

# ONE-SHOT CLUSTERING FOR CONTEXTUAL BANDITS WITH KNAPSACKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In this work, we study the problem of clustered linear contextual bandits with knapsack constraints, a setting that closely models real-world recommender systems. In such systems, the overwhelmed number of items makes it impractical to explore all options, and overexposing certain items can harm content diversity and fairness. To address these challenges, our algorithm clusters actions to enable knowledge transfer across similar items and incorporates global resource constraints to limit over-consumption. We provide a formal analysis showing that the algorithm achieves sublinear regret in the number of time periods, even without access to the full action set. Notably, we prove that it is sufficient to perform clustering once on a randomly selected subset of actions.

## 1 INTRODUCTION

In the contextual bandits problem (Langford & Zhang, 2007; Slivkins, 2011; Agarwal et al., 2014), an agent selects an action and collects a reward for that action over a sequence of rounds. At each round, the agent makes its choice based on the context for the current round and the feedback from previous rounds. The feedback only consists of the rewards for the chosen action, and the rewards of other actions remain unobserved. Online recommender systems (Tang et al., 2014) are often modeled as contextual bandit problems, where personalized recommendations are tailored for each user. The system receives feedback in the form of user interactions, such as clicks, likes, or comments, on the recommended item, but does not observe feedback for items that were not shown. Based on this partial feedback and the user’s profile, the system continuously updates its policy to recommend new items.

Contextual bandits problem theoretically works in the recommender system. However, real world recommender system often suffers from:

- The constant introduction of new items makes the environment highly dynamic. Additionally, there is sparse user interaction when the item set is large. Contextual bandits must explore each item individually, leading to high sample complexity and cold-start issues as feedback per item remains scarce.
- Risk of depleting limited recommendation capacity. Contextual bandits do not account for constraints such as user attention span, content diversity quotas, or exposure fairness limits. Without modeling these factors, the system may over-recommend certain items, leading to user fatigue, reduced content diversity, or unfair exposure.

To solve the first challenge on how to present new items to users and use all available user-item preference information gathered, previous works (Yang et al., 2020a; Nguyen & Lauw, 2014; Gentile et al., 2017) have proposed clustered contextual bandits. It performs clustering over user interests, items, or user–item preferences, enabling clustered contextual bandits to transfer knowledge across similar users and improve learning efficiency. To address the second challenge of depleting limited recommendation capacity, prior research has extended the clustered bandit framework by incorporating global resource constraints. This method is known as bandit with knapsacks (Ma et al., 2024; Jiang & Ye, 2024; Slivkins et al., 2023; Li et al., 2021b; Han et al., 2023). To our knowledge, no prior work has combined these two ideas into a single algorithmic framework. Thus a natural question arises:

054 *Can we design a contextual-bandit algorithm that simultaneously exploits cluster structure and*  
 055 *respects knapsack constraints, while retaining provable sublinear regret?*  
 056

057 In this work, we provide a positive answer to the above question. We introduce a new algorithm  
 058 for clustered contextual bandits with knapsack constraints that achieves sublinear regret in the time  
 059 horizon  $T$ . Arms are partitioned into unknown clusters. The arm’s reward and  $d$ -dimension con-  
 060 sumption are specified by a cluster-specific linear model in its observed context. At each round,  
 061 the agent observes the i.i.d. context vector of each arm, select an arm, and get both a reward and  
 062 a consumption vector. The goal of the algorithm is to maximize the sum of the rewards over pe-  
 063 riods. If the cumulative consumption of any resource exceeds its budget, the process terminates.  
 064 Our model extends that of (Agrawal & Devanur, 2016) by considering clusters. The approach we  
 065 employ builds on the algorithm of that paper. The main challenge lies in simultaneously learning the  
 066 unknown cluster structure of the arms while guaranteeing that the average regret remains vanishing  
 over time.

067 Our algorithm performs clustering only once. Initially, we sample a subset  $\mathcal{S}$  of the total  $K$  arms;  
 068 any arm not in  $\mathcal{S}$  is never played. The key intuition is that, to achieve vanishing average regret  
 069 as the number of periods grows, we must access arms from every cluster and accurately estimate  
 070 the fraction of the  $K$  arms belonging to each cluster. Once  $\mathcal{S}$  is fixed, we play each arm in  $\mathcal{S}$   
 071 a predetermined number of times to gather sufficient observations for accurate clustering. Playing an  
 072 arm yields  $d + 1$  observations—one reward and  $d$  resource-consumption values—each governed by  
 073 a different linear model. For clustering, we adopt the classifier-Lasso method from econometrics  
 074 (Su et al., 2016), which treats arm parameters as distinct (though not necessarily unequal) relative  
 075 to cluster parameters. After clustering, the algorithm selects arms according to the optimism-in-the-  
 076 face-of-uncertainty principle (Auer, 2002; Auer et al., 2002; Abbasi-Yadkori et al., 2011). In each  
 077 period, we form optimistic estimates of both reward and resource consumption for every arm in  
 078  $\mathcal{S}$ , based on the observed context and the cluster-specific linear models. Since clustering may have  
 079 some errors, some arms might be misgrouped and thus not share the true parameters of their assigned  
 080 cluster. To facilitate regret analysis, we treat all arms assigned to a cluster as if they belonged to the  
 081 same true cluster—but allow the linear-model error term to have a context-dependent expectation  
 082 (with probability tied to the clustering error) instead of zero. In other words, we model imperfect  
 083 clustering as measurement error in the context (Wansbeek & Meijer, 2001; Fuller, 2009). Finally, we  
 084 leverage results from convex online learning to account for the depletion of each of the  $d$  resources,  
 following the approach in (Srebro et al., 2011) and the algorithm of (Agrawal & Devanur, 2016).

085 In regard to the number of time periods  $T$ , the regret of our algorithm increases at the rate  $T^{1-\delta}$ , for  
 086  $\delta \in (0, \frac{1}{2})$ , requiring the budget for each resource to be  $B > \tilde{O}(T^{2\delta})$ . Thus, the regret is sublinear  
 087 in  $T$ , and it can approach the square root rate as the resource budget approaches  $T$ . However,  
 088 as it becomes clear in the next section, we must have  $B < T$ , since otherwise we could ignore the  
 089 constraints. The main drawback of our algorithm is that its regret depends on the number of arms  $K$ ,  
 090 as it operates only on a subset of them. Nevertheless, our results suggest that such dependence may  
 091 be inevitable. Overall, we believe that our approach points to a fruitful direction for the consideration  
 092 of online learning problems that involve clusters and constraints.

093 **Roadmap.** The rest of the paper is organized as follows. In Section 2, we provide the related works  
 094 on our algorithms. In Section 3, we provide notation and define the problem under consideration. In  
 095 Section 4, we derive optimistic estimates of the reward and consumption parameters and establish  
 096 confidence bounds for these estimate. Section 5 presents our `clusterLCBwK` algorithm along with  
 097 its regret analysis. In Section 6, we discuss the implications of our results and potential directions  
 098 for future work. Technical proofs and additional details are deferred to the appendix.  
 099

## 100 2 RELATED WORK

101  
 102 Some of the earliest works to consider the problem of bandits with knapsacks were (Badanidiyuru  
 103 et al., 2014) and an earlier version of (Badanidiyuru et al., 2018), with both of them considering  
 104 stochastic reward and resource consumption. The adversarial version of the problem was initially  
 105 considered in (Immorlica et al., 2019), and consequently in (Kesselheim & Singla, 2020). An algo-  
 106 rithm that achieves low expected regret for both the adversarial and the stochastic version is inves-  
 107 tigated in (Rangi et al., 2019), while in (Amani et al., 2019) bandit algorithms with constraints for  
 safety-critical systems are proposed. General forms for the correlation of the reward and the resource

consumption are considered in (Cayci et al., 2020). (Devanur et al., 2011; Agrawal & Devanur, 2014; Agrawal et al., 2016) study computationally efficient algorithms for online learning problems with constraints, including bandits with knapsacks. (Agrawal & Devanur, 2016; Wu et al., 2015) consider the case where there is a single resource. From a result of (Dani et al., 2008), it is known that the regret of a linear contextual bandits algorithm must be at least linear in the dimensionality of the context. The topic of contextual bandits with behavioral constraints is presented in (Balakrishnan et al., 2018). (Yang et al., 2020b) studies contextual bandits with a resource constraint in recommender systems (Carlsson et al., 2021). Recent works focus on establishing high-probability performance guarantees (Ma et al., 2024; Deb et al., 2024; Jiang & Ye, 2024; Slivkins et al., 2023; Sivakumar et al., 2022), exploring symmetry properties (Li et al., 2021b), deriving optimal algorithms (Han et al., 2023), and under non-stationary settings (Liu et al., 2022; Lyu & Cheung, 2023; Zhang & Cheung, 2024).

In works that model recommender systems as bandits with clusters of users (Nguyen & Lauw, 2014; Gentile et al., 2014; 2017; Li & Zhang, 2018; Ban & He, 2021; Li et al., 2021a), the users arrive at the system exogenously. Besides the consideration of constraints and the fact that the model we study is not specific to an application, we cannot adopt techniques from the literature of recommender systems because we consider clusters of arms. Thus, the units to perform clustering on are endogenously and not exogenously provided, which makes our problem harder.

Clustering under mixed linear models (like ours) has attracted a lot of the interest in the literature (Zhong et al., 2016; Li & Liang, 2018; Chen et al., 2020; Kong et al., 2020; Chen et al., 2021; Pal et al., 2023; Wang et al., 2024; Cheng et al., 2023; Yang et al., 2024), but with the results of these works typically requiring a lower bound on the number of units to be clustered. Interestingly, similar versions of this problem have been considered in the literature of econometrics (Lin & Ng, 2012; Ando & Bai, 2016; Su et al., 2016; 2019; Gu & Volgushev, 2019; Okui & Wang, 2021). Here, we exploit results from the latter literature, since the assumptions they require are more appropriate to our case. In particular, we make assumptions about the joint rate of the number of sampled arms and of the observations collected for each such arm, instead of imposing a lower bound on the former.

### 3 PROBLEM DEFINITION AND NOTATION

In this section, we first introduce notation that we use throughout the paper, and then we provide the definition of the problem we consider. For any positive integer  $N$ , we use  $[N]$  to denote the set  $\{1, 2, \dots, N\}$ . We use  $\mathbb{E}[X]$  to denote the expectation of the random variable  $X$ , and  $\Pr[E]$  to denote the probability of the event  $E$ . The notation  $\mathbb{1}[f]$  is used for the indicator variable which outputs 1 if  $f$  holds and 0 otherwise. Boldface lower case letters are reserved for vectors, and boldface upper case letters for matrices. The vector  $\mathbf{1}_d$  denotes the  $d$ -dimensional vector that has the value 1 in every dimension, and  $\mathbf{0}_d$  is defined equivalently. Calligraphic upper case letters, e.g.,  $\mathcal{S}$ , denote sets. The probability simplex in  $d$  dimensions is denoted as  $\Delta^d$ . For a square matrix  $\mathbf{A}$ , we define the matrix norm of the vector  $\mathbf{x}$  as  $\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^{\top} \mathbf{A} \mathbf{x}}$ .

The total number of arms is denoted by  $K$ . A subset of arms  $\mathcal{S} \subset [K]$  is sampled from  $[K]$ . Arms not in  $\mathcal{S}$  never being played. The number of clusters by  $C \in \mathbb{N}$ . Each arm belongs to one cluster, with  $c(a) \in [C]$  denoting the cluster of arm  $a \in [K]$ . The algorithm is given a budget  $B \in \mathbb{R}_+$  and the number of clusters  $C$ , but not the arm membership in the clusters. In each period  $t \in [T]$ , the algorithm observes the context of all the arms  $X_t = [x_t(1), x_t(2), \dots, x_t(K)] \in [0, 1]^{m \times K}$ .  $x_t(a)$  is the context vector of arm  $a$  in time  $t$ .  $m \in \mathbb{N}$  is the dimension of the context vector. Then, the agent chooses arm  $a_t \in [K]$ , and finally observes reward  $r_t(a_t) \in [0, 1]$  and consumption vector  $\mathbf{v}_t(a_t) \in [0, 1]^d$ . The algorithm can also play the “no-op” action that deterministically gives zero reward and consumption. The goal of the algorithm is to maximize the total reward  $\sum_{t \in [T]} r_t(a_t)$  under the constraints  $\sum_{t \in [T]} \mathbf{v}_t(a_t) \leq B \cdot \mathbf{1}_d$ . Here,  $\mathbf{v}_t(a_t)$  is a  $d$ -dimensional vector indicating how much of each resource is consumed by choosing arm  $a_t$  at time  $t$ . Notice that it is without loss of generality to consider the same budget for each of the  $d$  resources, since this can imposed by normalization.

Moreover, we assume that the reward and the consumption are generated by cluster-specific linear models, as indicated by the following assumptions.

**Assumption 3.1** (Linearity). *For each cluster  $c \in [C]$ , there is an unknown parameter vector  $\boldsymbol{\mu}_c \in [0, 1]^m$  and an unknown parameter matrix  $\mathbf{W}_c \in [0, 1]^{m \times d}$  such that for all rounds  $t \in [T]$  and all arms  $a \in [K]$ ,*

$$r_t(a) = \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}_t(a) + g_t(a) \quad (1)$$

$$\mathbf{v}_t(a) = \mathbf{W}_{c(a)}^\top \mathbf{x}_t(a) + \mathbf{q}_t(a) \quad (2)$$

where  $\boldsymbol{\mu}_{c(a)} = \boldsymbol{\mu}_c$  and  $\mathbf{W}_{c(a)} = \mathbf{W}_c$  if  $c(a) = c$ , and  $g_t(a) \in \mathbb{R}$  and  $\mathbf{q}_t(a) \in \mathbb{R}^d$  are additive error terms.

We make the following assumptions about the error terms.

**Assumption 3.2** (Zero conditional mean). *We assume  $\mathbb{E}[g_t(a)|\mathbf{x}_t(a)] = 0$ , and  $\mathbb{E}[\mathbf{q}_t(a)|\mathbf{x}_t(a)] = \mathbf{0}_d$ .*

**Assumption 3.3.** *We assume  $|g_t(a)| \leq 2R$ , and  $\|\mathbf{q}_t(a)\|_\infty \leq 2R$ .*

Furthermore, we make the following i.i.d. assumption about the context.

**Assumption 3.4** (i.i.d.). *The context  $\mathbf{x}_t(a)$  is i.i.d. across arms and periods.*

In the algorithm also rely on Assumption A.5, Assumption A.6, Assumption A.7, Assumption A.8. For brevity, we direct reader to Appendix A for more details.

We let the vector  $\mathbf{p} = (p_1, \dots, p_C) \in \Delta^C$  denote the proportions of arms that are in each cluster, i.e.,  $p_c = \frac{1}{K} \sum_{a \in [K]} \mathbb{1}[c(a) = c]$  for  $c \in [C]$ . The smallest proportion in  $\mathbf{p}$  is denoted as  $p_{\min} := \min_{c \in [C]} p_c$ . Even though the algorithm does not know the other proportions in  $\mathbf{p}$ , some knowledge about  $p_{\min}$  is required for the clustering, as described later. Also, notice that  $p_{\min} \leq 1/C$  by definition.

### 3.1 BENCHMARK

Our goal is to devise an algorithm that achieves sublinear regret in the number of time periods  $T$ . We employ as benchmark the expected reward of the optimal static policy which needs to satisfy the consumption constraints only in expectation. This benchmark policy knows the reward and consumption parameters for each cluster, as well as the cluster membership of each arm. It is known (Devanur et al., 2011; Agrawal & Devanur, 2016; Badanidiyuru et al., 2018) that the optimal static policy achieves the same expected reward to the optimal adaptive policy which knows the distribution of the context and needs to satisfy the constraints for the realizations of the resource consumption.

**Definition 3.5** (Optimal Static Policy). *Let  $\mathbf{X} = (\mathbf{x}(1), \dots, \mathbf{x}(K)) \in [0, 1]^{m \times K}$  be the matrix with the context of all arms in an arbitrary period, and  $\pi(i, \mathbf{X}) \in [0, 1]$  the probability with which the action  $i \in [K] \cup \{\text{“no-op”}\}$  is taken by the static policy  $\pi$  when the context is  $\mathbf{X}$ . The per-period expected reward and consumption vector of  $\pi$  are respectively defined as*

$$\begin{aligned} r(\pi) &:= \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in [K]} \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}(a) \pi(a, \mathbf{X}) \right], \\ \mathbf{v}(\pi) &:= \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in [K]} \mathbf{W}_{c(a)}^\top \mathbf{x}(a) \pi(a, \mathbf{X}) \right]. \end{aligned} \quad (3)$$

With  $\Pi$  denoting the set of all static policies, the optimal static policy is defined as

$$\pi^* := \arg \max_{\pi \in \Pi} r(\pi) \text{ subject to } \mathbf{v}(\pi) \leq \frac{B}{T} \cdot \mathbf{1}_d.$$

The expected total reward of  $\pi^*$  is defined as

$$\text{OPT} := T \cdot r(\pi^*).$$

Since the “no-op” action is allowed, the policy  $\pi^*$  is feasible and the definition of OPT is valid. Having specified the benchmark, we can now define the regret of our algorithm.

**Definition 3.6** (Regret). *The regret of the algorithm for  $T$  time periods is defined as*

$$\text{regret}(T) := \text{OPT} - \sum_{t=1}^T r_t(a_t). \quad (4)$$

## 4 OPTIMISM IN THE FACE OF UNCERTAINTY

In this section, we describe how optimistic estimates of the reward and consumption parameters are derived in time periods after the clustering is performed. Such estimates will consequently be utilized to devise a no-regret algorithm in the next section. At period  $t > N_S \cdot T_0$ , the parameters of cluster  $c \in [C]$  are estimated using all the observations prior to  $t$  that involve arms that have been assigned to  $c$ . We denote by  $t_c < t$  the number of periods in which arms estimated to be in  $c$  have been played prior to period  $t$ , i.e.,  $t_c := \sum_{i=1}^{t-1} \mathbb{1}[\hat{c}(a_i) = c]$ .

To formalize the estimation error that arises both from inherent stochasticity and from occasional clustering mistakes, we begin by introducing a general noisy-observation model. Specifically, we view each observed outcome as the true linear predictor plus an additive noise term that decomposes into a zero-mean fluctuation and a bias component due to mis-clustering.

**Definition 4.1.** *At time step  $t \in [T]$ , let  $\mathbf{x}_t \in [0, 1]^m$  denote the observed context vector, and let  $y_t \in [0, 1]$  be the corresponding observed outcome. We define  $y_t := \boldsymbol{\mu}^\top \mathbf{x}_t + \eta_t$ , where  $\boldsymbol{\mu} \in [0, 1]^m$  is an unknown parameter vector, and  $\eta_t$  is an error term composed of two components:  $\eta_t = u_t + h_t$ , with  $\mathbb{E}[u_t \mid \mathbf{x}_t] = 0$  (from Assumption 3.2) and  $|u_t| \leq 2R$  (from Assumption 3.3). The second component  $h_t$  accounts for clustering mismatch and is defined as*

$$h_t = \begin{cases} 0 & \text{w.p. } 1 - \epsilon, \\ \boldsymbol{\gamma}^\top \mathbf{x}_t & \text{w.p. } \epsilon, \end{cases}$$

where  $\boldsymbol{\gamma} \in [-1, 1]^m$  is perceived as the element-wise difference between parameters of different clusters.

Building on the noisy observation model of Definition 4.1, we now turn to estimating the parameter vector  $\boldsymbol{\mu}$  via regularized least squares. At each time step  $t$ , all past context–outcome pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^{t-1}$

are aggregated into a design matrix, and a ridge regression estimator is computed to balance data fitting against overfitting. Concretely, we form the regularized covariance matrix and then solve for  $\hat{\boldsymbol{\mu}}_t$  as follows.

**Definition 4.2.** *Let  $\lambda_2 > 0$  be the regularization parameter. At each time step  $i < t$ , the agent observes a context vector  $\mathbf{x}_i \in [0, 1]^m$  and its corresponding scalar outcome  $y_i \in [0, 1]$ , where  $y_i = \boldsymbol{\mu}^\top \mathbf{x}_i + \eta_i$ . Then the regularized matrix  $\mathbf{M}_t$  and the regression estimator  $\hat{\boldsymbol{\mu}}_t$  at time  $t$  are defined as:*

$$\begin{aligned} \mathbf{M}_t &:= \lambda_2 \mathbf{I}_m + \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^\top, \\ \hat{\boldsymbol{\mu}}_t &:= \mathbf{M}_t^{-1} \sum_{i=1}^{t-1} \mathbf{x}_i y_i. \end{aligned} \tag{5}$$

By substituting  $r_t(a)$ ,  $\boldsymbol{\mu}_c$  and  $\epsilon_c$  for  $y_t$ ,  $\boldsymbol{\mu}$  and  $\epsilon$  in Definition 4.1, where  $c = \hat{c}(a)$  for each  $a \in \mathcal{S}$ , the concentration results of Section A.4 apply to the reward parameter of cluster  $c$  using its  $t_c$  observations. Similarly, letting  $\mathbf{w}_{c,j}$  be the  $j$ th column of  $\mathbf{W}_c$  for  $j \in [d]$ , the same bounds hold for each consumption dimension. In particular, with  $R = \frac{1}{2}$  and  $\lambda_2 = 1$ , we define the cluster-specific design matrix

$$\mathbf{M}_{c,t} := \mathbf{I}_m + \sum_{i < t: \hat{c}(a_i) = c} \mathbf{x}_i(a_i) \mathbf{x}_i(a_i)^\top.$$

The parameters of cluster  $c \in [C]$  are estimated at period  $t$  as

$$\begin{aligned} \hat{\boldsymbol{\mu}}_{c,t} &:= \mathbf{M}_{c,t}^{-1} \sum_{i < t: \hat{c}(a_i) = c} \mathbf{x}_i(a_i) r_i(a_i) \\ \widehat{\mathbf{W}}_{c,t} &:= \mathbf{M}_{c,t}^{-1} \sum_{i < t: \hat{c}(a_i) = c} \mathbf{x}_i(a_i) \mathbf{v}_i(a_i)^\top. \end{aligned}$$

**Definition 4.3.** Let  $\zeta \in (0, 1)$  be a confidence parameter. At time step  $t$ , define the confidence ellipsoid  $\mathcal{C}_t \subset \mathbb{R}^m$  as the set of vectors  $\beta$  satisfying

$$\mathcal{C}_t := \{\beta \in \mathbb{R}^m : \|\beta - \hat{\mu}_t\|_{M_t} \leq \rho_t\},$$

where the radius  $\rho_t$  is given by

$$\rho_t := 2(R + 1)\sqrt{m \log\left(\frac{tm}{\lambda_2 \zeta}\right)} + \epsilon m \sqrt{t} + \sqrt{\lambda_2 m}.$$

For brevity, We refer the reader to Section A.4 for additional details.

With  $R = \frac{1}{2}$  and  $\lambda_2 = 1$ , we can define confidence ellipsoids for the parameters of each cluster and derive optimistic estimates. The confidence ellipsoid of the reward vector of cluster  $c$  at period  $t$  is defined as  $\mathcal{C}_{\mu,c,t} := \{\beta \in \mathbb{R}^m : \|\beta - \hat{\mu}_{c,t}\|_{M_{c,t}} \leq 3\sqrt{m \log(t_c m / \zeta)} + \epsilon m \sqrt{t_c} + \sqrt{m}\}$ , and the optimistic estimate of the reward parameter for arm  $a \in \mathcal{S}$  at period  $t$  is defined as

$$\tilde{\mu}_{a,t} := \arg \max_{\beta \in \mathcal{C}_{\mu,\hat{c}(a),t}} \mathbf{x}_t(a)^\top \beta. \quad (6)$$

We can similarly define the confidence ellipsoid for the vector of the consumption dimension  $j \in [d]$  for cluster  $c$  at period  $t$  as

$$\mathcal{C}_{w,c,t,j} := \{\beta \in \mathbb{R}^m : \|\beta - \hat{w}_{c,t,j}\|_{M_{c,t}} \leq 3\sqrt{m \log(dt_c m / \zeta)} + \epsilon m \sqrt{t_c} + \sqrt{m}\},$$

where  $\hat{w}_{c,t,j}$  is the  $j^{\text{th}}$  column of  $\widehat{W}_{c,t}$ . Given a vector  $\theta_t \in [0, 1]^d$ , we define the optimistic consumption estimate for arm  $a \in \mathcal{S}$  at time  $t$  by choosing the matrix in the Cartesian product of the  $d$  confidence sets that minimizes the weighted consumption:

$$\widetilde{W}_{a,t} := \arg \min_{W \in \times_{j=1}^d \mathcal{C}_{w,\hat{c}(a),t,j}} \mathbf{x}_t(a)^\top W \theta_t. \quad (7)$$

By Lemma A.11 and the union bound,  $W_c$  is in  $\times_{j=1}^d \mathcal{C}_{w,c,t,j}$  with probability at least  $1 - \zeta$ . The vector  $\theta_t$  will allow the algorithm in the next section to translate consumption into reward, so that arms that are expected to consume plenty of scarce resources will be relatively less appealing. For the optimistic estimate of a reward vector, the maximizer is picked in Eq. (6) since we want the choice of an arm to result in high reward. However, since large budget losses are undesirable, the optimistic estimate of a consumption matrix is taken to be the minimizer in Eq. (7). The following lemma relates the optimistic estimates to the true parameters.

**Lemma 4.4** (Informal version of Lemma C.1). *Given clustering  $\{\hat{c}(a)\}_{a \in \mathcal{S}}$  and vectors  $\{\theta_i\}_{i=N_S \cdot T_0 + 1}^t$ , where  $\theta_i \in [0, 1]^d$ , with probability at least  $1 - \zeta$  we have that for any  $a \in \mathcal{S}$ ,*

- a)  $\mathbf{x}_t(a)^\top (\tilde{\mu}_{a,t} - \mu_{\hat{c}(a)}) \geq 0$ ,
- b)  $\mathbf{x}_t(a)^\top (\widetilde{W}_{a,t} - W_{\hat{c}(a)}) \theta_t \leq 0$ ,
- c)  $|\sum_{i=N_S \cdot T_0 + 1}^t \mathbf{x}_i(a_i)^\top (\tilde{\mu}_{a_i,i} - \mu_{\hat{c}(a_i)})| \leq \rho$ ,
- d)  $\|\sum_{i=N_S \cdot T_0 + 1}^t \mathbf{x}_i(a_i)^\top (\widetilde{W}_{a_i,i} - W_{\hat{c}(a_i)})\|_\infty \leq \rho$ ,

where  $\rho$  is given by:

$$\rho := 4Cm\sqrt{t \log(tm/\zeta) \log(t)} + \epsilon_c m^{\frac{3}{2}} t \sqrt{\log(t)}.$$

For brevity, the proof of Lemma 4.4 is deferred to Section C.

## 5 ALGORITHM

In this section, we present our algorithm for the problem of clustered linear contextual bandits with knapsacks (clusterLCBwK) and results related to its regret. Initially, clustering is performed as specified in subsection A.2. We show that a small, randomly sampled subset of arms is sufficient to identify all underlying clusters with high probability.

**Lemma 5.1** (informal version of Lemma A.4). *For parameter  $\delta > 0$ , if the set  $\mathcal{S}$  is formed by sampling  $N_S = O(p_{\min}^{-1}(T^\delta + \log C))$  arms, where each arm is sampled with equal probability, then  $\mathcal{S}$  covers the clusters, i.e.,  $\cup_{a \in \mathcal{S}} \{c(a)\} = [C]$ , with probability at least  $1 - O(T^{-2\delta})$ .*

**Lemma 5.2** (Informal version of Lemma D.1). *If  $\max_{\pi \in \{\pi^*, \pi'\}} |r(\pi) - \rho(\pi)| \leq \epsilon'$  then  $|r(\pi^*) - \rho(\pi')| \leq \epsilon'$ .*

With Lemma 5.2, it suffices to prove the difference  $\max_{\pi \in \{\pi^*, \pi'\}} |r(\pi) - \rho(\pi)|$  is small. We first consider  $\pi^*$  and then  $\pi'$ .

**Lemma 5.3** (Informal version of Lemma D.2).  $|r(\pi^*) - \rho(\pi^*)| < o(1)$ .

**Lemma 5.4** (Informal version of Lemma D.3).  $|r(\pi') - \rho(\pi')| < o(1)$ .

Combining Lemmas 5.3 and 5.4, we conclude that  $\max_{\pi \in \{\pi^*, \pi'\}} |r(\pi) - \rho(\pi)|$  is vanishing. It remains to show that our empirical estimate of the optimal reward is accurate. Using the initial  $N_S T_0$  samples, we estimate the reward of the optimal static policy, as formalized in the following lemma.

**Lemma 5.5.** *Let  $\widehat{\text{OPT}} := T \cdot \widehat{r}$ , where  $\widehat{r}$  is the estimate of  $r(\pi^*)$  based on the initial  $N_S T_0$  random samples, defined as*

$$\begin{aligned} \widehat{r} &:= \max_{\pi} \frac{K}{N_S^2 T_0} \sum_{t=1}^{N_S T_0} \sum_{a \in \mathcal{S}} \widehat{\boldsymbol{\mu}}_{\widehat{c}(a), N_S T_0 + 1}^\top \mathbf{x}_t(a) \pi(a, \mathbf{X}_t) \\ \text{s.t. } &\frac{K}{N_S^2 T_0} \sum_{t=1}^{N_S T_0} \sum_{a \in \mathcal{S}} \widehat{\mathbf{W}}_{\widehat{c}(a), N_S T_0 + 1}^\top \mathbf{x}_t(a) \pi(a, \mathbf{X}_t) \leq \frac{B}{T} \mathbf{1}_d. \end{aligned}$$

Then,  $\widehat{\text{OPT}} - \text{OPT} = o(1)$  with high probability.

*Proof.* For notational convenience, let

$$\rho(\pi) := \frac{K}{N_S^2 T_0} \sum_{t=1}^{N_S T_0} \sum_{a \in \mathcal{S}} \widehat{\boldsymbol{\mu}}_{\widehat{c}(a), N_S T_0 + 1}^\top \mathbf{x}_t(a) \pi(a, \mathbf{X}_t).$$

Denote the maximizer of the program of the lemma by  $\pi'$ , i.e.,

$$\pi' := \arg \max_{\pi} \rho(\pi) \text{ s.t. } \frac{K}{N_S^2 T_0} \sum_{t=1}^{N_S T_0} \sum_{a \in \mathcal{S}} \widehat{\mathbf{W}}_{\widehat{c}(a), N_S T_0 + 1}^\top \mathbf{x}_t(a) \pi(a, \mathbf{X}_t) \leq \frac{B}{T} \mathbf{1}_d.$$

By the definition of the benchmark policy in Section 3, we have  $\pi^* = \arg \max_{\pi \in \Pi} r(\pi)$  subject to  $\mathbf{v}(\pi) \leq \frac{B}{T} \cdot \mathbf{1}_d$ . For ease of exposition, we will first ignore the constraints and account for them later. Notice that we want to prove that the difference  $|r(\pi^*) - \rho(\pi')|$  is small. The next lemma allows us to work with a more convenient term instead.  $\square$

The estimate  $\widehat{\text{OPT}}$  is then used to define the variable  $Z$  which will contribute to the choice made by the algorithm in periods  $t > N_S T_0$  by appropriately weighting the optimistic estimates of the consumption. In particular, this variable is defined as

$$Z := \frac{N_S \widehat{\text{OPT}}}{2KB'}, \quad (8)$$

where  $B' := B - N_S T_0$  and we also define  $T'$  similarly, i.e.,  $T' := T - N_S T_0$ . In periods  $t > N_S T_0$ , the choice of the algorithm is

$$a_t = \arg \max_{a \in \mathcal{S}} \mathbf{x}_t(a)^\top (\widetilde{\boldsymbol{\mu}}_{a,t} - Z \widetilde{\mathbf{W}}_{a,t} \boldsymbol{\theta}_t),$$

where  $\boldsymbol{\theta}_t \in [0, 1]^d$  is the choice of the online mirror descent algorithm when at the previous period the payoff  $\boldsymbol{\theta}_{t-1}^\top (\mathbf{v}(a_{t-1}) - \frac{B'}{T'} \mathbf{1}_d)$  is observed. Thus, in effect,  $Z$  and  $\boldsymbol{\theta}_t$  allow the resource consumption to be compared to the reward, so that arms estimated to consume a lot of a scarce resource can be avoided.

Our algorithm is an extension of the linear contextual bandits with knapsacks (linCBwK) algorithm (Agrawal & Devanur, 2016), but with the clustering step having been incorporated and accounted for in the derivation of the regret. Moreover, our estimation of OPT utilizes the initial randomly collected  $N_S T_0$  samples, and is different to the corresponding estimation in linCBwK which would lead to additional sampling and thus higher regret in our case.

---

**Algorithm 1** clusterLCBwK
 

---

- 1:  $N_S \leftarrow O(p_{\min}^{-1}(T^\delta + \log C))$
  - 2:  $\mathcal{S} \leftarrow$  random subset of  $[K]$  with size  $N_S$
  - 3:  $T_0 \leftarrow N_S$
  - 4:  $\forall a \in \mathcal{S}$ , play  $T_0$  times the arm  $a$
  - 5: Cluster the arms in  $\mathcal{S}$  per Eq. (11)
  - 6: Compute  $Z$  per Eq. (8)
  - 7: **for**  $t = N_S T_0 + 1, \dots, T$  **do**
  - 8:    $\forall a \in \mathcal{S}$ , obtain  $\tilde{\boldsymbol{\mu}}_{a,t}$  and  $\tilde{\mathbf{W}}_{a,t}$  per Eq. (6), (7)
  - 9:    $a_t \leftarrow \arg \max_{a \in \mathcal{S}} \mathbf{x}_t(a)^\top (\tilde{\boldsymbol{\mu}}_{a,t} - Z \tilde{\mathbf{W}}_{a,t} \boldsymbol{\theta}_t)$
  - 10:   Play the arm  $a_t$ , and observe  $r_t(a_t)$  and  $\mathbf{v}_t(a_t)$
  - 11:   **if**  $\exists j \in [d] : \sum_{i=1}^t \mathbf{v}_i(a_i)^\top \mathbf{e}_j \geq B$  **then** exit
  - 12:   Update  $\mathbf{M}_{c,t+1}$ ,  $\hat{\boldsymbol{\mu}}_{c,t+1}$ , and  $\widehat{\mathbf{W}}_{c,t+1}$ , for  $c = \hat{c}(a_t)$
  - 13:   Update  $\boldsymbol{\theta}_{t+1}$  with the online mirror descent algorithm for payoff  $\boldsymbol{\theta}_t^\top (\mathbf{v}_t(a_t) - \frac{B'}{T'} \mathbf{1}_d)$
  - 14: **end for**
- 

**Lemma 5.6** (Informal version of Lemma D.4). *With probability at least  $(1 - \zeta)^3$  we have:*

- a)  $|\sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \tilde{\boldsymbol{\mu}}_{a_t,t}| \leq R(T)$ ,
- b)  $\|\sum_{t=N_S T_0+1}^{T_\omega} \mathbf{v}_t(a_t) - \mathbf{x}_t(a_t)^\top \tilde{\mathbf{W}}_{a_t,t}\|_\infty \leq R(T)$ .

Now, let  $\mathcal{S}_1$  denote the subset of  $\mathcal{S}$  that contains correctly clustered arms,  $\mathcal{S}_1 := \{a \in \mathcal{S} : \hat{c}(a) = c(a)\}$ .

The following lemma provides a lower bound related to the choice of the algorithm.

**Lemma 5.7** (Informal version of Lemma D.5). *For  $t > N_S T_0$ , the following inequality holds with high probability:*

$$\mathbf{x}_t(a_t)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} - Z \tilde{\mathbf{W}}_{a_t,t} \boldsymbol{\theta}_t) \geq \frac{1}{\sum_{a' \in \mathcal{S}_1} \pi^*(a', \mathbf{X}_t)} \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\boldsymbol{\mu}_{c(a)} - Z \mathbf{W}_{c(a)} \boldsymbol{\theta}_t).$$

Lemma 5.7 implies the weaker condition where expectation is taken over  $\mathbf{X}_t$  only and it is conditional on the past realizations of the context.

**Lemma 5.8** (Informal version of Lemma D.6). *The following inequality holds with high probability:*

$$\begin{aligned} \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{E}_{\mathbf{X}_t} \left[ \mathbf{x}_t(a_t)^\top \tilde{\boldsymbol{\mu}}_{a_t,t} \right] &\geq O\left(\frac{N_S \cdot T_\omega}{K \cdot T} \text{OPT} \right. \\ &\quad \left. + Z \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{E}_{\mathbf{X}_t} \left[ \mathbf{x}_t(a_t)^\top \tilde{\mathbf{W}}_{a_t,t} - \frac{N_S}{K} \cdot \frac{B}{T} \mathbf{1}_d \right] \boldsymbol{\theta}_t \right). \end{aligned}$$

Since in the first  $N_S T_0$  periods the choices are made randomly, and so  $N_S T_0$  of the budget can potentially be consumed, the following lemma which is known from the literature is expressed in terms of  $B' = B - N_S T_0$  and  $T' = T - N_S T_0$  rather than  $B$  and  $T$ .

432 **Lemma 5.9** (Lemma 9 of (Agrawal & Devanur, 2016)).

$$433 \sum_{t=N_S T_0+1}^{T_\omega} \left( \mathbf{x}_t(a_t)^\top \widetilde{\mathbf{W}}_{a_t,t} - \frac{B'}{T'} \mathbf{1}_d \right) \boldsymbol{\theta}_t \geq B' \left( 1 - \frac{T_\omega - N_S T_0}{T'} \right) - R(T).$$

437 The following theorem is our main result, showing that the regret of Algorithm 1 is sublinear in  $T$ .

438 **Theorem 5.10** (Main Result, informal version of D.8). *For  $\delta \in (0, \frac{1}{2})$ , and  $B > N_S T_0$ , with high*

439 *probability*

$$440 \text{regret}(T) \leq O\left(R(T)\left(1 + \frac{N_S \text{OPT}}{KB'}\right) + \text{OPT}\left(1 - \frac{N_S}{K}\right)\right)$$

441 where

$$442 R(T) = O\left(C p_{\min}^{-1} m^{\frac{3}{2}} T^{1-\delta} \sqrt{\log(T)}\right).$$

443 Since  $N_S/K \leq 1$  and the exponent of  $T$  is at most  $1 - \delta < 1$ , the regret is sublinear in the number  
444 of time periods. In the special case where we sampled all arms, i.e.,  $N_S = K$ , the bound simplifies  
445 to  $O(R(T)(1 + \text{OPT}/B'))$  where  $B' = B - N_S T_0$ . This matches the structure of non-clustered  
446 knapsacks bound  $O(R'(T)(1 + \text{OPT}/B))$  (Agrawal & Devanur, 2016), with  $R'(T) = \tilde{O}(mT^{\frac{1}{2}})$   
447 Moreover, the difference between  $B' = B - N_S T_0$  and  $B$  in the division of  $\text{OPT}$  reflects the loss  
448 due to the initial  $N_S T_0$  periods of obtaining samples for the clustering.

449 The probability  $1 - \zeta$  that appears in Lemma 4.4 affects our regret only through  $R'(T)$ . Considering  
450 the case that is of greater interest here, where a subset of the arms is sampled, i.e.,  $N_S < K$ ,  
451 as  $N_S$  decreases, the regret due to  $R(T)$  decreases until  $\frac{N_S \text{OPT}}{KB'} < 1$ , while the term  $\text{OPT}\left(1 - \frac{N_S}{K}\right)$   
452 increases. These two terms express two different sources of regret. More specifically,  $R(T)$   
453 decreases in the number of sampled arms because it captures the difference in the performance  
454 between our algorithm and the optimal static policy with respect to the portion of the arms that are  
455 sampled. The second term,  $\text{OPT}\left(1 - \frac{N_S}{K}\right)$ , reflects the loss suffered due to the fact that as the  
456 number of sampled arms decreases, the options that are available to our algorithm become fewer.  
457 Thus, the algorithm will sometimes miss the opportunity to pick arms with favorable context due  
458 to such arms having been discarded. Therefore, the choice of  $N_S$  determines the impact of each  
459 of these two sources of regret. Understanding the consequences of this choice can be especially  
460 important in applications where for practical reasons operating on the set  $[K]$  is infeasible. Despite  
461 that our results suggest that the dependence of the regret on the number of arms  $K$  may be inevitable,  
462 the weakening of this dependence is an open question.

## 463 6 DISCUSSION

464 We have shown that regret sublinear in the number of time periods can be achieved for the problem  
465 of linear contextual bandits with knapsacks and clusters of arms. Our approach does not require  
466 access to all the available arms, and can thus be utilized in applications where the heterogeneity  
467 of the available choices can be meaningfully summarized with clusters, and where individually  
468 considering each possible choice is unrealistic, e.g., in online advertising campaigns. It is among  
469 our beliefs that the study of forms for summarizing choice heterogeneity in online learning can be a  
470 fruitful research direction.

471 Furthermore, our approach can be extended in a number of ways. For instance, the initially sampled  
472 arms can serve only as a basis for the clustering, and with additional arms being explored and clus-  
473 tered in later steps of the algorithm. Also, the results in (Su et al., 2016) suggest that the assumption  
474 about knowing the number of the clusters can be relaxed. The improvement of the dependency of  
475 the regret on the total number of arms  $K$  is a question we find interesting to be explored. Notice that  
476 we could have improved this dependency here by allowing the number of sampled arms  $N_S$  to be a  
477 function of  $K$ . However, the budget constraints would then depend on  $K$ .

486 ETHIC STATEMENT  
487

488 This paper does not involve human subjects, personally identifiable data, or sensitive applications.  
489 We do not foresee direct ethical risks. We follow the ICLR Code of Ethics and affirm that all aspects  
490 of this research comply with the principles of fairness, transparency, and integrity.  
491

492 REPRODUCIBILITY STATEMENT  
493

494 We ensure reproducibility of our theoretical results by including all formal assumptions, definitions,  
495 and complete proofs in the appendix. The main text states each theorem clearly and refers to the  
496 detailed proofs. No external data or software is required.  
497

498 REFERENCES  
499

- 500 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic  
501 bandits. *Advances in neural information processing systems*, 24:2312–2320, 2011.  
502
- 503 Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming  
504 the monster: A fast and simple algorithm for contextual bandits. In *International Conference on*  
505 *Machine Learning*, pp. 1638–1646. PMLR, 2014.
- 506 Shipra Agrawal and Nikhil Devanur. Linear contextual bandits with knapsacks. *Advances in Neural*  
507 *Information Processing Systems*, 29:3450–3458, 2016.  
508
- 509 Shipra Agrawal and Nikhil R Devanur. Bandits with concave rewards and convex knapsacks. In *Pro-*  
510 *ceedings of the fifteenth ACM conference on Economics and computation*, pp. 989–1006, 2014.
- 511 Shipra Agrawal, Nikhil R Devanur, and Lihong Li. An efficient algorithm for contextual bandits  
512 with knapsacks, and an extension to concave objectives. In *Conference on Learning Theory*, pp.  
513 4–18. PMLR, 2016.
- 514 Sanae Amani, Mahnoosh Alizadeh, and Christos Thrampoulidis. Linear stochastic bandits under  
515 safety constraints. *Advances in Neural Information Processing Systems*, 32:9256–9266, 2019.  
516
- 517 Tomohiro Ando and Jushan Bai. Panel data models with grouped factor structure under unknown  
518 group membership. *Journal of Applied Econometrics*, 31(1):163–191, 2016.
- 519 Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*  
520 *Learning Research*, 3(Nov):397–422, 2002.  
521
- 522 Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit  
523 problem. *Machine learning*, 47(2):235–256, 2002.  
524
- 525 Ashwinkumar Badanidiyuru, John Langford, and Aleksandrs Slivkins. Resourceful contextual ban-  
526 dits. In *Conference on Learning Theory*, pp. 1109–1134. PMLR, 2014.
- 527 Ashwinkumar Badanidiyuru, Robert Kleinberg, and Aleksandrs Slivkins. Bandits with knapsacks.  
528 *Journal of the ACM (JACM)*, 65(3):1–55, 2018.  
529
- 530 Avinash Balakrishnan, Djallel Bouneffouf, Nicholas Mattei, and Francesca Rossi. Using contextual  
531 bandits with behavioral constraints for constrained online movie recommendation. In *IJCAI*, pp.  
532 5802–5804, 2018.
- 533 Yikun Ban and Jingrui He. Local clustering in contextual multi-armed bandits. In *Proceedings of*  
534 *the Web Conference 2021*, pp. 2335–2346, 2021.
- 535 Emil Carlsson, Devdatt Dubhashi, and Fredrik D Johansson. Thompson sampling for bandits with  
536 clustered arms. *arXiv preprint arXiv:2109.01656*, 2021.  
537
- 538 Semih Cayci, Atilla Eryilmaz, and R Srikant. Budget-constrained bandits over general cost and  
539 reward distributions. In *International Conference on Artificial Intelligence and Statistics*, pp.  
4388–4398. PMLR, 2020.

- 540 Sitan Chen, Jerry Li, and Zhao Song. Learning mixtures of linear regressions in subexponential time  
541 via fourier moments. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of*  
542 *Computing*, pp. 587–600, 2020.
- 543 Yanxi Chen, Cong Ma, H Vincent Poor, and Yuxin Chena. Learning mixtures of low-rank models.  
544 *IEEE Transactions on Information Theory*, 2021.
- 546 Xiaotong Cheng, Cheng Pan, and Setareh Maghsudi. Parallel online clustering of bandits via hedonic  
547 game. In *International Conference on Machine Learning*, pp. 5485–5503. PMLR, 2023.
- 548 Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit  
549 feedback. In *Conference on Learning Theory*, 2008.
- 551 Rohan Deb, Aadirupa Saha, and Arindam Banerjee. Think before you duel: Understanding complexities  
552 of preference learning under constrained resources. In Sanjoy Dasgupta, Stephan Mandt,  
553 and Yingzhen Li (eds.), *Proceedings of The 27th International Conference on Artificial Intelligence*  
554 *and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pp. 4546–4554.  
555 PMLR, 02–04 May 2024. URL <https://proceedings.mlr.press/v238/deb24a.html>.
- 557 Nikhil R Devanur, Kamal Jain, Balasubramanian Sivan, and Christopher A Wilkens. Near optimal  
558 online algorithms and fast approximation algorithms for resource allocation problems. In  
559 *Proceedings of the 12th ACM conference on Electronic commerce*, pp. 29–38, 2011.
- 560 Philippe Flajolet, Daniele Gardy, and Loÿs Thimonier. Birthday paradox, coupon collectors, caching  
561 algorithms and self-organizing search. *Discrete Applied Mathematics*, 39(3):207–229, 1992.
- 563 Wayne A Fuller. *Measurement error models*, volume 305. John Wiley & Sons, 2009.
- 564 Claudio Gentile, Shuai Li, and Giovanni Zappella. Online clustering of bandits. In *International*  
565 *Conference on Machine Learning*, pp. 757–765. PMLR, 2014.
- 567 Claudio Gentile, Shuai Li, Purushottam Kar, Alexandros Karatzoglou, Giovanni Zappella, and  
568 Evans Etrue. On context-dependent clustering of bandits. In *International Conference on Machine*  
569 *Learning*, pp. 1253–1262. PMLR, 2017.
- 570 Jiaying Gu and Stanislav Volgushev. Panel data quantile regression with grouped fixed effects.  
571 *Journal of Econometrics*, 213(1):68–91, 2019.
- 573 Yuxuan Han, Jialin Zeng, Yang Wang, Yang Xiang, and Jiheng Zhang. Optimal contextual bandits  
574 with knapsacks under realizability via regression oracles. In Francisco Ruiz, Jennifer Dy,  
575 and Jan-Willem van de Meent (eds.), *Proceedings of The 26th International Conference on Artificial*  
576 *Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp.  
577 5011–5035. PMLR, 25–27 Apr 2023. URL <https://proceedings.mlr.press/v206/han23b.html>.
- 579 Nicole Immorlica, Karthik Abinav Sankararaman, Robert Schapire, and Aleksandrs Slivkins. Adversarial  
580 bandits with knapsacks. In *2019 IEEE 60th Annual Symposium on Foundations of Computer*  
581 *Science (FOCS)*, pp. 202–219. IEEE, 2019.
- 583 Jiashuo Jiang and Yinyu Ye. Achieving  $\tilde{O}(1/\epsilon)$  sample complexity for constrained  
584 markov decision process. In *The Thirty-eighth Annual Conference on Neural Information Processing*  
585 *Systems*, 2024. URL <https://openreview.net/forum?id=psG4LX1DNs>.
- 587 Thomas Kesselheim and Sahil Singla. Online learning with vector costs and bandits with knapsacks.  
588 In *Conference on Learning Theory*, pp. 2286–2305. PMLR, 2020.
- 589 Weihao Kong, Raghav Somani, Zhao Song, Sham Kakade, and Sewoong Oh. Meta-learning for  
590 mixed linear regression. In *International Conference on Machine Learning*, pp. 5394–5404.  
591 PMLR, 2020.
- 592 John Langford and Tong Zhang. Epoch-greedy algorithm for multi-armed bandits with side information.  
593 *Advances in Neural Information Processing Systems (NIPS 2007)*, 20:1, 2007.

- 594 Chuanhao Li, Qingyun Wu, and Hongning Wang. Unifying clustered and non-stationary bandits. In  
595 *International Conference on Artificial Intelligence and Statistics*, pp. 1063–1071. PMLR, 2021a.  
596
- 597 Shuai Li and Shengyu Zhang. Online clustering of contextual cascading bandits. In *Thirty-Second*  
598 *AAAI Conference on Artificial Intelligence*, 2018.
- 599 Xiaocheng Li, Chunlin Sun, and Yinyu Ye. The symmetry between arms and knapsacks: A primal-  
600 dual approach for bandits with knapsacks. In Marina Meila and Tong Zhang (eds.), *Proceedings of*  
601 *the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*  
602 *Learning Research*, pp. 6483–6492. PMLR, 18–24 Jul 2021b. URL <https://proceedings.mlr.press/v139/li21s.html>.  
603  
604
- 605 Yuanzhi Li and Yingyu Liang. Learning mixtures of linear regressions with nearly optimal com-  
606 plexity. In *Conference On Learning Theory*, pp. 1125–1144. PMLR, 2018.
- 607 Chang-Ching Lin and Serena Ng. Estimation of panel data models with parameter heterogeneity  
608 when group membership is unknown. *Journal of Econometric Methods*, 1(1):42–55, 2012.  
609
- 610 Shang Liu, Jiashuo Jiang, and Xiaocheng Li. Non-stationary bandits with knapsacks. In Ai-  
611 ce H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neu-*  
612 *ral Information Processing Systems*, 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=OVb3ZY0fzMk)  
613 [OVb3ZY0fzMk](https://openreview.net/forum?id=OVb3ZY0fzMk).
- 614 Lixing Lyu and Wang Chi Cheung. Bandits with knapsacks: Advice on time-varying demands.  
615 In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and  
616 Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*,  
617 volume 202 of *Proceedings of Machine Learning Research*, pp. 23212–23238. PMLR, 23–29 Jul  
618 2023. URL <https://proceedings.mlr.press/v202/lyu23a.html>.  
619
- 620 Wanteng Ma, Dong Xia, and Jiashuo Jiang. High-dimensional linear bandits with knapsacks. In  
621 Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scar-  
622 lett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine*  
623 *Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 34008–34037. PMLR,  
624 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/ma24p.html>.
- 625 Trong T Nguyen and Hady W Lauw. Dynamic clustering of contextual multi-armed bandits. In *Pro-*  
626 *ceedings of the 23rd ACM International Conference on Conference on Information and Knowl-*  
627 *edge Management*, pp. 1959–1962, 2014.
- 628 Ryo Okui and Wendun Wang. Heterogeneous structural breaks in panel data models. *Journal of*  
629 *Econometrics*, 220(2):447–473, 2021.
- 630 Soumyabrata Pal, Arun Sai Suggala, Karthikeyan Shanmugam, and Prateek Jain. Optimal algo-  
631 rithms for latent bandits with cluster structure. In *International Conference on Artificial Intelli-*  
632 *gence and Statistics*, pp. 7540–7577. PMLR, 2023.
- 633 Anshuka Rangi, Massimo Franceschetti, and Long Tran-Thanh. Unifying the stochastic and the  
634 adversarial bandits with knapsack. In *IJCAI*, 2019.
- 635 Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and trends*  
636 *in Machine Learning*, 4(2):107–194, 2011.
- 637 Vidyashankar Sivakumar, Shiliang Zuo, and Arindam Banerjee. Smoothed adversarial linear con-  
638 textual bandits with knapsacks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba  
639 Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Con-*  
640 *ference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,  
641 pp. 20253–20277. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v162/sivakumar22a.html)  
642 [v162/sivakumar22a.html](https://proceedings.mlr.press/v162/sivakumar22a.html).  
643  
644
- 645 Aleksandrs Slivkins. Contextual bandits with similarity information. In *Proceedings of the 24th*  
646 *annual Conference On Learning Theory*, pp. 679–702. JMLR Workshop and Conference Pro-  
647 ceedings, 2011.

- 648 Aleksandrs Slivkins, Karthik Abinav Sankararaman, and Dylan J Foster. Contextual bandits with  
649 packing and covering constraints: A modular lagrangian approach via regression. In Gergely Neu  
650 and Lorenzo Rosasco (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume  
651 195 of *Proceedings of Machine Learning Research*, pp. 4633–4656. PMLR, 12–15 Jul 2023. URL  
652 <https://proceedings.mlr.press/v195/foster23c.html>.
- 653 Nati Srebro, Karthik Sridharan, and Ambuj Tewari. On the universality of online mirror descent. In  
654 *Advances in neural information processing systems*, pp. 2645–2653, 2011.
- 655 Liangjun Su, Zhentao Shi, and Peter CB Phillips. Identifying latent structures in panel data. *Econo-*  
656 *metrica*, 84(6):2215–2264, 2016.
- 657 Liangjun Su, Xia Wang, and Sainan Jin. Sieve estimation of time-varying panel data models with  
658 latent structures. *Journal of Business & Economic Statistics*, 37(2):334–349, 2019.
- 659 Liang Tang, Yexi Jiang, Lei Li, and Tao Li. Ensemble contextual bandits for personalized recom-  
660 mendation. In *Proceedings of the 8th ACM Conference on Recommender Systems*, pp. 73–80,  
661 2014.
- 662 Zhiyong Wang, Jize Xie, Xutong Liu, Shuai Li, and John Lui. Online clustering of bandits with  
663 misspecified user models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 664 Tom Wansbeek and Erik Meijer. Measurement error and latent variables. *A companion to theoretical*  
665 *econometrics*, pp. 162–179, 2001.
- 666 Huasen Wu, R Srikant, Xin Liu, and Chong Jiang. Algorithms with logarithmic or sublinear regret  
667 for constrained contextual bandits. *Advances in Neural Information Processing Systems*, 28:433–  
668 441, 2015.
- 669 Junwen Yang, Zixin Zhong, and Vincent YF Tan. Optimal clustering with bandit feedback. *Journal*  
670 *of Machine Learning Research*, 25(186):1–54, 2024.
- 671 Liu Yang, Bo Liu, Leyu Lin, Feng Xia, Kai Chen, and Qiang Yang. Exploring clustering of bandits  
672 for online recommendation system. In *Proceedings of the 14th ACM Conference on Recommender*  
673 *Systems*, pp. 120–129, 2020a.
- 674 Mengyue Yang, Qingyang Li, Zhiwei Qin, and Jieping Ye. Hierarchical adaptive contextual bandits  
675 for resource constraint based recommendation. In *Proceedings of The Web Conference 2020*, pp.  
676 292–302, 2020b.
- 677 Xilin Zhang and Wang Chi Cheung. Piecewise-stationary bandits with knapsacks. In *The Thirty-*  
678 *eighth Annual Conference on Neural Information Processing Systems*, 2024. URL [https://](https://openreview.net/forum?id=haa457jwjw)  
679 [openreview.net/forum?id=haa457jwjw](https://openreview.net/forum?id=haa457jwjw).
- 680 Kai Zhong, Prateek Jain, and Inderjit S Dhillon. Mixed linear regression with multiple components.  
681 In *NIPS*, pp. 2190–2198, 2016.
- 682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

# Appendix

# Appendix

**Roadmap.** The appendix is organized as follows.

- In Section A we provide the preliminary of this work.
- In Section B, we present supplementary proofs for Section A.
- In Section C, we present supplementary Lemmas for Section 4.
- In Section D we present the proofs of the results from Section 5 related to the regret of the algorithm.

## A PRELIMINARIES

In this section we present results and technical ingredients that will allow us to derive the regret.

### A.1 PROBABILITY TOOLS

**Lemma A.1** (Cantelli’s inequality). *For random variable  $X$  and constant  $\xi > 0$ ,*

$$\Pr[X \geq \mathbb{E}[X] + \xi] \leq \frac{\text{Var}[X]}{\text{Var}[X] + \xi^2}$$

**Lemma A.2** (Azuma-Hoeffding inequality). *For a martingale  $\{X_i\}_{i=0,1,\dots}$  and constant  $\xi > 0$ , if  $|X_i - X_{i-1}| \leq c_i$  almost surely, then*

$$\Pr[|X_N - X_0| \geq \xi] \leq 2 \exp\left(-\frac{\xi^2}{2 \sum_{i=1}^N c_i^2}\right)$$

**Lemma A.3** ((Su et al., 2016), pp. 2250).  $\Pr[\hat{c}(a) = c(a)] = 1 - o(N_S^{-1})$

### A.2 CLUSTERING UNDER A MIXED LINEAR MODEL

In order to devise an algorithm that achieves low regret, it will be necessary to learn the parameters  $\mu_c$  and  $W_c$  of each cluster  $c \in [C]$ . Consequently, we need to assign arms to clusters, with the correctness of an arm’s assignment being defined only up to permutation of the cluster identities, as a cluster is defined by its members. However, clustering all of the arms requires obtaining at least  $K$  samples, which could be undesirable. Instead, our approach allows clustering and operating on a subset of the arms.

While clustering will be performed only for a subset of the  $K$  arms, it is important that this subset covers all the  $C$  clusters. We initially choose (without playing) at random a subset of the  $K$  arms, denoted by  $S \subseteq [K]$ , with  $|\mathcal{S}| = N_S$ . The set of arms  $[K] \setminus \mathcal{S}$  that is not chosen will be discarded for all the  $T$  time periods. Of course, the regret will be later derived given that the benchmark static policy has all the  $K$  arms in its availability. The following result indicates the size of the set  $\mathcal{S}$  needed in order to cover all of the clusters with a given probability. We use  $p_{\min}$  to denote the proportion of arms belonging to the smallest cluster.

**Lemma A.4** (formal version of Lemma 5.1). *For parameter  $\delta > 0$ , if the set  $\mathcal{S}$  is formed by sampling  $N_S = O(p_{\min}^{-1}(T^\delta + \log C))$  arms, where each arm is sampled with equal probability, then  $\mathcal{S}$  covers the clusters, i.e.,  $\cup_{a \in \mathcal{S}} \{c(a)\} = [C]$ , with probability at least  $1 - O(T^{-2\delta})$ .*

*Proof.* In this proof, we will treat the portions in  $\mathbf{p}$  as probabilities from a distribution by sampling an arm with replacement, instead of without replacement. Since putting a sampled arm back in the sampling distribution only decreases the probability of the next sample being from a non-sampled cluster, the result we will derive will imply a lower bound for our problem.

Consider the following two-step sampling process, repeated for  $j = 1, 2, 3, \dots$ , for the collection of arms: In the first step, a cluster  $c \sim \mathbf{p}$  is drawn and then the arm  $a_j$  with  $c(a_j) = c$  is considered.

In the second step, the arm  $a_j$  is kept in the set of collected arms with probability  $\frac{p_{\min}}{p_c}$ , and it is discarded with probability  $1 - \frac{p_{\min}}{p_c}$ . Let  $l_j$  denote the outcome of the  $j^{\text{th}}$  draw, so that  $l_j = 0$  if the sampled arm was discarded and  $l_j = c(a_j)$  if the arm was kept. Under this two-step sampling process, we have for  $c \in [C]$  that

$$\Pr[l_j = c] = p_c \frac{p_{\min}}{p_c} = p_{\min}$$

and so, the corresponding sampling distribution for a draw is

$$\begin{aligned} \mathbf{p}_0 &:= (\Pr[l_j = 0], \Pr[l_j = 1], \dots, \Pr[l_j = C]) \\ &= (1 - Cp_{\min}, p_{\min}, \dots, p_{\min}) \end{aligned}$$

Thus, with regard to  $[C]$ ,  $\mathbf{p}_0$  is a uniform distribution. Then, the probability of the event that a new draw gives a new cluster when  $i-1$  clusters have already been collected is  $\frac{C-i+1}{C}(1 - (1 - Cp_{\min})) = (C - i + 1)p_{\min}$ . Let  $L$  denote the number of draws needed to cover all of the clusters. From properties of the geometric distribution it follows that

$$\begin{aligned} \mathbb{E}[L] &= p_{\min}^{-1} \sum_{c \in [C]} \frac{1}{C - c + 1} \\ &= O(p_{\min}^{-1} \log C) \end{aligned}$$

and

$$\begin{aligned} \text{Var}[L] &< p_{\min}^{-2} \sum_{c \in [C]} \frac{1}{(C - c + 1)^2} \\ &< 2p_{\min}^{-2} \end{aligned}$$

By Cantelli's inequality we get

$$\begin{aligned} \Pr[L \leq \mathbb{E}[L] + p_{\min}^{-1} T^\delta] &\geq 1 - \frac{\text{Var}[L]}{\text{Var}[L] + p_{\min}^{-2} T^{2\delta}} \\ &> 1 - \frac{1}{1 + \frac{1}{2} T^{2\delta}} \end{aligned}$$

□

The proof of the above lemma relies on arguments similar to those that can be used for the coupon collector problem (Flajolet et al., 1992). Once the set  $\mathcal{S}$  is determined, in order to collect observations to perform clustering, each arm in  $\mathcal{S}$  is played  $T_0$  times, with the precise value of  $T_0$  being defined later in this subsection. Even though a consumption vector is also observed every time an arm is played, for the purposes of the clustering only, we will utilize just the context and the reward, but we could as well have chosen to utilize one of the  $d$  resources instead of the reward. In any case, the resource consumption due to these initial  $N_S \cdot T_0$  plays will still have to be subtracted from the budget.

The arms in  $\mathcal{S}$  are clustered using the classifier-Lasso method proposed by (Su et al., 2016), which relies on the following objective function

$$\begin{aligned} &Q((\boldsymbol{\mu}_a)_{a \in \mathcal{S}}, (\boldsymbol{\mu}_c)_{c \in [C]}) \\ &= \frac{1}{N_S \cdot T_0} \sum_{a \in \mathcal{S}} \sum_{t: a_t = a} \frac{1}{2} (r_t(a) - \boldsymbol{\mu}_a^\top \mathbf{x}_t(a))^2 \\ &+ \frac{\lambda_1}{N_S} \sum_{a \in \mathcal{S}} \prod_{c \in [C]} \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_c\| \end{aligned} \tag{9}$$

where  $\lambda_1 \in \mathbb{R}_+$  is a regularization parameter. It is valuable to point out that  $N_S + C$  vectors of parameters are estimated in total, since the method does not impose the parameter of an arm to be

810 corresponding to one of the cluster parameters. The estimation of the reward parameters follows by  
 811 minimizing this objective function,

$$812 \begin{aligned} & ((\hat{\boldsymbol{\mu}}_a)_{a \in \mathcal{S}}, (\hat{\boldsymbol{\mu}}_c)_{c \in [C]}) \\ 813 & := \arg \min_{((\boldsymbol{\mu}_a)_{a \in \mathcal{S}}, (\boldsymbol{\mu}_c)_{c \in [C]})} Q((\boldsymbol{\mu}_a)_{a \in \mathcal{S}}, (\boldsymbol{\mu}_c)_{c \in [C]}). \end{aligned} \quad (10)$$

814 Then, the arm  $a \in \mathcal{S}$  is clustered as

$$815 \hat{c}(a) := \sum_{c \in [C]} c \cdot \mathbb{1}[\hat{\boldsymbol{\mu}}_a = \hat{\boldsymbol{\mu}}_c] \quad (11)$$

816 Since under the classifier-Lasso method it is possible for an arm's estimated parameter to not be  
 817 equal to the estimated parameter of any cluster, an arm can be assigned to none of the clusters,  
 818 leading to the case where  $\hat{c}(a) = 0$ . However, as shown in the proof of the focal result of this  
 819 subsection, only few of the arms are not assigned to any cluster.

820 Besides the coverage of the clusters, we also want to achieve high clustering accuracy. We define  
 821 the clustering error for cluster  $c \in [C]$  as

$$822 \epsilon_c := \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c]} \quad (12)$$

823 Thus, the clustering error for  $c$  is the proportion of arms assigned to  $c$  that should have been assigned  
 824 to other clusters. We want to ensure that for each cluster  $c$ , the output of Eq. (11) results to low error  
 825  $\epsilon_c$ . Towards this end, we make the following assumptions.

826 **Assumption A.5** (Separation). *We assume  $\|\boldsymbol{\mu}_c - \boldsymbol{\mu}_{c'}\| \geq \xi_1 > 0$ , for  $c, c' \in [C], c \neq c'$ .*

827 **Assumption A.6.** *We assume the number of clusters  $C$  is fixed.*

828 **Assumption A.7.** *We assume*

$$829 T_0 \lambda_1^2 / (\log T_0)^{6+2\xi_2} \rightarrow \infty,$$

830 and

$$831 \lambda_1 (\log T_0)^{\xi_2} \rightarrow 0,$$

832 for  $\xi_2 > 0$  as  $(N_S, T_0) \rightarrow \infty$ .

833 **Assumption A.8.** *We assume*

$$834 N_S^{1/2} T_0^{-1} (\log T_0)^9 \rightarrow 0,$$

835 and

$$836 N_S^2 T_0^{1-\xi_3/2} \rightarrow \xi_4 < \infty,$$

837 for  $\xi_3 \geq 6$  as  $(N_S, T_0) \rightarrow \infty$ .

838 Assumption A.5 is needed to ensure that the clusters are distinguishable, Assumption A.6 disallows  
 839 the number of clusters to grow asymptotically. While Assumptions A.7 and A.8 impose appropriate  
 840 rates on  $\lambda_1, N_S$ , and  $T_0$ , the algorithm need to perform clustering based on one of the  $d$  resources  
 841 instead of the reward.

842 **Lemma A.9.** *Under Assumptions 3.1-3.4, A.5-A.8, the clustering error is  $\epsilon_c = o(p_{\min}^{-1} N_S^{-1})$  with  
 843 high probability, for any  $c \in [C]$ .*

844 *Proof.* In this proof, we will exploit the Azuma-Hoeffding inequality and a result from (Su et al.,  
 845 2016) that holds under the stated assumptions.

846 Now, for the clustering error  $\epsilon_c$  we have

$$847 \begin{aligned} \epsilon_c &= \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c]} \\ 848 &= \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{\sum_{a \in \mathcal{S}} \mathbb{1}[c(a) = c] - \sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) \neq c, c(a) = c] + \sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]} \end{aligned}$$

$$\begin{aligned}
& \leq \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{\sum_{a \in \mathcal{S}} \mathbb{1}[c(a) = c] - \sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) \neq c, c(a) = c]} \\
& = \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{\sum_{a: c(a)=c} 1 - \mathbb{1}[\hat{c}(a) \neq c]} \\
& \leq \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{O(1) \mathbb{E}[\sum_{a: c(a)=c} 1 - \mathbb{1}[\hat{c}(a) \neq c]]} \tag{13} \\
& = \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{O(1) \mathbb{E}[\sum_{a \in \mathcal{S}} \mathbb{1}[c(a) = c, \hat{c}(a) = c]]} \\
& = \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{O(1) \sum_{a \in \mathcal{S}} \Pr[c(a) = c] \Pr[\hat{c}(a) = c | c(a) = c]} \\
& = \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{O(1) N_S p_c (1 - o(N_S^{-1}))} \\
& \leq \frac{\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]}{O(1) N_S p_{\min} (1 - o(N_S^{-1}))} \\
& \leq O(1) \frac{\mathbb{E}[\sum_{a \in \mathcal{S}} \mathbb{1}[\hat{c}(a) = c, c(a) \neq c]]}{N_S p_{\min} (1 - o(N_S^{-1}))} \tag{14} \\
& = O(1) \frac{\sum_{a \in \mathcal{S}} \Pr[c(a) \neq c] \Pr[\hat{c}(a) = c | c(a) \neq c]}{N_S p_{\min} (1 - o(N_S^{-1}))} \\
& \leq O(1) \frac{\sum_{a \in \mathcal{S}} \Pr[c(a) \neq c] \Pr[\hat{c}(a) \neq c(a)]}{N_S p_{\min} (1 - o(N_S^{-1}))} \\
& = O(1) \frac{N_S (1 - p_c) o(N_S^{-1})}{N_S p_{\min} (1 - o(N_S^{-1}))} \\
& = O(1) \frac{(1 - p_c) o(N_S^{-1})}{p_{\min} (1 - o(N_S^{-1}))} \\
& \leq \frac{o(p_{\min}^{-1} N_S^{-1})}{1 - o(N_S^{-1})} \\
& \leq o(p_{\min}^{-1} N_S^{-1})
\end{aligned}$$

where Eq. (13) and (14) follow from the Azuma-Hoeffding inequality.  $\square$

In the proof of Lemma A.9 we exploit the fact that  $\Pr[\hat{c}(a) = c(a)] = 1 - o(N_S^{-1})$  (Su et al., 2016, pp. 2250). Since the clustering error is, in asymptotic terms, the same for all clusters, we shall use  $\epsilon_c$  to refer to this error for any cluster. Moreover, since  $N_S$  is given by Lemma A.4, the number of samples required from each arm in  $\mathcal{S}$  for the clustering can be specified by satisfying Assumption A.8.

**Claim A.10.** For  $T_0 = N_S$ , Assumption A.8 is satisfied.

### A.3 ONLINE LEARNING

A special case of the online convex optimization problem is the game where at period  $t \in [T]$  a learner chooses

$$\theta_t \in \{\theta \in [0, 1]^d : \|\theta\|_1 \leq 1\}$$

based on past observations, and the adversary chooses the learner's payoff to be the outcome of a function that is linear in  $\theta_t$ . It is known (Srebro et al., 2011; Shalev-Shwartz et al., 2011) that for  $\theta_t$  chosen based on the online mirror descent algorithm, the regret against the best fixed action of the learner in the online convex optimization problem is  $O(\sqrt{\log(d)T})$ . In the context of our problem, we utilize this result following the approach of (Agrawal & Devanur, 2016). In particular,  $\theta_t$  will allow the consideration of resources in the arm choices in periods  $t > N_S T_0$ , decreasing thus the probability that a choice will lead to depletion of one of the  $d$  resources and consequently to the

918 termination of the algorithm. To derive the regret for our problem (Eq. (4)), we will consider as the  
 919 payoff chosen by the adversary the value

$$921 \quad \boldsymbol{\theta}_t^\top (\mathbf{v}_t(a_t) - \frac{B - N_S T_0}{T - N_S T_0} \mathbf{1}_d),$$

922 as illustrated by the algorithm presented later.

#### 925 A.4 CONFIDENCE ELLIPSOID

926 In this subsection we present results about bounds on the parameter estimates that will be derived by  
 927 the algorithm after the clustering is conducted. The approach here is to express the clustering error  
 928 as violation of the zero conditional mean assumption about the error terms. Then, the technique of  
 929 confidence ellipsoids, which is standard in the literature of online learning, is employed to bound the  
 930 parameter estimates. The definition that follows considers a linear model for a response variable  $y_t$   
 931 that generalizes the response variables in the problem we consider, i.e.,  $r_t(a)$  and each dimension of  
 932  $\mathbf{v}_t(a)$ . Therefore, the results derived in this subsection will imply results about our actual problem.

933 This definition corresponds to a measurement error model (Wansbeek & Meijer, 2001; Fuller, 2009),  
 934 as the expected value of the error term is  $\mathbb{E}[\eta_t | \mathbf{x}_t] = \epsilon \boldsymbol{\gamma}^\top \mathbf{x}_t$ , which in the general case is not zero.  
 935 The error term  $u_t$  corresponds to an error term of the model of our actual problem, i.e.,  $u_t$  represents  
 936  $g_t(a)$  and the individual dimensions of  $\mathbf{q}_t(a)$ . Since the context is i.i.d. across arms and periods, the  
 937 probability  $\epsilon$  is perceived as the clustering error in Eq. (12). Thus,  $\boldsymbol{\gamma}$  is perceived as the element-  
 938 wise difference between parameters of different clusters. For instance, considering  $C = 2$  and  $y_t$   
 939 corresponding to  $r_t(a)$  for  $c(a) = 1$ , we have that  $\boldsymbol{\gamma}$  corresponds to  $\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1$ . Therefore,  $h_t$  represents  
 940 the part of the error term due to  $\mathbf{x}_t$  being endogenous, when zero clustering error is falsely assumed.  
 941 We let  $h_t$  have two branches instead of  $C$  because  $\boldsymbol{\gamma}$  can be perceived as a bound on parameter  
 942 differences.

943 This confidence ellipsoid captures the set of plausible values for the unknown parameter vector  $\boldsymbol{\mu}$ ,  
 944 given the observations up to time  $t$ . The radius  $\rho_t$  consists of three terms: a statistical concentration  
 945 term accounting for zero-mean stochastic noise  $u_t$  (as in Assumption 3.2), a bias term  $\epsilon m \sqrt{t}$  intro-  
 946 duced by clustering errors, and a regularization term  $\sqrt{\lambda_2 m}$  reflecting uncertainty due to the use of  
 947 ridge regression. With high probability (at least  $1 - \zeta$ ), the true parameter vector  $\boldsymbol{\mu}$  lies within  $\mathcal{C}_t$ .

948 The next two lemmas are the results derived in this subsection.

949 **Lemma A.11.** *Under the zero-mean and bounded-noise assumptions (Assumption 3.2 and 3.3) and  
 950 with regularization parameter  $\lambda_2$ , define the confidence radiue*

$$951 \quad \rho_t := 2(R + 1) \sqrt{m \log \left( \frac{tm}{\lambda_2 \zeta} \right)} + \epsilon m \sqrt{t} + \sqrt{\lambda_2 m}.$$

952 Then, for any time  $t \in [T]$  with probability at least  $1 - \zeta$ , the true parameter vector  $\boldsymbol{\mu}$  satisfies

$$953 \quad \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_t\|_{\mathbf{M}_t} \leq \rho_t$$

954 *Proof.* By lemmas B.4, B.5, B.6, and B.7, we have that for any  $\zeta \in (0, 1)$ , with probability at least  
 955  $1 - \zeta$ ,

$$956 \quad |\mathbf{x}^\top \hat{\boldsymbol{\mu}}_t - \mathbf{x}^\top \boldsymbol{\mu}| \leq \|\mathbf{x}\|_{\mathbf{M}_t^{-1}} \cdot A_1$$

957 where

$$958 \quad A_1 := \left( (R + 1) \sqrt{2 \log \left( \det(\mathbf{M}_t)^{1/2} \det(\lambda_2 \mathbf{I}_m)^{-1/2} / \zeta \right)} \right. \\ 959 \quad \left. + \epsilon \bar{\gamma} m \sqrt{t} + \sqrt{\lambda_2 m} \right)$$

960 Now, by letting  $\mathbf{x} = \mathbf{M}_t(\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu})$ , we have

$$961 \quad \|\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}\|_{\mathbf{M}_t}^2 \leq \|\mathbf{M}_t(\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu})\|_{\mathbf{M}_t^{-1}} \cdot A_1$$

$$= \|\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}\|_{\mathbf{M}_t} \cdot A_1$$

Dividing both sides by  $\|\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}\|_{\mathbf{M}_t}$  we get

$$\|\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}\|_{\mathbf{M}_t} \leq A_1 \quad (15)$$

Since  $\|\mathbf{x}_t\|_2 \leq \sqrt{m}$ , and  $\mathbf{M}_t$  and  $\lambda_2 \mathbf{I}_m$  are positive-definite matrices, we can upper bound the first term in  $A_1$  (ignoring the term  $R + 1$ ) as follows

$$\sqrt{2 \log \left( \frac{\det(\mathbf{M}_t)^{1/2}}{\det(\lambda_2 \mathbf{I}_m)^{1/2} / \zeta} \right)} \leq \sqrt{m \log \left( \frac{1 + tm/\lambda_2}{\zeta} \right)} \quad (16)$$

and thus

$$\begin{aligned} & \|\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}\|_{\mathbf{M}_t} \\ & \leq (R + 1) \sqrt{m \log \left( \frac{1 + tm/\lambda_2}{\zeta} \right)} + \epsilon \bar{\gamma} m \sqrt{t} + \sqrt{\lambda_2 m} \end{aligned} \quad (17)$$

$$\leq 2(R + 1) \sqrt{m \log \left( \frac{tm}{\lambda_2 \zeta} \right)} + \epsilon m \sqrt{t} + \sqrt{\lambda_2 m} \quad (18)$$

where  $\bar{\gamma} \leq 1$ .  $\square$

**Lemma A.12** (Sum of rewards). *Consider  $\tilde{\boldsymbol{\mu}}_t \in \mathcal{C}_t$ . For  $R = \frac{1}{2}$  and  $\lambda_2 = 1$ , with probability at least  $1 - \zeta$*

$$\begin{aligned} & \sum_{t=1}^T |\tilde{\boldsymbol{\mu}}_t^\top \mathbf{x}_t - \boldsymbol{\mu}^\top \mathbf{x}_t| \\ & \leq 4m \sqrt{T \log(Tm/\zeta) \log(T)} + \epsilon m^{\frac{3}{2}} T \sqrt{\log(T)}. \end{aligned}$$

Lemma A.11 implies that a parameter estimate lies within some fixed distance from the true parameter with some specified probability, while Lemma A.12 serves as the basis for employing optimistic parameter estimates, as described later.

## B USEFUL TOOLS FOR SECTION A

Here in this section, we provide some missing proofs for Section A. In Section B.1 we provide some lemmas for Lemma A.11. In Section B.2 we provide proof for Lemma A.12.

### B.1 PROOF OF LEMMA A.11

For notational convenience, let

**Definition B.1.**

$$\mathbf{X} := (\mathbf{x}_1^\top, \dots, \mathbf{x}_t^\top) \in [0, 1]^{t \times m}$$

$$\mathbf{y} := (y_1, \dots, y_t) \in [0, 1]^t$$

$$\boldsymbol{\eta} := (\eta_1, \dots, \eta_t) \in \mathbb{R}^t$$

$$\mathbf{s}_t := \sum_{i=1}^{t-1} \mathbf{x}_i (\eta_i - \mathbb{E}[\eta_i | \mathbf{x}_i])$$

$$\bar{\gamma} := \|\boldsymbol{\gamma}\|_\infty$$

$$\langle \mathbf{a}, \mathbf{b} \rangle_{\mathbf{M}} := \mathbf{a}^\top \mathbf{M} \mathbf{b}, \text{ for } \mathbf{a}, \mathbf{b} \in \mathbb{R}^m, \mathbf{M} \in \mathbb{R}^{m \times m}$$

1026 For the proof we will need the following two Facts and four Lemmas.

1027 **Fact B.2.** We can show that  $\hat{\boldsymbol{\mu}}_t = \boldsymbol{\mu} - \lambda_2 \mathbf{M}_t^{-1} \boldsymbol{\mu} + \mathbf{M}_t^{-1} \mathbf{X}^\top \boldsymbol{\eta}$

1028  
1029  
1030 *Proof.*

$$\begin{aligned}
1031 \hat{\boldsymbol{\mu}}_t &= \mathbf{M}_t^{-1} \mathbf{X}^\top \mathbf{y} \\
1032 &= \mathbf{M}_t^{-1} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\mu} + \boldsymbol{\eta}) \\
1033 &= \mathbf{M}_t^{-1} \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu} + \mathbf{M}_t^{-1} \mathbf{X}^\top \boldsymbol{\eta} \\
1034 &= \mathbf{M}_t^{-1} \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu} + \mathbf{M}_t^{-1} \lambda_2 \mathbf{I}_m \boldsymbol{\mu} - \mathbf{M}_t^{-1} \lambda_2 \mathbf{I}_m \boldsymbol{\mu} + \mathbf{M}_t^{-1} \mathbf{X}^\top \boldsymbol{\eta} \\
1035 &= \mathbf{M}_t^{-1} \mathbf{M}_t \boldsymbol{\mu} - \lambda_2 \mathbf{M}_t^{-1} \boldsymbol{\mu} + \mathbf{M}_t^{-1} \mathbf{X}^\top \boldsymbol{\eta} \\
1036 &= \boldsymbol{\mu} - \lambda_2 \mathbf{M}_t^{-1} \boldsymbol{\mu} + \mathbf{M}_t^{-1} \mathbf{X}^\top \boldsymbol{\eta}
\end{aligned}$$

1037 Thus, we complete the proof. □

1038 **Fact B.3.** For any  $\mathbf{x} \in \mathbb{R}^m$ , we have that  $\mathbf{x}^\top \hat{\boldsymbol{\mu}}_t - \mathbf{x}^\top \boldsymbol{\mu} = \langle \mathbf{x}, \mathbf{X}^\top \boldsymbol{\eta} \rangle_{\mathbf{M}_t^{-1}} - \lambda_2 \langle \mathbf{x}, \boldsymbol{\mu} \rangle_{\mathbf{M}_t^{-1}}$

1039  
1040 *Proof.* From Fact B.2 we have

$$\begin{aligned}
1041 \mathbf{x}^\top \hat{\boldsymbol{\mu}}_t - \mathbf{x}^\top \boldsymbol{\mu} &= \mathbf{x}^\top (\boldsymbol{\mu} - \lambda_2 \mathbf{M}_t^{-1} \boldsymbol{\mu} + \mathbf{M}_t^{-1} \mathbf{X}^\top \boldsymbol{\eta}) - \mathbf{x}^\top \boldsymbol{\mu} \\
1042 &= \mathbf{x}^\top \mathbf{M}_t^{-1} \mathbf{X}^\top \boldsymbol{\eta} - \lambda_2 \mathbf{x}^\top \mathbf{M}_t^{-1} \boldsymbol{\mu} \\
1043 &= \langle \mathbf{x}, \mathbf{X}^\top \boldsymbol{\eta} \rangle_{\mathbf{M}_t^{-1}} - \lambda_2 \langle \mathbf{x}, \boldsymbol{\mu} \rangle_{\mathbf{M}_t^{-1}}
\end{aligned}$$

1044 Thus, we complete the proof. □

1045 **Lemma B.4.** We can show that

$$1046 |\mathbf{x}^\top \hat{\boldsymbol{\mu}}_t - \mathbf{x}^\top \boldsymbol{\mu}| \leq \|\mathbf{x}\|_{\mathbf{M}_t^{-1}} \left( \|\mathbf{s}_t\|_{\mathbf{M}_t^{-1}} + \epsilon \|\mathbf{X}^\top \mathbf{X} \boldsymbol{\gamma}\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2 m} \right)$$

1047  
1048  
1049 *Proof.* Using the Cauchy-Schwarz inequality on Fact B.3, we have

$$1050 |\mathbf{x}^\top \hat{\boldsymbol{\mu}}_t - \mathbf{x}^\top \boldsymbol{\mu}| \leq \|\mathbf{x}\|_{\mathbf{M}_t^{-1}} \left( \|\mathbf{X}^\top \boldsymbol{\eta}\|_{\mathbf{M}_t^{-1}} + \lambda_2 \|\boldsymbol{\mu}\|_{\mathbf{M}_t^{-1}} \right)$$

1051 We can upper bound the second term of the above equation as follows:

$$\begin{aligned}
1052 \|\mathbf{X}^\top \boldsymbol{\eta}\|_{\mathbf{M}_t^{-1}} + \lambda_2 \|\boldsymbol{\mu}\|_{\mathbf{M}_t^{-1}} &\leq \|\mathbf{X}^\top \boldsymbol{\eta}\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2} \|\boldsymbol{\mu}\|_2 \\
1053 &\leq \|\mathbf{X}^\top \boldsymbol{\eta}\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2 m} \\
1054 &= \|\mathbf{X}^\top (\boldsymbol{\eta} + \mathbb{E}[\boldsymbol{\eta} | \mathbf{X}] - \mathbb{E}[\boldsymbol{\eta} | \mathbf{X}])\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2 m} \\
1055 &= \|\mathbf{s}_t + \mathbf{X}^\top \mathbb{E}[\boldsymbol{\eta} | \mathbf{X}]\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2 m} \\
1056 &= \|\mathbf{s}_t + \epsilon \mathbf{X}^\top \mathbf{X} \boldsymbol{\gamma}\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2 m} \\
1057 &\leq \|\mathbf{s}_t\|_{\mathbf{M}_t^{-1}} + \|\epsilon \mathbf{X}^\top \mathbf{X} \boldsymbol{\gamma}\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2 m} \\
1058 &= \|\mathbf{s}_t\|_{\mathbf{M}_t^{-1}} + \epsilon \|\mathbf{X}^\top \mathbf{X} \boldsymbol{\gamma}\|_{\mathbf{M}_t^{-1}} + \sqrt{\lambda_2 m}
\end{aligned}$$

1059 where the second inequality follows from that  $\|\boldsymbol{\mu}\|_{\mathbf{M}_t^{-1}}^2 \leq \frac{1}{\lambda_{\min}(\mathbf{M}_t) \|\boldsymbol{\mu}\|_2^2} \leq \frac{1}{\lambda_2 \|\boldsymbol{\mu}\|_2^2}$ . □

1060 **Lemma B.5.** Consider the  $\sigma$ -algebra  $F_t = \sigma(\mathbf{x}_1, \dots, \mathbf{x}_t, \eta_1, \dots, \eta_{t-1})$ , such that  $\{F_t\}_{t=1}^\infty$  is a filtration, and  $\eta_t - \mathbb{E}[\eta_t | \mathbf{x}_t]$  is  $F_t$ -measurable. Then,  $\eta_t - \mathbb{E}[\eta_t | \mathbf{x}_t]$  is  $(R + 1)$ -sub-Gaussian.

1080 *Proof.* It suffices to show that  $\eta_t - \mathbb{E}[\eta_t|\mathbf{x}_t]$  lies in an interval of length at most  $2(R+1)$ . We can  
 1081 upper bound  $|\eta_t - \mathbb{E}[\eta_t|\mathbf{x}_t]|$  as follows:  
 1082

$$\begin{aligned} 1083 \quad |\eta_t - \mathbb{E}[\eta_t|\mathbf{x}_t]| &= |u_t + h_t - \epsilon\boldsymbol{\gamma}^\top \mathbf{x}_t| \\ 1084 \quad &\leq |u_t| + |h_t| + |\epsilon\boldsymbol{\gamma}^\top \mathbf{x}_t| \\ 1085 \quad &\leq 2R + (1 + \epsilon)|\boldsymbol{\gamma}^\top \mathbf{x}_t| \\ 1086 \quad &\leq 2(R+1) \end{aligned}$$

□

1089 **Lemma B.6.** For any  $\zeta \in (0, 1)$ , with probability at least  $1 - \zeta$ , if  $\eta_t - \mathbb{E}[\eta_t|\mathbf{x}_t]$  is  $(R+1)$ -sub-  
 1090 Gaussian, then for all  $t > 0$ ,

$$1092 \quad \|\mathbf{s}_t\|_{\mathbf{M}_t^{-1}}^2 \leq 2(R+1)^2 \log(\det(\mathbf{M}_t)^{1/2} \det(\lambda_2 \mathbf{I}_m)^{-1/2} / \zeta)$$

1094 *Proof.* We apply Theorem 1 of (Abbasi-Yadkori et al., 2011) with  $\eta_t - \mathbb{E}[\eta_t|\mathbf{x}_t]$  in the place of  
 1095  $\eta_t$ . □

1097 **Lemma B.7.** We can show that

$$1098 \quad \|\mathbf{X}^\top \mathbf{X} \boldsymbol{\gamma}\|_{\mathbf{M}_t^{-1}} \leq \bar{\gamma} m \sqrt{t}$$

1101 *Proof.* Since  $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{m \times m}$  contains non-negative entries, we have

$$1102 \quad \|\mathbf{X}^\top \mathbf{X} \boldsymbol{\gamma}\|_{\mathbf{M}_t^{-1}} \leq \bar{\gamma} \cdot \|\mathbf{X}^\top \mathbf{X} \mathbf{1}_m\|_{\mathbf{M}_t^{-1}}$$

1104 Next, we just need to upper bound  $\|\mathbf{X}^\top \mathbf{X} \mathbf{1}_m\|_{\mathbf{M}_t^{-1}}^2$ . By the properties of PSD/PD matrices, we  
 1105 know that

$$1107 \quad \lambda_2 \mathbf{I}_m + \mathbf{X}^\top \mathbf{X} > \mathbf{X}^\top \mathbf{X} \geq 0$$

1108 which implies that

$$1110 \quad (\mathbf{X}^\top \mathbf{X})^{1/2} \cdot (\lambda_2 \mathbf{I}_m + \mathbf{X}^\top \mathbf{X})^{-1} \cdot (\mathbf{X}^\top \mathbf{X})^{1/2} < \mathbf{I}_m$$

1112 Thus, we have

$$\begin{aligned} 1113 \quad \mathbf{1}_m^\top \mathbf{X}^\top \mathbf{X} (\lambda_2 \mathbf{I}_m + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{1}_m &\leq \mathbf{1}_m^\top \mathbf{X}^\top \mathbf{X} \mathbf{1}_m \\ 1114 \quad &= \|\mathbf{X} \mathbf{1}_m\|_2^2 \\ 1115 \quad &\leq tm^2 \end{aligned}$$

1116 where the last step follows since each entry of  $\mathbf{X}$  is between 0 and 1. □

## 1120 B.2 PROOF OF LEMMA A.12

1121 We will need the following fact and lemma.

1123 **Fact B.8.** For any positive definite matrix  $\mathbf{M} \in \mathbb{R}^{m \times m}$  and any two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$ , it holds that  
 1124  $|\mathbf{a}^\top \mathbf{b}| \leq \|\mathbf{a}\|_{\mathbf{M}} \|\mathbf{b}\|_{\mathbf{M}^{-1}}$ .

1125 **Lemma B.9** (Lemma 3 of (Agrawal & Devanur, 2016)). For  $\mathbf{x}_i \in \mathbb{R}^m$  with  $\|\mathbf{x}_i\|_2 \leq \sqrt{m}$ , it holds  
 1126 that

$$1127 \quad \sum_{i=1}^t \|\mathbf{x}_i\|_{\mathbf{M}_i^{-1}} \leq \sqrt{mt \log(t)}$$

1131 Now, we derive the statement of the lemma as:

$$1132 \quad \sum_{t=1}^T |\tilde{\boldsymbol{\mu}}_t^\top \mathbf{x}_t - \boldsymbol{\mu}^\top \mathbf{x}_t| \leq \sum_{t=1}^T \|\mathbf{x}_t\|_{\mathbf{M}_t^{-1}} \|\tilde{\boldsymbol{\mu}}_t - \boldsymbol{\mu}\|_{\mathbf{M}_t}$$

$$\begin{aligned}
&\leq \left(3\sqrt{m \log(Tm/\zeta)} + \epsilon m\sqrt{T} + \sqrt{m}\right) \sum_{t=1}^T \|\mathbf{x}_t\|_{\mathbf{M}_t^{-1}} \\
&\leq \left(4\sqrt{m \log(Tm/\zeta)} + \epsilon m\sqrt{T}\right) \sum_{t=1}^T \|\mathbf{x}_t\|_{\mathbf{M}_t^{-1}} \\
&\leq \left(4\sqrt{m \log(Tm/\zeta)} + \epsilon m\sqrt{T}\right) \sqrt{mT \log(T)} \\
&= 4m\sqrt{T \log(Tm/\zeta) \log(T)} + \epsilon m^{\frac{3}{2}} T \sqrt{\log(T)}
\end{aligned}$$

where the first step follows from Fact B.8, the second from Lemma A.11, and the fourth from Lemma B.9.

## C USEFUL LEMMAS FOR SECTION 4

Here we provide the following lemmas.

**Lemma C.1** (Informal version of Lemma 4.4). *Given clustering  $\{\hat{c}(a)\}_{a \in \mathcal{S}}$  and vectors  $\{\boldsymbol{\theta}_i\}_{i=N_S \cdot T_0+1}^t$ , where  $\boldsymbol{\theta}_i \in [0, 1]^d$ , with probability at least  $1 - \zeta$  we have that for any  $a \in \mathcal{S}$ ,*

- a)  $\mathbf{x}_t(a)^\top (\tilde{\boldsymbol{\mu}}_{a,t} - \boldsymbol{\mu}_{\hat{c}(a)}) \geq 0$
- b)  $\mathbf{x}_t(a)^\top (\tilde{\mathbf{W}}_{a,t} - \mathbf{W}_{\hat{c}(a)}) \boldsymbol{\theta}_t \leq 0$
- c)  $|\sum_{i=N_S \cdot T_0+1}^t \mathbf{x}_i(a_i)^\top (\tilde{\boldsymbol{\mu}}_{a_i,i} - \boldsymbol{\mu}_{\hat{c}(a_i)})| \leq \rho$
- d)  $\|\sum_{i=N_S \cdot T_0+1}^t \mathbf{x}_i(a_i)^\top (\tilde{\mathbf{W}}_{a_i,i} - \mathbf{W}_{\hat{c}(a_i)})\|_\infty \leq \rho$

Where  $\rho$  is given by:

$$\rho := 4Cm\sqrt{t \log(tm/\zeta) \log(t)} + \epsilon_c m^{\frac{3}{2}} t \sqrt{\log(t)}$$

*Proof.* Statements a) and b) follow directly from Eq. (6) and (7). For statement c) we have:

$$\begin{aligned}
\left| \sum_{i=N_S \cdot T_0+1}^t \mathbf{x}_i(a_i)^\top (\tilde{\boldsymbol{\mu}}_{a_i,i} - \boldsymbol{\mu}_{\hat{c}(a_i)}) \right| &\leq \sum_{i=N_S \cdot T_0+1}^t |\mathbf{x}_i(a_i)^\top (\tilde{\boldsymbol{\mu}}_{a_i,i} - \boldsymbol{\mu}_{\hat{c}(a_i)})| \\
&= \sum_{c \in [C]} \sum_{\substack{i: \hat{c}(a_i)=c \\ i > N_S \cdot T_0}} |\mathbf{x}_i(a_i)^\top (\tilde{\boldsymbol{\mu}}_{a_i,i} - \boldsymbol{\mu}_c)| \\
&\leq \sum_{c \in [C]} 4m\sqrt{t_c \log(t_c m/\zeta) \log(t_c)} + \epsilon_c m^{\frac{3}{2}} t_c \sqrt{\log(t_c)} \\
&\leq 4Cm\sqrt{t \log(tm/\zeta) \log(t)} + \epsilon_c m^{\frac{3}{2}} t \sqrt{\log(t)}
\end{aligned}$$

where the first step follows from the triangle inequality, and the third from Lemma A.12. Statement d) follows similarly.  $\square$

## D USEFUL LEMMAS FOR SECTION 5

Here we provide the following results. In Section D.1 we provide the lemmas for Lemma 5.5. In Section D.2 we provide the lemmas for Theorem D.8.

### D.1 USEFUL LEMMAS FOR LEMMA 5.5

**Lemma D.1** (Formal version of Lemma 5.2). *If  $\max_{\pi \in \{\pi^*, \pi'\}} |r(\pi) - \rho(\pi)| \leq \epsilon'$  then  $|r(\pi^*) - \rho(\pi')| \leq \epsilon'$ .*

1188 *Proof.* Considering  $\pi^*$  we have

$$1189 \quad r(\pi^*) - \rho(\pi^*) \leq \epsilon'$$

1190 which implies that

$$1191 \quad r(\pi^*) - \rho(\pi') \leq \epsilon'$$

1192 and by considering  $\pi'$  we similarly get

$$1193 \quad \rho(\pi') - r(\pi^*) \leq \epsilon'$$

1194 Thus, it follows that  $|r(\pi^*) - \rho(\pi')| \leq \epsilon'$ .  $\square$

1197 Therefore, it suffices to prove that the difference  $\max_{\pi \in \{\pi^*, \pi'\}} |r(\pi) - \rho(\pi)|$  is small. We first  
1198 consider  $\pi^*$  and then  $\pi'$ .

1199 **Lemma D.2** (Formal version of Lemma 5.3).  $|r(\pi^*) - \rho(\pi^*)| < o(1)$

1202 *Proof.*

$$\begin{aligned}
1203 & |r(\pi^*) - \rho(\pi^*)| \\
1204 & \leq \left| \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in [K]} \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right] - \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right] \right| \\
1205 & + \left| \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right] - \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{\hat{c}(a)}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right] \right| \\
1206 & + \left| \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{\hat{c}(a)}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right] - \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \hat{\boldsymbol{\mu}}_{\hat{c}(a), N_S T_0 + 1}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right] \right| \\
1207 & + \left| \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \hat{\boldsymbol{\mu}}_{\hat{c}(a), N_S T_0 + 1}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right] - \frac{K}{N_S} \frac{1}{N_S T_0} \sum_{t=1}^{N_S T_0} \sum_{a \in \mathcal{S}} \hat{\boldsymbol{\mu}}_{\hat{c}(a), N_S T_0 + 1}^\top \mathbf{x}(a) \pi^*(a, \mathbf{X}) \right| \\
1208 & < o(1)
\end{aligned}$$

1212 The first difference is zero, as the set  $\mathcal{S}$  is sampled uniformly at random. The second difference is  
1213  $o(1)$  because the probability of clustering an arm incorrectly is  $o(N_S^{-1})$ . The third difference is  $o(1)$   
1214 because the regularized ordinary least squares estimator is consistent when there is no clustering  
1215 error, i.e.,  $\hat{\boldsymbol{\mu}}_{c, N_S T_0 + 1} \xrightarrow{P} \boldsymbol{\mu}_c$  as  $(N_S, T_0) \rightarrow \infty$  (and as long as the set  $[C]$  is covered by  $\mathcal{S}$ ), and the  
1216 clustering error vanishes asymptotically. The fourth difference is  $o(1)$  because it is the difference  
1217 between an expectation and its empirical counterpart.  $\square$

1223 **Lemma D.3** (Formal version of Lemma 5.4).  $|r(\pi') - \rho(\pi')| < o(1)$

1225 *Proof.*

$$\begin{aligned}
1226 & |r(\pi') - \rho(\pi')| \\
1227 & \leq \left| \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in [K]} \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right] - \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right] \right| \\
1228 & + \left| \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{c(a)}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right] - \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{\hat{c}(a)}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right] \right| \\
1229 & + \left| \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \boldsymbol{\mu}_{\hat{c}(a)}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right] - \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \hat{\boldsymbol{\mu}}_{\hat{c}(a), N_S T_0 + 1}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right] \right| \\
1230 & + \left| \frac{K}{N_S} \mathbb{E}_{\mathbf{X}} \left[ \sum_{a \in \mathcal{S}} \hat{\boldsymbol{\mu}}_{\hat{c}(a), N_S T_0 + 1}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right] - \frac{K}{N_S} \frac{1}{N_S T_0} \sum_{t=1}^{N_S T_0} \sum_{a \in \mathcal{S}} \hat{\boldsymbol{\mu}}_{\hat{c}(a), N_S T_0 + 1}^\top \mathbf{x}(a) \pi'(a, \mathbf{X}) \right| \\
1231 & < o(1)
\end{aligned}$$

1232 Each difference is  $o(1)$  following arguments similar to those in the proof of Lemma D.2.  $\square$

1241 The result follows by considering the Karush-Kuhn-Tucker conditions to incorporate the constraints.

## D.2 USEFUL RESULTS FOR THEOREM D.8

**Lemma D.4** (Formal version of Lemma 5.6). *With probability at least  $(1 - \zeta)^3$  we have:*

$$\begin{aligned} a) & \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \tilde{\boldsymbol{\mu}}_{a_t,t} \right| \leq R(T) \\ b) & \left\| \sum_{t=N_S T_0+1}^{T_\omega} \mathbf{v}_t(a_t) - \mathbf{x}_t(a_t)^\top \tilde{\mathbf{W}}_{a_t,t} \right\|_\infty \leq R(T) \end{aligned}$$

*Proof.* Considering part a), we start with the probability of the event of interest, and we utilize Lemma 4.4 c) in order to make the first  $(1 - \zeta)$  term show up, and the Azuma-Hoeffding inequality for the remaining two  $(1 - \zeta)$  terms.

$$\begin{aligned} & \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \tilde{\boldsymbol{\mu}}_{a_t,t} \right| \leq R(T) \right] \\ &= \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} + \boldsymbol{\mu}_{\hat{c}(a_t)} - \boldsymbol{\mu}_{\hat{c}(a_t)}) \right| \leq R(T) \right] \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{\hat{c}(a_t)} \right| + \left| \sum_{t=N_S T_0+1}^{T_\omega} \mathbf{x}_t(a_t)^\top (\boldsymbol{\mu}_{\hat{c}(a_t)} - \tilde{\boldsymbol{\mu}}_{a_t,t}) \right| \leq R(T) \right] \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{\hat{c}(a_t)} \right| \leq \frac{R(T)}{2} \right] \\ &\quad \cdot \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} \mathbf{x}_t(a_t)^\top (\boldsymbol{\mu}_{\hat{c}(a_t)} - \tilde{\boldsymbol{\mu}}_{a_t,t}) \right| \leq \frac{R(T)}{2} \right] \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{\hat{c}(a_t)} \right| \leq \frac{R(T)}{2} \right] \cdot (1 - \zeta) \\ &= \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top (\boldsymbol{\mu}_{\hat{c}(a_t)} + \boldsymbol{\mu}_{c(a_t)} - \boldsymbol{\mu}_{c(a_t)}) \right| \leq \frac{R(T)}{2} \right] \cdot (1 - \zeta) \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{c(a_t)} \right| + \left| \sum_{t=N_S T_0+1}^{T_\omega} \mathbf{x}_t(a_t)^\top (\boldsymbol{\mu}_{c(a_t)} - \boldsymbol{\mu}_{\hat{c}(a_t)}) \right| \leq \frac{R(T)}{2} \right] \\ &\quad \cdot (1 - \zeta) \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{c(a_t)} \right| + \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{1}[\hat{c}(a_t) \neq c(a_t)] \leq \frac{R(T)}{2} \right] \cdot (1 - \zeta) \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{c(a_t)} \right| \leq \frac{R(T)}{4} \right] \\ &\quad \cdot \Pr \left[ \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{1}[\hat{c}(a_t) \neq c(a_t)] \leq \frac{R(T)}{4} \right] \cdot (1 - \zeta) \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{c(a_t)} \right| \leq \frac{R(T)}{4} \right] \cdot \left( 1 - 2 \exp \left( - \frac{(R(T)/4)^2}{2(T_\omega - N_S T_0)} \right) \right) \\ &\quad \cdot (1 - \zeta) \\ &\geq \Pr \left[ \left| \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) - \mathbf{x}_t(a_t)^\top \boldsymbol{\mu}_{c(a_t)} \right| \leq \frac{R(T)}{4} \right] \cdot (1 - \zeta)^2 \\ &\geq \left( 1 - 2 \exp \left( - \frac{(R(T)/4)^2}{2(T_\omega - N_S T_0)} \right) \right) \cdot (1 - \zeta)^2 \end{aligned}$$

$$\geq (1 - \zeta)^3$$

The proof for part b) follows the same steps.  $\square$

Now, let  $\mathcal{S}_1$  denote the subset of  $\mathcal{S}$  that contains correctly clustered arms,

$$\mathcal{S}_1 := \{a \in \mathcal{S} : \hat{c}(a) = c(a)\}$$

The following lemma provides a lower bound related to the choice of the algorithm.

**Lemma D.5** (Formal version of Lemma 5.7). *For  $t > N_S T_0$ , the following inequality holds with high probability:*

$$\mathbf{x}_t(a_t)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} - Z\tilde{\mathbf{W}}_{a_t,t}\boldsymbol{\theta}_t) \geq \frac{1}{\sum_{a' \in \mathcal{S}_1} \pi^*(a', \mathbf{X}_t)} \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\boldsymbol{\mu}_{c(a)} - Z\mathbf{W}_{c(a)}\boldsymbol{\theta}_t)$$

*Proof.* Let  $\Pi^{\mathcal{S}}$  denote the set of static policies that assign non-zero probability only to arms in  $\mathcal{S}$ , and let  $\Pi^{\mathcal{S}_1}$  be defined equivalently. Then, for  $t > N_S T_0$  we have

$$\begin{aligned} & \mathbf{x}_t(a_t)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} - Z\tilde{\mathbf{W}}_{a_t,t}\boldsymbol{\theta}_t) \\ & \geq \max_{\pi^{\mathcal{S}} \in \Pi^{\mathcal{S}}} \sum_{a \in \mathcal{S}} \pi^{\mathcal{S}}(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} - Z\tilde{\mathbf{W}}_{a_t,t}\boldsymbol{\theta}_t) \\ & \geq \max_{\pi^{\mathcal{S}_1} \in \Pi^{\mathcal{S}_1}} \sum_{a \in \mathcal{S}_1} \pi^{\mathcal{S}}(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} - Z\tilde{\mathbf{W}}_{a_t,t}\boldsymbol{\theta}_t) \\ & \geq \frac{1}{\sum_{a' \in \mathcal{S}_1} \pi^*(a', \mathbf{X}_t)} \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} - Z\tilde{\mathbf{W}}_{a_t,t}\boldsymbol{\theta}_t) \\ & \geq \frac{1}{\sum_{a' \in \mathcal{S}_1} \pi^*(a', \mathbf{X}_t)} \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\boldsymbol{\mu}_{\hat{c}(a)} - Z\mathbf{W}_{\hat{c}(a)}\boldsymbol{\theta}_t) \\ & = \frac{1}{\sum_{a' \in \mathcal{S}_1} \pi^*(a', \mathbf{X}_t)} \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\boldsymbol{\mu}_{c(a)} - Z\mathbf{W}_{c(a)}\boldsymbol{\theta}_t) \end{aligned}$$

where the first inequality follows from the choice of the algorithm, the second from restricting the set of arms to  $\mathcal{S}_1$ , the third from considering the normalization of  $\pi^*$  as a policy in  $\Pi^{\mathcal{S}_1}$ , the fourth from Lemma 4.4 a) and b), and the last equality from the definition of  $\mathcal{S}_1$ .  $\square$

**Lemma D.6** (Formal version of Lemma 5.8). *The following inequality holds with high probability:*

$$\begin{aligned} \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{E}_{\mathbf{X}_t} \left[ \mathbf{x}_t(a_t)^\top \tilde{\boldsymbol{\mu}}_{a_t,t} \right] & \geq O\left(\frac{N_S \cdot T_\omega}{K \cdot T} \text{OPT}\right) \\ & \quad + Z \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{E}_{\mathbf{X}_t} \left[ \mathbf{x}_t(a_t)^\top \tilde{\mathbf{W}}_{a_t,t} - \frac{N_S}{K} \cdot \frac{B}{T} \mathbf{1}_d \right] \boldsymbol{\theta}_t \end{aligned}$$

*Proof.* Lemma D.5 implies the weaker condition where expectation is taken over  $\mathbf{X}_t$  only and it is conditional on the past realizations of the context:

$$\begin{aligned} & \mathbb{E}_{\mathbf{X}_t} \left[ \mathbf{x}_t(a_t)^\top (\tilde{\boldsymbol{\mu}}_{a_t,t} - Z\tilde{\mathbf{W}}_{a_t,t}\boldsymbol{\theta}_t) \right] \\ & \geq \mathbb{E}_{\mathbf{X}_t} \left[ \frac{1}{\sum_{a' \in \mathcal{S}_1} \pi^*(a', \mathbf{X}_t)} \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\boldsymbol{\mu}_{c(a)} - Z\mathbf{W}_{c(a)}\boldsymbol{\theta}_t) \right] \end{aligned}$$

$$\begin{aligned}
1350 & \geq \mathbb{E}_{\mathbf{X}_t} \left[ \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top (\boldsymbol{\mu}_{c(a)} - ZW_{c(a)}\boldsymbol{\theta}_t) \right] \\
1351 & \\
1352 & = \mathbb{E}_{\mathbf{X}_t} \left[ \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top \boldsymbol{\mu}_{c(a)} \right] - Z \mathbb{E}_{\mathbf{X}_t} \left[ \sum_{a \in \mathcal{S}_1} \pi^*(a, \mathbf{X}_t) \cdot \mathbf{x}_t(a)^\top W_{c(a)} \right] \boldsymbol{\theta}_t \\
1353 & \\
1354 & = O\left(\frac{N_S}{K}\right) \cdot \left(\frac{\text{OPT}}{T} - Z \frac{B}{T} \mathbf{1}_d \boldsymbol{\theta}_t\right) \\
1355 & \\
1356 &
\end{aligned}$$

1357 where the second step holds since the probability-normalization term is greater than one. From  
1358 Lemma A.3 the size of  $\mathcal{S}_1$  is  $N_S(1 - o(N_S^{-1})) = O(N_S)$ . We get the statement of the lemma by  
1359 summing from period  $N_S T_0 + 1$  to period  $T_\omega$ .  $\square$

1360

1361 Since in the first  $N_S T_0$  periods the choices are made randomly, and so  $N_S T_0$  of the budget can  
1362 potentially be consumed, the following lemma which is known from the literature is expressed in  
1363 terms of  $B' = B - N_S T_0$  and  $T' = T - N_S T_0$  rather than  $B$  and  $T$ .

1364 **Lemma D.7.**  $\frac{N_S \cdot B}{K \cdot T} < 2 \frac{B'}{T'}$

1365

1366 *Proof.*  $\frac{N_S B T'}{K T B'} \leq \frac{N_S B}{K B'} \leq \frac{B}{B'} = \frac{1}{1 - \frac{N_S T_0}{B}} < 2$ , since  $B/2 > N_S T_0$ .  $\square$

1367

1368 **Theorem D.8** (Main Result, formal version of 5.10). *For  $\delta \in (0, \frac{1}{2})$ , and  $B > N_S T_0$ , with high  
1369 probability*

1370

1371

$$1372 \text{regret}(T) \leq O\left(R(T)\left(1 + \frac{N_S \text{OPT}}{K B'}\right) + \text{OPT}\left(1 - \frac{N_S}{K}\right)\right)$$

1373

1374 where

1375

$$1376 R(T) = O\left(C p_{\min}^{-1} m^{\frac{3}{2}} T^{1-\delta} \sqrt{\log(T)}\right).$$

1377

1378 *Proof.* Let  $T_\omega \leq T$  be the stopping time of the algorithm. Starting from the definition of regret we  
1379 get:

1380

1381

$$\begin{aligned}
1382 \text{regret}(T) &= \text{OPT} - \sum_{t=1}^T r_t(a_t) \\
1383 & \\
1384 &= \text{OPT} - \sum_{t=1}^{T_\omega} r_t(a_t) \\
1385 & \\
1386 &= \text{OPT} - \sum_{t=1}^{N_S T_0} r_t(a_t) - \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) \\
1387 & \\
1388 &\leq \text{OPT} - \sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t) \\
1389 & \\
1390 & \\
1391 & \\
1392 & \\
1393 &
\end{aligned}$$

1394 where the inequality follows since  $r_t(a_t) \in [0, 1]$ . Now, let

1395

$$1396 R(T) := O\left(C m \sqrt{T \log(d T m / \zeta)} \log(T) + C \epsilon_c m^{\frac{3}{2}} T \sqrt{\log(T)}\right)$$

1397

1398 The proof now proceeds with first stating Lemma D.4 that allows us to work with the optimistic  
1399 estimates instead of the actual realizations of the reward and the consumption.  $\square$

1400

1401 Thus, we have

1402

$$1403 \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{E}_{\mathbf{X}_t} \left[ \mathbf{x}_t(a_t)^\top \tilde{\boldsymbol{\mu}}_{a_t, t} \right]$$

$$\begin{aligned}
&\geq O\left(\frac{N_S T_\omega}{KT} \text{OPT} + Z \sum_{t=N_S T_0+1}^{T_\omega} \mathbb{E}_{\mathbf{X}_t} \left[ \mathbf{x}_t(a_t)^\top \widetilde{\mathbf{W}}_{a_t,t} - 2\frac{B'}{T'} \mathbf{1}_d \right] \boldsymbol{\theta}_t\right) \quad (19) \\
&\geq O\left(\frac{N_S T_\omega}{KT} \text{OPT} + Z \left(2B' \left(1 - \frac{T_\omega - N_S T_0}{T'}\right) - R(T)\right)\right) \\
&\geq O\left(\frac{N_S T_\omega}{KT} \text{OPT} + \frac{N_S \text{OPT}}{2KB'} \left(2B' \left(1 - \frac{T_\omega - N_S T_0}{T'}\right) - R(T)\right)\right) \\
&\geq O\left(\frac{N_S}{K} \text{OPT} \left(\frac{T_\omega}{T} + 1 - \frac{T_\omega - N_S T_0}{T'} - \frac{R(T)}{2B'}\right)\right) \\
&\geq O\left(\frac{N_S}{K} \text{OPT} \left(1 - \frac{R(T)}{B'}\right)\right) \quad (20)
\end{aligned}$$

where the first step follows from Lemmas D.6 and D.7, the second from Lemma 5.9, and the third from Lemma 5.5.

From Lemma D.4, the bound in Eq. (20) applies to  $\sum_{t=N_S T_0+1}^{T_\omega} \mathbb{E}_{\mathbf{X}_t} [r_t(a_t)]$  by adding  $R(T)$ , and an application of the Azuma-Hoeffding inequality to the realized reward  $\sum_{t=N_S T_0+1}^{T_\omega} r_t(a_t)$  then gives:

$$\begin{aligned}
\text{regret}(T) &\leq O\left(\text{OPT} - \frac{N_S}{K} \text{OPT} \left(1 - \frac{R(T)}{B'}\right) + R(T)\right) \\
&= O\left(\text{OPT} \left(1 - \frac{N_S}{K}\right) + R(T) \left(1 + \frac{N_S \text{OPT}}{KB'}\right)\right).
\end{aligned}$$

## LLM USAGE DISCLOSURE

LLMs were used only to polish language, such as grammar and wording. These models did not contribute to idea creation or writing, and the authors take full responsibility for this paper's content.