
Interactive Recommendations for Optimal Allocations in Markets with Constraints

Yigit Efe Erginbas¹*, Soham Phade²*, Kannan Ramchandran¹

¹ Department of Electrical Engineering and Computer Science
University of California, Berkeley

² Salesforce Research
{erginbas, soham_phade, kannanr}@berkeley.edu

Abstract

Recommendation systems when employed in markets play a dual role: they assist users in selecting their most desired items from a large pool and they help in allocating a limited number of items to the users who desire them the most. Despite the prevalence of capacity constraints on allocations in many real-world recommendation settings, a principled way of incorporating them in the design of these systems has been lacking. Motivated by this, we propose an interactive framework where the system provider can enhance the quality of recommendations to the users by opportunistically exploring allocations that maximize user rewards and respect the capacity constraints using appropriate pricing mechanisms. We model the problem as an instance of a low-rank combinatorial multi-armed bandit problem with selection constraints on the arms. We employ an integrated approach using techniques from collaborative filtering, combinatorial bandits, and optimal resource allocation to provide an algorithm that provably achieves sub-linear regret, namely $\tilde{O}(\sqrt{NM(N+M)RT})$ in T rounds for a problem with N users, M items and rank R mean reward matrix. Empirical studies on synthetic and real-world data also demonstrate the effectiveness and performance of our approach.

1 Introduction

Online recommendation systems have become an integral part of our socioeconomic life with the rapid increases in online services that help users discover options matching their preferences. Despite providing efficient ways to discover information about the preferences of users, they have played a largely complementary role to searching and browsing with little consideration of the accompanying *markets* within which recommended items are allocated to the users. Indeed, in many real-world scenarios, recommendations bring about the *allocation* of the corresponding items in a market that has possibly intrinsic constraints. In particular, recommendations of candidate items that have associated notions of limited *capacities* naturally give rise to a market setting where users compete for the allocation of the recommended items.

Allocation constraints are common in recommendation contexts. A few interesting examples include: (1) Point-of-Interest (PoI) recommendation systems (e.g., restaurants, theme parks, hotels), where the PoI can only accommodate limited number of visitors, (2) book recommendation systems employed by libraries, where the books recommended to the borrowers have limited copies, (3) route recommendation systems which aim to suggest the optimal road for travelling while avoiding traffic congestion, (4) course recommendation systems for universities, where each recommended course has limited number of seats. As similar systems become more ubiquitous and impactful in the broader aspects of daily life, there is a huge application drive and potential for delivering recommendations that respect the requirements of the market. Therefore, it is crucial to consider capacity-aware recommendation systems to maximize the user experience.

*equal contribution

Main Challenges: We model the user preferences as rewards that users obtain by consuming different items, while the social welfare is the aggregate reward over the entire system comprising multiple users with heterogeneous preferences, and a provider who continually recommends items to the users and receives interactive reward feedback from them. The provider aims to maximize the social welfare while respecting the *time-varying* allocation constraints: indeed we consider system *dynamics* in terms of user demands and item capacities to be an important aspect of our problem. In the process of identifying the best match between users and target items, the provider encounters two challenges: The first challenge relates to the element of *recommendation* as the provider needs to make recommendations without exact knowledge of the user preferences ahead of time, and hence has to continue exploring user preferences while continually making recommendations. The second challenge relates to the *allocation* aspect of the problem induced by the market constraints. Note that even if matching the users with their most preferred items would result in high rewards, such an allocation may not respect the constraints of the market. For example, in a restaurant recommendation setting, if there is a hugely popular restaurant that most people love, a naive recommender would send many users to the same restaurant, causing overcrowding and considerable user dissatisfaction.

The key to overcoming the (first) challenge of making accurate recommendations is to learn the user preferences from the reward feedback. Since the preferences of different users for different items are highly correlated, it is natural to employ collaborative filtering techniques that have been widely applied in recommender systems [1, 2, 3, 4]. In order to learn the user preferences efficiently, previous works have established interactive collaborative filtering systems that query the users with well-chosen recommendations [5, 6]. Typically, these works consider a setting where a single user arrives to the system at each round and the system makes a recommendation that will match the user’s preferences. However, this assumption no longer holds in applications having an associated market structure, as recommendations made to different users in the same time period must also respect the constraints of the market.

The common strategy to tackling the (second) allocation challenge is through pricing mechanisms that ensure social optimality. Such mechanisms have been studied in economics for two-sided (supply and demand) markets and are called Walrasian auctions [7]. In the networking literature, Kelly has also used similar mechanisms to do optimal bandwidth allocation over a network [8]. In pricing-based mechanisms, the users choose the items based on their preferences as well as the posted prices. The provider meanwhile successively adjusts these prices in response to the user’s demand for the items, so that capacity constraints are satisfied in equilibrium. The equilibrium prices ensure that the limited number of items are allocated to users that are expected to obtain the largest reward. However, these mechanisms still require the users to know and evaluate their preferences for *all* possible items and respond constantly with their updated bids/demands for each item. This is definitely not a scalable solution for the large-scale system (comprising large numbers of users and items) that we target. Furthermore, this framework assumes that users already know their preferences for all items, which is clearly not true in our setting, where users report their preferences through feedback *after* being targeted with their recommended items. For this reason, the provider must *learn* the user preferences in its quest to perform optimal capacity-constrained allocations.

Hence, as depicted in Figure 1, the goal of the provider is twofold: (1) to learn the user preferences and make recommendations that will guide the users to choose the items that they are likely to obtain high rewards, (2) to achieve allocations that will satisfy the capacity constraints. To achieve these goals, we envision developing the following market-aware recommendation mechanism for the provider. By recommending items, the provider helps the users to narrow down their options so that users can comprehend and evaluate their preferences among a smaller number of offered items. In addition, being aware of the market structure, the provider carefully determines the item prices that play the role of an intermediary for satisfying the constraints of the market. We believe that this is an important and practically-relevant question to be resolved because it allows for the analysis of many interesting real-world interactive recommendation settings with market constraints. In its full generality, this framework requires us to model the user decisions in a way that will capture the effects of the recommendations and prices that they are presented. In order to avoid the complications intro-

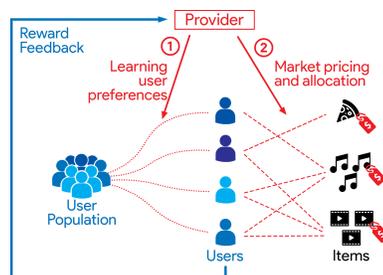


Figure 1: The provider interactively learns the user preferences to achieve socially optimal capacity-constrained allocations.

duced by this modelling challenge and to obtain a profound understanding of fundamental aspects of the problem, we begin with focusing our attention on a central question whose solution will be key to making progress towards our longer-term goal of developing a complete framework.

Specifically, we focus our study on these essential aspects of the problem: recommending and allocating the items while interactively learning the user preferences, which to the best of our knowledge has not been addressed in the literature. In essence, we analyze a special case of the mechanism introduced above, by assuming that the provider makes recommendations such that the number of presented choices matches with the number of items the user is willing to consume, so that the users obtain all of the recommended items regardless of their prices. Then, the provider’s task reduces to deciding on high-reward allocations while satisfying the constraints by allocating each item to at most certain number of users.

Structured Combinatorial Multi-Armed Bandits: The provider seeks to choose high-reward allocations subject to the constraints, while actively learning the user preferences by making queries that will give rise to the most informative responses. Therefore, it encounters the well-known *exploration-exploitation* dilemma. In essence, there exists a trade-off between two competing goals: maximizing social welfare using the historical feedback data, and gathering new information to improve the performance in the future. In the literature of interactive collaborative filtering, this dilemma is typically formulated as a multi-armed bandit problem where each arm corresponds to allocation of an item to a user [6, 9, 10]. When an item is allocated to a user, a random reward is obtained and the reward information is fed back to the provider to improve its future allocation strategies. However, in contrast to prior works, our setting further requires that a collection of actions taken for different users satisfy the constraints of the market.

We formulate our problem as a bandit problem with arms having correlated means, and call it Structured Combinatorial Bandit. Based on the standard OFU (Optimism in Face of Uncertainty) principle for linear bandits [11, 12], we devise a procedure that learns the mean reward values opportunistically so as to solve the system problem of optimal allocation with minimum regret. The estimation method benefits from both the combinatorial nature of the feedback and the dependencies induced by the low-rank structure of collaborative filtering setting. Moreover, using matrix factorization techniques, the algorithm is efficient even at scale in settings with a large number of users and items. As is standard with OFU-based methods, our algorithm maintains a confidence set of the mean rewards for all user-item pairs. If it has less data about some user-item allocation pair, the confidence set becomes wider in the corresponding direction. Then, due to optimism, the algorithm becomes more inclined to attempt the corresponding allocation pairs to explore and collect more information.

Our contributions:

- We formulate the problem of making recommendations that will facilitate socially optimal allocation of items with constraints. Our formulation further allows for the analysis of problem settings with dynamic (i.e., time-varying) item capacities and user demands.
- We pose the Structured Combinatorial Bandit problem under generic structural assumptions (not only low-rank) and propose an algorithm that achieves sublinear regret bounds in terms of parameters that depend on the problem-specific structure of the arms.
- For the recommendation setting, we specialize our results to low-rank structures and obtain a Low-Rank Combinatorial Bandit (LR-COMB) algorithm that achieves $\tilde{O}(\sqrt{NM(N+M)RT})$ regret in T rounds for a problem with N users, M items and rank R mean reward matrix.

Experiments: We run experiments both on synthetic and real-world datasets to show the efficacy of the proposed algorithms. Results show that proposed algorithm can obtain significant improvements over naive approaches in solving the problem of recommendation and allocation with constraints.

Related work:

- **Combinatorial Multi-Armed Bandits (CMAB) and Semi-Bandits:** The frameworks of CMAB [13, 14] and semi-bandits [15] model multi-armed bandit problems where the player chooses a subset of arms in each round and observes individual outcomes of the played arms. However, they do not incorporate any structural assumptions about the rewards obtained from the arms. However, in a collaborative filtering setting like ours, the main promise is to leverage the intrinsic structure between different user-item pairs. To close this gap, we pose the problem of Structured Combinatorial Bandits and devise an algorithm that makes use of the structure of the arms as well. Additionally, CMAB framework assumes availability of an oracle that takes the means rewards for the arms and outputs the optimum subset of arms subject to the selection constraints. Due to the

combinatorial nature of the problem, this oracle may not be readily available in general CMAB settings. In our case, due to the special structure of the capacity constraints, we can efficiently solve for the optimum allocations given the mean rewards of the allocation pairs.

- **Structured Linear Bandits:** Our formulation also shows parallelism with the frameworks of structured linear bandits [16, 17] and low-rank linear bandits [18]. However, it is distinct from them by having the additional ability to capture the combinatorial nature of the problem. In linear bandits, the player only observes the final total reward, but no outcome of any individual arm. Our setup differs from their case because the player (provider) is able to observe individual outcomes of all played arms. Due to this richer nature of the observation model, we can achieve lower regret guarantees than what is available in the literature of structured linear bandits.
- **Recommendation with Capacity Constraints:** There have been a few works using the notion of constrained resources to model and solve the problem of recommendation with capacity constraints [19, 20]. However, these works only consider optimizing the recommendation accuracy subject to item usage constraints without any consideration of the interactive mechanisms that discover user preferences through recommendations.
- **Competing Bandits in Matching Markets:** One other related line of literature studies the stable matching problem in two-sided markets [21]. The model assumes that entities on each side of the market has preference orderings for the other side of the market and the allocations are driven by these preference orderings rather than the prices. In contrast to our work, these mechanisms necessitate at least the entities on one side of the market know their preferences over all the entities on the other side of the market. However, in many real-world settings of optimum recommendation and allocation, like the examples given above, the explicit preferences are not known ahead of time and can only be discovered through interactions. Furthermore, the matching markets only model one-to-one matches, meaning that they do not allow for the items to be allocated for multiple users.

2 Problem setting

We use bold font for vectors \mathbf{x} and matrices \mathbf{X} , and calligraphic font \mathcal{X} for sets. We denote by $[K]$ the set $\{1, 2, \dots, K\}$. For a vector \mathbf{x} , we denote its i -th entry by x_i and for a matrix \mathbf{X} , we denote its (i, j) -th entry by x_{ij} . We denote the Frobenius inner product of two matrices by $\langle \mathbf{A}, \mathbf{B} \rangle = \text{trace}(\mathbf{A}^T \mathbf{B})$, and the Frobenius norm of a matrix \mathbf{A} by $\|\mathbf{A}\|_F$.

Suppose the system has N users and M items in record. The items are *allocated* to the users in multiple *rounds* (or *periods*) denoted by $t \in \mathbb{N}$. Allocation of an item $i \in [M]$ to a user $u \in [N]$ results in a random *reward* that has a distribution unknown to the system provider. The expected reward obtained from allocating item i to user u is denoted by θ_{ui}^* and these values are collected into the mean reward matrix $\Theta^* \in \mathbb{R}^{N \times M}$.

We assume that each item has (time-varying) capacity that corresponds to the maximum number of different users it can be allocated to. We denote the capacity of item $i \in [M]$ by $c_{t,i}$, and collect these values into vectors $\mathbf{c}_t \in \mathbb{R}^M$. Similarly, each user has a (time-varying) demand that corresponds to the maximum number of different items it can get allocated. We denote the demand of user $u \in [N]$ by $d_{t,u}$, and collect these values into vectors $\mathbf{d}_t \in \mathbb{R}^N$. Therefore, each item can only be allocated to at most $c_{t,i}$ different users, while each user can only get allocated at most $d_{t,u}$ different types of items in the period t . We shall call these the *allocation constraints*. One can consider the special case where $d_{t,u}$ parameters only take values from $\{0, 1\}$ so that each *active* user gets at most one allocation while the *inactive* users do not get any allocations.

Let \mathbf{X}_t denote the *allocation matrix* for round t where the (u, i) -th entry is one if user u is allocated item i at round t , and zero otherwise. Due to the allocation constraints, any valid \mathbf{X}_t must belong to the set of valid allocation matrices $\mathcal{X}_t \subseteq \{0, 1\}^{N \times M}$ defined as:

$$\mathcal{X}_t = \{\mathbf{X} \in \{0, 1\}^{N \times M} : \mathbf{X}\mathbf{1}_M \leq \mathbf{d}_t \text{ and } \mathbf{X}^T\mathbf{1}_N \leq \mathbf{c}_t\}$$

where the inequalities are entry-wise and $\mathbf{1}_p$ denotes the all-ones vector of size p .

2.1 Optimal allocations

Given the knowledge about the mean reward matrix Θ^* , the optimal allocation \mathbf{X}_t^* at time t can be obtained by solving the integer program:

$$\mathbf{X}_t^* \in \arg \max_{\mathbf{X} \in \mathcal{X}_t} \langle \mathbf{X}, \Theta^* \rangle \quad (2.1)$$

This integer program can be relaxed to a linear program by dropping the integral constraints (setting $0 \leq x_{ui} \leq 1$). In Appendix B, we show that the integrality gap of this problem is zero.² Hence, any integer solution found for the relaxed problem is also a solution for the allocation problem.

When the provider does not have direct knowledge of the mean rewards associated with user-item allocation pairs, one standard approach is to employ pricing mechanisms [7, 8]. The idea is to apply dual decomposition on the (partial) Lagrangian function $L(\mathbf{X}, \boldsymbol{\lambda}) = \langle \mathbf{X}, \boldsymbol{\Theta}^* \rangle + \boldsymbol{\lambda}^\top (\mathbf{c}_t - \mathbf{X}^\top \mathbf{1}_N)$ where $\boldsymbol{\lambda} \geq 0$ are the Lagrange multipliers (item prices) associated with the capacity constraints. Then, the allocation problem is decomposed into one problem for each user and one problem for the provider where the item prices mediate between the subsidiary problems. Each user calculates its demand by maximizing the corresponding component of the Lagrangian for a given set of prices. On the other side, the provider iteratively updates the prices based on users demands to achieve the optimal pricing. At the end of many consecutive updates from users and the provider, the equilibrium ensures that the limited items are allocated to users that are expected to obtain the largest reward. (See [22] for further details.)

However, as discussed in the introduction, this pricing mechanism has limitations in many real-world applications. Most importantly, it requires the user to solve a problem that involves the valuations even for the items that the user has no prior experience with. However, in many real-world scenarios, it is infeasible to request the users to choose among all the items in the system. Secondly, in the process of price discovery, the mechanism asks the users to repetitively respond to the prices by recomputing their demands. However, since it might take many iterations until convergence to the optimal pricing, asking the users to respond many times would be a burden for them. Furthermore, the final prices found by this iterative mechanism are only guaranteed to be optimal for the problem defined by the capacity \mathbf{c}_t and demand \mathbf{d}_t parameters at round t . If the capacities and demands vary with time, the optimal pricing and allocation for the next allocation round $t + 1$ will be different and will be needed to be rediscovered.

2.2 Learning the optimal allocations

To address the issues discussed in the previous section, we need mechanisms that can find the optimum allocations using fewer and simpler interactions. One resolution is to recommend a subset of items along with prices intelligently chosen by the provider. This way, the users will be able to easily evaluate their preference on the small number of recommended items and decide on their demand without requiring to consider all items in the system. The provider will decide on well-chosen offerings with correct prices so that it can satisfy the capacity constraints. However, as the provider does not have the complete knowledge of the user preferences, it needs to learn the unknown preference parameters $\boldsymbol{\Theta}^*$ from the user feedback so that it can determine better recommendations as well as the correct prices. Based on the examples of applications provided in the introduction, we believe that design of such system dynamics is a practically-relevant question to be resolved.

As a first step in this direction, we decide to restrict our attention to a setting that itself has interesting interactions between learning the user preferences and allocating the items. In order to facilitate our analysis, we consider that the number of choices presented to each user u at round t is limited exactly by their demand $d_{t,u}$ and users are allocated with all of the items that they are recommended. Therefore, the problem essentially reduces to an allocation problem in which users get allocated a set of items directly by the provider instead of users choosing between the offerings. Then, after each round of allocation, users provide feedback about the items that they have been allocated so that the provider can enhance its performance in the following rounds. Hence, whilst the users get allocated sequentially, the predictions are constantly refined using the reward feedback.

The provider determines the allocations according to an *optimistic* estimate of the true mean reward matrix $\boldsymbol{\Theta}^*$. It solves the allocation problem (2.1) assuming that the estimated parameter is the underlying reward parameter and obtains an estimate for the optimum allocation at each round t . Even though these allocations can be suboptimal due to estimation errors, our analysis shows that the cumulative regret obtained from these sequential allocations can only grow sublinearly with the time horizon. Using this approach, we are coupling the general principle of optimism in the face of uncertainty (OFU) along with capacity aware resource allocation. In the experiments section, we show the importance of this connection by comparing our strategy with algorithms that only focus on one aspect of the problem: a non-OFU algorithm that only aims for achieving momentary performance and an OFU-based algorithm that is unaware of the capacities.

²The integrality gap is the difference between optimal values of the integer program and its linear relaxation.

Remark 2.1. When the allocation problem (2.1) is solved with the estimated parameters, the Lagrange multipliers for the capacity constraints give estimates for the optimum prices of the items. As long as the user preferences are estimated well enough, these prices emerging from provider's problem are such that users who are aware of their preference for all items would still choose the recommended items. Hence, when the user preferences are learned, the mechanism is able to achieve high-reward allocations that complies with the user incentives under the optimal pricing.

2.3 Problem formulation

In this section, we formulate the provider's problem and its objective. At each time period t , the provider chooses multiple user-item allocation pairs collected into a set $\mathcal{A}_t \subseteq [N] \times [M]$. Then, the provider observes a random reward $R_{t,u,i}$ if user u is allocated with item i at round t . The total reward is the sum of rewards obtained from the system at all rounds during a time horizon T . The task is to repeatedly allocate the items to the users in multiple rounds so that the total expected reward of the system is as close to the reward of the optimal allocation as possible.

Letting $\mathbf{E}_{u,i} \in \mathbb{R}^{N \times M}$ denote the zero-one matrix with a single one at the (u, i) entry, we can write the indicator matrix for the allocation at time t as $\mathbf{X}_t = \sum_{(u,i) \in \mathcal{A}_t} \mathbf{E}_{u,i}$. Consequently, \mathbf{X}_t becomes a zero-one matrix with ones at entries \mathcal{A}_t and zeros everywhere else. Note that there is a one-to-one relation between the matrix \mathbf{X}_t and the set \mathcal{A}_t .

We denote by H_t the history $\{\mathbf{X}_\tau, (R_{\tau,u,i})_{(u,i) \in \mathcal{A}_\tau}\}_{\tau=1}^{t-1}$ of observations available to the provider when choosing the next allocation \mathbf{X}_t . The allocator employs a policy $\pi = \{\pi_t | t \in \mathbb{N}\}$, which is a sequence of functions, each mapping the history H_t to an action \mathbf{X}_t . Then, the T period cumulative regret of a policy π is the random variable

$$\mathcal{R}(T, \pi) = \sum_{t=1}^T [\langle \mathbf{X}_t^*, \Theta^* \rangle - \langle \mathbf{X}_t, \Theta^* \rangle]$$

where $\mathbf{X}_t^* \in \arg \max_{\mathbf{X} \in \mathcal{X}_t} \langle \mathbf{X}, \Theta^* \rangle$ denotes optimum allocation at time t .

3 Methodology

In order to facilitate our analysis, we start by making the following assumptions that are standard in the multi-armed bandits literature.

Assumption 1. For all $u \in [N]$, $i \in [M]$ and $t \in \mathbb{N}$, the rewards $R_{t,u,i}$ are independent and η -sub-Gaussian with mean $\theta_{u,i}^* \in [0, B]$.

To model the dependency between the mean rewards obtained from different user-item pairs, we employ the following assumption. We first present our algorithm and theoretical results under the general setting given by this assumption, and specialize for the setting of the collaborative filtering in following sections.

Assumption 2. The mean reward matrix Θ^* belongs to a known structure set $\mathcal{L} \subseteq \mathbb{R}^{N \times M}$.

In order to make use of initial historical data possibly available to the provider, we assume that the algorithm has access to an initial rough estimate $\bar{\Theta}$ that satisfies $\|\bar{\Theta} - \Theta^*\|_F \leq G$. Such an estimate can be constructed using an off-the-shelf low-rank matrix completion algorithm on the initialization data. If such observations are not readily available at the time of initialization, they can be obtained by randomly sampling some of the user-item allocation pairs once. It is worth to note that one can also set $\bar{\Theta} = \mathbf{0}$ and let G be some number satisfying $\|\Theta^*\|_F \leq G$.

Algorithm 1 Structured Combinatorial Multi-Armed Bandit

Require: horizon T , initial estimate $\bar{\Theta} \in \mathbb{R}^{N \times M}$ with $\|\bar{\Theta} - \Theta^*\|_F \leq G$.

for $t = 1, 2, \dots, T$ **do**

Find the regularized least squares estimate $\hat{\Theta}_t = \arg \min_{\Theta \in \mathcal{L}} \{L_{2,t}(\Theta) + \gamma \|\Theta - \bar{\Theta}\|_2^2\}$

Construct the confidence set $\mathcal{C}_t = \{\Theta \in \mathcal{L} : \|\Theta - \hat{\Theta}_t\|_{2, E_t} \leq \sqrt{\beta_t^*}(\delta, \alpha, \gamma)\}$

Compute the action vector $\mathbf{X}_t = \arg \max_{\mathbf{X} \in \mathcal{X}_t} \max_{\Theta \in \mathcal{C}_t} \langle \mathbf{X}, \Theta \rangle$

Play the arms \mathcal{A}_t according to \mathbf{X}_t

Observe $R_{t,u,i}$ for all $(u, i) \in \mathcal{A}_t$

end for

Our method summarized in Algorithm 1 follows the standard OFU (Optimism in Face of Uncertainty) principle [12]. It maintains a confidence set \mathcal{C}_t which contains the true parameter Θ^* with high probability and chooses the allocation \mathbf{X}_t according to

$$\mathbf{X}_t = \arg \max_{\mathbf{X} \in \mathcal{X}_t} \left\{ \max_{\Theta \in \mathcal{C}_t} \langle \mathbf{X}, \Theta \rangle \right\} \quad (3.1)$$

Typically, the faster the confidence set \mathcal{C}_t shrinks, the lower regret we have. However, the main difficulty is to construct a series of \mathcal{C}_t that leverage the combinatorial observation model as well as the structure of the parameter so that we have low regret bounds. In this work, we consider constructing confidence sets that are centered around the regularized least square estimates. We let the cumulative squared prediction error at time t be

$$L_{2,t}(\Theta) = \sum_{\tau=1}^{t-1} \sum_{(u,i) \in \mathcal{A}_\tau} (\theta_{ui} - R_{\tau,u,i})^2,$$

and define the regularized least squares estimate at time t as

$$\widehat{\Theta}_t = \arg \min_{\Theta \in \mathcal{L}} \{L_{2,t}(\Theta) + \gamma \|\Theta - \overline{\Theta}\|_2^2\}. \quad (3.2)$$

Then, the confidence sets take the form $\mathcal{C}_t := \{\Theta \in \mathcal{L} : \|\Theta - \widehat{\Theta}_t\|_{2,E_t} \leq \sqrt{\beta_t}\}$ where β_t is an appropriately chosen confidence parameter, and the regularized empirical 2-norm $\|\cdot\|_{2,E_t}$ is

$$\|\Delta\|_{2,E_t}^2 := \sum_{u=1}^N \sum_{i=1}^M (n_{t,u,i} + \gamma) (\Delta_{ui})^2,$$

where $n_{t,u,i} := \sum_{\tau=1}^{t-1} \mathbb{1}\{(u,i) \in \mathcal{A}_\tau\}$ is the number of times item i has been allocated to user u before time t (excluding time t). Hence, the empirical 2-norm is a measure of discrepancy that weighs the entries depending on how much they have been explored. Roughly speaking, since the confidence ellipsoid constructed using the 2-norm is wider in directions that are not yet well-explored, the OFU step described in 3.1 is more inclined to make allocations that include the corresponding user-item pairs. In order to obtain low-regret guarantees for the allocations, the first step is to choose correct β_t parameter such that \mathcal{C}_t will contain the true parameter Θ^* for all t with high probability. In order to take advantage of the structure of the arms, we let $\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_F)$ denote the α -covering number of \mathcal{F} in the Frobenious-norm $\|\cdot\|_F$, and let

$$\beta_t^*(\delta, \alpha, \gamma) := 8\eta^2 \log(\mathcal{N}(\mathcal{L}, \alpha, \|\cdot\|_F)/\delta) + 2\alpha t N M \left[8B + \sqrt{8\eta^2 \log(4NMt^2/\delta)} \right] + 4\gamma G^2.$$

Then, the following Lemma establishes that if we set $\beta_t = \beta_t^*(\delta, \alpha, \gamma)$, the resulting confidence sets have the desired properties.

Lemma 3.1. *For any $\delta > 0$, $\alpha > 0$, $\gamma > 0$, let $\widehat{\Theta}_t$ be the regularized least squares estimate given in 3.2. If the confidence sets are given as*

$$\mathcal{C}_t := \{\Theta \in \mathcal{L} : \|\Theta - \widehat{\Theta}_t\|_{2,E_t} \leq \sqrt{\beta_t^*(\delta, \alpha, \gamma)}\}, \quad (3.3)$$

then with probability at least $1 - 2\delta$, $\mathcal{C}_t \ni \Theta^$, for all $t \in \mathbb{N}$.*

Finally, we show that if the structured combinatorial bandits algorithm follows the OFU allocations given in (3.1) while constructing the confidence sets according to (3.3), it obtains the following overall regret guarantee:

Theorem 3.2. *Under Assumptions 1 and 2, for any $\delta > 0$, $\alpha > 0$, $\gamma \geq 1$, with probability $1 - 2\delta$, the cumulative regret of Algorithm 1 is bounded by*

$$\mathcal{R}(T, \pi) \leq \sqrt{8NM\beta_T^*(\delta, \alpha, \gamma)T \log(1 + T/\gamma)}.$$

3.1 Low-Rank COMbinatorial Bandits (LR-COMB)

As common in collaborative filtering settings, the correlation between users and arms can be captured through a matrix factorization model that leads to a low-rank mean reward matrix. Each user

u (item i) is associated with a feature vector $\mathbf{p}_u(\mathbf{q}_i)$ in a shared R -dimensional space (typically $R \ll M, N$), and the mean reward of each user-item allocation pair is given by $\theta_{ui}^* = \mathbf{p}_u^T \mathbf{q}_i$. Consequently, the mean reward matrix satisfies the factorization $\Theta^* = \mathbf{P}\mathbf{Q}^T$ for some $\mathbf{P} \in \mathbb{R}^{N \times R}$ and $\mathbf{Q} \in \mathbb{R}^{M \times R}$. Based on this observation and the boundedness condition given in Assumption 1, we can choose the structure set \mathcal{L} as

$$\mathcal{L} = \{\Theta \in \mathbb{R}^{N \times M} : \text{rank}(\Theta) \leq R, \theta_{ui} \in [0, B], \forall u, i\}. \quad (3.4)$$

Then, Lemma F.1 in the appendix shows that the covering number for \mathcal{L} given in equation (3.4) is upper bounded by $\log \mathcal{N}(\mathcal{L}, \alpha, \|\cdot\|_F) \leq (N + M + 1)R \log(9B\sqrt{NM}/\alpha)$. Therefore, the regret guarantee for a setting with low-rank mean reward matrix becomes:

Theorem 3.3 (Regret of LR-COMB). *Under Assumption 1 and Assumption 2 with \mathcal{L} given in (3.4), the Algorithm 1 achieves cumulative regret*

$$\mathcal{R}(T, \pi) = \tilde{O}\left(\sqrt{NM(N+M)RT}\right), \quad (3.5)$$

where \tilde{O} is the big-O notation, ignoring the poly-logarithmic factors of N, M, T, R .

In comparison, if we were to ignore the low-rank structure between the mean rewards obtained from user-item allocation pairs and apply the standard combinatorial bandit algorithms (e.g., CUCB [13]), we would suffer $\tilde{O}(NM\sqrt{T})$ regret [14]. Since $R \ll M, N$ in many applications of collaborative filtering, our algorithm significantly outperforms this naive approach. As common in the literature of combinatorial bandits, one possible approach to improve upon our theoretical analysis might be by assuming a problem setting where at most K of the arms can be played in each round. However, our current analysis techniques do not allow us to incorporate and leverage such an assumption together with the low-rank structure of collaborative filtering.

Implementation via Matrix Factorization: In order to efficiently solve optimization problems (3.1) and (3.2) in large scales, we take advantage of the matrix factorization model. As a result, we factorize $\Theta = \mathbf{P}\mathbf{Q}^T$ where $\mathbf{P} \in \mathbb{R}^{N \times R}$ and $\mathbf{Q} \in \mathbb{R}^{M \times R}$, and solve the problems by optimizing over \mathbf{P} and \mathbf{Q} rather than directly optimizing over Θ . Even if the problem (3.2) is not convex in the joint variable (\mathbf{P}, \mathbf{Q}) , it is convex in \mathbf{P} for fixed \mathbf{Q} and it is convex in \mathbf{Q} for fixed \mathbf{P} . Therefore, an alternating minimization algorithm becomes a feasible choice to find a reasonable solution for the least squares problem. Similarly, an alternating minimization approach is also useful to solve the problem (3.1). We can fix an allocation \mathbf{X} and minimize over \mathbf{P} and \mathbf{Q} . Then, for fixed \mathbf{P} and \mathbf{Q} , the allocation \mathbf{X} is determined through the dual decomposition mechanism described in the section 2.1. We call the resulting algorithm LR-COMB with Matrix Factorization and present it as Algorithm 2 in the Appendix.

4 Experiments

In this section, we demonstrate the efficiency of our proposed algorithm by conducting an experimental study over both synthetic and real-world datasets. The goal of our experimental evaluation is twofold: (i) evaluate our algorithm for making online recommendations and allocations in various market settings and (ii) understand the qualitative performance and intuition of our algorithm.

Baseline algorithms: We demonstrate the performance of our method by comparing it with baseline algorithms. To the best of our knowledge, there are no current approaches specifically designed to make interactive recommendations and allocations considering the capacity constraints. Therefore, based on currently available algorithms, we construct our baseline with methods that are designed for similar goals:

1. **ACF:** (Allocations with Collaborative Filtering) It solves for the least squares problem (3.2) to estimate the mean rewards obtained from user-item allocation pairs. Then, makes the best allocation with respect to the estimated parameters at each round.
2. **CUCB:** It runs the Combinatorial-UCB algorithm [13] to decide on allocations without assuming any low-rank structure between the users and items. It views the user-item allocation pairs as arms that have no correlation in between. In every round, it pulls some subset of the arms according to the capacity and demand constraints.
3. **ICF:** It runs the Interactive Collaborative Filtering algorithm with linear UCB [6] without considering the capacity constraints. For each user, the algorithm recommend the items that it estimates

the users will obtain the most reward. Since this method does not consider the capacities, the recommendations do not necessarily satisfy the capacity constraints. Therefore, if an item is recommended to more users than its capacity, we assume that only a randomly chosen subset of the assigned users are able to get the item. The users that are not able to get the item do not send any reward feedback to the system.

4. **ICF2:** It is the same as ICF method described above, except that the algorithm observes a zero reward ($R_{t,u,i} = 0$) for all the user-item allocations that were not successfully achieved. As a result of the low rewards obtained from allocations that lead to capacity violations, the algorithm learns to avoid violating the capacities.

Experimental setup and datasets: We use a synthetic dataset and two real world datasets to evaluate our approach. For the synthetic data, we generate an (approximately) low-rank random matrix $\Theta^* \in \mathbb{R}^{N \times M}$ with entries from $[0, B]$. For the real-world data, we consider the following publicly available data sets: Movielens 100k [23] which includes ratings from 943 users on 1682 movies and the RC (Restaurant and Consumer) dataset [24] which includes ratings from 138 users on 130 restaurants. As the information of capacities are not given in the considered data, and to the best of our knowledge to any of the publicly available recommendation datasets, we consider instantiating random capacities for all items as described shortly. We consider settings with static and time-varying capacities/demands. For the static case, we assume that all users request one item at all iterations, and the capacity of each each remains unchanged with time. In the dynamic setting, we allow both the demands \mathbf{d}_t and capacities \mathbf{c}_t vary with time t . At each allocation round, we consider that each entry of \mathbf{d}_t is independently sampled from a fixed probability distribution over $\{0, 1\}$. Therefore, while active users (with demand 1) are allocated at most one item, the inactive users (with demand 0) do not get allocated any item. Similarly, each entry of \mathbf{c}_t is independently sampled from a uniform distribution over $\{0, 1, \dots, C_{\max}\}$. At each round t , if user u is allocated the item i , the system observes a reward with normal distribution $\mathcal{N}(\theta_{ui}^*, \eta^2)$.

Results: We summarize our results in Figure 2. Further experimental details and results are left to Appendix G. The observations can be summed up into following points: (1) LR-COMB (our proposed approach) is able to achieve lower regret than all other baseline methods in all experimental settings. (2) Even though the ACF method performs slightly better than LR-COMB in the initial rounds, it often gets stuck at high-regret allocations, and hence cannot achieve *no-regret*. It suffers from large regrets in the long-term because it tries to directly exploit the information it acquired so far without making any deliberate explorations. Therefore, we observe the significance of employing a bandit-based approach in achieving a no-regret algorithm. (3) Since CUCB does not leverage the low-rank structure of the parameters, it needs to sample and learn about each user-item allocation pair separately. Hence, it takes much longer for it to learn the optimum allocations. (4) Since ICF does not consider the capacities while making the allocations, it ends up incurring very large regrets. Even if it is able to identify the high-reward allocation pairs via collaborative filtering, the recommendations exceed the respective capacities of the items and we cannot obtain high rewards. (5) One possible ad-hoc approach to mitigate the issues with ICF is to use ICF2 which can indirectly capture the effects of the capacities since it receives zero rewards when the items are not successfully allocated. Nevertheless, ICF2 still does not directly use the knowledge of the capacities and hence it is still quite suboptimal. Even though it is able to show decent performance in static settings, its performance significantly degrades when the capacities dynamically change with time.

5 Conclusion and future directions

In this paper, we have studied the setting of interactive recommendations that achieve socially optimal allocations under capacity constraints. We have formulated the problem as a low-rank combinatorial multi-armed bandit and proposed an algorithm that enjoys low regret. Building on the ideas founded in this work, we aim to pursue joint recommendation and pricing mechanisms that will achieve optimal allocations in the general problem setting with users actively reacting to the recommendations based on the prices determined by the provider. We believe that this is a practically-relevant question to be resolved because it allows for design of many interesting real-world recommendation applications for settings with associated markets.

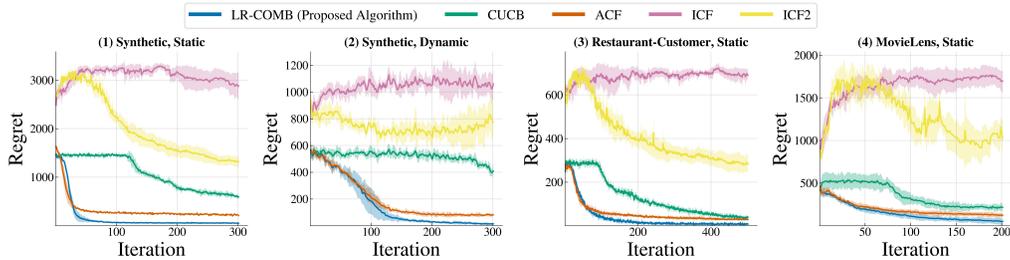


Figure 2: Instantaneous regret incurred in each round in different experimental settings. From left to right: (1) synthetic data in a static setting with $N = 800$, $M = 400$, $R = 20$, (2) synthetic data in a dynamic setting with $N = 1000$, $M = 150$, $R = 20$, probability of user activity 0.2, (3) Restaurant-Customer data in a static setting, (4) MovieLens 100k data in a static setting. In all settings, the experiments are run on 10 problem instances and means are reported together with error regions that indicate one standard deviation of uncertainty.

References

- [1] J. B. Schafer, D. Frankowski, J. Herlocker, and S. Sen, “Collaborative filtering recommender systems,” *The Adaptive Web*, vol. 4321, pp. 291–324, 2007.
- [2] J. Bennett and S. Lanning, “The netflix prize,” *Proc. KDD Cup Workshop*, pp. 3–6, Aug. 2007.
- [3] Y. Koren, R. Bell, and C. Volinsky, “Matrix factorization techniques for recommender systems,” *Computer*, vol. 42, no. 8, pp. 30–37, Aug. 2009.
- [4] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, “Item-based collaborative filtering recommendation algorithms,” *Proc. 10th International Conference on World Wide Web*, pp. 285–295, May 2001.
- [5] J. Kawale, H. H. Bui, B. Kveton, L. Tran-Thanh, and S. Chawla, “Efficient thompson sampling for online matrix-factorization recommendation,” *Advances in Neural Information Processing Systems*, vol. 28, pp. 1297–1305, Dec. 2015.
- [6] X. Zhao, W. Zhang, and J. Wang, “Interactive collaborative filtering,” *Proceedings of the 22nd ACM International Conference on Information & Knowledge Management*, pp. 1411–1420, Oct. 2013.
- [7] V. L. Smith, “Experimental auction markets and the walrasian hypothesis,” *Journal of Political Economy*, vol. 73, no. 4, pp. 64–70, 1965.
- [8] F. Kelly, “Charging and rate control for elastic traffic,” *European Transactions on Telecommunications*, vol. 8, no. 1, pp. 33–37, Jan. 1997.
- [9] A. Barraza-Urbina, “The exploration-exploitation trade-off in interactive recommender systems,” *Proceedings of the Eleventh ACM Conference on Recommender Systems*, pp. 431–435, Aug. 2017.
- [10] Q. Wang, C. Zeng, W. Zhou, T. Li, S. S. Iyengar, L. Shwartz, and G. Y. Grabarnik, “Online interactive collaborative filtering using multi-armed bandit with dependent arms,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 8, pp. 1569–1580, Aug. 2019.
- [11] V. Dani, T. Hayes, and S. M. Kakade, “Stochastic linear optimization under bandit feedback,” *21st Annual Conference on Learning Theory - COLT 2008, Helsinki, Finland*, pp. 355–366, Jul. 2008.
- [12] Y. Abbasi-Yadkori, “Improved algorithms for linear stochastic bandits,” *Advances in Neural Information Processing Systems*, pp. 2312–2320, Dec. 2011.
- [13] W. Chen, Y. Wang, and Y. Yuan, “Combinatorial multi-armed bandit: General framework and applications,” *Proceedings of the 30th International Conference on Machine Learning*, vol. 28, no. 1, pp. 151–159, Jun. 2013.
- [14] B. Kveton, Z. Wen, A. Ashkan, and C. Szepesvari, “Tight Regret Bounds for Stochastic Combinatorial Semi-Bandits,” *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, vol. 38, pp. 535–543, May 2015.
- [15] J.-Y. Audibert, S. Bubeck, and G. Lugosi, “Minimax policies for combinatorial prediction games,” *Proceedings of the 24th Annual Conference on Learning Theory*, vol. 19, pp. 107–132, Jun. 2011.
- [16] N. Johnson, V. Sivakumar, and A. Banerjee, “Structured stochastic linear bandits,” *arXiv preprint arXiv:1606.05693*, 2016.
- [17] R. Combes, S. Magureanu, and A. Proutiere, “Minimal exploration in structured stochastic bandits,” *Advances in Neural Information Processing Systems*, vol. 30, pp. 1761–1769, Dec. 2017.
- [18] Y. Lu, A. Meisami, and A. Tewari, “Low-rank generalized linear bandit problems,” *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics*, vol. 130, pp. 460–468, Apr. 2021.
- [19] K. Christakopoulou, J. Kawale, and A. Banerjee, “Recommendation with capacity constraints,” *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pp. 1439–1448, Nov. 2017.
- [20] R. Makhijani, S. Chakrabarti, D. Struble, and Y. Liu, “Lore: A large-scale offer recommendation engine with eligibility and capacity constraints,” *Proceedings of the 13th ACM Conference on Recommender Systems*, pp. 160–168, Sep. 2019.

- [21] L. Liu, H. Mania, and M. I. Jordan, “Competing bandits in matching markets,” *Proceedings of the Twenty-Third Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1618–1628, Aug. 2020.
- [22] D. Palomar and M. Chiang, “A tutorial on decomposition methods for network utility maximization,” *IEEE Journal on Selected Areas in Communications*, vol. 24, no. 8, pp. 1439–1451, Jul. 2006.
- [23] F. M. Harper and J. A. Konstan, “The movielens datasets: History and context,” *ACM Trans. Interact. Intell. Syst.*, vol. 5, no. 4, Dec. 2015.
- [24] V. Blanca, G. Gabriel, and P. Rafael, “Effects of relevant contextual features in the performance of a restaurant recommender system,” *3rd Workshop on Context-Aware Recsys*, Oct. 2011.
- [25] D. P. Bertsekas, *Linear network optimization: algorithms and codes*. Mit Press, 1991.
- [26] I. Heller and C. B. Tompkins, “An extension of a theorem of Dantzig’s,” *Linear Inequalities and Related Systems, Annals of Mathematics Studies*, vol. 38, pp. 247–254, 1957.
- [27] E. J. Candes and Y. Plan, “Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements,” *IEEE Transactions on Information Theory*, vol. 57, no. 4, pp. 2342–2359, Mar. 2011.

A Implementation via Matrix Factorization

The following algorithm describes an efficient implementation of our Low-Rank Combinatorial Bandit algorithm using matrix factorization. Note that converged $\hat{\Theta}_t$ and \mathbf{X}_t are not necessarily the optimum solution for problems (3.1) and (3.2) since the problems are not convex. However, the alternating optimization algorithm guarantees that, in each iteration, the objective value only decreases for (3.2). Similarly, the objective value for (3.1) increases in each iteration of the alternating optimization.

Algorithm 2 LR-COMB with Matrix Factorization

Require: horizon T , initial estimate $\bar{\Theta} \in \mathbb{R}^d$ with $\|\bar{\Theta} - \Theta^*\|_F \leq G$, parameters $\delta, \alpha > 0, \gamma \geq 1$.

for $t = 1, 2, \dots, T$ **do**

randomly initialize $\hat{\mathbf{P}}$ and $\hat{\mathbf{Q}}$

while convergence criterion not satisfied **do**

$\hat{\mathbf{P}} \leftarrow \arg \min_{\mathbf{P} \in \mathbb{R}^{N \times R}} \left\{ \sum_{\tau=1}^{t-1} \sum_{(u,i) \in \mathcal{A}_\tau} (\mathbf{p}_u^T \mathbf{q}_i - R_{\tau,u,i})^2 + \gamma \|\mathbf{P}\mathbf{Q}^T - \bar{\Theta}\|_F^2 \right\}$

$\hat{\mathbf{Q}} \leftarrow \arg \min_{\mathbf{Q} \in \mathbb{R}^{M \times R}} \left\{ \sum_{\tau=1}^{t-1} \sum_{(u,i) \in \mathcal{A}_\tau} (\mathbf{p}_u^T \mathbf{q}_i - R_{\tau,u,i})^2 + \gamma \|\mathbf{P}\mathbf{Q}^T - \bar{\Theta}\|_F^2 \right\}$

end while

$\hat{\Theta}_t \leftarrow \hat{\mathbf{P}}\hat{\mathbf{Q}}^T$

$\mathbf{X} \leftarrow \mathbb{1}_{N \times M}, \mathbf{P} \leftarrow \hat{\mathbf{P}}, \mathbf{Q} \leftarrow \hat{\mathbf{Q}}$

while convergence criterion not satisfied **do**

while convergence criterion not satisfied **do**

$\mathbf{P} \leftarrow \arg \max_{\mathbf{P} \in \mathbb{R}^{N \times R}} \langle \mathbf{X}, \mathbf{P}\mathbf{Q}^T \rangle$ s.t. $\|\mathbf{P}\mathbf{Q}^T - \hat{\Theta}_t\|_{2, E_t} \leq \sqrt{\beta_t^*(\delta, \alpha, \gamma)}$

$\mathbf{Q} \leftarrow \arg \max_{\mathbf{Q} \in \mathbb{R}^{M \times R}} \langle \mathbf{X}, \mathbf{P}\mathbf{Q}^T \rangle$ s.t. $\|\mathbf{P}\mathbf{Q}^T - \hat{\Theta}_t\|_{2, E_t} \leq \sqrt{\beta_t^*(\delta, \alpha, \gamma)}$

end while

$\Theta \leftarrow \mathbf{P}\mathbf{Q}^T$

while convergence criterion not satisfied **do**

for $u \in [N]$ **do**

$\mathbf{x}_u \leftarrow \arg \max_{\mathbf{x}} \{ \mathbf{x}^T (\theta_u - \lambda) \mid \mathbf{x} \in \{0, 1\}^M, \mathbf{x}^T \mathbb{1}_M \leq d_{t,u} \}$

end for

$\lambda \leftarrow \left[\lambda - \alpha \left(\mathbf{c}_t - \sum_{u=1}^N \mathbf{x}_u \right) \right]^+$

end while

$\mathbf{X} \leftarrow [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$

end while

$\mathbf{X}_t \leftarrow \mathbf{X}$

Play the arms \mathcal{A}_t according to \mathbf{X}_t

Observe $R_{t,u,i}$ for all $(u, i) \in \mathcal{A}_t$

end for

B Relaxation of Integer Program

A traditional linear integer program (IP) in matrix form is formulated as

$$\begin{aligned}
 & \max_{\mathbf{x}} \mathbf{t}^T \mathbf{x} \\
 & \text{s.t. } \mathbf{A}\mathbf{x} \leq \mathbf{b} \\
 & \mathbf{x} \in \mathbb{Z}_+^d
 \end{aligned} \tag{B.1}$$

This problem can be relaxed to a linear program by dropping the integral constraints (setting $\mathbf{x} \in \mathbb{R}_+^d$). The integrality gap of an integer program is defined as the difference between the optimal values of the integer program in (IP) and its relaxed linear program. When the vector \mathbf{b} is integral and the matrix \mathbf{A} is totally unimodular (all entries are 1, 0, or -1 and every square sub-minor has determinant of +1 or -1) then the integrality gap is zero and the solution of the relaxed linear program is integer valued [25].

Hence, we can solve (B.1) by instead solving the following relaxed linear program:

$$\begin{aligned} \max_{\mathbf{x}} \quad & \mathbf{t}^\top \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \in \mathbb{R}_+^d \end{aligned} \tag{B.2}$$

For a matrix \mathbf{A} whose rows can be partitioned into two disjoint sets \mathcal{C} and \mathcal{D} , the following four conditions together are sufficient for \mathbf{A} to be totally unimodular [26]:

1. Every entry in \mathbf{A} is 0, +1, or -1.
2. Every column of \mathbf{A} contains at most two non-zero entries.
3. If two non-zero entries in a column of \mathbf{A} have the same sign, then the row of one is in \mathcal{C} , and the other in \mathcal{D} .
4. If two non-zero entries in a column of \mathbf{A} have opposite signs, then the rows of both are in \mathcal{C} , or both in \mathcal{D} .

In the setting of resource allocation, we can write problem (2.1) equivalently as problem (B.1) where $\mathbf{x} = \text{vec}(\mathbf{X})$, $\mathbf{t} = \text{vec}(\Theta^*)$, \mathbf{A} and \mathbf{b} are given as

$$\mathbf{A} = \begin{bmatrix} \mathbf{1}_N^\top \otimes \mathbf{I}_M \\ \mathbf{I}_N \otimes \mathbf{1}_M^\top \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} \mathbf{c}_t \\ \mathbf{d}_t \end{bmatrix} \tag{B.3}$$

For matrix \mathbf{A} given in (B.3), we can set \mathcal{C} to be the set of first M rows corresponding to the capacity constraints, and \mathcal{D} to be the set of remaining rows corresponding to the demand constraints. Since this \mathbf{A} matrix satisfies the conditions of the proposition for sets \mathcal{C} and \mathcal{D} , we obtain that \mathbf{A} is totally unimodular. Finally, since the vector \mathbf{b} is integral and the matrix \mathbf{A} is totally unimodular, the integrality gap is zero.

C Structured combinatorial multi-armed bandits

For the ease of exposition, we present our proofs in the following setting with d structured arms.

We consider a CMAB (Combinatorial Multi Armed Bandit) problem setting with d arms associated with a set of independent random rewards $R_{t,i}$ for $i \in [d]$ and $t \in \mathbb{N}$. Assume that the set of rewards $\{R_{t,i} | t \in \mathbb{N}\}$ associated with arm i are η -sub-Gaussian with mean $\theta_i^* \in [0, B]$. Let $\theta^* = (\theta_1, \theta_2, \dots, \theta_d)$ be the vector of expectations of all arms and assume we know that it belongs to a structure set $\mathcal{F} \subseteq \mathbb{R}^d$. We further assume that we have access to an initial rough estimate $\bar{\theta}$ that satisfies $\|\bar{\theta} - \theta^*\|_2 \leq G$ (one can also set $\bar{\theta} = \mathbf{0}$ and let G be some number satisfying $\|\theta^*\|_2 \leq G$).

At each round t , a subset of arms $\mathcal{A}_t \subseteq [d]$ are played and the individual outcomes of arms in \mathcal{A}_t are revealed. The total reward at round t is the sum of the rewards obtained from all arms in \mathcal{A}_t . Letting $\mathbf{e}_i \in \mathbb{R}^d$ denote the zero-one vector with a single one at the i -th entry, define the action vector for time t as $\mathbf{x}_t = \sum_{i \in \mathcal{A}_t} \mathbf{e}_i$. Consequently, \mathbf{x}_t becomes a zero-one vector with ones at entries \mathcal{A}_t and zeros everywhere else. The problem contains a constraint that any valid action \mathbf{x}_t must belong to a (time-varying) constraint set $\mathcal{X}_t \subseteq \{0, 1\}^d$.

The optimum allocation \mathbf{x}_t^* at time t is given by $\mathbf{x}_t^* \in \arg \max_{\mathbf{x} \in \mathcal{X}_t} \langle \mathbf{x}, \theta^* \rangle$. We denote by H_t the history $\{\mathbf{x}_\tau, (R_{\tau,i})_{i \in \mathcal{A}_\tau}\}_{\tau=1}^{t-1}$ of observations available when choosing the next action \mathbf{x}_t . Let π be a policy which takes the action \mathbf{x}_t using the history H_t . Then, the T period regret of a policy π is the random variable $\mathcal{R}(T, \pi) = \sum_{t=1}^T [\langle \mathbf{x}_t^*, \theta^* \rangle - \langle \mathbf{x}_t, \theta^* \rangle]$.

C.1 OFU for Structured Combinatorial Bandits

We present our algorithm in the setting where the mean reward vector θ^* belongs to a structure set $\mathcal{F} \subseteq \mathbb{R}^d$. Then, we analyze the algorithm to establish performance guarantees.

The algorithm maintains a confidence set \mathcal{C}_t that contains the true parameter θ^* with high probability and chooses the action \mathbf{x}_t according to

$$(\mathbf{x}_t, \tilde{\theta}_t) = \arg \max_{(\mathbf{x}, \theta) \in \mathcal{X}_t \times \mathcal{C}_t} \langle \mathbf{x}, \theta \rangle \tag{C.1}$$

Algorithm 3 OFU for Structured Combinatorial Bandits

Require: horizon T , initial estimate $\bar{\boldsymbol{\theta}} \in \mathbb{R}^d$ with $\|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 \leq G$, parameters $\delta, \alpha > 0, \gamma \geq 1$.

for $t = 1, 2, \dots, T$ **do**

Find the least squares estimate $\hat{\boldsymbol{\theta}}_t = \arg \min_{\boldsymbol{\theta} \in \mathcal{F}} \{L_{2,t}(\boldsymbol{\theta}) + \gamma \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2\}$

Construct the confidence set $\mathcal{C}_t = \{\boldsymbol{\theta} \in \mathcal{F} : \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t\|_{2,E_t} \leq \sqrt{\beta_t^*(\delta, \alpha, \gamma)}\}$

Compute the action vector $\mathbf{x}_t = \arg \max_{\mathbf{x} \in \mathcal{X}_t} \max_{\boldsymbol{\theta} \in \mathcal{C}_t} \langle \mathbf{x}, \boldsymbol{\theta} \rangle$

Play the arms \mathcal{A}_t according to \mathbf{x}_t

Observe $(R_{\tau,i})_{i \in \mathcal{A}_\tau}$

end for

The confidence sets that we construct are centered around the regularized least square estimates defined next. We first let the cumulative squared prediction error at time t be

$$L_{2,t}(\boldsymbol{\theta}) = \sum_{\tau=1}^{t-1} \sum_{i \in \mathcal{A}_\tau} (\theta_i - R_{\tau,i})^2 \quad (\text{C.2})$$

and define the regularized least squares estimate at time t as

$$\hat{\boldsymbol{\theta}}_t = \arg \min_{\boldsymbol{\theta} \in \mathcal{F}} \{L_{2,t}(\boldsymbol{\theta}) + \gamma \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2\} \quad (\text{C.3})$$

Then, the confidence sets take the form $\mathcal{C}_t := \{\boldsymbol{\theta} \in \mathcal{F} : \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t\|_{2,E_t} \leq \sqrt{\beta_t}\}$ where β_t is an appropriately chosen confidence parameter, and the regularized empirical 2-norm $\|\cdot\|_{2,E_t}$ is defined by

$$\|\boldsymbol{\Delta}\|_{2,E_t}^2 := \sum_{\tau=1}^{t-1} \sum_{i \in \mathcal{A}_\tau} \langle \boldsymbol{\Delta}, \mathbf{e}_i \rangle^2 + \gamma \|\boldsymbol{\Delta}\|_2^2 = \sum_{i=1}^d (n_{t,i} + \gamma) (\Delta_i)^2$$

where $n_{t,i} := \sum_{\tau=1}^{t-1} \mathbb{1}\{i \in \mathcal{A}_\tau\}$ denotes the number of times arm i has been pulled before time t (excluding time t). For future reference, we also define the (non-regularized) empirical 2-norm $\|\cdot\|_{2,\tilde{E}_t}$ by

$$\|\boldsymbol{\Delta}\|_{2,\tilde{E}_t}^2 := \sum_{\tau=1}^{t-1} \sum_{i \in \mathcal{A}_\tau} \langle \boldsymbol{\Delta}, \mathbf{e}_i \rangle^2 = \sum_{i=1}^d n_{t,i} (\Delta_i)^2$$

Note that the regularized empirical 2-norm is related to (non-regularized) empirical 2-norm as

$$\|\boldsymbol{\Delta}\|_{2,E_t}^2 = \|\boldsymbol{\Delta}\|_{2,\tilde{E}_t}^2 + \gamma \|\boldsymbol{\Delta}\|_2^2$$

By Lemma D.3, we establish that for any $\boldsymbol{\theta} \in \mathbb{R}^d$,

$$\mathbb{P} \left(L_{2,t}(\boldsymbol{\theta}) \geq L_{2,t}(\boldsymbol{\theta}^*) + \frac{1}{2} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\tilde{E}_t}^2 - 4\eta^2 \log(1/\delta) \quad , \forall t \in \mathbb{N} \right) \geq 1 - \delta \quad (\text{C.4})$$

Hence, with high probability, $\boldsymbol{\theta}$ can achieve lower squared error than $\boldsymbol{\theta}^*$ only if the empirical deviation $\|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\tilde{E}_t}^2$ is less than $8\eta^2 \log(1/\delta)$.

In order to make this property hold uniformly for all $\boldsymbol{\theta}$ in a subset \mathcal{C}_t of \mathcal{F} , we discretize \mathcal{C}_t at some discretization scale α and apply a union bound for this finite discretization set. Let $\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_2)$ denote the α -covering number of \mathcal{F} in the 2-norm $\|\cdot\|_2$, and let

$$\beta_t^*(\delta, \alpha, \gamma) := 8\eta^2 \log(\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_2)/\delta) + 2\alpha t \sqrt{d} \left[8B + \sqrt{8\eta^2 \log(4dt^2/\delta)} \right] + 4\gamma G^2 \quad (\text{C.5})$$

Then, Lemma D.5 shows that if we set $\beta_t = \beta_t^*(\delta, \alpha)$, the confidence sets \mathcal{C}_t contain the true parameter $\boldsymbol{\theta}^*$ for all t with high probability. Following the construction of the confidence sets, the next step is to obtain the overall regret guarantee. As given in Corollary E.10, we find that the regret of Algorithm 3 satisfies

$$\mathcal{R}(T, \pi) = \tilde{\mathcal{O}} \left(\sqrt{\eta^2 d \log(\mathcal{N}(\mathcal{F}, T^{-1}, \|\cdot\|_2)) T} \right) \quad (\text{C.6})$$

Remark C.1. In the literature of linear bandits, the typical observation model is such that each action \mathbf{x}_t results in a single reward feedback with mean $\langle \mathbf{x}_t, \boldsymbol{\theta}^* \rangle$ and sub-Gaussianity parameter $d\eta^2$ (since each arm has a η^2 -sub-Gaussian reward). Therefore, that observation model can only obtain $\tilde{\mathcal{O}}(\sqrt{\eta^2 d^2 \log(\mathcal{N}(\mathcal{F}, T^{-1}, \|\cdot\|_2)) T})$ regret guarantee. However, in our setting, the observations are sets of independent rewards $\{R_{t,i}\}_{i \in \mathcal{A}_t}$ where each element is η^2 -sub-Gaussian. Due to this richer nature of the observation model, we are able to achieve lower regret guarantees than only observing a single cumulative reward.

D Proofs for Confidence Sets

D.1 Martingale Exponential Inequalities

We start with preliminary results on martingale exponential inequalities.

Consider a sequence of random variables $(Z_n)_{n \in \mathbb{N}}$ adapted to the filtration $(\mathcal{H}_n)_{n \in \mathbb{N}}$. Assume $\mathbb{E}[\exp(\lambda Z_i)]$ is finite for all λ . Define the conditional mean $\mu_i = \mathbb{E}[Z_i | \mathcal{H}_{i-1}]$, and define the conditional cumulant generating function of the centered random variable $[Z_i - \mu_i]$ by $\psi_i(\lambda) := \log \mathbb{E}[\exp(\lambda[Z_i - \mu_i]) | \mathcal{H}_{i-1}]$. Let

$$M_n(\lambda) = \exp \left\{ \sum_{i=1}^n \lambda[Z_i - \mu_i] - \psi_i(\lambda) \right\}$$

Lemma D.1. $(M_n(\lambda))_{n \in \mathbb{N}}$ is a martingale with respect to the filtration $(\mathcal{H}_n)_{n \in \mathbb{N}}$, and $\mathbb{E}[M_n(\lambda)] = 1$.

Proof. By definition, we have

$$\mathbb{E}[M_1(\lambda) | \mathcal{H}_0] = \mathbb{E}[\exp\{\lambda[Z_1 - \mu_1] - \psi_1(\lambda)\} | \mathcal{H}_0] = 1$$

Then, for any $n \geq 2$,

$$\begin{aligned} \mathbb{E}[M_n(\lambda) | \mathcal{H}_{n-1}] &= \mathbb{E}[M_{n-1}(\lambda) \exp\{\lambda[Z_n - \mu_n] - \psi_n(\lambda)\} | \mathcal{H}_{n-1}] \\ &= M_{n-1}(\lambda) \mathbb{E}[\exp\{\lambda[Z_n - \mu_n] - \psi_n(\lambda)\} | \mathcal{H}_{n-1}] \\ &= M_{n-1}(\lambda) \end{aligned}$$

since $M_{n-1}(\lambda)$ is a measurable function of the filtration \mathcal{H}_{n-1} . □

Lemma D.2. For all $x \geq 0$ and $\lambda \geq 0$,

$$\mathbb{P} \left(\sum_{i=1}^n \lambda Z_i \leq x + \sum_{i=1}^n [\lambda \mu_i + \psi_i(\lambda)] \quad , \forall t \in \mathbb{N} \right) \geq 1 - e^{-x}$$

Proof. For any λ , $(M_n(\lambda))_{n \in \mathbb{N}}$ is a martingale with respect to $(\mathcal{H}_n)_{n \in \mathbb{N}}$ and $\mathbb{E}[M_n(\lambda)] = 1$ by Lemma D.1. For arbitrary $x \geq 0$, define $\tau_x = \inf\{n \geq 0 | M_n(\lambda) \geq x\}$ and note that τ_x is a stopping time corresponding to the first time M_n crosses the boundary x . Since τ is a stopping time with respect to $(\mathcal{H}_n)_{n \in \mathbb{N}}$, we have $\mathbb{E}[M_{\tau_x \wedge n}(\lambda)] = 1$. Then, by Markov's inequality

$$x \mathbb{P}(M_{\tau_x \wedge n}(\lambda) \geq x) \leq \mathbb{E}[M_{\tau_x \wedge n}(\lambda)] = 1$$

Noting that the event $\{M_{\tau_x \wedge n}(\lambda) \geq x\} = \bigcup_{k=1}^n \{M_k(\lambda) \geq x\}$, we have

$$\mathbb{P} \left(\bigcup_{k=1}^n \{M_k(\lambda) \geq x\} \right) \leq \frac{1}{x}$$

Taking the limit as $n \rightarrow \infty$, and applying monotone convergence theorem shows that $\mathbb{P}(\bigcup_{k=1}^{\infty} \{M_k(\lambda) \geq x\}) \leq \frac{1}{x}$ or $\mathbb{P}(\bigcup_{k=1}^{\infty} \{M_k(\lambda) \geq e^x\}) \leq e^{-x}$. Then, by definition of $M_k(\lambda)$, we conclude

$$\mathbb{P} \left(\bigcup_{k=1}^{\infty} \left\{ \sum_{i=1}^k \lambda[Z_i - \mu_i] - \psi_i(\lambda) \geq x \right\} \right) \leq e^{-x}$$

□

D.2 Proofs for the construction of confidence sets

Lemma D.3. For any $\delta > 0$ and $\boldsymbol{\theta} \in \mathbb{R}^d$,

$$\mathbb{P}\left(L_{2,t}(\boldsymbol{\theta}) \geq L_{2,t}(\boldsymbol{\theta}^*) + \frac{1}{2}\|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\tilde{E}_t}^2 - 4\eta^2 \log(1/\delta), \forall t \in \mathbb{N}\right) \geq 1 - \delta \quad (\text{D.1})$$

Proof. Let \mathcal{H}_{t-1} be the σ -algebra generated by (H_t, \mathcal{A}_t) and let $\mathcal{H}_0 = \sigma(\emptyset, \Omega)$. Then, define $\epsilon_{t,i} := R_{t,i} - \langle \mathbf{X}_{t,i}, \boldsymbol{\theta}^* \rangle$ for all $t \in \mathbb{N}$ and $i \in \mathcal{A}_t$. By previous assumptions, $\mathbb{E}[\epsilon_{t,i} | \mathcal{H}_{t-1}] = 0$ and $\mathbb{E}[\exp(\lambda \epsilon_{t,i}) | \mathcal{H}_{t-1}] \leq \exp\left(\frac{\lambda^2 \eta^2}{2}\right)$ for all t .

Define $Z_{t,i} = (R_{t,i} - \langle \mathbf{X}_{t,i}, \boldsymbol{\theta}^* \rangle)^2 - (R_{t,i} - \langle \mathbf{X}_{t,i}, \boldsymbol{\theta} \rangle)^2$. Then, we have

$$\begin{aligned} Z_{t,i} &= -(\langle \mathbf{X}_{t,i}, \boldsymbol{\theta} \rangle - \langle \mathbf{X}_{t,i}, \boldsymbol{\theta}^* \rangle)^2 + 2\epsilon_{t,i}(\langle \mathbf{X}_{t,i}, \boldsymbol{\theta} \rangle - \langle \mathbf{X}_{t,i}, \boldsymbol{\theta}^* \rangle) \\ &= -\langle \mathbf{X}_{t,i}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle^2 + 2\epsilon_{t,i} \langle \mathbf{X}_{t,i}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle \end{aligned}$$

Therefore, the conditional mean and conditional cumulant generating function satisfy

$$\begin{aligned} \mu_{t,i} &:= \mathbb{E}[Z_{t,i} | \mathcal{H}_{t-1}] = -\langle \mathbf{X}_{t,i}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle^2 \\ \psi_t(\lambda) &:= \log \mathbb{E}[\exp(\lambda Z_{t,i} - \mu_{t,i}) | \mathcal{H}_{t-1}] \\ &= \log \mathbb{E}[\exp(2\lambda \langle \mathbf{X}_{t,i}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle \epsilon_{t,i}) | \mathcal{H}_{t-1}] \\ &\leq \frac{(2\lambda \langle \mathbf{X}_{t,i}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle)^2 \eta^2}{2} \end{aligned}$$

Applying Lemma D.2 shows that for all $x \geq 0$ and $\lambda \geq 0$,

$$\mathbb{P}\left(\sum_{\tau=1}^{t-1} \sum_{u \in \mathcal{A}_\tau} Z_{\tau,i} \leq \frac{x}{\lambda} + \sum_{\tau=1}^{t-1} \sum_{u \in \mathcal{A}_\tau} \langle \mathbf{X}_{\tau,i}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle^2 (2\lambda \eta^2 - 1), \forall t \in \mathbb{N}\right) \geq 1 - e^{-x}$$

Note that we have $\sum_{\tau=1}^{t-1} \sum_{u \in \mathcal{A}_\tau} Z_{\tau,i} = L_{2,t}(\boldsymbol{\theta}^*) - L_{2,t}(\boldsymbol{\theta})$,

and $\sum_{\tau=1}^{t-1} \sum_{u \in \mathcal{A}_\tau} \langle \mathbf{X}_{\tau,i}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle^2 = \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\tilde{E}_t}^2$.

Then, choosing $\lambda = \frac{1}{4\eta^2}$ and $x = \log \frac{1}{\delta}$ gives

$$\mathbb{P}\left(L_{2,t}(\boldsymbol{\theta}) \geq L_{2,t}(\boldsymbol{\theta}^*) + \frac{1}{2}\|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\tilde{E}_t}^2 - 4\eta^2 \log(1/\delta), \forall t \in \mathbb{N}\right) \geq 1 - \delta$$

□

Lemma D.4. If $\boldsymbol{\theta}^\alpha \in \mathcal{F}^\alpha$ satisfies $\|\boldsymbol{\theta} - \boldsymbol{\theta}^\alpha\|_2 \leq \alpha$, then with probability at least $1 - \delta$,

$$\left| \frac{1}{2}\|\boldsymbol{\theta}^* - \boldsymbol{\theta}^\alpha\|_{2,\tilde{E}_t}^2 - \frac{1}{2}\|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\tilde{E}_t}^2 + L_{2,t}(\boldsymbol{\theta}) - L_{2,t}(\boldsymbol{\theta}^\alpha) \right| \leq \alpha t d \left[8B + \sqrt{8\eta^2 \log(4dt^2/\delta)} \right] \quad (\text{D.2})$$

Proof. Since any two $\boldsymbol{\theta}, \boldsymbol{\theta}^\alpha \in \mathcal{F}$ satisfy $\|\boldsymbol{\theta} - \boldsymbol{\theta}^\alpha\|_2 \leq \sqrt{d}B$, it is enough to consider $\alpha \leq \sqrt{d}B$. We find

$$\begin{aligned} \sum_{i=1}^d |\langle \boldsymbol{\theta}, \mathbf{e}_i \rangle^2 - \langle \boldsymbol{\theta}^\alpha, \mathbf{e}_i \rangle^2| &\leq \max_{\|\boldsymbol{\Delta}\|_2 \leq \alpha} \left\{ \sum_{i=1}^d |\theta_i^2 - (\theta_i + \Delta_i)^2| \right\} \\ &= \max_{\|\boldsymbol{\Delta}\|_2 \leq \alpha} \left\{ \sum_{i=1}^d |2\theta_i \Delta_i + \Delta_i^2| \right\} \\ &\leq \max_{\|\boldsymbol{\Delta}\|_2 \leq \alpha} \left\{ 2 \sum_{i=1}^d |\theta_i \Delta_i| + \sum_{i=1}^d \Delta_i^2 \right\} \\ &\leq \max_{\|\boldsymbol{\Delta}\|_2 \leq \alpha} \{ 2B \|\boldsymbol{\Delta}\|_1 + \|\boldsymbol{\Delta}\|_2^2 \} \\ &\leq 2B\sqrt{d}\alpha + \alpha^2 \end{aligned}$$

Therefore, it implies

$$\begin{aligned} \sum_{i=1}^d |\langle \boldsymbol{\theta} - \boldsymbol{\theta}^*, \mathbf{e}_i \rangle^2 - \langle \boldsymbol{\theta}^\alpha - \boldsymbol{\theta}^*, \mathbf{e}_i \rangle^2| &= \sum_{i=1}^d |\langle \boldsymbol{\theta}, \mathbf{e}_i \rangle^2 - \langle \boldsymbol{\theta}^\alpha, \mathbf{e}_i \rangle^2 + 2\langle \boldsymbol{\theta}^*, \mathbf{e}_i \rangle \langle \boldsymbol{\theta}^\alpha - \boldsymbol{\theta}, \mathbf{e}_i \rangle| \\ &\leq 2B\sqrt{d}\alpha + \alpha^2 + 2B\|\boldsymbol{\theta} - \boldsymbol{\theta}^\alpha\|_1 \\ &\leq 4B\sqrt{d}\alpha + \alpha^2 \end{aligned}$$

Similarly, for any t , we have

$$\begin{aligned} \sum_{i=1}^d |(R_{t,i} - \langle \boldsymbol{\theta}, \mathbf{e}_i \rangle)^2 - (R_{t,i} - \langle \boldsymbol{\theta}^\alpha, \mathbf{e}_i \rangle)^2| &= \sum_{i=1}^d |2R_{t,i}\langle \boldsymbol{\theta}^\alpha - \boldsymbol{\theta}, \mathbf{e}_i \rangle + \langle \boldsymbol{\theta}, \mathbf{e}_i \rangle^2 - \langle \boldsymbol{\theta}^\alpha, \mathbf{e}_i \rangle^2| \\ &\leq 2 \sum_{i=1}^d |R_{t,i}| |\langle \boldsymbol{\theta}^\alpha - \boldsymbol{\theta}, \mathbf{e}_i \rangle| + 2B\sqrt{d}\alpha + \alpha^2 \\ &\leq 2\|\boldsymbol{\theta}^\alpha - \boldsymbol{\theta}\|_\infty \sum_{i=1}^d |R_{t,i}| + 2B\sqrt{d}\alpha + \alpha^2 \\ &\leq 2\alpha \sum_{i=1}^d |R_{t,i}| + 2B\sqrt{d}\alpha + \alpha^2 \end{aligned}$$

Summing over t and noting that $\mathcal{A}_t \subseteq [d]$, the left hand side of (D.2) is bounded by

$$\sum_{\tau=1}^{t-1} \left(\frac{1}{2} [4B\sqrt{d}\alpha + \alpha^2] + 2\alpha \sum_{i=1}^d |R_{\tau,i}| + 2B\sqrt{d}\alpha + \alpha^2 \right) \leq \alpha \sum_{\tau=1}^{t-1} \left(6B\sqrt{d} + 2 \sum_{i=1}^d |R_{\tau,i}| \right)$$

Because $\epsilon_{\tau,i}$ is η -sub-Gaussian, $\mathbb{P}\left(|\epsilon_{\tau,i}| > \sqrt{2\eta^2 \log(2/\delta)}\right) \leq \delta$. By a union bound, $\mathbb{P}\left(\exists \tau, i \text{ s.t. } |\epsilon_{\tau,i}| > \sqrt{2\eta^2 \log(4d\tau^2/\delta)}\right) \leq \frac{\delta d}{2} \sum_{\tau=1}^{\infty} \frac{1}{d\tau^2} \leq \delta$. Since $|R_{\tau,i}| \leq B + |\epsilon_{\tau,i}|$, we have $|R_{\tau,i}| \leq B + \sqrt{2\eta^2 \log(4d\tau^2/\delta)}$ with probability at least $1 - \delta$. Consequently, the bound for the discretization error becomes

$$\alpha t d \left[8B + 2\sqrt{2\eta^2 \log(4dt^2/\delta)} \right]$$

□

Lemma D.5. For any $\delta > 0$, $\alpha > 0$ and $\gamma > 0$, if

$$\mathcal{C}_t = \{\boldsymbol{\theta} \in \mathcal{F} : \|\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_t\|_{2, E_t} \leq \sqrt{\beta_t^*(\delta, \alpha, \gamma)}\} \quad (\text{D.3})$$

for all $t \in \mathbb{N}$, then

$$\mathbb{P}(\boldsymbol{\theta}^* \in \mathcal{C}_t, \forall t \in \mathbb{N}) \geq 1 - 2\delta \quad (\text{D.4})$$

Proof. Let $\mathcal{F}^\alpha \subset \mathcal{F}$ be an α -cover of \mathcal{F} in the 2-norm so that for any $\boldsymbol{\theta} \in \mathcal{F}$, there exists $\boldsymbol{\theta}^\alpha \in \mathcal{F}^\alpha$ such that $\|\boldsymbol{\theta} - \boldsymbol{\theta}^\alpha\|_2 \leq \alpha$. By a union bound applied to Lemma D.3, with probability at least $1 - \delta$,

$$L_{2,t}(\boldsymbol{\theta}^\alpha) - L_{2,t}(\boldsymbol{\theta}^*) \geq \frac{1}{2} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}^\alpha\|_{2, \tilde{E}_t}^2 - 4\eta^2 \log(|\mathcal{F}^\alpha|/\delta), \quad \forall \boldsymbol{\theta}^\alpha \in \mathcal{F}^\alpha, t \in \mathbb{N}$$

Therefore, with probability at least $1 - \delta$, for all $\boldsymbol{\theta} \in \mathcal{F}$, $t \in \mathbb{N}$,

$$\begin{aligned} L_{2,t}(\boldsymbol{\theta}) - L_{2,t}(\boldsymbol{\theta}^*) &\geq \frac{1}{2} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2, \tilde{E}_t}^2 - 4\eta^2 \log(|\mathcal{F}^\alpha|/\delta) \\ &\quad + \min_{\boldsymbol{\theta}^\alpha \in \mathcal{F}^\alpha} \left\{ \frac{1}{2} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}^\alpha\|_{2, \tilde{E}_t}^2 - \frac{1}{2} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2, \tilde{E}_t}^2 + L_{2,t}(\boldsymbol{\theta}) - L_{2,t}(\boldsymbol{\theta}^\alpha) \right\} \end{aligned}$$

By Lemma D.4, with probability at least $1 - 2\delta$,

$$L_{2,t}(\boldsymbol{\theta}) - L_{2,t}(\boldsymbol{\theta}^*) \geq \frac{1}{2} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\bar{E}_t}^2 - D_t$$

where $D_t := 4\eta^2 \log(|\mathcal{F}^\alpha|/\delta) + \alpha t d \left[8B + \sqrt{8\eta^2 \log(4dt^2/\delta)} \right]$.

Adding the regularization terms to both sides, we obtain

$$L_{2,t}(\boldsymbol{\theta}) + \gamma \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2 - L_{2,t}(\boldsymbol{\theta}^*) - \gamma \|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2 \geq \frac{1}{2} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_{2,\bar{E}_t}^2 + \gamma \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2 - D_t - \gamma \|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2$$

Note the definition of the least square estimate $\hat{\boldsymbol{\theta}}_t = \arg \min_{\boldsymbol{\theta} \in \mathcal{F}} \{L_{2,t}(\boldsymbol{\theta}) + \gamma \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2\}$. By letting $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_t$, the left hand side becomes non-positive, and hence

$$\frac{1}{2} \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_t\|_{2,\bar{E}_t}^2 \leq D_t + \gamma \left(\|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2 - \|\hat{\boldsymbol{\theta}}_t - \bar{\boldsymbol{\theta}}\|_2^2 \right)$$

Then,

$$\frac{1}{2} \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_t\|_{2,\bar{E}_t}^2 + \gamma \left(\|\hat{\boldsymbol{\theta}}_t - \bar{\boldsymbol{\theta}}\|_2^2 + \|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2 \right) \leq D_t + 2\gamma \|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2$$

By triangle inequality we have $\|\hat{\boldsymbol{\theta}}_t - \bar{\boldsymbol{\theta}}\|_2 + \|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2 \geq \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_t\|_2$. Taking squares on both sides, we obtain $\|\hat{\boldsymbol{\theta}}_t - \bar{\boldsymbol{\theta}}\|_2^2 + \|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2 \geq \frac{1}{2} \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_t\|_2^2$. Then, noting that $\|\boldsymbol{\Delta}\|_{2,E_t}^2 = \|\boldsymbol{\Delta}\|_{2,\bar{E}_t}^2 + \gamma \|\boldsymbol{\Delta}\|_2^2$, we have

$$\frac{1}{2} \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_t\|_{2,E_t}^2 \leq D_t + 2\gamma \|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2$$

Lastly, using the inequality $\|\boldsymbol{\theta}^* - \bar{\boldsymbol{\theta}}\|_2^2 \leq G^2$,

$$\|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_t\|_{2,E_t}^2 \leq 8\eta^2 \log(|\mathcal{F}^\alpha|/\delta) + 2\alpha t d \left[8B + \sqrt{8\eta^2 \log(4dt^2/\delta)} \right] + 4\gamma G^2$$

Taking the infimum over the size of α -covers, we obtain the final result. □

E Proofs for Regret Bounds

Throughout this section we will use the shorthand $\beta_t = \beta_t^*(\delta, \alpha, \gamma)$.

We start with the definitions of weighted inner product and norms.

Definition E.1. For a symmetric positive definite matrix $\mathbf{W} \in \mathbb{R}^{d \times d}$, define

- \mathbf{W} -inner product of two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ as $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{W}} := \langle \mathbf{W}\mathbf{x}, \mathbf{y} \rangle$
- \mathbf{W} -norm of a vector $\mathbf{x} \in \mathbb{R}^d$ as $\|\mathbf{x}\|_{\mathbf{W}} := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_{\mathbf{W}}}$.

Then, the regularized empirical 2-norm of a vector $\mathbf{z} \in \mathbb{R}^d$ can be written as

$$\|\mathbf{z}\|_{2,E_t} = \|\mathbf{z}\|_{\mathbf{A}_t} \tag{E.1}$$

where \mathbf{A}_t is a diagonal matrix with diagonal entries $\{(n_{t,1} + \gamma), (n_{t,2} + \gamma), \dots, (n_{t,d} + \gamma)\}$.

Recall that $n_{t,i} = \sum_{\tau=1}^{t-1} \mathbb{1}\{i \in \mathcal{A}_\tau\}$ denotes the number of times arm i has been pulled before time t (excluding time t).

Lemma E.2. Let $x \in \mathcal{X}$ and $\Theta \in \mathcal{F}_t$. Then,

$$|\langle \boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t^{LS}, \mathbf{x} \rangle| \leq w \sqrt{\beta_t} \tag{E.2}$$

where $w = \|\mathbf{x}\|_{\mathbf{A}_t^{-1}}$ is the "confidence width" of an action \mathbf{x} at time t .

Proof. Let $\Delta = \theta - \hat{\theta}_t^{\text{LS}}$. Then,

$$\begin{aligned}
|\langle \Delta, \mathbf{x} \rangle| &= |\Delta^\top \mathbf{x}| \\
&= |\Delta^\top \mathbf{A}_t^{1/2} \mathbf{A}_t^{-1/2} \mathbf{x}| \\
&= |(\mathbf{A}_t^{1/2} \Delta)^\top \mathbf{A}_t^{-1/2} \mathbf{x}| \\
&\leq \|\mathbf{A}_t^{1/2} \Delta\| \|\mathbf{A}_t^{-1/2} \mathbf{x}\| \\
&= \|\Delta\|_{\mathbf{A}_t} \|\mathbf{x}\|_{\mathbf{A}_t^{-1}} \\
&= w \|\Delta\|_{2, E_t} \\
&\leq w \sqrt{\beta_t}
\end{aligned}$$

□

Define the widths of the allocations as

$$w_t := \|\mathbf{x}_t\|_{\mathbf{A}_t^{-1}} \quad \text{and} \quad w_{t,i} := \|\mathbf{e}_i\|_{\mathbf{A}_t^{-1}} \quad (\text{E.3})$$

Lemma E.3. For any $t \in \mathbb{N}$, we have the identity

$$w_t^2 = \sum_{i \in \mathcal{A}_t} w_{t,i}^2$$

Proof.

$$\begin{aligned}
w_t^2 &= \langle \mathbf{x}_t, \mathbf{x}_t \rangle_{\mathbf{A}_t^{-1}} \\
&= \left\langle \mathbf{A}_t^{-1} \sum_{i \in \mathcal{A}_t} \mathbf{e}_i, \sum_{i \in \mathcal{A}_t} \mathbf{e}_i \right\rangle \\
&= \sum_{i \in \mathcal{A}_t} \sum_{j \in \mathcal{A}_t} \langle \mathbf{A}_t^{-1} \mathbf{e}_i, \mathbf{e}_j \rangle \\
&= \sum_{i \in \mathcal{A}_t} \langle \mathbf{A}_t^{-1} \mathbf{e}_i, \mathbf{e}_i \rangle \\
&= \sum_{i \in \mathcal{A}_t} w_{t,i}^2
\end{aligned}$$

where the penultimate step follows because $\langle \mathbf{A}_t^{-1} \mathbf{e}_i, \mathbf{e}_j \rangle = 0$ for $i \neq j$. □

Lemma E.4. Let the regret at time t be $r_t = \langle \mathbf{x}_t^*, \boldsymbol{\theta}^* \rangle - \langle \mathbf{x}_t, \boldsymbol{\theta}^* \rangle$. If $\boldsymbol{\theta}^* \in \mathcal{C}_t$, then

$$r_t \leq 2w_t \sqrt{\beta_t}$$

Proof. By the choice of $(\mathbf{x}_t, \tilde{\boldsymbol{\theta}}_t)$, we have

$$\langle \mathbf{x}_t, \tilde{\boldsymbol{\theta}}_t \rangle = \max_{(\mathbf{x}, \boldsymbol{\theta}) \in \mathcal{X}_t \times \mathcal{C}_t} \langle \mathbf{x}, \boldsymbol{\theta} \rangle \geq \langle \mathbf{x}_t^*, \boldsymbol{\theta}^* \rangle$$

where the inequality uses $\boldsymbol{\theta}^* \in \mathcal{C}_t$. Hence,

$$\begin{aligned}
r_t &= \langle \mathbf{x}_t^*, \boldsymbol{\theta}^* \rangle - \langle \mathbf{x}_t, \boldsymbol{\theta}^* \rangle \\
&\leq \langle \mathbf{x}_t, \tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^* \rangle \\
&= \langle \mathbf{x}_t, \tilde{\boldsymbol{\theta}}_t - \hat{\boldsymbol{\theta}}_t^{\text{LS}} \rangle + \langle \mathbf{x}_t, \hat{\boldsymbol{\theta}}_t^{\text{LS}} - \boldsymbol{\theta}^* \rangle \\
&\leq 2w_t \sqrt{\beta_t}
\end{aligned}$$

where the last step follows from Lemma E.2. □

Next, we show that the confidence widths do not grow too fast.

Lemma E.5. For every t ,

$$\log \det \mathbf{A}_{t+1} = d \log \gamma + \sum_{\tau=1}^t \sum_{i \in \mathcal{A}_\tau} \log(1 + w_{\tau,i}^2) \quad (\text{E.4})$$

Proof. By the definition of \mathbf{A}_t , we have

$$\begin{aligned} \det \mathbf{A}_{t+1} &= \det \left(\mathbf{A}_t + \sum_{i \in \mathcal{A}_t} \mathbf{e}_i \mathbf{e}_i^\top \right) \\ &= \det \left(\mathbf{A}_t^{1/2} \left(\mathbf{I} + \mathbf{A}_t^{-1/2} \left(\sum_{i \in \mathcal{A}_t} \mathbf{e}_i \mathbf{e}_i^\top \right) \mathbf{A}_t^{-1/2} \right) \mathbf{A}_t^{1/2} \right) \\ &= \det(\mathbf{A}_t) \det \left(\mathbf{I} + \sum_{i \in \mathcal{A}_t} \mathbf{A}_t^{-1/2} \mathbf{e}_i \mathbf{e}_i^\top \mathbf{A}_t^{-1/2} \right) \end{aligned}$$

Each $\mathbf{A}_t^{-1/2} \mathbf{e}_i \mathbf{e}_i^\top \mathbf{A}_t^{-1/2}$ term has zeros everywhere except one entry on the diagonal and that non-zero entry is equal to $w_{t,i}^2$. Furthermore, the location of the non-zero entry is different in for each term. Hence,

$$\det \left(\mathbf{I} + \sum_{i \in \mathcal{A}_t} \mathbf{A}_t^{-1/2} \mathbf{e}_i \mathbf{e}_i^\top \mathbf{A}_t^{-1/2} \right) = \prod_{i \in \mathcal{A}_t} (1 + w_{t,i}^2)$$

Therefore, we have

$$\log \det \mathbf{A}_{t+1} = \log \det \mathbf{A}_t + \sum_{i \in \mathcal{A}_t} \log(1 + w_{t,i}^2)$$

Since $\mathbf{A}_1 = \gamma \mathbf{I}$, we have $\log \det \mathbf{A}_1 = d \log \gamma$ and the result follows by induction. \square

Lemma E.6. For all t , $\log \det \mathbf{A}_t \leq d \log(t + \gamma - 1)$.

Proof. Noting that \mathbf{A}_t is a diagonal matrix with diagonals $(n_{t,i} + \gamma)$,

$$\begin{aligned} \text{trace } \mathbf{A}_t &= d\gamma + \sum_{i=1}^d n_{t,i} \\ &= d\gamma + d(t-1) \\ &= d(t + \gamma - 1) \end{aligned}$$

Now, recall that $\text{trace } \mathbf{A}_t$ equals the sum of the eigenvalues of \mathbf{A}_t . On the other hand, $\det(\mathbf{A}_t)$ equals the product of the eigenvalues. Since \mathbf{A}_t is positive definite, its eigenvalues are all positive. Subject to these constraints, $\det(\mathbf{A}_t)$ is maximized when all the eigenvalues are equal; the desired bound follows. \square

Lemma E.7. Let $\gamma \geq 1$. Then, for all t , we have

$$\sum_{\tau=1}^t \sum_{i \in \mathcal{A}_\tau} w_{\tau,i}^2 \leq 2d \log \left(1 + \frac{t}{\gamma} \right)$$

Proof. Note that $0 \leq w_{\tau,i}^2 \leq 1$, if $\gamma \geq 1$. Using the inequality $y \leq 2 \log(1 + y)$ for $0 \leq y \leq 1$, we have

$$\begin{aligned} \sum_{\tau=1}^t \sum_{i \in \mathcal{A}_\tau} w_{\tau,i}^2 &\leq 2 \sum_{\tau=1}^t \sum_{i \in \mathcal{A}_\tau} \log(1 + w_{\tau,i}^2) \\ &= 2 \log \det \mathbf{A}_{t+1} - 2d \log \gamma \\ &\leq 2d \log \left(1 + \frac{t}{\gamma} \right) \end{aligned}$$

where the last two lines follow from Lemmas E.5 and E.6 respectively. \square

Lemma E.8. Let the instantaneous regret at time t be $r_t = \langle \mathbf{x}_t^*, \boldsymbol{\theta}^* \rangle - \langle \mathbf{x}_t, \boldsymbol{\theta}^* \rangle$. If $\gamma \geq 1$ and $\boldsymbol{\theta}^* \in \mathcal{C}_t$ for all $t \leq T$, then

$$\sum_{t=1}^T r_t^2 \leq 8\beta_T d \log \left(1 + \frac{T}{\gamma} \right)$$

Proof. Assuming that $\boldsymbol{\theta}^* \in \mathcal{C}_t$ for all $t \leq T$,

$$\begin{aligned} \sum_{t=1}^T r_t^2 &\leq \sum_{t=1}^T 4w_t^2 \beta_t \\ &\leq 4\beta_T \sum_{t=1}^T w_t^2 \\ &= 4\beta_T \sum_{t=1}^T \sum_{i \in \mathcal{A}_t} w_{t,i}^2 \\ &\leq 8\beta_T d \log \left(1 + \frac{T}{\gamma} \right) \end{aligned}$$

where the first step follows from Lemma E.4, second step follows from the definition of β_t , third step uses the identity given in Lemma E.3 and the last step is due to Lemma E.7. \square

Theorem E.9. Let $\gamma \geq 1$. Then, with probability at least $1 - 2\delta$, the T period regret is bounded by

$$\mathcal{R}(T, \pi) \leq \sqrt{8d\beta_T^*(\delta, \alpha, \gamma)T \log \left(1 + \frac{T}{\gamma} \right)}$$

where

$$\beta_T^*(\delta, \alpha, \gamma) = 8\eta^2 \log(\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_2)/\delta) + 2\alpha dT \left(8B + \sqrt{8\eta^2 \log(4dT^2/\delta)} \right) + 4\gamma G^2 \quad (\text{E.5})$$

Proof. Assuming that $\boldsymbol{\theta}^* \in \mathcal{C}_t$ for all $t \leq T$,

$$\begin{aligned} \mathcal{R}(T, \pi) &= \sum_{t=1}^T r_t \\ &\leq \left(T \sum_{t=1}^T r_t^2 \right)^{1/2} \\ &\leq \left(8d\beta_T T \log \left(1 + \frac{T}{\gamma} \right) \right)^{1/2} \end{aligned}$$

where the last step follows from Lemma E.8. Then, by Lemma D.5, $\boldsymbol{\theta}^* \in \mathcal{C}_t$ for all $t \leq T$ with probability at least $1 - 2\delta$. Therefore, the bound holds true with probability at least $1 - 2\delta$. \square

Corollary E.10. Letting $\delta = \mathcal{O}((dT)^{-1})$, $\alpha = \mathcal{O}((dT)^{-1})$ and $\gamma = 1$ results in a regret bound that satisfies

$$\mathcal{R}(T, \pi) = \tilde{\mathcal{O}} \left(\sqrt{\eta^2 d \log(\mathcal{N}(\mathcal{F}, T^{-1}, \|\cdot\|_2)) T} \right) \quad (\text{E.6})$$

Proof. By Theorem E.9, with probability 1,

$$\mathcal{R}(T, \pi) \leq (1 - \delta) \sqrt{8d\beta_T^*(\delta, \alpha, \gamma)T \log \left(1 + \frac{T}{\gamma} \right)} + 2\delta BdT$$

Letting $\delta = \mathcal{O}(T^{-1})$, $\alpha = \mathcal{O} = (T^{-1})$ and $\gamma = 1$,

$$\mathcal{R}(T, \pi) = \tilde{\mathcal{O}} \left(\sqrt{d\beta_T^*(T, T^{-1}, 1)T} \right)$$

Noting that $\beta_T^*(T, T^{-1}, 1) = \tilde{\mathcal{O}}(\eta^2 \log(\mathcal{N}(\mathcal{F}, T^{-1}, \|\cdot\|_2)))$, the proof is complete. \square

F Proofs for OFU-based Allocations

F.1 Proof of Theorem 3.2

In the allocation problem, the mean reward of the arms are given in the matrix $\Theta \in \mathbb{R}^{N \times M}$. Consider setting $\theta = \text{vec}(\Theta)$ as the mean reward vector for the problem described in Appendix Section C. Noting that $d = NM$ and $\|\cdot\|_F = \|\text{vec}(\cdot)\|_2$, the proof becomes a direct extension of Theorem E.9.

F.2 Proof of Theorem 3.3

The proof is direct extension of Corollary E.10 where the covering number is replaced with the following upper bound given for the choice of \mathcal{L} defined in equation (3.4). We provide the upper bound for the covering number of \mathcal{L} in the following lemma.

Lemma F.1 (Covering number for low-rank matrices). *The covering number of \mathcal{L} given in (3.4) obeys*

$$\log \mathcal{N}(\mathcal{L}, \alpha, \|\cdot\|_F) \leq (N + M + 1)R \log \left(\frac{9B\sqrt{NM}}{\alpha} \right) \quad (\text{F.1})$$

Proof. This proof is modified from [27]. Let $\mathcal{S} = \{\Theta \in \mathbb{R}^{N \times M} : \text{rank}(\Theta) \leq R, \|\Theta\|_F \leq 1\}$. We will first show that there exists an ϵ -net \mathcal{S}^ϵ for the Frobenious norm obeying

$$|\mathcal{S}^\epsilon| \leq (9/\epsilon)^{(N+M+1)R}$$

For any $\Theta \in \mathcal{S}$, singular value decomposition gives $\Theta = \mathbf{U}\Sigma\mathbf{V}^T$, where $\|\Sigma\|_F \leq 1$. We will construct an ϵ -net for \mathcal{S} by covering the set of permissible \mathbf{U} , Σ and \mathbf{V} . Let \mathcal{D} be the set of diagonal matrices with nonnegative diagonal entries and Frobenious norm less than or equal to one. We take $\mathcal{D}^{\epsilon/3}$ be an $\epsilon/3$ -net for \mathcal{D} with $|\mathcal{D}^{\epsilon/3}| \leq (9/\epsilon)^R$. Next, let $\mathcal{O}_{N,R} = \{\mathbf{U} \in \mathbb{R}^{N \times R} : \mathbf{U}^T \mathbf{U} = \mathbf{I}\}$. To cover $\mathcal{O}_{N,R}$, we use the $\|\cdot\|_{1,2}$ norm defined as

$$\|\mathbf{U}\|_{1,2} = \max_i \|\mathbf{u}_i\|_{\ell_2}$$

where \mathbf{u}_i denotes the i th column of \mathbf{U} . Let $\mathcal{Q}_{N,R} = \{\mathbf{U} \in \mathbb{R}^{N,R} : \|\mathbf{U}\|_{1,2} \leq 1\}$. It is easy to see that $\mathcal{O}_{N,R} \subset \mathcal{Q}_{N,R}$ since the columns of an orthogonal matrix are unit normed. We see that there is an $\epsilon/3$ -net $\mathcal{O}_{N,R}^{\epsilon/3}$ for $\mathcal{O}_{N,R}$ obeying $|\mathcal{O}_{N,R}^{\epsilon/3}| \leq (9/\epsilon)^{NR}$. Similarly, let $\mathcal{P}_{M,R} = \{\mathbf{V} \in \mathbb{R}^{M \times R} : \mathbf{V}^T \mathbf{V} = \mathbf{I}\}$. By the same argument, there is an $\epsilon/3$ -net $\mathcal{P}_{M,R}^{\epsilon/3}$ for $\mathcal{P}_{M,R}$ obeying $|\mathcal{P}_{M,R}^{\epsilon/3}| \leq (9/\epsilon)^{MR}$. We now let $\mathcal{S}^\epsilon = \{\bar{\mathbf{U}}\bar{\Sigma}\bar{\mathbf{V}}^T : \bar{\mathbf{U}} \in \mathcal{O}_{N,R}^{\epsilon/3}, \bar{\mathbf{V}} \in \mathcal{P}_{M,R}^{\epsilon/3}, \bar{\Sigma} \in \mathcal{D}^{\epsilon/3}\}$, and remark $|\mathcal{S}^\epsilon| \leq |\mathcal{O}_{N,R}^{\epsilon/3}| |\mathcal{P}_{M,R}^{\epsilon/3}| |\mathcal{D}^{\epsilon/3}| \leq (9/\epsilon)^{(N+M+1)R}$. It remains to show that for all $\Theta \in \mathcal{S}$, there exists $\bar{\Theta} \in \mathcal{S}^\epsilon$ with $\|\Theta - \bar{\Theta}\|_F \leq \epsilon$.

Fix $\Theta \in \mathcal{S}$ and decompose it as $\Theta = \mathbf{U}\Sigma\mathbf{V}^T$. Then, there exists $\bar{\Theta} = \bar{\mathbf{U}}\bar{\Sigma}\bar{\mathbf{V}}^T \in \mathcal{S}^\epsilon$ with $\bar{\mathbf{U}} \in \mathcal{O}_{N,R}^{\epsilon/3}$, $\bar{\mathbf{V}} \in \mathcal{P}_{M,R}^{\epsilon/3}$, $\bar{\Sigma} \in \mathcal{D}^{\epsilon/3}$ satisfying $\|\mathbf{U} - \bar{\mathbf{U}}\|_{1,2} \leq \epsilon/3$, $\|\mathbf{V} - \bar{\mathbf{V}}\|_{1,2} \leq \epsilon/3$ and $\|\Sigma - \bar{\Sigma}\|_F \leq \epsilon/3$. This gives

$$\begin{aligned} \|\Theta - \bar{\Theta}\|_F &= \|\mathbf{U}\Sigma\mathbf{V}^T - \bar{\mathbf{U}}\bar{\Sigma}\bar{\mathbf{V}}^T\|_F \\ &= \|\mathbf{U}\Sigma\mathbf{V}^T - \bar{\mathbf{U}}\Sigma\mathbf{V}^T + \bar{\mathbf{U}}\Sigma\mathbf{V}^T - \bar{\mathbf{U}}\bar{\Sigma}\mathbf{V}^T + \bar{\mathbf{U}}\bar{\Sigma}\mathbf{V}^T - \bar{\mathbf{U}}\bar{\Sigma}\bar{\mathbf{V}}^T\|_F \\ &\leq \|(\mathbf{U} - \bar{\mathbf{U}})\Sigma\mathbf{V}^T\|_F + \|\bar{\mathbf{U}}(\Sigma - \bar{\Sigma})\mathbf{V}^T\|_F + \|\bar{\mathbf{U}}\bar{\Sigma}(\mathbf{V} - \bar{\mathbf{V}})^T\|_F \end{aligned}$$

For the first term, since \mathbf{V} is an orthogonal matrix,

$$\begin{aligned} \|(\mathbf{U} - \bar{\mathbf{U}})\boldsymbol{\Sigma}\mathbf{V}^T\|_F^2 &= \|(\mathbf{U} - \bar{\mathbf{U}})\boldsymbol{\Sigma}\|_F^2 \\ &\leq \|\boldsymbol{\Sigma}\|_F^2 \|\mathbf{U} - \bar{\mathbf{U}}\|_{1,2}^2 \leq (\epsilon/3)^2 \end{aligned}$$

By the same argument, $\|\bar{\mathbf{U}}\bar{\boldsymbol{\Sigma}}(\mathbf{V} - \bar{\mathbf{V}})^T\|_F \leq \epsilon/3$ as well. Lastly, $\|\bar{\mathbf{U}}(\boldsymbol{\Sigma} - \bar{\boldsymbol{\Sigma}})\mathbf{V}^T\|_F = \|\boldsymbol{\Sigma} - \bar{\boldsymbol{\Sigma}}\|_F \leq \epsilon/3$. Therefore, $\|\boldsymbol{\Theta} - \bar{\boldsymbol{\Theta}}\|_F \leq \epsilon$, showing that \mathcal{S}^ϵ is an ϵ -net for \mathcal{S} with respect to the Frobenious norm.

Next, we will construct an α -net for \mathcal{L} given in equation 3.4. Let $\kappa = B\sqrt{NM}$. We start by noting that for all $\boldsymbol{\Theta} \in \mathcal{L}$, the Frobenious norm obeys $\|\boldsymbol{\Theta}\|_F \leq \kappa$. Then, define $\mathbf{X} = \frac{1}{\kappa}\boldsymbol{\Theta} \in \mathcal{S}$ and $\mathcal{L}^\alpha := \{\kappa\bar{\mathbf{X}} : \bar{\mathbf{X}} \in \mathcal{S}^\epsilon\}$. We previously showed that for any $\mathbf{X} \in \mathcal{S}$, there exists $\bar{\mathbf{X}} \in \mathcal{S}^\epsilon$ such that $\|\mathbf{X} - \bar{\mathbf{X}}\|_F \leq \epsilon$. Therefore, for any $\boldsymbol{\Theta} \in \mathcal{L}$, there exists $\bar{\boldsymbol{\Theta}} = \kappa\bar{\mathbf{X}} \in \mathcal{L}^\alpha$ such that $\|\boldsymbol{\Theta} - \bar{\boldsymbol{\Theta}}\|_F \leq \kappa\epsilon$. Setting $\epsilon = \alpha/\kappa$, we obtain that \mathcal{L}^α is an α -net for \mathcal{L} with respect to the Frobenious norm. Finally, the size of \mathcal{L}^α obeys

$$|\mathcal{L}^\alpha| = |\mathcal{S}^{\alpha/\kappa}| \leq (9\kappa/\alpha)^{(N+M+1)R}$$

This completes the proof. \square

G Additional Experimental Details

All experiments are implemented in Python and carried out on a machine with 2.3GHz 8-core Intel Core i9 CPU and 16GB of RAM. We solve the allocation integer program (2.1) using large-scale mixed integer programming (MIP) solver packages to have efficient computations.

Parameter setup:

- In synthetic data, $\boldsymbol{\Theta}^*$ is scaled such that $B = 10$.
- Standard deviation of the rewards: $\eta = 1$.
- In the static setting, $d_{t,u} = 1$ for all $u \in [N]$.
- In the dynamic setting, $d_{t,u} = 1$ with probability 0.2, 0 otherwise, independently for each $u \in [N]$.
- $C_{\max} = \text{ceil}\left(\frac{3}{M} \sum_{u=1}^N d_{t,u}\right)$.
- $c_{t,i}$ are uniformly sampled over $\{1, \dots, C_{\max}\}$ independently for each $t \in [T]$ and $i \in [M]$.

To complement our discussion on the importance of capacity-aware recommendations, Figure 3 shows the effects of having stricter capacity constraints. When C_{\max} is low (the capacity constraints are stricter), we see that the performance gap between our proposed method and ICF2 is larger. Therefore, it is more crucial to employ capacity-aware mechanisms in settings with tighter capacity constraints.

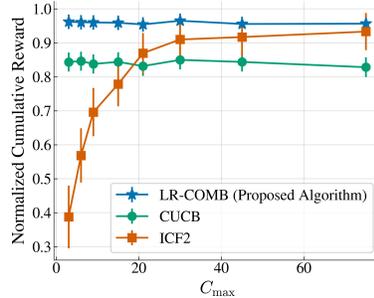


Figure 3: Normalized cumulative reward obtained in $T = 300$ rounds for different choices of C_{\max} (normalized by the cumulative reward of optimal allocations). Synthetic data in a static setting with $N = 400$, $M = 200$, $R = 10$. For each data point, the experiments are run on 20 problem instances and means are reported together with error regions that indicate one standard deviation of uncertainty.

In Figures 4, 5, 6 and 7, we provide detailed results for different experimental settings described in Section 4. Reward indicates the instantaneous reward obtained in each iteration, regret is the gap between the reward of the optimum allocation and the allocation achieved by the algorithm. Cumulative regret (defined in (2.3)) is the cumulative sum of instantaneous regrets up to iteration t . The average cumulative regret is obtained by normalizing the cumulative regret with $1/t$.

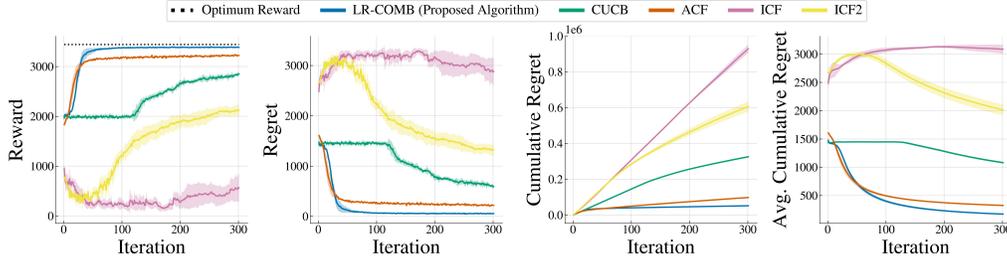


Figure 4: Experimental results for synthetic data in a static setting with $N = 800$, $M = 400$, $R = 20$. The experiments are run on 10 problem instances and means are reported together with error regions that indicate one standard deviation of uncertainty.

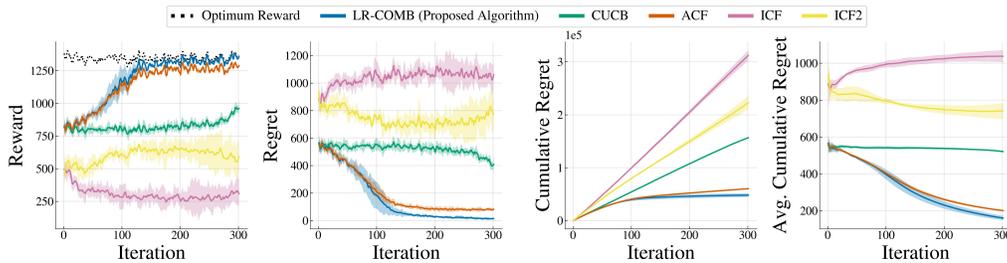


Figure 5: Experimental results for synthetic data in a dynamic setting with $N = 1000$, $M = 150$, $R = 20$, probability of activity 0.2. The experiments are run on 10 problem instances and means are reported together with error regions that indicate one standard deviation of uncertainty.

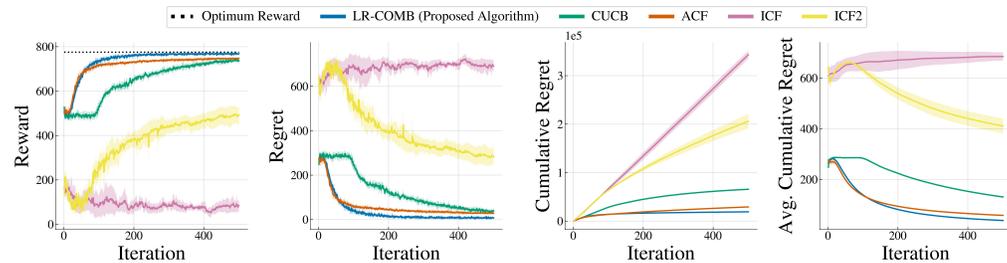


Figure 6: Experimental results for Restaurant-Customer data in a static setting. The experiments are run on 10 problem instances and means are reported together with error regions that indicate one standard deviation of uncertainty.

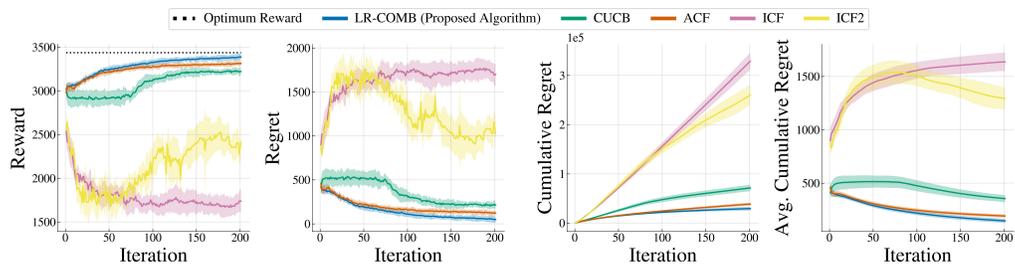


Figure 7: Experimental results for MovieLens 100k data in a static setting. The experiments are run on 10 problem instances and means are reported together with error regions that indicate one standard deviation of uncertainty.