$t \geq 8, \quad (12a)$$

$$t^{\frac{1}{3}} \geq 12 \times 14 \left(1 + \frac{1}{h}\right) (2d_x + 2d_a + 6 \log t + 6 \log L)^2 \log t, \quad (12b)$$

$$\lambda_0 \geq \frac{8h\sigma}{t^{\frac{1}{2}}} \sqrt{\log(tL)} \sqrt{(d_x + 3 \log(Lt))(d_a + 3 \log t)(\log(d_x + d_a) + 2 \log t) + 2h\sigma t^{-1/6} \log t \sqrt{\frac{\max\{d_a, d_x\} \log(d_a + d_x)}{L}}}. \quad (12c)$$

Clearly, B_{init} is well defined because there is a constant C_{h,L,λ_0} depending on L, h and λ_0 such that for

$$t \geq C_{h,L,\lambda_0} (d_x + d_a)^6 (\log(d_x + d_a))^3,$$

Inequalities (12) holds. Therefore,

$$B_{\text{init}} \leq C_{h,L,\lambda_0} (d_x + d_a)^6 (\log(d_x + d_a))^3.$$

546 Then we have for $t \geq B_{\text{init}}$,

$$\|\hat{\Theta}_t - \Theta^*\|_F \leq \frac{9t^{\frac{1}{3}}(2 + \sqrt{t})(d_x + 2 \log tL)\sigma}{th^2} + 6\lambda_0 \frac{\sqrt{2r}}{h^2 t^{\frac{1}{6}}}. \quad (13)$$

Next we will proceed to bound the regret. Denote the event that equation (13) holds to be \mathcal{E}_t and its complement as \mathcal{E}_t^c . Then $\mathbb{P}(\mathcal{E}_t^c) \leq \frac{3}{t} + \frac{2}{t^2} + \frac{3}{Lt} + \frac{2}{L^3 t^3} + \frac{1}{L^2 t^2}$ for $t \geq B_{\text{init}}$, and $(\mathcal{E}_t^c, \hat{\Theta}_t) \perp\!\!\!\perp \mathbf{b}_{t+1}$.

Let the oracle optimal action at time t be \mathbf{a}_t^* and $\mathbf{b}_t = \sum_{l=1}^L \mathbf{x}_{t,l}$. Apparently, the cumulative regret $\mathcal{TR}^\pi(T)$ is the sum of two parts. 1. The cumulative regret until the time B_{init} :

$$B_T = \mathbb{E} \left(\sum_{t=0}^{\min(B_{\text{init}}-1, T-1)} \sum_{l=1}^L (\mathbf{a}_{t+1}^{*\top} \Theta^* \mathbf{x}_{t+1,l} - \mathbf{a}_{t+1}^\top \Theta^* \mathbf{x}_{t+1,l}) \right),$$

and 2. the cumulative regret after the time B_{init} :

$$D_T = \mathbb{E} \left(\sum_{t=\min(B_{\text{init}}, T)}^{T-1} \sum_{l=1}^L (\mathbf{a}_{t+1}^{*\top} \Theta^* \mathbf{x}_{t+1,l} - \mathbf{a}_{t+1}^\top \Theta^* \mathbf{x}_{t+1,l}) \right).$$

547 We start by bounding D_T .

548 By definition, we have

$$\mathbf{a}_{t+1}^{*\top} = \frac{\Theta^* \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2}, \quad (14)$$

549 and

$$\mathbb{E}(\mathbf{a}_{t+1}^\top \mid \{\mathbf{x}_{j,l}\}_{j \leq t+1, l \leq L}, \{\mathbf{a}_j\}_{j \leq t}, \{y_{j,l}\}_{j \leq t, l \leq L}) = \frac{\hat{\Theta}_t \mathbf{b}_{t+1}}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2}. \quad (15)$$

550 Therefore, by taking conditional expectation on $\{\mathbf{x}_{t+1,l}\}_{l \leq L}$ for D_T first and plugging in Equations (14) (15), we have

$$D_T = \mathbb{E} \left(\sum_{t=\min(B_{\text{init}}, T)}^{T-1} \left\langle \frac{\Theta^* \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2} - \frac{\hat{\Theta}_t \mathbf{b}_{t+1}}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2}, \Theta^* \mathbf{b}_{t+1} \right\rangle \right). \quad (16)$$

552 Further split the each term in the summation of D_T into two parts based on \mathcal{E}_t and \mathcal{E}_t^c gives

$$D_T = \sum_{t=B_{\text{init}} \wedge T}^{T-1} \mathbb{E} \left(\left\langle \frac{\Theta^* \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2} - \frac{\hat{\Theta}_t \mathbf{b}_{t+1}}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2}, \Theta^* \mathbf{b}_{t+1} \right\rangle \mathbb{1}\{\mathcal{E}_t\} \right) + \mathbb{E} \left(\left\langle \frac{\Theta^* \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2} - \frac{\hat{\Theta}_t \mathbf{b}_{t+1}}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2}, \Theta^* \mathbf{b}_{t+1} \right\rangle \mathbb{1}\{\mathcal{E}_t^c\} \right). \quad (17)$$

553 We analyze the two parts separately.

554 For the first term, we have

$$\sum_{t=B_{\text{init}} \wedge T}^{T-1} \mathbb{E} \left(\left\langle \frac{\Theta^* \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2} - \frac{\hat{\Theta}_t \mathbf{b}_{t+1}}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2}, \Theta^* \mathbf{b}_{t+1} \right\rangle \mathbb{1}\{\mathcal{E}_t\} \right) \quad (18)$$

$$= \sum_{t=(B_{\text{init}} \wedge T)}^{T-1} \left(\mathbb{E} \left(\left\langle \frac{(\Theta^* - \hat{\Theta}_t) \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2} + \frac{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2 - \|\Theta^* \mathbf{b}_{t+1}\|_2}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2 \|\Theta^* \mathbf{b}_{t+1}\|_2} \hat{\Theta}_t \mathbf{b}_{t+1}, \Theta^* \mathbf{b}_{t+1} \right\rangle \mathbb{1}\{\mathcal{E}_t\} \right) \right) \quad (19)$$

$$\leq \sum_{t=(B_{\text{init}} \wedge T)}^{T-1} 2\mathbb{E}(\|\Delta\|_{op} \|\mathbf{b}_{t+1}\| \mathbb{1}\{\mathcal{E}_t\}) \leq \sum_{t=(B_{\text{init}} \wedge T)}^{T-1} 2\mathbb{E}(\|\Delta\|_F \|\mathbf{b}_{t+1}\| \mathbb{1}\{\mathcal{E}_t\}) \quad (20)$$

$$\leq 2 \sum_{t=(B_{\text{init}} \wedge T)}^{T-1} \left(\frac{9t^{\frac{1}{3}}(2 + \sqrt{t})(d_x + 2 \log tL)\sigma}{th^2} + 6\lambda_0 \frac{\sqrt{2r}}{h^2 t^{\frac{1}{6}}} \right) \sqrt{Ld_x} \quad (21)$$

$$\leq T^{\frac{5}{6}} \log T \left(1 + (d_x + 2 \log L) \left(\frac{1 + 5T^{-1/2}}{2 \log T} \right) + 5 \frac{1}{\sqrt{T}} \right) \frac{72\sigma}{5h^2} + \frac{72}{5} \lambda_0 \frac{\sqrt{2rd_x L}}{h^2} T^{\frac{5}{6}}. \quad (22)$$

555 For the second term, we have

$$\sum_{t=B_{\text{init}} \wedge T}^{T-1} \mathbb{E} \left(\left\langle \frac{\Theta^* \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2} - \frac{\hat{\Theta}_t \mathbf{b}_{t+1}}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2}, \Theta^* \mathbf{b}_{t+1} \right\rangle \mathbb{1}\{\mathcal{E}_t^c\} \right) \quad (23)$$

$$\stackrel{(i)}{\leq} \sum_{t=B_{\text{init}} \wedge T}^{T-1} 2\mathbb{E}(\|\Theta^* \mathbf{b}_{t+1}\|_2 \mathbb{1}\{\mathcal{E}_t^c\}) \stackrel{(ii)}{=} \sum_{t=B_{\text{init}} \wedge T}^{T-1} 2\mathbb{E}(\|\Theta^* \mathbf{b}_{t+1}\|_2) \mathbb{E}(\mathbb{1}\{\mathcal{E}_t^c\}) \quad (24)$$

$$\stackrel{(iii)}{\leq} \sum_{t=B_{\text{init}} \wedge T}^{T-1} 2\|\Theta^*\|_{op} \sqrt{\mathbb{E}(\|\mathbf{b}_{t+1}\|^2)} \left(\frac{3}{t} + \frac{2}{t^2} + \frac{3}{Lt} + \frac{2}{L^3 t^3} + \frac{1}{Lt^2} \right) \quad (25)$$

$$\stackrel{(iv)}{<} 7\|\Theta^*\|_{op} \sqrt{Ld_x} \log T, \quad (26)$$

556 where step (i) is by $\left\| \frac{\Theta^* \mathbf{b}_{t+1}}{\|\Theta^* \mathbf{b}_{t+1}\|_2} - \frac{\hat{\Theta}_t \mathbf{b}_{t+1}}{\|\hat{\Theta}_t \mathbf{b}_{t+1}\|_2} \right\|_2 \leq 2$, step (ii) follows from $(\mathcal{E}_t^c, \hat{\Theta}_t) \perp \mathbf{b}_{t+1}$, step
 557 (iii) uses Cauchy-Schwarz Inequality and $\mathbb{P}(\mathcal{E}_t^c) \leq \frac{3}{t} + \frac{2}{t^2} + \frac{3}{Lt} + \frac{2}{L^3 t^3} + \frac{1}{Lt^2}$ for $t \geq B_{\text{init}}$, step (iv)
 558 follows from elementary calculation.

559 Therefore,

$$D_T < T^{\frac{5}{6}} \log T \left(\left(1 + (d_x + 2 \log L) \left(\frac{1 + 5T^{-1/2}}{2 \log T} \right) + 5 \frac{1}{\sqrt{T}} \right) \sigma + \frac{\lambda_0 \sqrt{2rd_x L}}{\log T} \right) \frac{72}{5h^2} + 7\|\Theta^*\|_{op} \sqrt{Ld_x} \log T \quad (27)$$

560 Similar arguments also give

$$B_T = \mathbb{E} \left(\sum_{t=0}^{(B_{\text{init}}-1) \wedge (T-1)} \sum_{l=1}^L \mathbb{E}(\mathbf{a}_t^\top \Theta^* \mathbf{x}_{t,l} - \mathbf{a}_t^\top \Theta^* \mathbf{x}_{t,l}) \right) \leq B_{\text{init}} \times 2\sqrt{Ld_x} \|\Theta^*\|_{op}. \quad (28)$$

Therefore,

$$\begin{aligned} \mathcal{R}^\pi(T) &= \frac{B_T + D_T}{T} \leq 2\sqrt{Ld_x} \|\Theta^*\|_{op} B_{\text{init}} T^{-1} + 7\sqrt{Ld_x} \|\Theta^*\|_{op} \frac{\log T}{T} \\ &\quad + \frac{\log T}{T^{\frac{1}{6}}} \left(\left(1 + (d_x + 2\log L) \left(\frac{1 + 5T^{-1/2}}{2\log T} \right) + 5\frac{1}{\sqrt{T}} \right) \sigma + \frac{\lambda_0 \sqrt{2rd_x L}}{\log T} \right) \frac{72}{5h^2}. \end{aligned} \quad (29)$$

A.1 Proof of Lemma A.2

For simplicity of notation, we prove that for any $T \geq 2$, with probability at least $1 - \frac{1}{T^2 L} - \frac{2}{L^3 T^3} - \frac{2}{LT} - \frac{1}{T} - \frac{1}{T^2}$, the following holds for all $\Delta \in \mathbb{R}^{d_a \times d_x}$:

$$\begin{aligned} |\langle \nabla \mathcal{L}_T(\Theta^*), \Delta \rangle| &\leq \|\Delta\|_F \frac{\sigma(6 + 3\sqrt{T})(d_x + 2\log TL)}{T} + \left(2h\sigma T^{-2/3} \log T \sqrt{\frac{\max\{d_a, d_x\} \log(d_a + d_x)}{L}} + \right. \\ &\quad \left. \frac{8h\sigma}{T} \sqrt{\log(TL)} \sqrt{(d_x + 3\log(LT))(d_a + 3\log T)(\log(d_x + d_a) + 2\log T)} \right) \|\Delta\|_*. \end{aligned} \quad (30)$$

Recall that $y_{t,l} = \mathbf{a}_t^T \Theta^* \mathbf{x}_{t,l} + \sigma \varepsilon_{t,l}$, then

$$\begin{aligned} \nabla \mathcal{L}_T(\Theta^*) &= \frac{\sigma}{LT} \sum_{t=1}^T \sum_{l=1}^L -\varepsilon_{t,l} \mathbf{x}_{t,l} \mathbf{a}_t^\top \\ &= \frac{\sigma}{LT} \sum_{l=1}^L -\varepsilon_{1,l} \mathbf{x}_{1,l} \mathbf{a}_1^\top + \frac{\sigma}{LT} \sum_{t=2}^T \sum_{l=1}^L (-\varepsilon_{t,l} \mathbf{x}_{t,l} \hat{\mathbf{a}}_t^\top - \varepsilon_{t,l} \mathbf{x}_{t,l} \boldsymbol{\delta}_t^\top). \end{aligned} \quad (31)$$

Now we consider the terms in (31) separately. Let

$$\begin{aligned} S_2 &= \frac{\sigma}{LT} \sum_{l=1}^L -\varepsilon_{1,l} \mathbf{x}_{1,l} \mathbf{a}_1^\top + \frac{\sigma}{LT} \sum_{t=2}^T \sum_{l=1}^L -\varepsilon_{t,l} \mathbf{x}_{t,l} \hat{\mathbf{a}}_t^\top. \\ S_3 &= \frac{\sigma}{LT} \sum_{t=2}^T \sum_{l=1}^L -\varepsilon_{t,l} \mathbf{x}_{t,l} \boldsymbol{\delta}_t^\top. \end{aligned} \quad (32)$$

Define the following event

$$J_T = \{\|\mathbf{x}_{t,l}\|_2^2 \leq d_x + 2\sqrt{2d_x \log TL} + 4\log TL\} \text{ for } 1 \leq t \leq T, 1 \leq l \leq L. \quad (33)$$

Since $\|\mathbf{x}_{t,l}\|_2^2$ follows chi-square distribution with degree of freedom d_x , we have

$$P(J_T) \geq 1 - LT \frac{1}{T^2 L^2} \geq 1 - \frac{1}{TL}. \quad (34)$$

For simplicity of expression, denote

$$\mathfrak{J} = d_x + 2\sqrt{2d_x \log TL} + 4\log TL \quad (35)$$

Define

$$W(l; T) = \sum_{t=2}^T -\varepsilon_{t,l} \mathbf{x}_{t,l} \hat{\mathbf{a}}_t^\top. \quad (36)$$

Clearly

$$\|S_2\|_F \leq \frac{\sigma}{LT} \sum_{l=1}^L |\varepsilon_{1,l}| \cdot \|\mathbf{x}_{1,l}\|_F + \frac{\sigma}{LT} \sum_{l=1}^L \|W(l; T)\|_F. \quad (37)$$

We start with bounding the first term in Equation (37). Set event

$$\tilde{J}_T = \{|\varepsilon_{1,l}| \leq 2\sqrt{\log TL} : 1 \leq l \leq L\}.$$

Clearly on event $J_T \cap \tilde{J}_T$,

$$\frac{\sigma}{LT} \sum_{l=1}^L |\varepsilon_{1,l}| \cdot \|\mathbf{x}_{1,l}\|_F \leq \frac{\sigma}{T} 6(d_x + 2 \log TL).$$

Elementary calculation shows that $P(\tilde{J}_T) \geq 1 - \frac{1}{T^2 L}$.

Now we proceed with the second term in Equation (37). Denote the history to time T as

$$H_T = (\mathbf{x}_{t,l}, \mathbf{a}_{t,l}, y_{t,l} : 1 \leq t \leq T, 1 \leq l \leq L). \quad (38)$$

Then clearly $\{\varepsilon_{T,l} : 1 \leq l \leq L\} \perp H_{T-1}$. For $\lambda > 0$, elementary calculation show that

$$\mathbb{E} [\exp(\lambda \|W(l; T)\|_F^2) \mathbb{1}\{J_T\}] \quad (39)$$

$$= \mathbb{E} \left[\mathbb{E} \left[e^{\lambda \|W(l; T-1)\|_F^2 - 2\lambda \hat{\mathbf{a}}_T^\top W(l; T-1)^\top \mathbf{x}_{T,l} \varepsilon_{T,l} + \lambda \|\mathbf{x}_{T,l}\|_F^2 \varepsilon_{T,l}^2} \mathbb{1}\{J_T\} \mid H_{T-1}, \mathbf{x}_{T,l} \right] \right] \quad (40)$$

$$= \mathbb{E} \left[\frac{1}{1 - 2\lambda \|\mathbf{x}_{T,l}\|_F^2} \exp \left(\lambda \|W(l; T-1)\|_F^2 + \frac{(-2\lambda \hat{\mathbf{a}}_T^\top W(l; T-1)^\top \mathbf{x}_{T,l})^2}{2(1 - 2\lambda \|\mathbf{x}_{T,l}\|_F^2)} \right) \mathbb{1}\{J_T\} \right] \quad (41)$$

$$\leq \frac{1}{1 - 2\lambda \mathfrak{J}} \mathbb{E} \left[\exp \left(\lambda \|W(l; T-1)\|_F^2 + 2\lambda^2 \frac{\|W(l; T-1)\|_2^2 \mathfrak{J}}{1 - 2\lambda \mathfrak{J}} \right) \mathbb{1}\{J_T\} \right] \quad (42)$$

$$\leq \frac{1}{1 - 2\lambda \mathfrak{J}} \mathbb{E} \left[\exp \left(\frac{\lambda}{1 - 2\lambda \mathfrak{J}} \|W(l; T-1)\|_F^2 \right) \mathbb{1}\{J_T\} \right]. \quad (43)$$

Set $\lambda = \frac{1}{2T\mathfrak{J}}$, recursively using the above arguments, elementary calculation show that

$$\mathbb{E} \left(\exp \left(\frac{1}{2T\mathfrak{J}} \|W(l; T)\|_F^2 \right) \mathbb{1}\{J_T\} \right) \leq \prod_{t=2}^T \frac{1}{1 - 1/t} = T. \quad (44)$$

Therefore

$$P(\|W(l; T)\|_F^2 \mathbb{1}\{J_T\} \geq \mathfrak{w}^2) \leq T \exp(-\frac{1}{2T\mathfrak{J}} \mathfrak{w}^2). \quad (45)$$

Set $\mathfrak{w} = \sqrt{4T\mathfrak{J} \log TL} = \sqrt{4T(d_x + 2\sqrt{2d_x} \log TL + 4 \log TL) \log TL}$, we know that

$$P(\{\|W(l; T)\|_F \leq \mathfrak{w} \text{ for all } 1 \leq l \leq L\}^c \cap J_T) \leq L \frac{1}{TL^2} \quad (46)$$

Note that $\mathfrak{w} \leq \sqrt{T}(3d_x + 6 \log TL)$, we have the following holds with probability at least $1 - \frac{1}{TL} -$

$P(J_T^c) - \frac{1}{T^2 L}$

$$\|S_2\|_F \leq \frac{\sigma}{T} (6 + 3\sqrt{T}) (d_x + 2 \log TL). \quad (47)$$

For S_3 , let G be an event defined as

$$\begin{aligned} G = & \left\{ \max\{|\varepsilon_{t,l}| : 1 \leq t \leq T, 1 \leq l \leq L\} \leq 3\sqrt{\log TL}, \right. \\ & \max\{\|\mathbf{x}_{t,l}\|^2 : 1 \leq t \leq T, 1 \leq l \leq L\} \leq 2d_x + 6 \log LT, \\ & \left. \max\{\|\boldsymbol{\delta}_t/h\|_2^2 : 1 \leq t \leq T\} \leq 2d_a + 6 \log T \right\}. \end{aligned} \quad (48)$$

Then elementary calculation shows that

$$P(G^c) \leq \frac{2}{T^3 L^3} + \frac{1}{LT} + \frac{1}{T}. \quad (49)$$

Clearly $G \subset J_T$.

Using Matrix Bernstein Inequality [48] on event G , we have the operator norm of S_3 on G is bounded as follows

$$P \left(\left\{ \left\| \frac{LT}{\sigma} S_3 \right\|_2 \geq \alpha \right\} \cap G \right) \leq (d_x + d_a) \exp \left(\frac{-\alpha^2}{2\sigma_{S_3}^2 + 2D\alpha/3} \right), \quad (50)$$

585 where

$$\sigma_{S_3}^2 \geq \max \left\{ \left\| \sum_{t=1}^T \mathbb{E} \left(\left(\sum_{l=1}^L \varepsilon_{t,l} \mathbf{x}_{t,l} \boldsymbol{\delta}_t^\top \right) \left(\sum_{l=1}^L \varepsilon_{t,l} \mathbf{x}_{t,l} \boldsymbol{\delta}_t^\top \right)^\top \right) \right\|_2, \right. \\ \left. \left\| \sum_{t=1}^T \mathbb{E} \left(\left(\sum_{l=1}^L \varepsilon_{t,l} \mathbf{x}_{t,l} \boldsymbol{\delta}_t^\top \right)^\top \left(\sum_{l=1}^L \varepsilon_{t,l} \mathbf{x}_{t,l} \boldsymbol{\delta}_t^\top \right) \right) \right\|_2 \right\}, \quad (51)$$

586 and

$$D = \max_t \sup_{\text{event } G \text{ holds}} \left\| \sum_{l=1}^L -\varepsilon_{t,l} \mathbf{x}_{t,l} \boldsymbol{\delta}_t^\top \right\|_2 \leq 6Lh \sqrt{\log TL} \sqrt{(d_x + 3 \log LT)(d_a + 3 \log T)}. \quad (52)$$

587 Elementary calculation shows that taking

$$\sigma_{S_3}^2 = h^2 \lfloor T^{\frac{2}{3}} \rfloor L \max\{d_a, d_x\} \quad (53)$$

588 satisfies equation (51).

589 Taking

$$\alpha = 2hT^{\frac{1}{3}} \log T \sqrt{L \max\{d_a, d_x\} \log(d_a + d_x)} + \\ 8hL \sqrt{\log TL} \sqrt{(d_x + 3 \log(LT))(d_a + 3 \log T)(\log(d_x + d_a) + 2 \log T)} \quad (54)$$

590

$$P(\{\| \frac{LT}{\sigma} S_3 \|_2 \geq \alpha\} \cap G) \leq \frac{1}{T^2}. \quad (55)$$

591 Therefore, we have

$$P\left(\|S_3\|_2 \leq 2h\sigma T^{-2/3} \log T \sqrt{\frac{\max\{d_a, d_x\} \log(d_a + d_x)}{L}} + \right. \\ \left. \frac{8h\sigma}{T} \sqrt{\log(TL)} \sqrt{(d_x + 3 \log(LT))(d_a + 3 \log T)(\log(d_x + d_a) + 2 \log T)} \right) \\ \geq 1 - \frac{2}{L^3 T^3} - \frac{1}{LT} - \frac{1}{T} - \frac{1}{T^2} \quad (56)$$

592 Recalling that

$$|\langle \nabla \mathcal{L}_T(\Theta^*), \Delta \rangle| = |\langle S_2, \Delta \rangle + \langle S_3, \Delta \rangle| \leq \|S_2\|_F \|\Delta\|_F + \|S_3\|_2 \|\Delta\|_*, \quad (57)$$

593 we get the statement of the lemma.

594 A.2 Proof of Lemma A.1

To simplify the notation in the proof, we prove for any $T \geq 2$, with probability at least $1 - \frac{1}{LT} - \frac{2}{T} - \frac{1}{T^2}$,

$$\mathcal{E}_T(\Delta) \geq \frac{\lfloor T^{\frac{2}{3}} \rfloor}{2T} h^2 \|\Delta\|_F^2 - 14T^{-\frac{2}{3}} (h + h^2) (2d_x + 2d_a + 6 \log T + 6 \log L)^2 \log T \|\Delta\|_2^2.$$

595 Let $\mathbf{b}_t = \sum_{l=1}^L \mathbf{x}_{t,l}$. Let $\boldsymbol{\delta}_t = \mathbf{0}$ for exploitation rounds.

596 Then we know that

$$\mathcal{E}_T(\Delta) = \frac{1}{2LT} \sum_{t=1}^T \sum_{l=1}^L (\mathbf{a}_{t,l}^\top \Delta \mathbf{x}_{t,l})^2 \\ = \frac{1}{2LT} \sum_{t=1}^T \sum_{l=1}^L \left(\left(\frac{\mathbf{b}_t^\top \hat{\boldsymbol{\Theta}}_{t-1}^\top}{\|\mathbf{b}_t^\top \hat{\boldsymbol{\Theta}}_{t-1}^\top\|} + \boldsymbol{\delta}_t^\top \right) \Delta \mathbf{x}_{t,l} \right)^2 \quad (58)$$

597 Define

$$\begin{aligned}
\mathcal{D}_T(\Delta) &= \frac{1}{2LT} \sum_{t=1}^T \sum_{l=1}^L \left(\left(\frac{\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top}{\|\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top\|} \Delta \mathbf{x}_{t,l} \right)^2 + (\delta_t^\top \Delta \mathbf{x}_{t,l})^2 \right), \\
\mathcal{D}_{1,T}(\Delta) &= \frac{1}{2LT} \sum_{t=1}^T \sum_{l=1}^L \left(\frac{\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top}{\|\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top\|} \Delta \mathbf{x}_{t,l} \right)^2 \\
\mathcal{D}_{2,T}(\Delta) &= \frac{1}{2LT} \sum_{t=1}^T \sum_{l=1}^L (\delta_t^\top \Delta \mathbf{x}_{t,l})^2
\end{aligned} \tag{59}$$

598 Then

$$\mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta) = \frac{1}{LT} \sum_{t=1}^T \sum_{l=1}^L \left(\frac{\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top}{\|\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top\|} \Delta \mathbf{x}_{t,l} \right) (\delta_t^\top \Delta \mathbf{x}_{t,l}) \tag{60}$$

599 Elementary calculation shows that

$$\mathbb{E}(\mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta)) = 0, \tag{61}$$

600 and

$$\mathbb{E}(\mathcal{D}_{2,T}(\Delta)) \geq \frac{\lfloor T^{\frac{2}{3}} \rfloor}{2T} h^2 \|\Delta\|_F^2. \tag{62}$$

601 Now we proceed with proving that the following two bounds hold with high probability:

$$\begin{aligned}
\inf_{\|\Delta\|_2 > 0} \frac{\mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta)}{\|\Delta\|_2^2} &\geq -7T^{-\frac{2}{3}} (h + h^2) (2d_x + 2d_a + 6 \log T + 6 \log L)^2 \log T \\
\inf_{\|\Delta\|_2 > 0} \frac{\mathcal{D}_{2,T}(\Delta) - \mathbb{E}(\mathcal{D}_{2,T}(\Delta))}{\|\Delta\|_2^2} &\geq -7T^{-\frac{2}{3}} (h + h^2) (2d_x + 2d_a + 6 \log T + 6 \log L)^2 \log T.
\end{aligned} \tag{63}$$

602 Note that $\|\mathbf{x}_{t,l}\|_2^2 \sim \chi_{d_x}^2$, $\|\delta_t/h\|_2^2 \sim \chi_{d_a}^2$. Therefore, we have that

$$\begin{aligned}
P(\sup_{t,l} \|\mathbf{x}_{t,l}\|_2^2 \leq d_x + 2\epsilon_1 + 2\sqrt{\epsilon_1 d_x}, \sup_t \|\delta_t/h\|_2^2 \leq d_a + 2\epsilon_2 + 2\sqrt{\epsilon_2 d_a}) \\
\geq 1 - (LT \exp(-\epsilon_1) + T \exp(-\epsilon_2)).
\end{aligned} \tag{64}$$

603 Let $\epsilon_1 = 2 \log LT$, $\epsilon_2 = 2 \log T$.

604 Denote

$$U_1 = d_x + 2\epsilon_1 + 2\sqrt{\epsilon_1 d_x}, U_2 = d_a + 2\epsilon_2 + 2\sqrt{\epsilon_2 d_a}. \tag{65}$$

605 And let the event O be

$$O = \{\sup_{t,l} \|\mathbf{x}_{t,l}\|_2^2 \leq U_1, \sup_t \|\delta_t/h\|_2^2 \leq U_2\}. \tag{66}$$

606 For the following, we restrict our attention to event O .

607 Note that

$$\inf_{U_0/1.1 \leq \|\Delta\|_2 \leq U_0} \mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta) \geq \inf_{U_0/1.1 \leq \|\Delta\|_2 \leq U_0, \hat{\Theta}_{t-1}^\top \neq \mathbf{0} \text{ for } 1 \leq t \leq T} \mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta), \tag{67}$$

608 also at most $\lfloor T^{\frac{2}{3}} \rfloor$ terms in the sum of $\mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta)$ are not zero, and for any term

$$\begin{aligned}
609 \text{ in the exploration round } &\left(\sup_{U_0/1.1 \leq \|\Delta\|_2 \leq U_0, \hat{\Theta}_{t-1}^\top \neq \mathbf{0} \text{ for } 1 \leq t \leq T} \left(\frac{\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top}{\|\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top\|} \Delta \mathbf{x}_{t,l} \right) (\delta_t^\top \Delta \mathbf{x}_{t,l}) \right) - \\
610 &\left(\inf_{U_0/1.1 \leq \|\Delta\|_2 \leq U_0, \hat{\Theta}_{t-1}^\top \neq \mathbf{0} \text{ for } 1 \leq t \leq T} \left(\frac{\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top}{\|\mathbf{b}_t^\top \hat{\Theta}_{t-1}^\top\|} \Delta \mathbf{x}_{t,l} \right) (\delta_t^\top \Delta \mathbf{x}_{t,l}) \right) \leq 2U_1 \sqrt{U_2} h U_0^2
\end{aligned}$$

Therefore, through Functional Hoeffding theorem (Theorem 3.26 in Wainwright [51]), we have

$$P(\mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta) \leq -\gamma_1 | O) \leq \exp \left(-\frac{\frac{T^2}{\lfloor T^{\frac{2}{3}} \rfloor} \gamma_1^2}{16U_1^2 U_2^2 h^2 U_0^4} \right) \quad (68)$$

for $\gamma_1 > 0$.

Similarly, for the exploration rounds in $\mathcal{D}_{2,T}(\Delta)$, we have

$$\left(\sup_{U_0/1.1 \leq \|\Delta\| \leq U_0} (\delta_t^\top \Delta \mathbf{x}_{t,l})^2 \right) - \left(\inf_{U_0/1.1 \leq \|\Delta\| \leq U_0} (\delta_t^\top \Delta \mathbf{x}_{t,l})^2 \right) \leq U_1 U_2 U_0^2 h^2. \quad (69)$$

Again, according to Functional Hoeffding theorem, we have

$$P(\mathcal{D}_{2,T}(\Delta) - \mathbb{E}(\mathcal{D}_{2,T}(\Delta)) \leq -\gamma_2 | O) \leq \exp \left(-\frac{\frac{T^2}{\lfloor T^{\frac{2}{3}} \rfloor} \gamma_2^2}{4U_1^2 U_2^2 U_0^4 h^4} \right) \quad (70)$$

Take $\gamma_1 = \gamma_2 = 7T^{-\frac{2}{3}}(h + h^2)(2d_x + 2d_a + 6\log T + 6\log L)^2 \|\Delta\|_2^2 \log T$.

Therefore,

$$\begin{aligned} & P(\mathcal{E}_T(\Delta) - \mathbb{E}(\mathcal{D}_{2,T}) \leq -14T^{-\frac{2}{3}}(h + h^2)(2d_x + 2d_a + 6\log T + 6\log L)^2 \|\Delta\|_2^2 \log T) \\ & \leq P(\mathcal{E}_T(\Delta) - \mathcal{D}_T(\Delta) \leq -7T^{-\frac{2}{3}}(h + h^2)(2d_x + 2d_a + 6\log T + 6\log L)^2 \|\Delta\|_2^2 \log T | O) \\ & + P(\mathcal{D}_{1,T}(\Delta) \leq 0 | O) \\ & + P(\mathcal{D}_{2,T} - \mathbb{E}(\mathcal{D}_{2,T}) \leq -7T^{-\frac{2}{3}}(h + h^2)(2d_x + 2d_a + 6\log T + 6\log L)^2 \|\Delta\|_2^2 \log T | O) \\ & + P(O^c) \\ & \leq \frac{1}{LT} + \frac{2}{T} + \frac{1}{T^2} \end{aligned} \quad (71)$$

Hence with probability at least $1 - \frac{1}{LT} - \frac{2}{T} - \frac{1}{T^2}$,

$$\mathcal{E}_T(\Delta) \geq \frac{\lfloor T^{\frac{2}{3}} \rfloor}{2T} h^2 \|\Delta\|_F^2 - 14T^{-\frac{2}{3}}(h + h^2)(2d_x + 2d_a + 6\log T + 6\log L)^2 \|\Delta\|_2^2 \log T. \quad (72)$$

B Details on the simulation study

In this section, we detail the tuning parameters of each algorithm we used for the simulation study. We ran our simulation on servers with Intel(R) Xeon(R) CPU E5-2680 v3 @ 2.50GHz.

Hi-CCAB. There are three tuning parameters for Hi-CCAB: we set the steps for initialization $t_1 = 100$, the initialized penalization parameter $\lambda_0 = \|\frac{1}{2t_1 L} \sum_{i=1}^{t_1} \sum_{j=1}^L |\mathbf{a}_i^\top \hat{\Theta}_{t_1} \mathbf{x}_{i,j} - y_{i,j}| \mathbf{x}_{i,j} \mathbf{a}_i^\top\|_2$, and the exploration parameter $h = .1$.

LinUCB [33]. We apply the LinUCB with disjoint linear models and set multiplier for the upper confidence bound $\alpha = 1 + \sqrt{\ln(2/\delta)/2}$ with $\delta = .05$ as suggested in the paper.

Lasso Bandit [6]. There are a couple of tuning parameters in the original algorithm including h for the set of “near-optimal arms”, q for the force-sample set, and λ_1 and $\lambda_{2,0}$ as the regularization parameters for the “forced sample estimate” and “all-sample estimate”. We follow the original paper and set $h = 5$, $\lambda_1 = \lambda_{2,0} = 0.05$. We set $q = 2$ so that the size of initialized forced sample set is close to that we used for Hi-CCAB.

NeuralUCB [55]. The tuning parameters of NeuralUCB include the confidence parameter as in all UCB-based algorithm, the size of neural network, as well as the step size, regularization parameter for gradient descent to train the neural network. We adapted the code from <https://github.com/uclaml/NeuralUCB> and used the default settings.

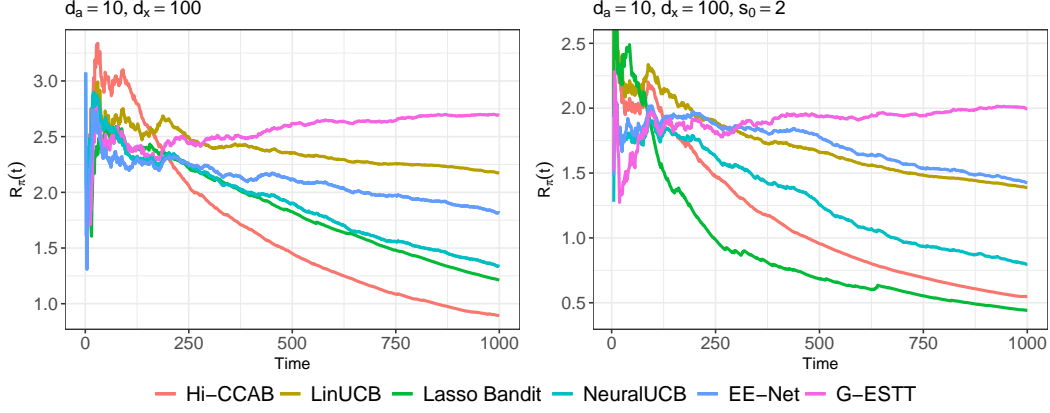


Figure 2: Expected average regret.

EE-Net [5]. EE-Net involves tuning parameters for gradient descent to train the exploitation network, exploration network, and the decision-maker network. We adapted the code from <https://github.com/banyikun/EE-Net-ICLR-2022> and used the default settings.

G-ESTT [22]. We implement the algorithm proposed in Appendix H of [22], which is claimed to be an extension of their main algorithm to contextual setting. We set the initialization steps T_1 according to Theorem 4.3 in [22], which claims good performance for their main algorithm. However, since we focus on $T = 1000$ and our moderate dimensions are already relatively large for their algorithm, we instead set $T_1 = \sqrt{rT \log((d_1 + d_2)/\delta)/D_{rr}}$ where $d_1 = d_x$, $d_2 = d_a$, $\delta = 0.01$ as in their setup and $D_{rr} = .5$ is the smallest non-zero singular value.

Figure 2 shows the expected average regret comparing Hi-CCAB, LinUCB [33], Lasso Bandit [6], NeuralUCB [55], EE-Net [5] and G-ESTT [22] under the same settings of the first left column in Figure 1a. The purpose is to compare a potential extension of G-ESTT in Appendix H of [22].

In [22], the main algorithm with theoretical guarantees is designed for low-rank matrix bandit but not for contextual bandit. In Appendix H, the author sketched an extension to the contextual setting, but their theory and numerical validation only apply to non-contextual bandits. We applied the modified G-ESTT to the contextual bandit setting where $d_a = 10, d_x = 100$ and the sparse setting with $d_a = 10, d_x = 100, s_0 = 2$ as in Section 4. Figure 2 shows that the modified G-ESTT does not perform well compared to other contextual bandit algorithms. One reason can be the following. They adopt the explore-then-commit algorithm and their initialization step is required to be of the order $\sqrt{d_1 d_2 r T}$. In our simulation setting, the dimension of covariate $d_2 = d_x$ is high and therefore their algorithm will require a large T_1 to perform well. Therefore, the modified G-ESTT does not perform well in high-dimensional settings (note that the dimensions in their simulations is of the order 10), especially when T is relatively small. In addition, the computational complexity of the modified G-ESTT is high compared to other methods. We tried to run the modified G-ESTT for other settings as in Section 4 with larger d_x but it would have taken too long and we do not expect a better performance of the modified G-ESTT in higher dimension settings given the above-mentioned claim.

C More details on the case study and additional numerical results

In this section, we provides more background information on the case study and additional interpretations of the representation matrix Θ and numerical results.

Figure 3a shows the daily sales by product and each color represents one product (only products that appeared more than 95% of the days are colored; the rest are colored as grey). The days corresponding to the vertical dashed grey lines are days with promotion. The two red vertical lines correspond to the annual sales events. The variation between products was large and one product dominated the rest most of the time. The sales were also driven by the promotion – the sales went up when there is a promotion. Figure 3b shows the median unit price across time with the 25th and 75th quantiles as the boundaries of the grey area. The median unit price was around 3.2 RMB and there were variations in

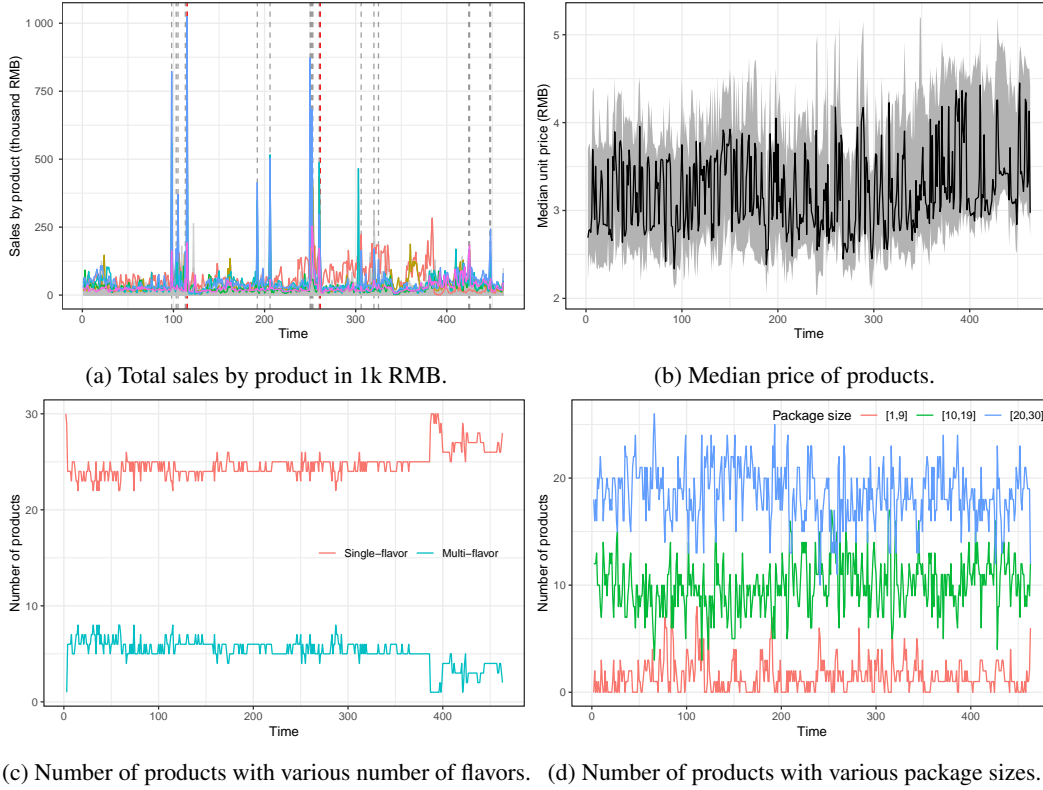


Figure 3: Summary of the products.

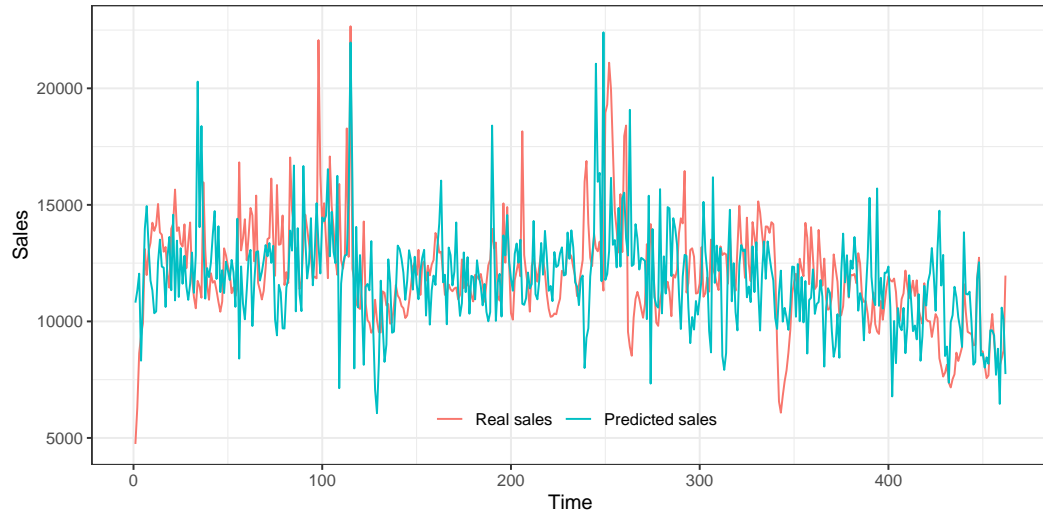


Figure 4: Real sales vs predicted sales.

unit price among products. Figure 3c shows the number of single-flavor and multi-flavor products. Three-quarters of the products were single-flavored. Note that products with the same flavor can have different package sizes. Figure 3d shows the number of products with different package sizes. The package size of about 60% of the products is larger than 20 with 30% having package sizes between 10 and 20 and the rest less than 10.

To check our model assumption (1) on the data, Figure 4 shows the hold-out-sample prediction of the sales versus the real sales. The predicted sale at each time point t is the predicted total sales across

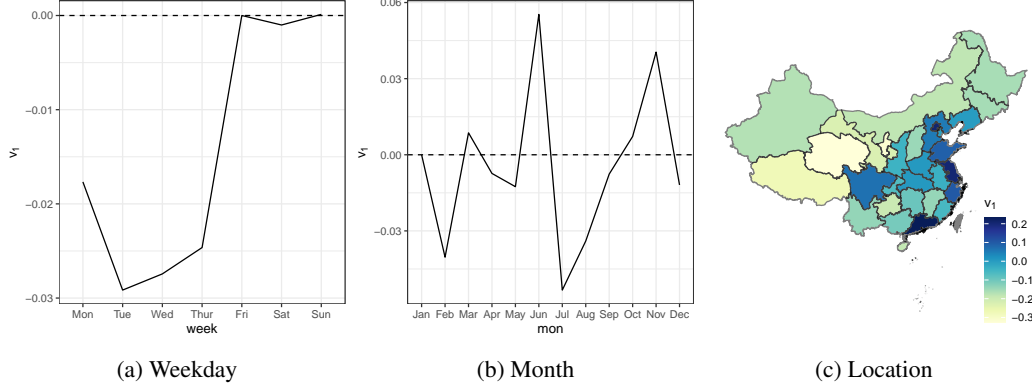


Figure 5: Loadings of the leading right singular vectors for the covariates.

$L = 31$ locations based on $\hat{\Theta}_{-t}$ estimated from all the data except for data at time t , i.e.,

$$\hat{y}_{t,j} = \mathbf{a}_t^\top \hat{\Theta}_{-t} \mathbf{x}_{t,j}$$

where $\hat{\Theta}_{-t} = \arg \min_{\Theta} \sum_{i=1, i \neq t}^T \sum_{j=1}^L (\mathbf{a}_i^\top \Theta \mathbf{x}_{i,j} - r_{i,j})^2 + \lambda \|\Theta\|_*$. As shown in Figure 4, the real sales and the out-of-sample predicted sales follow quite closely across time, which indicates that both our model and estimation are reasonable.

Structure of the representation matrix Θ . One advantage of model is the interpretability which allows us to gain insights from the representation matrix Θ . Specifically, our model is able to discover the underlying factors of the effect of both arms and covariates on the reward. In the following, we will examine the pseudo ground truth Θ we obtained using all the data.

The rank of Θ is 5 with the singular values being (2.5, 0.3, 0.2, 0.02, 0.002). The leading singular value dominates the rest and thus the leading left and right singular vectors are the most important ones in explaining the effect on the reward and we focus on the leading singular vectors in what follows.

Figure 5 shows the loadings for different covariates (i.e., the leading right singular vector) and our algorithm is able to learn interpretable patterns of the effects on the reward – for weekday, the effects are drastically different during the weekend and during the weekend; for months, the effects show different patterns during the promotion month (June and November) from other months; for location, the effects of the coastal provinces are different from the rest, which exactly corresponds to the levels of economic development of different regions in China. In sum, our model can exploit the underlying structure of the covariates and provide insights into purchasing behavior and seasonality.

On the other hand, Table 1 explores the loadings for the arm on May 29th 2022, the last Sunday in our data (i.e., the leading left singular vectors multiplied with $\langle \mathbf{v}_1, \bar{\mathbf{x}} \rangle$ where $\bar{\mathbf{x}}$ is the average of \mathbf{x}_j for $j = 1, \dots, L$ on May 29th 2022). Specifically, we investigate the effect of flavors on the reward given the context. We take the average of the loadings of the linear and quadratic terms for each flavor in all 30 products and compare with the total sales of each flavor across all Sundays in Mays. For ease of comparison, we further scale the sales and the loadings by their corresponding largest numbers. The loadings and sales are closely related to each other.² As in Table 1, on May 29th 2022, flavor 1 (F1) has the largest effect, followed by flavor 10, 13, 7, 9 and 11. Therefore, our model learns the values of the flavors (per unit).

	F1	F2	F3	F4	F5	F6	F7	F8	F9	F10	F11	F12	F13
Sales	1.00	0.05	0.00	0.00	0.00	0.03	0.19	0.00	0.08	0.19	0.18	0.00	0.38
\tilde{u}_1 (linear)	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
\tilde{u}_1 (quadratic)	1.00	0.12	-0.00	0.00	0.00	0.03	0.19	0.00	0.15	0.39	0.16	0.03	0.33

Table 1: Total sales and loadings of the linear and quadratic terms (scaled) of the 13 flavors.

²The correlation of sales and the linear-term loadings is 0.91 and that of the quadratic-term loadings is 0.97.

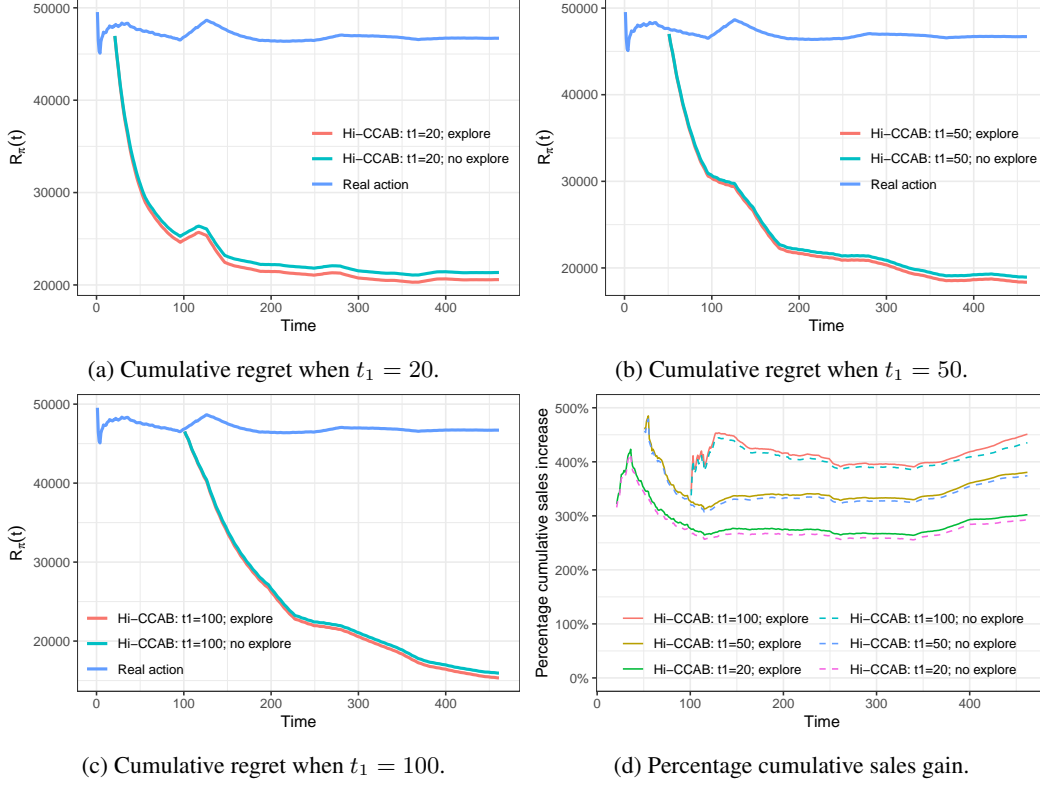


Figure 6. Performance of Hi-CCAB with different initialization times t_1 and with exploration and without exploration.

More on simulation with additional numerical results. We first detail how we ran the simulation and then provide more simulation results. To be specific, we first use $t_1 = 100$ for the initialization step to estimate $\hat{\Theta}_{t_1}$; and then at each time $t = t_1 + 1, \dots, T$, we follow Algorithm 1 to decide on the action a_t for assortment and pricing. After determining a_t , we generate the sales r_t according to (1) using the pseudo true Θ and σ . We further compare the performance of the assortment-pricing policy with exploration and without exploration and with different initialization time t_1 . Each setup is simulated 100 times.

Figures 6a-6c show cumulative regret and Figures 6d show percentage gain in cumulative sales when $t_1 = 20, 50, 100$ with exploration and without exploration. Hi-CCAB with exploration performs better than without exploration. As expected, longer initialization steps provide a better initial estimation of the Θ and thus helps with the performance in a short time windows. As time goes by, all of the expected regrets converge to zero and the percentage gain in cumulative sales should converge.