
Online Clustering of Bandits with Misspecified User Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

The contextual linear bandit is an important online learning problem where given arm features, a learning agent selects an arm at each round to maximize the cumulative rewards in the long run. A line of works, called the clustering of bandits (CB), utilize the collaborative effect over user preferences and have shown significant improvements over classic linear bandit algorithms. However, existing CB algorithms require well-specified linear user models and can fail when this critical assumption does not hold. Whether robust CB algorithms can be designed for more practical scenarios with misspecified user models remains an open problem. In this paper, we are the first to present the important problem of clustering of bandits with misspecified user models (CBMUM), where the expected rewards in user models can be perturbed away from perfect linear models. We devise two robust CB algorithms, RCLUMB and RSCLUMB (representing the learned clustering structure with dynamic graph and sets, respectively), that can accommodate the inaccurate user preference estimations and erroneous clustering caused by model misspecifications. We prove regret upper bounds of $O(\epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T)$ for our algorithms under milder assumptions than previous CB works, which match the lower bound asymptotically in T up to logarithmic factors, and also match the state-of-the-art results in several degenerate cases. Our regret analysis is novel and different from the typical proof flow of previous CB works. The techniques in proving the regret caused by misclustering users are quite general and may be of independent interest. Experiments on both synthetic and real-world data show our outperformance over previous algorithms.

1 Introduction

Stochastic multi-armed bandit (MAB) [2, 4, 20] is an online sequential decision-making problem, where the learning agent selects an action and receives a corresponding reward at each round, so as to maximize the cumulative reward in the long run. MAB algorithms have been widely applied in recommendation systems to handle the exploration and exploitation trade-off [19, 36].

To deal with large-scale applications, the contextual linear bandits [22, 8, 1] have been studied, where the expected reward of each arm is assumed to be perfectly linear in their features. Leveraging the contextual side information about the user and arms, linear bandits can provide more personalized recommendations [15]. Classical linear bandit approaches, however, ignore the often useful tool of collaborative filtering. To utilize the relationships among users, the problem of clustering of bandits (CB) has been proposed [11]. Specifically, CB algorithms adaptively partition users into clusters and utilize the collaborative effect of users to enhance learning performance.

Although existing CB algorithms have shown great success in improving recommendation qualities, there exist two major limitations. First, all previous works on CB [11, 23, 25] assume that for each user, the expected rewards follow a *perfectly linear* model with respect to the user preference vector and arms' feature vectors. In many real-world scenarios, due to feature noises or uncertainty [14], the reward may not necessarily conform to a perfectly linear function, or even deviates a lot from linearity

[13]. Second, previous CB works assume that for users within the same cluster, their preferences are exactly the same. Due to the heterogeneity in users' personalities and interests, similar users may not have identical preferences, invalidating this strong assumption.

To address these issues, we propose a novel problem of clustering of bandits with misspecified user models (CBMUM). In CBMUM, the expected reward model of each user does not follow a perfectly linear function but with possible additive deviations. We assume users in the same underlying cluster share a common preference vector, meaning they have the same linear part in reward models, but the deviation parts are allowed to be different, better reflecting the varieties of user personalities.

The relaxation of perfect linearity and the reward homogeneity within the same cluster bring many challenges to the CBMUM problem. In CBMUM, we not only need to handle the uncertainty from the *unknown* user preference vectors, but also have to tackle the additional uncertainty from model misspecifications. Due to such uncertainties, it becomes highly challenging to design a robust algorithm that can cluster the users appropriately and utilize the clustered information judiciously. On the one hand, the algorithm needs to be more tolerant in the face of misspecifications so that more similar users can be clustered together to utilize the collaborative effect. On the other hand, it has to be more selective to rule out the possibility of *misclustering* users with large preference gaps.

1.1 Our Contributions

This paper makes the following four contributions.

New Model Formulation. We are the first to formulate the clustering of bandits with misspecified user models (CBMUM) problem, which is more practical by removing the perfect linearity assumption in previous CB works.

Novel Algorithm Designs. We design two novel algorithms, RCLUMB and RSCLUMB, which robustly learn the clustering structure and utilize this collaborative information for faster user preference elicitation. Specifically, RCLUMB keeps updating a dynamic graph over all users, where users connected directly by edges are supposed to be in the same cluster. RCLUMB adaptively removes edges and recommends items based on historical interactions. RSCLUMB represents the clustering structure with sets, which are dynamically merged and split during the learning process. Due to the page limit, we only illustrate the RCLUMB algorithm in the main paper. We leave the exposition, illustration, and regret analysis of the RSCLUMB algorithm in Appendix K.

To overcome the challenges brought by model misspecifications, we do the following key steps in the RCLUMB algorithm. (i) To ensure that with high probability, similar users will not be partitioned apart, we design a more tolerant edge deletion rule by taking model misspecifications into consideration. (ii) Due to inaccurate user preference estimations caused by model misspecifications, trivially following previous CB works [11, 23, 26] to directly use connected components in the maintained graph as clusters would *miscluster* users with big preference gaps, causing a large regret. To be discriminative in cluster assignments, we filter users directly linked with the current user in the graph to form the cluster used in this round. With these careful designs of (i) and (ii), we can guarantee that with high probability, information of all similar users can be leveraged, and only users with close enough preferences might be *misclustered*, which will only mildly impair the learning accuracy. Additionally: (iii) we design an enlarged confidence radius to incorporate both the exploration bonus and the additional uncertainty from misspecifications when recommending arms. The design of RSCLUMB follows similar ideas, which we leave in the Appendix K due to page limit.

Theoretical Analysis with Milder Assumptions. We prove regret upper bounds for our algorithms of $O(\epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T)$ in CBMUM under much milder and practical assumptions (in arm generation distribution) than previous CB works, which match the state-of-the-art results in degenerate cases. Our proof is quite different from the typical proof flow of previous CB works (details in Appendix C). One key challenge is to bound the regret caused by *misclustering* users with close but not the same preference vectors and use the inaccurate cluster-based information to recommend arms. To handle the challenge, we prove a key lemma (Lemma 5.7) to bound this part of regret. We defer its details in Section 5 and Appendix G. The techniques and results for bounding this part are quite general and may be of independent interest. We also give a regret lower bound of $\Omega(\epsilon_* T \sqrt{d})$ for CBMUM, showing that our upper bounds are asymptotically tight with respect to T up to logarithmic factors. We leave proving a tighter lower bound for CBMUM as an open problem.

Good Experimental Performance. Extensive experiments on both synthetic and real-world data show the advantages of our proposed algorithms over the existing algorithms.

2 Related Work

Our work is closely related to two lines of research: online clustering of bandits (CB) and misspecified linear bandits (MLB). More discussions on related works can be found in Appendix A.

The paper [11] first formulates the CB problem and proposes a graph-based algorithm. The work [24] further considers leveraging the collaborative effects on items to guide the clustering of users. The work [23] considers the CB problem in the cascading bandits setting with random prefix feedback. The paper [25] also considers users with different arrival frequencies. A recent work [26] proposes the setting of clustering of federated bandits, considering both privacy protection and communication requirements. However, all these works assume that the reward model for each user follows a perfectly linear model, which is unrealistic in many real-world applications. To the best of our knowledge, this paper is the first work to consider user model misspecifications in the CB problem.

The work [13] first proposes the misspecified linear bandits (MLB) problem, shows the vulnerability of linear bandit algorithms under deviations, and designs an algorithm RLB that is only robust to non-sparse deviations. The work [21] proposes two algorithms to handle general deviations, which are modifications of the phased elimination algorithm [20] and LinUCB [1]. Some recent works [27, 10] use model selection methods to deal with unknown exact maximum model misspecification level. Note that the work [10] has an additional assumption on the access to an online regression oracle, and the paper [27] still needs to know an upper bound of the unknown exact maximum model deviation level. None of them consider the CB setting with multiple users, thus differing from ours.

We are the first to initialize the study of the important CBMUM problem, and propose a general framework for dealing with model misspecifications in CB problems. Our study is based on fundamental models on CB [11, 25] and MLB [21], the algorithm design ideas and theoretical analysis are pretty general. We leave incorporating the model selection methods [27, 10] into our framework to address the unknown exact maximum model misspecification level as an interesting future work.

3 Problem Setup

This section formulates the problem of “clustering of bandits with misspecified user models” (CB-MUM). We use boldface **lowercase** and boldface **CAPITALIZED** letters for vectors and matrices. We use $|\mathcal{A}|$ to denote the number of elements in \mathcal{A} , $[m]$ to denote $\{1, \dots, m\}$, and $\|x\|_M = \sqrt{x^\top M x}$ to denote the matrix norm of vector x regarding the positive semi-definite (PSD) matrix M .

In CBMUM, there are u users denoted by $\mathcal{U} = \{1, 2, \dots, u\}$. Each user $i \in \mathcal{U}$ is associated with an *unknown* preference vector $\theta_i \in \mathbb{R}^d$, with $\|\theta_i\|_2 \leq 1$. We assume there is an *unknown* underlying clustering structure over users representing the similarity of their behaviors. Specifically, \mathcal{U} can be partitioned into a small number m (i.e., $m \ll u$) clusters, V_1, V_2, \dots, V_m , where $\cup_{j \in [m]} V_j = \mathcal{U}$, and $V_j \cap V_{j'} = \emptyset$, for $j \neq j'$. We call these clusters *ground-truth clusters* and use $\mathcal{V} = \{V_1, V_2, \dots, V_m\}$ to denote the set of these clusters. Users in the same *ground-truth cluster* share the same preference vector, while users from different *ground-truth clusters* have different preference vectors. Let θ^j denote the common preference vector for V_j and $j(i) \in [m]$ denote the index of the *ground-truth cluster* that user i belongs to. For any $\ell \in \mathcal{U}$, if $\ell \in V_{j(i)}$, then $\theta_\ell = \theta_i = \theta^{j(i)}$.

At each round $t \in [T]$, a user $i_t \in \mathcal{U}$ comes to be served. The learning agent receives a finite arm set $\mathcal{A}_t \subseteq \mathcal{A}$ to choose from (with $|\mathcal{A}_t| \leq C, \forall t$), where each arm $a \in \mathcal{A}$ is associated with a feature vector $x_a \in \mathbb{R}^d$, and $\|x_a\|_2 \leq 1$. The agent assigns an appropriate cluster \bar{V}_t for user i_t and recommends an item $a_t \in \mathcal{A}_t$ based on the aggregated historical information gathered from cluster \bar{V}_t . After receiving the recommended item a_t , user i_t gives a random reward $r_t \in [0, 1]$ to the agent. To better model real-world scenarios, we assume that the reward r_t follows a misspecified linear function of the item feature vector x_{a_t} and the *unknown* user preference vector θ_{i_t} . Formally,

$$r_t = x_{a_t}^\top \theta_{i_t} + \epsilon_{a_t}^{i_t, t} + \eta_t, \quad (1)$$

where $\epsilon^{i_t, t} = [\epsilon_1^{i_t, t}, \epsilon_2^{i_t, t}, \dots, \epsilon_{|\mathcal{A}_t|}^{i_t, t}]^\top \in \mathbb{R}^{|\mathcal{A}_t|}$ denotes the *unknown* deviation in the expected rewards of arms in \mathcal{A}_t from linearity for user i_t at t , and η_t is the 1-sub-Gaussian noise. We allow the deviation vectors for users in the same *ground-truth cluster* to be different.

We assume the clusters, users, items, and model misspecifications satisfy the following assumptions.

Assumption 3.1 (Gap between different clusters). The gap between any two preference vectors for different *ground-truth clusters* is at least an *unknown* positive constant γ

$$\|\theta^j - \theta^{j'}\|_2 \geq \gamma > 0, \forall j, j' \in [m], j \neq j'.$$

146 **Assumption 3.2** (Uniform arrival of users). At each round t , a user i_t comes uniformly at random
 147 from \mathcal{U} with probability $1/u$, independent of the past rounds.

148 **Assumption 3.3** (Item regularity). At each time step t , the feature vector \mathbf{x}_a of each arm $a \in \mathcal{A}_t$
 149 is drawn independently from a fixed but unknown distribution ρ over $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 \leq 1\}$, where
 150 $\mathbb{E}_{\mathbf{x} \sim \rho}[\mathbf{x}\mathbf{x}^\top]$ is full rank with minimal eigenvalue $\lambda_x > 0$. Additionally, at any time t , for any fixed
 151 unit vector $\boldsymbol{\theta} \in \mathbb{R}^d$, $(\boldsymbol{\theta}^\top \mathbf{x})^2$ has sub-Gaussian tail with variance upper bounded by σ^2 .

152 **Assumption 3.4** (Bounded misspecification level). We assume that there is a pre-specified maximum
 153 misspecification level parameter ϵ_* such that $\|\epsilon^{i,t}\|_\infty \leq \epsilon_*$, $\forall i \in \mathcal{U}, t \in [T]$.

154 **Remark 1.** All these assumptions basically follow previous works on CB [11, 12, 23, 3, 26] and MLB
 155 [21]. Note that Assumption 3.3 is less stringent and more practical than previous CB works which also
 156 put restrictions on the variance upper bound σ^2 . For Assumption 3.2, our results can easily generalize
 157 to the case where the user arrival follows any distributions with minimum arrival probability greater
 158 than p_{min} . For Assumption 3.4, note that ϵ_* can be an upper bound on the maximum misspecification
 159 level, not the exact maximum itself. In real-world applications, the deviations are usually small [13],
 160 and we can set a relatively big ϵ_* as an upper bound. For more discussions please refer to Appendix B

161 Let $a_t^* \in \arg \max_{a \in \mathcal{A}_t} \mathbf{x}_a^\top \boldsymbol{\theta}_{i_t} + \epsilon_a^{i_t,t}$ denote an optimal arm which gives the highest expected reward
 162 at t . The goal of the agent is to minimize the expected cumulative regret

$$R(T) = \mathbb{E}[\sum_{t=1}^T (\mathbf{x}_{a_t^*}^\top \boldsymbol{\theta}_{i_t} + \epsilon_{a_t^*}^{i_t,t} - \mathbf{x}_{a_t}^\top \boldsymbol{\theta}_{i_t} - \epsilon_{a_t}^{i_t,t})]. \quad (2)$$

163 4 Algorithm

Algorithm 1 Robust Clustering of Misspecified Bandits Algorithm (RCLUMB)

- 1: **Input:** Deletion parameter $\alpha_1, \alpha_2 > 0$, $f(T) = \sqrt{\frac{1+\ln(1+T)}{1+T}}$, $\lambda, \beta, \epsilon_* > 0$.
- 2: **Initialization:** $\mathbf{M}_{i,0} = 0_{d \times d}$, $\mathbf{b}_{i,0} = 0_{d \times 1}$, $T_{i,0} = 0$, $\forall i \in \mathcal{U}$; a complete Graph $G_0 = (\mathcal{U}, E_0)$ over \mathcal{U} .
- 3: **for all** $t = 1, 2, \dots, T$ **do**
- 4: Receive the index of the current user $i_t \in \mathcal{U}$, and the current feasible arm set \mathcal{A}_t ;
- 5: Filter user i_t and users $i \in \mathcal{U}$ that are *directly* connected with user i_t via edge $(i, i_t) \in E_{t-1}$, to form the cluster \bar{V}_t ;
- 6: Compute the estimated statistics for cluster \bar{V}_t

$$\bar{\mathbf{M}}_{\bar{V}_t, t-1} = \lambda \mathbf{I} + \sum_{i \in \bar{V}_t} \mathbf{M}_{i, t-1}, \bar{\mathbf{b}}_{\bar{V}_t, t-1} = \sum_{i \in \bar{V}_t} \mathbf{b}_{i, t-1}, \hat{\boldsymbol{\theta}}_{\bar{V}_t, t-1} = \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \bar{\mathbf{b}}_{\bar{V}_t, t-1};$$
- 7: Recommend an arm a_t with the largest UCB index (Eq.(5)), and receive the reward $r_t \in [0, 1]$;
- 8: Update the statistics for user i_t $\mathbf{M}_{i_t, t} = \mathbf{M}_{i_t, t-1} + \mathbf{x}_{a_t} \mathbf{x}_{a_t}^\top$, $\mathbf{b}_{i_t, t} = \mathbf{b}_{i_t, t-1} + r_t \mathbf{x}_{a_t}$, $T_{i_t, t} = T_{i_t, t-1} + 1$, $\hat{\boldsymbol{\theta}}_{i_t, t} = (\lambda \mathbf{I} + \mathbf{M}_{i_t, t})^{-1} \mathbf{b}_{i_t, t}$;
- 9: Keep the statistics of other users unchanged
 $\mathbf{M}_{\ell, t} = \mathbf{M}_{\ell, t-1}$, $\mathbf{b}_{\ell, t} = \mathbf{b}_{\ell, t-1}$, $T_{\ell, t} = T_{\ell, t-1}$, $\hat{\boldsymbol{\theta}}_{\ell, t} = \hat{\boldsymbol{\theta}}_{\ell, t-1}$, for all $\ell \in \mathcal{U}, \ell \neq i_t$;
- 10: Delete the edge $(i_t, \ell) \in E_{t-1}$, if

$$\|\hat{\boldsymbol{\theta}}_{i_t, t} - \hat{\boldsymbol{\theta}}_{\ell, t}\|_2 \geq \alpha_1 \left(f(T_{i_t, t}) + f(T_{\ell, t}) \right) + \alpha_2 \epsilon_*,$$

and get an updated graph $G_t = (\mathcal{U}, E_t)$;

164 This section introduces our algorithm called ‘‘Robust CLustering of Misspecified Bandits’’
 165 (RCLUMB) (Algo.1). RCLUMB is a graph-based algorithm. The ideas and techniques of RCLUMB
 166 can be easily generalized to set-based algorithms. To illustrate this generalizability, we also design a
 167 set-based algorithm RSCLUMB. We leave the exposition and analysis of RSCLUMB in Appendix K.

168 For ease of interpretation, we define the coefficient

$$\zeta \triangleq 2\epsilon_* \sqrt{\frac{2}{\tilde{\lambda}_x}}, \quad (3)$$

169 where $\tilde{\lambda}_x \triangleq \int_0^{\lambda_x} (1 - e^{-\frac{(\lambda_x - x)^2}{2\sigma^2}})^C dx$. ζ is theoretically the minimum gap between two users’
 170 preference vectors that an algorithm can distinguish with high probability, as supported by Eq.(50) in

the proof of Lemma H.1 in Appendix H. Note that the algorithm does not require knowledge of ζ . We also make the following definition for illustration.

Definition 4.1 (ζ -close users and ζ -good clusters). Two users $i, i' \in \mathcal{U}$ are ζ -close if $\|\theta_i - \theta_{i'}\|_2 \leq \zeta$. Cluster \bar{V} is a ζ -good cluster at time t , if $\forall i \in \bar{V}$, user i and the coming user i_t are ζ -close.

We also say that two *ground-truth clusters* are “ ζ -close” if their preference vectors’ gap is less than ζ .

Now we introduce the process and intuitions of RCLUMB (Algo.1). The algorithm maintains an undirected user graph $G_t = (\mathcal{U}, E_t)$, where users are connected with edges if they are inferred to be in the same cluster. We denote the connected component in G_{t-1} containing user i_t at round t as \tilde{V}_t .

Cluster Detection. G_0 is initialized to be a complete graph, and will be updated adaptively based on the interactive information. At round t , user $i_t \in \mathcal{U}$ comes to be served with a feasible arm set \mathcal{A}_t (Line 4). Due to model misspecifications, it is impossible to cluster users with exactly the same preference vector θ , but similar users whose preference vectors are within the distance of ζ . According to the proof of Lemma H.1, after a sufficient time, with high probability, any pair of users directly connected by an edge in E_{t-1} are ζ -close. However, if we trivially follow previous CB works [11, 23, 26] to directly use the connected component \tilde{V}_t as the inferred cluster for user i_t at round t , it will cause a large regret. The reason is that in the worst case, the preference vector θ of the user in \tilde{V}_t who is h -hop away from user i_t could deviate by $h\zeta$ from θ_{i_t} , where h can be as large as $|\tilde{V}_t|$. Based on this reasoning, our key point is to select the cluster \bar{V}_t as the users at most 1-hop away from i_t in the graph. In other words, after some interactions, \bar{V}_t forms a ζ -good cluster with high probability; thus, RCLUMB can avoid using misleading information from dissimilar users for recommendations.

Cluster-based Recommendation. After finding the appropriate cluster \bar{V}_t for i_t , the agent estimates the common user preference vector based on the historical information associated with cluster \bar{V}_t by

$$\hat{\theta}_{\bar{V}_t, t-1} = \arg \min_{\theta \in \mathbb{R}^d} \sum_{s \in [t-1]} \sum_{i_s \in \bar{V}_t} (r_s - \mathbf{x}_{a_s}^\top \theta)^2 + \lambda \|\theta\|_2^2, \quad (4)$$

where $\lambda > 0$ is a regularization coefficient. Its closed-form solution is $\hat{\theta}_{\bar{V}_t, t-1} = \bar{M}_{\bar{V}_t, t-1}^{-1} \bar{\mathbf{b}}_{\bar{V}_t, t-1}$, where $\bar{M}_{\bar{V}_t, t-1} = \lambda \mathbf{I} + \sum_{s \in [t-1]} \sum_{i_s \in \bar{V}_t} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top$, $\bar{\mathbf{b}}_{\bar{V}_t, t-1} = \sum_{s \in [t-1]} \sum_{i_s \in \bar{V}_t} r_{a_s} \mathbf{x}_{a_s}$.

Based on this estimation, in Line 7, the agent recommends an arm using the UCB strategy

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}_t} \min \left\{ \underbrace{1, \mathbf{x}_a^\top \hat{\theta}_{\bar{V}_t, t-1}}_{\hat{R}_{a,t}} + \underbrace{\beta \|\mathbf{x}_a\|_{\bar{M}_{\bar{V}_t, t-1}^{-1}} + \epsilon_* \sum_{s \in [t-1]} \sum_{i_s \in \bar{V}_t} \left| \mathbf{x}_a^\top \bar{M}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \right|}_{C_{a,t}} \right\}, \quad (5)$$

where $\beta = \sqrt{\lambda} + \sqrt{2 \log(\frac{1}{\delta}) + d \log(1 + \frac{T}{\lambda d})}$, $\hat{R}_{a,t}$ denotes the estimated reward of arm a at t , $C_{a,t}$ denotes the confidence radius of arm a at round t .

Due to deviations from linearity, the estimation $\hat{R}_{a,t}$ computed by a linear function is no longer accurate. To handle the estimation uncertainty of model misspecifications, we design an enlarged confidence radius $C_{a,t}$. The first term of $C_{a,t}$ in Eq.(5) captures the uncertainty of online learning for the linear part, and the second term related to ϵ_* reflects the additional uncertainty from deviations from linearity. The design of $C_{a,t}$ theoretically relies on Lemma 5.6 which will be given in Section 5.

Update User Statistics. Based the feedback r_t , in Line 8 and 9, the agent updates the statistics for user i_t . Specifically, the agent estimates the preference vector θ_{i_t} by

$$\hat{\theta}_{i_t, t} = \arg \min_{\theta \in \mathbb{R}^d} \sum_{s \in [t]} \sum_{i_s = i_t} (r_s - \mathbf{x}_{a_s}^\top \theta)^2 + \lambda \|\theta\|_2^2, \quad (6)$$

with solution $\hat{\theta}_{i_t, t} = (\lambda \mathbf{I} + M_{i_t, t})^{-1} \mathbf{b}_{i_t, t}$, where $M_{i_t, t} = \sum_{s \in [t]} \sum_{i_s = i_t} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top$, $\mathbf{b}_{i_t, t} = \sum_{s \in [t]} \sum_{i_s = i_t} r_{a_s} \mathbf{x}_{a_s}$.

Update the Graph G_t . Finally, in Line 10, the agent verifies whether the similarities between user i_t and other users are still true based on the updated estimation $\hat{\theta}_{i_t, t}$. For every user $\ell \in \mathcal{U}$ connected with user i_t via edge $(i_t, \ell) \in E_{t-1}$, if the gap between her estimated preference vector $\hat{\theta}_{\ell, t}$ and $\hat{\theta}_{i_t, t}$ is larger than a threshold supported by Lemma H.1, the agent will delete the edge (i_t, ℓ) to split them apart. The threshold in Line 10 is carefully designed, taking both estimation uncertainty in a linear model and deviations from linearity into consideration. As shown in the proof of Lemma

H.1 (in Appendix H), using this threshold, with high probability, edges between users in the same *ground-truth clusters* will not be deleted, and edges between users that are not ζ -close will always be deleted. Together with the filtering step in Line 5, with high probability, the algorithm will leverage all the collaborative information of similar users and avoid misusing the information of dissimilar users. The updated graph G_t will be used in the next round.

5 Theoretical Analysis

In this section, we theoretically analyze the performance of the RCLUMB algorithm by giving an upper bound of the expected regret defined in Eq.(2). Due to the space limitation, we only show the main result (Theorem 5.3), key lemmas, and a sketched proof for Theorem 5.3. Detailed proofs, other technical lemmas, and the regret analysis of the RSLUMB algorithm can be found in the Appendix.

To state our main result, we first give two definitions as follows. The first definition is about the minimum separable gap constant γ_1 of a CBMUM problem instance.

Definition 5.1 (Minimum separable gap γ_1). The minimum separable gap constant γ_1 of a CBMUM problem instance is the minimum gap over the gaps among users that are greater than ζ (Eq. (3))

$$\gamma_1 = \min\{\|\theta_i - \theta_\ell\|_2 : \|\theta_i - \theta_\ell\|_2 > \zeta, \forall i, \ell \in \mathcal{U}\}, \text{ with } \min \emptyset = \infty.$$

Remark 2. In CBMUM, the role of $\gamma_1 - \zeta$ is similar to that of γ (given in Assumption 3.1) in the previous CB problem with perfectly linear models, quantifying the hardness of performing clustering on the problem instance. Intuitively, users are easier to cluster if γ_1 is larger, and the deduction of ζ shows the additional difficulty due to model deviations. If there are no misspecifications, i.e., $\zeta = 2\epsilon_*\sqrt{\frac{2}{\lambda_x}} = 0$, then $\gamma_1 = \gamma$, recovering the minimum separable gap between clusters in the classic CB problem [11, 23] without model misspecifications.

The second definition is about the number of “hard-to-cluster users” \tilde{u} .

Definition 5.2 (Number of “hard-to-cluster users” \tilde{u}). The number of “hard-to-cluster users” \tilde{u} is the number of users in the *ground-truth clusters* which are ζ -close to some other *ground-truth clusters*

$$\tilde{u} = \sum_{j \in [m]} |V_j| \times \mathbb{I}\{\exists j' \in [m], j' \neq j : \|\theta^{j'} - \theta^j\|_2 \leq \zeta\},$$

where $\mathbb{I}\{\cdot\}$ denotes the indicator function of the argument, $|V_j|$ denotes the number of users in V_j .

Remark 3. \tilde{u} captures the number of users who belong to different *ground-truth clusters* but their gaps are less than ζ . These users may be merged into one cluster by mistake and cause certain regret.

The following theorem gives an upper bound on the expected regret achieved by RCLUMB.

Theorem 5.3 (Main result on regret bound). Suppose that the assumptions in Section 3 are satisfied. Then the expected regret of the RCLUMB algorithm for T rounds satisfies

$$R(T) \leq O\left(u \left(\frac{d}{\bar{\lambda}_x(\gamma_1 - \zeta)^2} + \frac{1}{\bar{\lambda}_x^2} \right) \log T + \frac{\tilde{u} \epsilon_* \sqrt{dT}}{u \bar{\lambda}_x^{1.5}} + \epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T\right) \quad (7)$$

$$\leq O(\epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T), \quad (8)$$

where γ_1 is defined in Definition 5.1, and \tilde{u} is defined in Definition 5.2).

Discussion and Comparison. The bound in Eq.(7) has four terms. The first term is the time needed to gather enough information to assign appropriate clusters for users. The second term is the regret caused by *misclustering* ζ -close but not precisely similar users together, which is unavoidable with model misspecifications. The third term is from the preference estimation errors caused by model deviations. The last term is the usual term in CB with perfectly linear models [11, 23, 25].

Let us discuss how the parameters affect this regret bound.

- If $\gamma_1 - \zeta$ is large, the gaps between clusters that are not “ ζ -close” are much greater than the minimum gap ζ for the algorithm to distinguish, the first term in Eq.(7) will be small as it is easy to identify their dissimilarities. The role of $\gamma_1 - \zeta$ in CBMUM is similar to that of γ in the previous CB.
 - If \tilde{u} is small, indicating that few *ground-truth clusters* are “ ζ -close”, RCLUMB will hardly *miscluster* different *ground-truth clusters* together thus the second term in Eq.(7) will be small.
 - If the deviation level ϵ_* is small, the user models are close to linearity and the misspecifications will not affect the estimations much, then both the second and third term in Eq.(7) will be small.
- The following theorem gives a regret lower bound of the CBMUM problem.

256 **Theorem 5.4** (Regret lower bound for CBMUM). *There exists a problem instance for the CBMUM*
 257 *problem such that for any algorithm $R(T) \geq \Omega(\epsilon_* T \sqrt{d})$.*

258 The proof can be found in Appendix F. The upper bounds in Theorem 5.3 asymptotically match this
 259 lower bound with respect to T up to logarithmic factors (and a constant factor of \sqrt{m} where m is
 260 typically small in real-applications), showing the tightness of our theoretical results. Additionally, we
 261 conjecture the gap for the m factor is due to the strong assumption that cluster structures are known
 262 to prove this lower bound, and whether there exists a tighter lower bound is left for future work.

263 We then compare our results with two degenerate cases. First, when $m = 1$ (indicating $\tilde{u} = 0$), our
 264 setting degenerates to the MLB problem where all users share the same preference vector. In this
 265 case, our regret bound is $O(\epsilon_* T \sqrt{d \log T} + d \sqrt{T} \log T)$, exactly matching the current best bound of
 266 MLB [21]. Second, when $\epsilon_* = 0$, our setting reduces to the CB problem with perfectly linear user
 267 models and our bounds become $O(d \sqrt{m T} \log T)$, also perfectly match the existing best bound of
 268 the CB problem [23, 25]. The above discussions and comparisons show the tightness of our regret
 269 bounds. Additionally, we also provide detailed discussions on why trivially combining existing works
 270 on CB and MLB would not get any non-vacuous regret upper bound in Appendix D.

271 We define the following “good partition” for ease of interpretation.

272 **Definition 5.5** (Good partition). RCLUMB does a “good partition” at t , if the cluster \bar{V}_t assigned to
 273 i_t is a ζ -good cluster, and it contains all the users in the same *ground-truth cluster* as i_t , i.e.,

$$\|\theta_{i_t} - \theta_\ell\|_2 \leq \zeta, \forall \ell \in \bar{V}_t, \text{ and } V_{j(i_t)} \subseteq \bar{V}_t. \quad (9)$$

274 Note that when the algorithm does a “good partition” at t , \bar{V}_t will contain all the users in the same
 275 *ground-truth cluster* as i_t and may only contain some other ζ -close users with respect to i_t , which
 276 means the gathered information associated with \bar{V}_t can be used to infer user i_t ’s preference with high
 277 accuracy. Also, it is obvious that under a “good partition”, if $\bar{V}_t \in \mathcal{V}$, then $\bar{V}_t = V_{j(i_t)}$ by definition.

278 Next, we give a sketched proof for Theorem 5.3.

279 *Proof.* [Sketch for Theorem 5.3] The proof mainly contains two parts. First, we prove there is a
 280 sufficient time T_0 for RCLUMB to get a “good partition” with high probability. Second, we prove the
 281 regret upper bound for RCLUMB after maintaining a “good partition”. The most challenging part is
 282 to bound the regret caused by *misclustering* ζ -close users after getting a “good partition”.

283 **1. Sufficient time to maintain a “good partition”.** With the item regularity (Assumption 3.3),
 284 we can prove after some T_0 (defined in Lemma H.1 in Appendix H), RCLUMB will always have a
 285 “good partition”. Specifically, after $t \geq O\left(u \left(\frac{d}{\lambda_x(\gamma_1 - \zeta)^2} + \frac{1}{\lambda_x^2} \right) \log T\right)$, for any user $i \in \mathcal{U}$, the gap
 286 between the estimated $\hat{\theta}_{i,t}$ and the ground-truth $\theta^{j(i)}$ is less than $\frac{\gamma_1}{4}$ with high probability. With this,
 287 we can get: for any two users i and ℓ , if their gap is greater than ζ , it will trigger the deletion of the
 288 edge (i, ℓ) (Line 10 of Algo.1) with high probability; on the other hand, when the deletion condition
 289 of the edge (i, ℓ) is satisfied, then $\|\theta^{j(i)} - \theta^{j(\ell)}\|_2 > 0$, which means user i and ℓ belong to different
 290 *ground-truth clusters* by Assumption 3.1 with high probability. Therefore, we can get that with high
 291 probability, all those users in the same *ground-truth cluster* as i_t will be directly connected with i_t ,
 292 and users directly connected with i_t must be ζ -close to i_t . By filtering users directly linked with i_t as
 293 the cluster \bar{V}_t (Algo.1 Line 5) and the definition of “good partition”, we can ensure that RCLUMB
 294 will keep a “good partition” afterward with high probability.

295 **2. Bounding the regret after getting a “good partition”.** After T_0 , with the “good partition”, we
 296 can prove the following lemma that gives a bound of the difference between $\hat{\theta}_{\bar{V}_t, t-1}$ and ground-truth
 297 θ_{i_t} in direction of action vector \mathbf{x}_a , and supports the design of the confidence radius $C_{a,t}$ in Eq.(5).

298 **Lemma 5.6.** *With probability at least $1 - 5\delta$ for some $\delta \in (0, \frac{1}{5})$, $\forall t \geq T_0$*

$$\left| \mathbf{x}_a^\top (\theta_{i_t} - \hat{\theta}_{\bar{V}_t, t-1}) \right| \leq \frac{\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}} \mathbb{I}\{\bar{V}_t \notin \mathcal{V}\} + \epsilon_* \sum_{s \in [t-1]} \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \right| + \beta \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}}.$$

299 To prove this lemma, we consider the following two situations.

300 **(i) Assigning a perfect cluster for i_t .** In this case, $\bar{V}_t \in \mathcal{V}$, meaning the cluster assigned for user i_t
 301 is the same as her *ground-truth cluster*, i.e., $\bar{V}_t = V_{j(i_t)}$. Therefore, we have that $\forall \ell \in \bar{V}_t, \theta_\ell = \theta_{i_t}$.
 302 With careful analysis, we can bound $\left| \mathbf{x}_a^\top (\theta_{i_t} - \hat{\theta}_{\bar{V}_t, t-1}) \right|$ by $C_{a,t}$ (defined in Eq.(5)).

(ii) **Bounding the term of *misclustering* i_t 's ζ -close users.** In this case, $\bar{V}_t \notin \mathcal{V}$, meaning the algorithm *misclusters* user i_t , i.e., $\bar{V}_t \neq V_{j(i_t)}$. Thus, we do not have $\forall \ell \in \bar{V}_t, \theta_\ell = \theta_{i_t}$ anymore, but we have all the users in \bar{V}_t are ζ -close to i_t (by “good partition”), i.e., $\|\theta_{i_s} - \theta_{i_t}\|_2 \leq \zeta, \forall \ell \in \bar{V}_t$. Then an additional term can be caused by using the information of i_t 's ζ -close users in \bar{V}_t lying in different *ground-truth clusters* from i_t to estimate θ_{i_t} . It is highly challenging to bound this part.

We will get an extra term $\left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\theta_{i_s} - \theta_{i_t}) \right|$ when bounding the regret in this case, where $\|\theta_\ell - \theta_{i_t}\|_2 \leq \zeta, \forall \ell \in \bar{V}_t$. It is an easy-to-be-made mistake to directly drag $\|\theta_{i_s} - \theta_{i_t}\|_2$ out to bound it by $\left\| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top \right\|_2 \times \zeta$. With subtle analysis, we propose the following lemma to bound the above term.

Lemma 5.7 (Bound of error caused by *misclustering*). $\forall t \geq T_0$, if the current partition by RCLUMB is a “good partition”, and $\bar{V}_t \notin \mathcal{V}$, then for all $\mathbf{x}_a \in \mathbb{R}^d, \|\mathbf{x}_a\|_2 \leq 1$, with probability at least $1 - \delta$:

$$\left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\theta_{i_s} - \theta_{i_t}) \right| \leq \frac{\epsilon_* \sqrt{2d}}{\bar{\lambda}_x^{\frac{3}{2}}}$$

This lemma is quite general. Please see Appendix G for details about its proof.

The expected occurrences of $\{\bar{V}_t \notin \mathcal{V}\}$ is bounded by $\frac{\bar{u}}{u}T$ with Assumption 3.2, Definition 5.2 and 5.5. The result follows by bounding the expected sum of the bounds for the instantaneous regret using Lemma 5.6 with delicate analysis due to the time-varying clustering structure kept by RCLUMB. \square

6 Experiments

This section compares RCLUMB and RSCLUMB with CLUB [11], SCLUB [25], LinUCB with a single estimated vector for all users, LinUCB-Ind with separate estimated vectors for each user, and two modifications of LinUCB in [21] which we name as RLinUCB and RLinUCB-Ind. We use averaged reward as the evaluation metric, where the average is taken over ten independent trials.

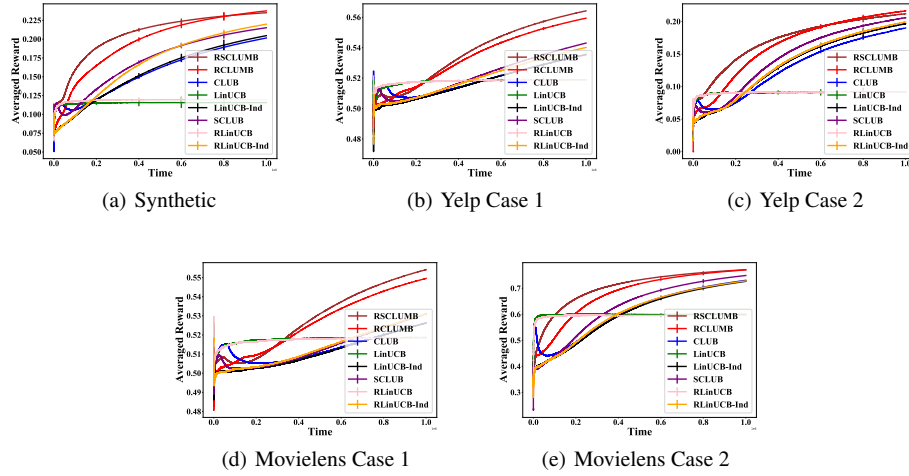


Figure 1: The figures compare RCLUMB and RSCLUMB with the baselines. (a) shows the result on synthetic data, (b) and (c) show the results on Yelp dataset, (d) and (e) show the results on Movielens dataset. All experiments are under the setting of $u = 1,000$ users, $m = 10$ clusters, and $d = 50$. All results are averaged under 10 random trials. The error bars are standard deviations divided by $\sqrt{10}$.

6.1 Synthetic Experiments

We consider a setting with $u = 1,000$ users, $m = 10$ clusters and $T = 10^6$ rounds. The preference and feature vectors are in $d = 50$ dimension with each entry drawn from a standard Gaussian distribution, and are normalized to vectors with $\|\cdot\|_2 = 1$ [25]. We fix an arm set with $|\mathcal{A}| = 1000$ items, at each round t , 20 items are randomly selected to form a set \mathcal{A}_t for the user to choose from. We construct a matrix $\epsilon \in \mathbb{R}^{1,000 \times 1,000}$ in which each element $\epsilon(i, j)$ is drawn uniformly from the range $(-0.2, 0.2)$ to represent the deviation. At t , for user i_t and the item a_t , $\epsilon(i_t, a_t)$ will be added to the feedback as the deviation, which corresponds to the $\epsilon_{a_t}^{i_t, t}$ defined in Eq.(1).

The result is provided in Figure 1(a), showing that our algorithms have clear advantages: RCLUMB improves over CLUB by 21.9%, LinUCB by 194.8%, LinUCB-Ind by 20.1%, SCLUB by 12.0%, RLinUCB by 185.2% and RLinUCB-Ind by 10.6%. The performance difference between RCLUMB and RSCLUMB is very small as expected. RLinUCB performs better than LinUCB; RLinUCB-Ind performs better than LinUCB-Ind and CLUB, showing that the modification of the recommendation policy is effective. The set-based RSCLUMB and SCLUB can separate clusters quicker and have advantages in the early period, but eventually RCLUMB catches up with RSCLUMB, and SCLUB is surpassed by RLinUCB-Ind because it does not consider misspecifications. RCLUMB and RSCLUMB perform better than RLinUCB-Ind, which shows the advantage of the clustering. So it can be concluded that both the modification for misspecification and the clustering structure are critical to improving the algorithm’s performance. We also have done some ablation experiments on different scales of ϵ^* in Appendix P, and we can notice that under different ϵ^* , our algorithms always outperform the baselines, and some baselines will perform worse as ϵ^* increases.

6.2 Experiments on Real-world Datasets

We conduct experiments on the Yelp data and the 20m MovieLens data [16]. For both data, we have two cases due to the different methods for generating feedback. For case 1, we extract 1,000 items with most ratings and 1,000 users who rate most; then we construct a binary matrix $\mathbf{H}^{1,000 \times 1,000}$ based on the user rating [34, 37]: if the user rating is greater than 3, the feedback is 1; otherwise, the feedback is 0. Then we use this binary matrix to generate the preference and feature vectors by singular-value decomposition (SVD) [25, 23, 34]. Similar to the synthetic experiment, we construct a matrix $\epsilon \in \mathbb{R}^{1,000 \times 1,000}$ in which each element is drawn uniformly from the range $(-0.2, 0.2)$. For case 2, we extract 1,100 users who rate most and 1000 items with most ratings. We construct a binary feedback matrix $\mathbf{H}^{1,100 \times 1,000}$ based on the same rule as case 1. Then we select the first 100 rows $\mathbf{H}_1^{100 \times 1,000}$ to generate the feature vectors by SVD. The remaining 1,000 rows $\mathbf{F}^{1,000 \times 1,000}$ is used as the feedback matrix, meaning user i receives $\mathbf{F}(i, j)$ as feedback while choosing item j . In both cases, at time t , we randomly select 20 items for the algorithms to choose from. In case 1, the feedback is computed by the preference and feature vector with misspecification, in case 2, the feedback is from the feedback matrix.

The results on Yelp are shown in Fig 1(b) and Fig 1(c). In case 1, RCLUMB improves CLUB by 45.1%, SCLUB by 53.4%, LinUCB-One by 170.1%, LinUCB-Ind by 46.2%, RLinUCB by 171.0% and RLinUCB-Ind by 21.5%. In case 2, RCLUMB improves over CLUB by 13.9%, SCLUB by 5.1%, LinUCB-One by 135.6%, LinUCB-Ind by 10.1%, RLinUCB by 138.6% and RLinUCB-Ind by 8.5%. It is notable that our modeling assumption 3.4 is violated in case 2 since the misspecification range is unknown. We set $\epsilon_* = 0.2$ following our synthetic dataset and it can still perform better than other algorithms. When the misspecification level is known as in case 1, our algorithms’ improvement is significantly enlarged, e.g., RCLUMB improves over SCLUB from 5.1% to 53.4%.

The results on Movielens are shown in Fig 1(d) and 1(e). In case 1, RCLUMB improves CLUB by 58.8%, SCLUB by 92.1%, LinUCB-One by 107.7%, LinUCB-Ind by 61.5%, RLinUCB by 109.5%, and RLinUCB-Ind by 21.3%. In case 2, RCLUMB improves over CLUB by 5.5%, SCLUB by 2.9%, LinUCB-One by 28.5%, LinUCB-Ind by 6.1%, RLinUCB by 29.3% and RLinUCB-Ind by 5.8%. The results are consistent with the Yelp data, confirming our superior performance.

7 Conclusion

We present a new problem of clustering of bandits with misspecified user models (CBMUM), where the agent has to adaptively assign appropriate clusters for users under the disturbance of model misspecifications. We propose two robust CB algorithms, RCLUMB and RSCLUMB. We prove the regret bounds of our algorithms, which match the lower bound asymptotically in T up to logarithmic factors, and match the state-of-the-art results in several degenerate cases. It is highly challenging to bound the regret caused by *misclustering* users with close but not the same preference vectors and use inaccurate cluster-based information to select arms. Our analysis to bound this part of the regret is quite general and may be of independent interest. Experiments on synthetic and real-world data demonstrate the advantage of our algorithms. We would like to state that there are some interesting future works: (1) Prove a tighter regret lower bound for CBMUM, (2) Incorporate recent model selection methods into our fundamental framework to design robust algorithms for CBMUM with unknown exact maximum model misspecification level, and (3) Consider the setting with misspecifications in the underlying user clustering structure rather than user models.

References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- [2] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- [3] Yikun Ban and Jingrui He. Local clustering in contextual multi-armed bandits. In *Proceedings of the Web Conference 2021*, pages 2335–2346, 2021.
- [4] Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- [5] Leonardo Cella and Massimiliano Pontil. Multi-task and meta-learning with sparse linear bandits. In *Uncertainty in Artificial Intelligence*, pages 1692–1702. PMLR, 2021.
- [6] Leonardo Cella, Alessandro Lazaric, and Massimiliano Pontil. Meta-learning with stochastic linear bandits. In *International Conference on Machine Learning*, pages 1360–1370. PMLR, 2020.
- [7] Leonardo Cella, Karim Lounici, Grégoire Pacreau, and Massimiliano Pontil. Multi-task representation learning with stochastic linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 4822–4847. PMLR, 2023.
- [8] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214. JMLR Workshop and Conference Proceedings, 2011.
- [9] Aniket Anand Deshmukh, Urun Dogan, and Clay Scott. Multi-task learning for contextual bandits. *Advances in neural information processing systems*, 30, 2017.
- [10] Dylan J Foster, Claudio Gentile, Mehryar Mohri, and Julian Zimmert. Adapting to misspecification in contextual bandits. *Advances in Neural Information Processing Systems*, 33: 11478–11489, 2020.
- [11] Claudio Gentile, Shuai Li, and Giovanni Zappella. Online clustering of bandits. In *International Conference on Machine Learning*, pages 757–765. PMLR, 2014.
- [12] Claudio Gentile, Shuai Li, Purushottam Kar, Alexandros Karatzoglou, Giovanni Zappella, and Evans Etrue. On context-dependent clustering of bandits. In *International Conference on machine learning*, pages 1253–1262. PMLR, 2017.
- [13] Avishek Ghosh, Sayak Ray Chowdhury, and Aditya Gopalan. Misspecified linear bandits. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [14] Jens Hainmueller and Chad Hazlett. Kernel regularized least squares: Reducing misspecification bias with a flexible and interpretable machine learning approach. *Political Analysis*, 22(2): 143–168, 2014.
- [15] Negar Hariri, Bamshad Mobasher, and Robin Burke. Context adaptation in interactive recommender systems. In *Proceedings of the 8th ACM Conference on Recommender Systems*, pages 41–48, 2014.
- [16] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- [17] Joey Hong, Branislav Kveton, Manzil Zaheer, and Mohammad Ghavamzadeh. Hierarchical bayesian bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 7724–7741. PMLR, 2022.
- [18] Ruiquan Huang, Weiqiang Wu, Jing Yang, and Cong Shen. Federated linear contextual bandits. *Advances in neural information processing systems*, 34:27057–27068, 2021.

- [19] Pushmeet Kohli, Mahyar Salek, and Greg Stoddard. A fast bandit algorithm for recommendation to users with heterogenous tastes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 27, pages 1135–1141, 2013.
- [20] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [21] Tor Lattimore, Csaba Szepesvari, and Gellert Weisz. Learning with good feature representations in bandits and in rl with a generative model. In *International Conference on Machine Learning*, pages 5662–5670. PMLR, 2020.
- [22] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- [23] Shuai Li and Shengyu Zhang. Online clustering of contextual cascading bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [24] Shuai Li, Alexandros Karatzoglou, and Claudio Gentile. Collaborative filtering bandits. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 539–548, 2016.
- [25] Shuai Li, Wei Chen, Shuai Li, and Kwong-Sak Leung. Improved algorithm on online clustering of bandits. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence, IJCAI’19*, page 2923–2929. AAAI Press, 2019. ISBN 9780999241141.
- [26] Xutong Liu, Haoru Zhao, Tong Yu, Shuai Li, and John Lui. Federated online clustering of bandits. In *The 38th Conference on Uncertainty in Artificial Intelligence*, 2022.
- [27] Aldo Pacchiano, My Phan, Yasin Abbasi Yadkori, Anup Rao, Julian Zimmert, Tor Lattimore, and Csaba Szepesvari. Model selection in contextual stochastic bandit problems. *Advances in Neural Information Processing Systems*, 33:10328–10337, 2020.
- [28] Chengshuai Shi and Cong Shen. Federated multi-armed bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9603–9611, 2021.
- [29] Marta Soare, Ouais Alsharif, Alessandro Lazaric, and Joelle Pineau. Multi-task linear bandits. In *NIPS2014 workshop on transfer and multi-task learning: theory meets practice*, 2014.
- [30] Runzhe Wan, Lin Ge, and Rui Song. Metadata-based multi-task bandits with bayesian hierarchical models. *Advances in Neural Information Processing Systems*, 34:29655–29668, 2021.
- [31] Runzhe Wan, Lin Ge, and Rui Song. Towards scalable and robust structured bandits: A meta-learning framework. In *International Conference on Artificial Intelligence and Statistics*, pages 1144–1173. PMLR, 2023.
- [32] Zhi Wang, Chicheng Zhang, Manish Kumar Singh, Laurel Riek, and Kamalika Chaudhuri. Multitask bandit learning through heterogeneous feedback aggregation. In *International Conference on Artificial Intelligence and Statistics*, pages 1531–1539. PMLR, 2021.
- [33] Zhi Wang, Chicheng Zhang, and Kamalika Chaudhuri. Thompson sampling for robust transfer in multi-task bandits. *arXiv preprint arXiv:2206.08556*, 2022.
- [34] Junda Wu, Canzhe Zhao, Tong Yu, Jingyang Li, and Shuai Li. Clustering of conversational bandits for user preference learning and elicitation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 2129–2139, 2021.
- [35] Qingyun Wu, Huazheng Wang, Quanquan Gu, and Hongning Wang. Contextual bandits in a collaborative environment. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 529–538, 2016.
- [36] Chunqiu Zeng, Qing Wang, Shekoofeh Mokhtari, and Tao Li. Online context-aware recommendation with time varying multi-armed bandit. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 2025–2034, 2016.

479 [37] Shi Zong, Hao Ni, Kenny Sung, Nan Rosemary Ke, Zheng Wen, and Branislav Kveton. Cas-
480 cading bandits for large-scale recommendation problems. *arXiv preprint arXiv:1603.05359*,
481 2016.

482 Appendix

483 A More Discussions on Related Work

484 In this section, we will give more comparisons and discussions on some previous works that are
485 related to our work to some extent.

486 There are some other works on bandits leveraging user (or task) relations, which have some relations
487 with the clustering of bandits (CB) works to some extent, but are in different lines of research from
488 CB, and are quite different from our work. First, besides CB, the work [35] also leverages user
489 relations. Specifically, it utilizes a *known* user adjacency graph to share context and payoffs among
490 neighbors, whereas in CB, the user relations are *unknown* and need to be learnt, thus the setting
491 differs a lot from CB. Second, there are lines of works on multi-task learning [5, 9, 29, 7, 33, 32],
492 meta-learning [31, 17, 6] and federated learning [28, 18], where multiple different tasks are solved
493 jointly and share information. Note that all of these works do not assume an underlying *unknown*
494 user clustering structure which needs to be inferred by the agent to speed up learning. For works
495 on multi-task learning [5, 9, 29, 7, 33, 32], they assume the tasks are related but no user clustering
496 structures, and to the best of our knowledge, none of them consider model misspecifications, thus
497 differing a lot from ours. For some recent works on meta-learning [31, 17, 30], they propose general
498 Bayesian hierarchical models to share knowledge across tasks, and design Thompson-Sampling-
499 based algorithms to optimize the Bayes regret, which are quite different from the line of CB works,
500 and differ a lot from ours. And additionally, as supported by the discussions in the works [6, 32],
501 multi-task learning and meta-learning are different lines of research from CB. For the works on
502 federated learning [28, 18], they consider the privacy and communication costs among multiple
503 servers, whose setting is also very different from the previous CB works and our work.

504 **Remark.** Again, we emphasize that the goal of this work is to initialize the study of the important
505 CBMUM problem, and propose general design ideas for dealing with model misspecifications in
506 CB problems. Therefore, our study is based on fundamental models on CB [11, 25] and MLB [21],
507 and the algorithm design ideas and theoretical analysis are pretty general. We leave incorporating
508 the more recent model selection methods [27, 10] into our framework to address the unknown exact
509 maximum model misspecification level as an interesting future work. It would also be interesting to
510 consider incorporating our methods and ideas of tackling model misspecifications into the studies of
511 multi-task learning, meta learning and federated learning.

512 B More Discussions on Assumptions

513 All the assumptions (Assumptions 3.1,3.2,3.3,3.4) in this work are natural and basically follow (or
514 less stringent than) previous works on CB and MLB [11, 23, 25, 26, 21].

515 B.1 Less Stringent Assumption on on the Generating Distribution of Arm Vectors

516 We also make some contributions to relax a widely-used but stringent assumption on the generating
517 distribution of arm vectors. Specifically, our Assumption 3.3 on item regularity relaxes the previous
518 one used in previous CB works [11, 23, 25, 26] by removing the condition that the variance should
519 be upper bounded by $\frac{\lambda^2}{8 \log(4|\mathcal{A}_t|)}$. For technical details on this, please refer to the theoretical analysis
520 and discussions in Appendix J.

521 B.2 Discussions on Assumption 3.4 about Bounded Misspecification Level

522 This assumption follows [21]. Note that this ϵ_* can be an upper bound on the maximum misspeci-
523 fication level, not the exact maximum itself. In real-world applications, the deviations are usually
524 small [13], and we can set a relatively big ϵ_* (e.g., 0.2) to be the upper bound. Our experimental
525 results support this claim. As shown in our experimental results on real-data case 2, even when ϵ_* is
526 unknown, our algorithms still perform well by setting $\epsilon_* = 0.2$. Some recent studies [27, 10] use
527 model selection methods to theoretically deal with unknown exact maximum misspecification level in
528 the single-user case, which is not the emphasis of this work. Additionally, the work [10] assumes that
529 the learning agent has access to a regression oracle. And for the work [27], though their regret bound

is dependent on the exact maximum misspecification level that needs not to be known by the agent, an upper bound of the exact maximum misspecification level is still needed. We leave incorporating their methods to deal with unknown exact maximum misspecification level as an interesting future work.

B.3 Discussions on Assumption 3.2 about the Theoretical Results under General User Arrival Distributions

The uniform arrival in Assumption 3.2 follows previous CB works [11, 23, 26], it only affects the T_0 term, which is the time after which the algorithm maintains a “good partition” and is of $O(u \log T)$. For an arbitrary arrival distribution, T_0 becomes $O(1/p_{\min} \log T)$, where p_{\min} is the minimal arrival probability of a user. And since it is a lower-order term (of $O(\log T)$), it will not affect the main order of our regret upper bound which is of $O(\epsilon_* T \sqrt{md \log T} + d\sqrt{mT} \log T)$. The work [25] studies arbitrary arrivals and aims to remove the $1/p_{\min}$ factor in this term, but their setting is different. They make an additional assumption that users in the same cluster not only have the same preference vector, but also the same arrival probability, which is different from our setting and other classic CB works [11, 23, 26] where we only assume users in the same cluster share the same preference vector.

C Highlight of the Theoretical Analysis

Our proof flow and methodologies are novel in clustering of bandits (CB), which are expected to inspire future works on model misspecifications and CB. The main challenge of the regret analysis in CBMUM is that due to the estimation inaccuracy caused by misspecifications, it is impossible to cluster all users exactly correctly, and it is highly non-trivial to bound the regret caused by “misclustering” ζ -close users.

To the best of our knowledge, the common proof flow of previous CB works (e.g., [11, 23, 26]) can be summarized in two steps: The first is to prove a sufficient time T'_0 after which the algorithms can cluster all users **exactly correctly** with high probability. Note that the inferred clustering structure remains static after T'_0 , making the analysis easy. Second, after the **correct static clustering**, the regret can be trivially bounded by bounding m (number of underlying clusters) independent linear bandit algorithms, resulting in a $O(d\sqrt{mT} \log T)$ regret.

The above common proof flow is straightforward in CB with perfectly linear models, but it would fail to get a non-vacuous regret bound for CBMUM. In CBMUM, it is impossible to learn an exactly correct static clustering structure with model misspecifications. In particular, we prove that we can only expect the algorithm to cluster ζ -close users together rather than cluster all users exactly correctly. Therefore, the previous flow can not be applied to the more challenging CBMUM problem.

We do the following to address the challenges in obtaining a tight regret bound for CBMUM. With the carefully-designed novel key components of RCLUMB, we can prove a sufficient time T_0 after which RCLUMB can get a “good partition” (Definition 5.5) with high probability, which means the cluster \bar{V}_t assigned to i_t contains all users in the same ground-truth cluster as i_t , and possibly some other i_t ’s ζ -close users. Intuitively, after T_0 , the algorithm can leverage all the information from the users’ ground-truth clusters but may misuse some information from other ζ -close users with preference gaps up to ζ , causing a regret of “misclustering” ζ -close users. It is highly non-trivial to bound this part of regret, and the proof methods would be beneficial for future studies in CB in challenging cases when it is impossible to cluster all users exactly correctly. For details, please refer to the discussions “(ii) Bounding the term of misclustering it’s ζ -close users” in Section 5, the key Lemma 5.7 (Bound of error caused by misclustering), its proof and tightness discussion in Appendix G. Also, a more subtle analysis is needed to handle the time-varying inferred clustering structure since the “good partition” may change over time, whereas in the previous CB works, the clustering structure remains static after T'_0 . For theoretical details on this, please refer to Appendix E.

D Discussions on why Trivially Combining Existing CB and MLB Works Could Not Achieve a Non-vacuous Regret Upper Bound

We consider discussing regret upper bounds for CB without considering misspecifications for three cases: (1) neither the clustering process nor the decision process considers misspecifications (previous

CB algorithms); (2) the decision process does not consider misspecifications; (3) the clustering process does not consider misspecifications.

For cases (1) and (2), the decision process could contribute to the leading regret. We consider the case where there are m underlying clusters, with each cluster's arrival being T/m , and the agent knows the underlying clustering structure. For this case, there exist some instances where the regret upper bound $R(T)$ is strictly larger than $\epsilon_* T \sqrt{m \log T}$ asymptotically in T . Formally, in the discussion of "Failure of unmodified algorithm" in Appendix E in [21], they give an example to show that in the single-user case, the regret $R_1(T)$ of the classic linear bandit algorithms without considering misspecifications

will have: $\lim_{T \rightarrow +\infty} \frac{R_1(T)}{\epsilon_* T \sqrt{m \log T}} = +\infty$. In our problem with multiple users and m underlying clusters, even if we know the underlying clustering structure and keep m independent linear bandit algorithms with T_i for the cluster $i \in [m]$ to leverage the common information of clusters, the best we can get is $R_2(T) = \sum_{i \in [m]} R_1(T_i)$. By the above results, if the decision process does not consider misspecifications, we have $\lim_{T \rightarrow +\infty} \frac{R_2(T)}{\epsilon_* T \sqrt{m \log T}} = \lim_{T \rightarrow +\infty} \frac{m R_1(T/m)}{\epsilon_* T \sqrt{m \log T}} = +\infty$. Recall that the regret upper bound $R(T)$ of our proposed algorithms is of $O(\epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T)$ (thus, we have $\lim_{T \rightarrow +\infty} \frac{R(T)}{\epsilon_* T \sqrt{m \log T}} < +\infty$), which gives a proof that that the regret upper bound of our proposed algorithms is asymptotically much better than CB algorithms in cases (1)(2).

For case (3), if the clustering process does not use the more tolerant deletion rule in Line 10 of Algo.1, the gap between users linked by edges would possibly exceed ζ ($\zeta = 2\epsilon_* \sqrt{\frac{2}{\lambda_x}}$) even after T_0 , which will result in a regret upper bound no better than $O(\epsilon_* u \sqrt{dT})$. As the number of users u is usually huge in practice, this result is vacuous. The reasons for getting the above claim are as follows. Even if the clustering process further uses our deletion rule considering misspecifications, and the users linked by edges are within ζ distance, failing to extract 1-hop users (Line 5 in Algo.1) would cause the leading $O(\epsilon_* u \sqrt{dT})$ regret term, as in the worst case, the preference vector θ of the user in \bar{V}_t who is h -hop away from user i_t could deviate by $h\zeta$ from θ_{i_t} , where h can be as large as u , and it would make the second term in Eq.(8) a $O(\epsilon_* u \sqrt{dT})$ term. If we completely do not consider the misspecifications in the clustering process, the above user gap between users linked by edges would possibly exceed ζ , which will cause a regret upper bound worse than $O(\epsilon_* u \sqrt{dT})$.

607 E Proof of Theorem 5.3

We first prove the result in the case when γ_1 defined in Definition 5.1 is not infinity, i.e., $4\epsilon_* \sqrt{\frac{2}{\lambda_x}} < \gamma_1 < \infty$. The proof of the special case when $\gamma_1 = \infty$ will directly follow the proof of this case.

For the instantaneous regret R_t at round t , with probability at least $1 - 5\delta$ for some $\delta \in (0, \frac{1}{5})$, at $\forall t \geq T_0$:

$$\begin{aligned} R_t &= (\mathbf{x}_{a_t^*}^\top \boldsymbol{\theta}_{i_t} + \boldsymbol{\epsilon}_{a_t^*}^{i_t, t}) - (\mathbf{x}_{a_t}^\top \boldsymbol{\theta}_{i_t} + \boldsymbol{\epsilon}_{a_t}^{i_t, t}) \\ &= \mathbf{x}_{a_t^*}^\top (\boldsymbol{\theta}_{i_t} - \hat{\boldsymbol{\theta}}_{\bar{V}_t, t-1}) + (\mathbf{x}_{a_t^*}^\top \hat{\boldsymbol{\theta}}_{\bar{V}_t, t-1} + C_{a_t^*, t}) - (\mathbf{x}_{a_t}^\top \hat{\boldsymbol{\theta}}_{\bar{V}_t, t-1} + C_{a_t, t}) \\ &\quad + \mathbf{x}_{a_t}^\top (\hat{\boldsymbol{\theta}}_{\bar{V}_t, t-1} - \boldsymbol{\theta}_{i_t}) + C_{a_t, t} - C_{a_t^*, t} + (\boldsymbol{\epsilon}_{a_t^*}^{i_t, t} - \boldsymbol{\epsilon}_{a_t}^{i_t, t}) \\ &\leq 2C_{a_t, t} + \frac{2\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}} \mathbb{I}\{\bar{V}_t \notin \mathcal{V}\} + 2\epsilon_*, \end{aligned} \tag{10}$$

where the last inequality holds by the UCB arm selection strategy in Eq.(5), the concentration bound given in Lemma 5.6, and the fact that $\|\boldsymbol{\epsilon}^{i, t}\|_\infty \leq \epsilon_*, \forall i \in \mathcal{U}, \forall t$.

We define the following events. Let

$$\begin{aligned} \mathcal{E}_0 &= \{R_t \leq 2C_{a_t, t} + \frac{2\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}} \mathbb{I}\{\bar{V}_t \notin \mathcal{V}\} + 2\epsilon_*, \text{ for all } \{t : t \geq T_0, \text{ and the algorithm maintains a "good partition" at } t\}\}, \\ \mathcal{E}_1 &= \{\text{the algorithm maintains a "good partition" for all } t \geq T_0\}, \\ \mathcal{E} &= \mathcal{E}_0 \cap \mathcal{E}_1. \end{aligned}$$

615 $\mathbb{P}(\mathcal{E}_0) \geq 1 - 2\delta$. According to Lemma H.1, $\mathbb{P}(\mathcal{E}_1) \geq 1 - 3\delta$. Thus, $\mathbb{P}(\mathcal{E}) \geq 1 - 5\delta$ for some
 616 $\delta \in (0, \frac{1}{5})$. Take $\delta = \frac{1}{T}$, we can get that

$$\begin{aligned}\mathbb{E}[R(T)] &= \mathbb{P}(\mathcal{E})\mathbb{I}\{\mathcal{E}\}R(T) + \mathbb{P}(\bar{\mathcal{E}})\mathbb{I}\{\bar{\mathcal{E}}\}R(T) \\ &\leq \mathbb{I}\{\mathcal{E}\}R(T) + 5 \times \frac{1}{T} \times T \\ &= \mathbb{I}\{\mathcal{E}\}R(T) + 5,\end{aligned}\tag{11}$$

617 where $\bar{\mathcal{E}}$ denotes the complementary event of \mathcal{E} , $\mathbb{I}\{\mathcal{E}\}R(T)$ denotes $R(T)$ under event \mathcal{E} , $\mathbb{I}\{\bar{\mathcal{E}}\}R(T)$
 618 denotes $R(T)$ under event $\bar{\mathcal{E}}$, and we use $R(T) \leq T$ to bound $R(T)$ under event $\bar{\mathcal{E}}$.

619 Then it remains to bound $\mathbb{I}\{\mathcal{E}\}R(T)$:

$$\begin{aligned}\mathbb{I}\{\mathcal{E}\}R(T) &\leq R(T_0) + \mathbb{E}[\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T R_t] \\ &\leq T_0 + 2\mathbb{E}[\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T C_{a_t,t}] + \frac{2\epsilon_*\sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}} \sum_{t=T_0+1}^T \mathbb{E}[\mathbb{I}\{\mathcal{E}, \bar{V}_t \notin \mathcal{V}\}] + 2\epsilon_*T \tag{12} \\ &= T_0 + 2\mathbb{E}[\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T C_{a_t,t}] + \frac{2\epsilon_*\sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}} \sum_{t=T_0+1}^T \mathbb{P}(\mathbb{I}\{\mathcal{E}, \bar{V}_t \notin \mathcal{V}\}) + 2\epsilon_*T \\ &\leq T_0 + 2\mathbb{E}[\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T C_{a_t,t}] + \frac{2\epsilon_*\sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}} \times \frac{\tilde{u}}{u}T + 2\epsilon_*T,\end{aligned}\tag{13}$$

620 where Eq.(12) follows from Eq.(10). Eq.(13) holds since under Assumption 3.2 about user arrival
 621 uniformness and by Definition 5.5 of “good partition”, $\mathbb{P}(\mathbb{I}\{\mathcal{E}, \bar{V}_t \notin \mathcal{V}\}) \leq \frac{\tilde{u}}{u}, \forall t \geq T_0$, where \tilde{u} is
 622 defined in Definition 5.2.

623 Then we need to bound $\mathbb{E}[\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T C_{a_t,t}]$:

$$\begin{aligned}\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T C_{a_t,t} &= \left(\sqrt{\lambda} + \sqrt{2\log(\frac{1}{\delta}) + d\log(1 + \frac{T}{\lambda d})}\right) \mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T \|\mathbf{x}_{a_t}\|_{\bar{\mathbf{M}}_{\bar{V}_t,t-1}^{-1}} \\ &\quad + \mathbb{I}\{\mathcal{E}\}\epsilon_* \sum_{t=T_0+1}^T \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \left| \mathbf{x}_{a_t}^\top \bar{\mathbf{M}}_{\bar{V}_t,t-1}^{-1} \mathbf{x}_{a_s} \right|.\end{aligned}\tag{14}$$

624 Next, we bound the $\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T \|\mathbf{x}_{a_t}\|_{\bar{\mathbf{M}}_{\bar{V}_t,t-1}^{-1}}$ term in Eq.(14):

$$\begin{aligned}\mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T \|\mathbf{x}_{a_t}\|_{\bar{\mathbf{M}}_{\bar{V}_t,t-1}^{-1}} &= \mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T \sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \|\mathbf{x}_{a_t}\|_{\bar{\mathbf{M}}_{\tilde{V}'_{t,k},t-1}^{-1}} \\ &\leq \mathbb{I}\{\mathcal{E}\} \sum_{t=T_0+1}^T \sum_{j=1}^m \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\bar{\mathbf{M}}_{V_j,t-1}^{-1}}\end{aligned}\tag{15}$$

$$\begin{aligned}&\leq \mathbb{I}\{\mathcal{E}\} \sum_{j=1}^m \sqrt{\sum_{t=T_0+1}^T \mathbb{I}\{i_t \in V_j\} \sum_{t=T_0+1}^T \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\bar{\mathbf{M}}_{V_j,t-1}^{-1}}^2} \\ &\tag{16}\end{aligned}$$

$$\leq \mathbb{I}\{\mathcal{E}\} \sum_{j=1}^m \sqrt{2T_{V_j,T} d \log(1 + \frac{T}{\lambda d})}\tag{17}$$

$$\begin{aligned}&\leq \mathbb{I}\{\mathcal{E}\} \sqrt{2 \sum_{j=1}^m 1 \sum_{j=1}^m T_{V_j,T} d \log(1 + \frac{T}{\lambda d})} = \mathbb{I}\{\mathcal{E}\} \sqrt{2mdT \log(1 + \frac{T}{\lambda d})}, \\ &\tag{18}\end{aligned}$$

where we use m_t to denote the number of connected components partitioned by the algorithm at t , $\tilde{V}'_{t,k}, k \in [m_t]$ to denote the connected components partitioned by the algorithm at t , $\bar{V}'_{t,k} \subseteq \tilde{V}'_{t,k}$ to denote the subset extracted to be the cluster \bar{V}_t for i_t from $\tilde{V}'_{t,k}$ conditioned on $i_t \in \tilde{V}'_{t,k}$, and $T_{V_j,T}$ to denote the number of times that the served users lie in the *ground-truth cluster* V_j up to time T , i.e., $T_{V_j,T} = \sum_{t \in [T]} \mathbb{I}\{i_t \in V_j\}$.

The reasons for having Eq.(15) are as follows. Under event \mathcal{E} , the algorithm will always have a “good partition” after T_0 . By Definition 5.5 and the proof process of Lemma H.1 about the edge deletion conditions, we can get $m_t \leq m$ and if $i_t \in \tilde{V}'_{t,k}, i_t \in V_j$, then $V_j \subseteq \bar{V}'_{t,k}$ since $\bar{V}'_{t,k}$ contains V_j and possibly other *ground-truth clusters* $V_n, n \in [m]$, whose preference vectors are ζ -close to θ^j . Therefore, by the definition of the regularized Gramian matrix, we can get $M_{\bar{V}'_{t,k},t-1} \succeq M_{V_j,t-1}, \forall t \geq T_0 + 1$. Thus by the above reasoning, $\sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \|\mathbf{x}_{a_t}\|_{\bar{M}_{\bar{V}'_{t,k},t-1}^{-1}} \leq \sum_{j=1}^m \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\bar{M}_{V_j,t-1}^{-1}}, \forall t \geq T_0 + 1$. Eq.(16) holds by the Cauchy–Schwarz inequality; Eq.(17) follows by the following technical Lemma J.2. Eq.(18) is from the Cauchy–Schwarz inequality and the fact that $\sum_{j=1}^m T_{V_j,T} = T$.

We then bound the last term in Eq.(14):

$$\begin{aligned} \mathbb{I}\{\mathcal{E}\} \epsilon_* \sum_{t=T_0+1}^T \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \left| \mathbf{x}_{a_t}^\top \bar{M}_{\bar{V}_t,t-1}^{-1} \mathbf{x}_{a_s} \right| &= \mathbb{I}\{\mathcal{E}\} \epsilon_* \sum_{t=T_0+1}^T \sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}'_{t,k}}} \left| \mathbf{x}_{a_t}^\top \bar{M}_{\bar{V}'_{t,k},t-1}^{-1} \mathbf{x}_{a_s} \right| \\ &\leq \mathbb{I}\{\mathcal{E}\} \epsilon_* \sum_{t=T_0+1}^T \sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \sqrt{\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}'_{t,k}}} 1 \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}'_{t,k}}} \left| \mathbf{x}_{a_t}^\top \bar{M}_{\bar{V}'_{t,k},t-1}^{-1} \mathbf{x}_{a_s} \right|^2} \end{aligned} \quad (19)$$

$$\leq \mathbb{I}\{\mathcal{E}\} \epsilon_* \sum_{t=T_0+1}^T \sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \sqrt{T_{\bar{V}'_{t,k},t-1} \|\mathbf{x}_{a_t}\|_{\bar{M}_{\bar{V}'_{t,k},t-1}^{-1}}^2} \quad (20)$$

$$\leq \mathbb{I}\{\mathcal{E}\} \epsilon_* \sum_{t=T_0+1}^T \sqrt{\sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} T_{\bar{V}'_{t,k},t-1} \|\mathbf{x}_{a_t}\|_{\bar{M}_{\bar{V}'_{t,k},t-1}^{-1}}^2} \quad (21)$$

$$\leq \mathbb{I}\{\mathcal{E}\} \epsilon_* \sqrt{T} \sum_{t=T_0+1}^T \sqrt{\sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \|\mathbf{x}_{a_t}\|_{\bar{M}_{\bar{V}'_{t,k},t-1}^{-1}}^2} \quad (22)$$

$$\leq \mathbb{I}\{\mathcal{E}\} \epsilon_* \sqrt{T} \sqrt{\sum_{t=T_0+1}^T 1 \sum_{t=T_0+1}^T \sum_{k=1}^{m_t} \mathbb{I}\{i_t \in \tilde{V}'_{t,k}\} \|\mathbf{x}_{a_t}\|_{\bar{M}_{\bar{V}'_{t,k},t-1}^{-1}}^2} \quad (23)$$

$$\leq \mathbb{I}\{\mathcal{E}\} \epsilon_* \sqrt{T} \sqrt{T \sum_{t=T_0+1}^T \sum_{j=1}^m \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\bar{M}_{V_j,t-1}^{-1}}^2} \quad (24)$$

$$\begin{aligned} &= \mathbb{I}\{\mathcal{E}\} \epsilon_* T \sqrt{\sum_{j=1}^m \sum_{t=T_0+1}^T \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\bar{M}_{V_j,t-1}^{-1}}^2} \\ &\leq \mathbb{I}\{\mathcal{E}\} \epsilon_* T \sqrt{2md \log(1 + \frac{T}{\lambda d})}, \end{aligned} \quad (25)$$

where Eq.(19), Eq.(21) and Eq.(23) hold because of the Cauchy–Schwarz inequality, Eq.(20) holds since $\overline{M}_{\overline{V}'_{t,k},t-1} \succeq \sum_{\substack{s \in [t-1] \\ i_s \in \overline{V}'_{t,k}}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top$, Eq.(22) is because $T_{\overline{V}'_{t,k},t-1} \leq T$, Eq. (24) follows from the same reasoning as Eq.(15), and Eq.(25) comes from the following technical Lemma J.2. Finally, plugging Eq.(18) and Eq.(25) into Eq.(14), take expectation and plug it into Eq.(13), we can get:

$$R(T) \leq 5 + T_0 + \frac{\tilde{u}}{u} \times \frac{2\epsilon_* \sqrt{2d}T}{\tilde{\lambda}_x^{\frac{3}{2}}} + 2\epsilon_* T \left(1 + \sqrt{2md \log(1 + \frac{T}{\lambda d})} \right) + 2 \left(\sqrt{\lambda} + \sqrt{2 \log(T) + d \log(1 + \frac{T}{\lambda d})} \right) \times \sqrt{2mdT \log(1 + \frac{T}{\lambda d})}, \quad (26)$$

where

$$T_0 = 16u \log(\frac{u}{\delta}) + 4u \max \max \left\{ \frac{8d}{\tilde{\lambda}_x(\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2\lambda_x}})^2} \log(\frac{u}{\delta}), \frac{16}{\tilde{\lambda}_x^2} \log(\frac{8d}{\tilde{\lambda}_x^2 \delta}) \right\}$$

is given in the following Lemma H.1 in Appendix H.

F Proof and Discussions of Theorem 5.4

In the work [21], they give a lower bound for misspecified linear bandits with a single user. The lower bound of $R(T)$ is given by: $R_3(T) \geq \epsilon_* T \sqrt{d}$. Therefore, suppose our problem with multiple users and m underlying clusters where the arrival times are T_i for each cluster, then for any algorithms, even if they know the underlying clustering structure and keep m independent linear bandit algorithms to leverage the common information of clusters, the best they can get is $R(T) = \sum_{i \in [m]} R_3(T_i) \geq \epsilon_* \sum_{i \in [m]} T_i \sqrt{d} = \epsilon_* T \sqrt{d}$, which gives a lower bound of $O(\epsilon_* T \sqrt{d})$ for the CBMUM problem. Recall that the regret upper bound of our algorithms is of $O(\epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T)$, asymptotically matching this lower bound with respect to T up to logarithmic factors and with respect to m up to $O(\sqrt{m})$ factors, showing the tightness of our theoretical results (where m are typically very small for real-applications).

We conjecture that the gap for the m factor is due to the strong assumption that cluster structures are known to prove our lower bound, and whether there exists a tighter lower bound will be left for future work.

G Proof of the key Lemma 5.7

In Lemma 5.7, we want to bound the term $\left| \mathbf{x}_a^\top \overline{M}_{\overline{V}_t,t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \overline{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\boldsymbol{\theta}_{i_s} - \boldsymbol{\theta}_{i_t}) \right|$. By the definition of “good partition”, we have $\|\boldsymbol{\theta}_{i_s} - \boldsymbol{\theta}_{i_t}\|_2 \leq \zeta, \forall i_s \in \overline{V}_t$. It is an easy-to-be-made mistake to directly drag $\|\boldsymbol{\theta}_{i_s} - \boldsymbol{\theta}_{i_t}\|_2$ out to upper bound it by $\left\| \mathbf{x}_a^\top \overline{M}_{\overline{V}_t,t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \overline{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top \right\|_2 \times \zeta$ and then proceed. We need more careful analysis.

We first prove the following general lemma.

Lemma G.1. For vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k \in \mathbb{R}^d, \|\mathbf{x}_i\|_2 \leq 1, \forall i \in [k]$, and vectors $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_k \in \mathbb{R}^d, \|\boldsymbol{\theta}_i\|_2 \leq C, \forall i \in [k]$, where $C > 0$ is a constant, we have:

$$\left\| \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\theta}_i \right\|_2 \leq C \sqrt{d} \left\| \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top \right\|_2.$$

669 *Proof.* Let $\mathbf{X} \in \mathbb{R}^{d \times k}$ be a matrix such that it has \mathbf{x}_i s as its columns, i.e., $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_k] =$

$$670 \begin{bmatrix} \mathbf{x}_{11} & x_{21} & \cdots & \mathbf{x}_{k1} \\ \mathbf{x}_{12} & x_{22} & \cdots & \mathbf{x}_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_{1d} & x_{2d} & \cdots & \mathbf{x}_{kd} \end{bmatrix}.$$

671 Let $\mathbf{y} \in \mathbb{R}^{k \times 1}$ be a vector that has $\mathbf{x}_i^\top \boldsymbol{\theta}_i$ s as its elements, i.e., $\mathbf{y} = [\mathbf{x}_1^\top \boldsymbol{\theta}_1, \dots, \mathbf{x}_k^\top \boldsymbol{\theta}_k]^\top$. Then we
672 have:

$$\left\| \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\theta}_i \right\|_2^2 = \|\mathbf{X} \mathbf{y}\|_2^2 \leq \|\mathbf{X}\|_2^2 \|\mathbf{y}\|_2^2 \quad (27)$$

$$\begin{aligned} &= \|\mathbf{X}\|_2^2 \sum_{i=1}^k (\mathbf{x}_i^\top \boldsymbol{\theta}_i)^2 \\ &\leq \|\mathbf{X}\|_2^2 \sum_{i=1}^k \|\mathbf{x}_i\|_2^2 \|\boldsymbol{\theta}_i\|_2^2 \end{aligned} \quad (28)$$

$$\begin{aligned} &\leq C^2 \|\mathbf{X}\|_2^2 \sum_{i=1}^k \|\mathbf{x}_i\|_2^2 \\ &= C^2 \|\mathbf{X}\|_2^2 \|\mathbf{X}\|_F^2 \\ &\leq C^2 d \|\mathbf{X}\|_2^4 \end{aligned} \quad (29)$$

$$= C^2 d \left\| \mathbf{X} \mathbf{X}^\top \right\|_2^2 \quad (30)$$

$$= C^2 d \left\| \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top \right\|_2^2, \quad (31)$$

673 where Eq. (27) follows by the matrix operator norm inequality, Eq. (28) follows by the
674 Cauchy–Schwarz inequality, Eq. (29) follows by $\|\mathbf{X}\|_F \leq \sqrt{d} \|\mathbf{X}\|_2$, Eq. (30) follows from
675 $\|\mathbf{X}\|_2^2 = \left\| \mathbf{X} \mathbf{X}^\top \right\|_2$. \square

676 The above result is tight. We can show that the lower bound of $\left\| \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\theta}_i \right\|_2$ under the conditions
677 in the lemma is exactly $C\sqrt{d} \left\| \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top \right\|_2$. Specifically, let $k = 2$, $C = 1$, $d = 2$, $\mathbf{x}_1 = [0, 1]^\top$,
678 $\mathbf{x}_2 = [1, 0]^\top$, $\boldsymbol{\theta}_1 = [1, 0]^\top$, $\boldsymbol{\theta}_2 = [0, 1]^\top$, then we have $\left\| \sum_{i=1}^2 \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\theta}_i \right\|_2 = \|[1, 1]^\top\|_2 = \sqrt{2}$, and
679 $C\sqrt{d} \left\| \sum_{i=1}^2 \mathbf{x}_i \mathbf{x}_i^\top \right\|_2 = 1 \times \sqrt{2} \times \left\| \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right\|_2 = \sqrt{2}$. Therefore, we have that the upper bound
680 given in Lemma G.1 matches the lower bound.

681 We are now ready to prove the key Lemma 5.7 with the above Lemma G.1.

At any $t \geq T_0$, if the current partition is a “good partition”, and $\bar{V}_t \notin \mathcal{V}$, then for all $\mathbf{x}_a \in \mathbb{R}^d$, $\|\mathbf{x}_a\|_2 \leq 1$, with probability at least $1 - \delta$:

$$\left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\boldsymbol{\theta}_{i_s} - \boldsymbol{\theta}_{i_t}) \right| \leq \|\mathbf{x}_a\|_2 \left\| \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\boldsymbol{\theta}_{i_s} - \boldsymbol{\theta}_{i_t}) \right\|_2 \quad (32)$$

$$\leq \left\| \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \right\|_2 \left\| \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\boldsymbol{\theta}_{i_s} - \boldsymbol{\theta}_{i_t}) \right\|_2 \quad (33)$$

$$\leq 2\epsilon_* \sqrt{\frac{2d}{\tilde{\lambda}_x}} \times \left\| \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \right\|_2 \left\| \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top \right\|_2 \quad (34)$$

$$\begin{aligned} &\leq 2\epsilon_* \sqrt{\frac{2d}{\tilde{\lambda}_x}} \times \frac{\lambda_{\max}(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top)}{\lambda_{\min}(\bar{\mathbf{M}}_{\bar{V}_t, t-1})} \\ &\leq 2\epsilon_* \sqrt{\frac{2d}{\tilde{\lambda}_x}} \times \frac{T_{\bar{V}_t, t-1}}{2T_{\bar{V}_t, t-1} \tilde{\lambda}_x + \lambda} \\ &\leq \frac{\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}}, \end{aligned} \quad (35)$$

where Eq.(32) follows by the Cauchy–Schwarz inequality, Eq.(33) follows from the inequality of matrix’s operator norm, Eq.(34) follows from the fact that in a “good partition”, $\|\boldsymbol{\theta}_{i_t} - \boldsymbol{\theta}_t\|_2 \leq 2\epsilon_* \sqrt{\frac{2}{\tilde{\lambda}_x}}$, $\forall t \in \bar{V}_t$ and Lemma G.1, Eq.(35) follows by Eq.(47) with probability $\geq 1 - \delta$.

H Lemma H.1 of the sufficient time T_0 and its proof

The following lemma gives a sufficient time T_0 for the algorithm to get a “good partition”.

Lemma H.1. *With the carefully designed edge deletion rule, after*

$$\begin{aligned} T_0 &\triangleq 16u \log\left(\frac{u}{\delta}\right) + 4u \max \max \left\{ \frac{8d}{\tilde{\lambda}_x(\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}})^2} \log\left(\frac{u}{\delta}\right), \frac{16}{\tilde{\lambda}_x^2} \log\left(\frac{8d}{\tilde{\lambda}_x^2 \delta}\right) \right\} \\ &= O\left(u \left(\frac{d}{\tilde{\lambda}_x(\gamma_1 - \zeta)^2} + \frac{1}{\tilde{\lambda}_x^2} \right) \log \frac{1}{\delta}\right) \end{aligned}$$

rounds, with probability at least $1 - 3\delta$ for some $\delta \in (0, \frac{1}{3})$, RCLUMB can always get a “good partition”.

Below is the detailed proof of Lemma H.1.

Proof. We first prove the following result:

With probability at least $1 - \delta$ for some $\delta \in (0, 1)$, at any $t \in [T]$:

$$\left\| \hat{\boldsymbol{\theta}}_{i,t} - \boldsymbol{\theta}^{j(i)} \right\|_2 \leq \frac{\beta(T_{i,t}, \frac{\delta}{u}) + \epsilon_* \sqrt{T_{i,t}}}{\sqrt{\lambda + \lambda_{\min}(\mathbf{M}_{i,t})}}, \forall i \in \mathcal{U}, \quad (36)$$

695 where $\beta(T_{i,t}, \frac{\delta}{u}) \triangleq \sqrt{\lambda} + \sqrt{2 \log(\frac{u}{\delta}) + d \log(1 + \frac{T_{i,t}}{\lambda d})}$.

$$\begin{aligned}
\hat{\theta}_{i,t} - \theta^{j(i)} &= \left(\sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} (\mathbf{x}_{a_s}^\top \theta^{j(i)} + \epsilon_{a_s}^{i_s, s} + \eta_s) \right) - \theta^{j(i)} \\
&= \left(\sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left[\left(\sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right) \theta^{j(i)} - \lambda \theta^{j(i)} + \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \eta_s \right] - \theta^{j(i)} \\
&= -\lambda \tilde{\mathbf{M}}_{i,t}^{-1} \theta^{j(i)} + \tilde{\mathbf{M}}_{i,t}^{-1} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \tilde{\mathbf{M}}_{i,t}^{-1} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \eta_s,
\end{aligned} \tag{37}$$

696 where we denote $\tilde{\mathbf{M}}_{i,t} = \mathbf{M}_{i,t} + \lambda \mathbf{I}$, and Eq.(37) holds by definition.

697 Therefore,

$$\left\| \hat{\theta}_{i,t} - \theta^{j(i)} \right\|_2 \leq \lambda \left\| \tilde{\mathbf{M}}_{i,t}^{-1} \theta^{j(i)} \right\|_2 + \left\| \tilde{\mathbf{M}}_{i,t}^{-1} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} \right\|_2 + \left\| \tilde{\mathbf{M}}_{i,t}^{-1} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \eta_s \right\|_2. \tag{38}$$

698 We then bound the three terms in Eq.(38) one by one. For the first term:

$$\lambda \left\| \tilde{\mathbf{M}}_{i,t}^{-1} \theta^{j(i)} \right\|_2 \leq \lambda \left\| \tilde{\mathbf{M}}_{i,t}^{-\frac{1}{2}} \right\|_2^2 \left\| \theta^{j(i)} \right\|_2 \leq \frac{\sqrt{\lambda}}{\sqrt{\lambda_{\min}(\tilde{\mathbf{M}}_{i,t})}}, \tag{39}$$

699 where we use the Cauchy–Schwarz inequality, the inequality for the operator norm of matrices, and
700 the fact that $\lambda_{\min}(\tilde{\mathbf{M}}_{i,t}) \geq \lambda$.

701 For the second term in Eq.(38):

$$\begin{aligned}
\left\| \tilde{\mathbf{M}}_{i,t}^{-1} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} \right\|_2 &= \max_{\mathbf{x} \in S^{d-1}} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}^\top \tilde{\mathbf{M}}_{i,t}^{-1} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} \\
&\leq \max_{\mathbf{x} \in S^{d-1}} \sum_{\substack{s \in [t] \\ i_s = i}} \left| \mathbf{x}^\top \tilde{\mathbf{M}}_{i,t}^{-1} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} \right| \\
&\leq \max_{\mathbf{x} \in S^{d-1}} \sum_{\substack{s \in [t] \\ i_s = i}} \left| \mathbf{x}^\top \tilde{\mathbf{M}}_{i,t}^{-1} \mathbf{x}_{a_s} \right| \left\| \epsilon_{a_s}^{i_s, s} \right\|_\infty \\
&\leq \epsilon_* \max_{\mathbf{x} \in S^{d-1}} \sum_{\substack{s \in [t] \\ i_s = i}} \left| \mathbf{x}^\top \tilde{\mathbf{M}}_{i,t}^{-1} \mathbf{x}_{a_s} \right|
\end{aligned} \tag{40}$$

$$\leq \epsilon_* \max_{\mathbf{x} \in S^{d-1}} \sqrt{\sum_{\substack{s \in [t] \\ i_s = i}} 1 \sum_{\substack{s \in [t] \\ i_s = i}} \left| \mathbf{x}^\top \tilde{\mathbf{M}}_{i,t}^{-1} \mathbf{x}_{a_s} \right|^2} \tag{41}$$

$$\leq \epsilon_* \sqrt{T_{i,t}} \sqrt{\max_{\mathbf{x} \in S^{d-1}} \mathbf{x}^\top \tilde{\mathbf{M}}_{i,t}^{-1} \mathbf{x}} \tag{42}$$

$$= \frac{\epsilon_* \sqrt{T_{i,t}}}{\sqrt{\lambda_{\min}(\tilde{\mathbf{M}}_{i,t})}}, \tag{43}$$

702 where we denote $S^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$, Eq.(40) follows from Holder’s inequality, Eq.(41)
703 follows by the Cauchy–Schwarz inequality, Eq.(42) holds because $\tilde{\mathbf{M}}_{i,t} \succeq \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top$, Eq.(43)
704 follows from the Courant-Fischer theorem.

705 For the last term in Eq.(38)

$$\left\| \tilde{\mathbf{M}}_{i,t}^{-1} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \eta_s \right\|_2 \leq \left\| \tilde{\mathbf{M}}_{i,t}^{-\frac{1}{2}} \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \eta_s \right\|_2 \left\| \tilde{\mathbf{M}}_{i,t}^{-\frac{1}{2}} \right\|_2 \quad (44)$$

$$= \frac{\left\| \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \eta_s \right\|_{\tilde{\mathbf{M}}_{i,t}^{-1}}}{\sqrt{\lambda_{\min}(\tilde{\mathbf{M}}_{i,t})}}, \quad (45)$$

706 where Eq.(44) follows by the Cauchy-Schwarz inequality and the inequality for the operator norm of
707 matrices, and Eq.(45) follows by the Courant-Fischer theorem.

708 Following Theorem 1 in [1], with probability at least $1 - \delta$ for some $\delta \in (0, 1)$, for any $i \in \mathcal{U}$, we
709 have:

$$\begin{aligned} \left\| \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \eta_s \right\|_{\tilde{\mathbf{M}}_{i,t}^{-1}} &\leq \sqrt{2 \log\left(\frac{u}{\delta}\right) + \log\left(\frac{\det(\tilde{\mathbf{M}}_{i,t})}{\det(\lambda \mathbf{I})}\right)} \\ &\leq \sqrt{2 \log\left(\frac{u}{\delta}\right) + d \log\left(1 + \frac{T_{i,t}}{\lambda d}\right)}, \end{aligned} \quad (46)$$

710 where $\det(\mathbf{M})$ denotes the determinant of matrix \mathbf{M} , Eq.(46) is because $\det(\tilde{\mathbf{M}}_{i,t}) \leq$
711 $\left(\frac{\text{trace}(\lambda \mathbf{I} + \sum_{\substack{s \in [t] \\ i_s = i}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top)}{d}\right)^d \leq \left(\frac{\lambda d + T_{i,t}}{d}\right)^d$, and $\det(\lambda \mathbf{I}) = \lambda^d$.

712 Plugging Eq.(46) into Eq. (45), then plugging Eq. (39), Eq.(43) and Eq.(45) into Eq.(38), we can get
713 that Eq.(73) holds with probability $\geq 1 - \delta$.

714 Then, with the item regularity assumption stated in Assumption 3.3, the technical Lemma J.1,
715 together with Lemma 7 in [23], with probability at least $1 - \delta$, for a particular user i , at any t such
716 that $T_{i,t} \geq \frac{16}{\tilde{\lambda}_x^2} \log\left(\frac{8d}{\tilde{\lambda}_x^2 \delta}\right)$, we have:

$$\lambda_{\min}(\tilde{\mathbf{M}}_{i,t}) \geq 2\tilde{\lambda}_x T_{i,t} + \lambda. \quad (47)$$

717 Based on the above reasoning, we have: if $T_{i,t} \geq \frac{16}{\tilde{\lambda}_x^2} \log\left(\frac{8d}{\tilde{\lambda}_x^2 \delta}\right)$, then with probability $\geq 1 - 2\delta$, we
718 have:

$$\begin{aligned} \left\| \hat{\boldsymbol{\theta}}_{i,t} - \boldsymbol{\theta}^{j(i)} \right\|_2 &\leq \frac{\beta(T_{i,t}, \frac{\delta}{u}) + \epsilon_* \sqrt{T_{i,t}}}{\sqrt{\lambda_{\min}(\tilde{\mathbf{M}}_{i,t})}} \\ &\leq \frac{\beta(T_{i,t}, \frac{\delta}{u}) + \epsilon_* \sqrt{T_{i,t}}}{\sqrt{2\tilde{\lambda}_x T_{i,t} + \lambda}} \\ &\leq \frac{\sqrt{\lambda} + \sqrt{2 \log\left(\frac{u}{\delta}\right) + d \log\left(1 + \frac{T_{i,t}}{\lambda d}\right)}}{\sqrt{2\tilde{\lambda}_x T_{i,t} + \lambda}} + \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}}, \end{aligned} \quad (48)$$

719 for any $i \in \mathcal{U}$.

720 Let

$$\frac{\sqrt{\lambda} + \sqrt{2 \log\left(\frac{u}{\delta}\right) + d \log\left(1 + \frac{T_{i,t}}{\lambda d}\right)}}{\sqrt{2\tilde{\lambda}_x T_{i,t} + \lambda}} + \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} < \frac{\gamma_1}{4}, \quad (49)$$

721 which is equivalent to

$$\frac{\sqrt{\lambda} + \sqrt{2 \log\left(\frac{u}{\delta}\right) + d \log\left(1 + \frac{T_{i,t}}{\lambda d}\right)}}{\sqrt{2\tilde{\lambda}_x T_{i,t} + \lambda}} < \frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}}, \quad (50)$$

722 where γ_1 is given in Definition 5.1.

723 Assume $\lambda \leq 2 \log(\frac{u}{\delta}) + d \log(1 + \frac{T_{i,t}}{\lambda d})$, which is typically held, then a sufficient condition for Eq.
724 (50) is:

$$\frac{2 \log(\frac{u}{\delta}) + d \log(1 + \frac{T_{i,t}}{\lambda d})}{2 \tilde{\lambda}_x T_{i,t}} < \frac{1}{4} (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2. \quad (51)$$

725 To satisfy the condition in Eq.(51), it is sufficient to show

$$\frac{2 \log(\frac{u}{\delta})}{2 \tilde{\lambda}_x T_{i,t}} < \frac{1}{8} (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2 \quad (52)$$

726 and

$$\frac{d \log(1 + \frac{T_{i,t}}{\lambda d})}{2 \tilde{\lambda}_x T_{i,t}} < \frac{1}{8} (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2. \quad (53)$$

727 From Eq.(52), we can get:

$$T_{i,t} \geq \frac{8 \log(\frac{u}{\delta})}{\tilde{\lambda}_x (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2}. \quad (54)$$

728 Following Lemma 9 in [23], we can get the following sufficient condition for Eq.(53):

$$T_{i,t} \geq \frac{8d \log(\frac{4}{\lambda \tilde{\lambda}_x (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2})}{\tilde{\lambda}_x (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2}. \quad (55)$$

729 Assume $\frac{u}{\delta} \geq \frac{4}{\lambda \tilde{\lambda}_x (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2}$, which is typically held, we can get that

$$T_{i,t} \geq \frac{8d}{\tilde{\lambda}_x (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2} \log(\frac{u}{\delta}) \quad (56)$$

730 is a sufficient condition for Eq.(49). Together with the condition that $T_{i,t} \geq \frac{16}{\tilde{\lambda}_x^2} \log(\frac{8d}{\tilde{\lambda}_x^2 \delta})$, we can get
731 that if

$$T_{i,t} \geq \max\{\frac{8d}{\tilde{\lambda}_x (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2} \log(\frac{u}{\delta}), \frac{16}{\tilde{\lambda}_x^2} \log(\frac{8d}{\tilde{\lambda}_x^2 \delta})\}, \forall i \in \mathcal{U}, \quad (57)$$

732 then with probability $\geq 1 - 2\delta$:

$$\|\hat{\theta}_{i,t} - \theta^{j(i)}\|_2 < \frac{\gamma_1}{4}, \forall i \in \mathcal{U}.$$

733 By Lemma 8 in [23], and Assumption 3.2 of user arrival uniformness, we have that for all

$$t \geq T_0 \triangleq 16u \log(\frac{u}{\delta}) + 4u \max\{\frac{8d}{\tilde{\lambda}_x (\frac{\gamma_1}{4} - \epsilon_* \sqrt{\frac{1}{2 \tilde{\lambda}_x}})^2} \log(\frac{u}{\delta}), \frac{16}{\tilde{\lambda}_x^2} \log(\frac{8d}{\tilde{\lambda}_x^2 \delta})\}, \quad (58)$$

734 with probability at least $1 - \delta$, condition in Eq.(57) is satisfied.

735 Therefore we have that for all $t \geq T_0$, with probability $\geq 1 - 3\delta$:

$$\|\hat{\theta}_{i,t} - \theta^{j(i)}\|_2 < \frac{\gamma_1}{4}, \forall i \in \mathcal{U}. \quad (59)$$

736 Next, we show that with Eq.(59), we can get that the RCLUMB keeps a “good partition”. First,
737 if we delete the edge (i, l) , then user i and user j belong to different *ground-truth clusters*, i.e.,
738 $\|\theta_i - \theta_l\|_2 > 0$. This is because by the deletion rule of the algorithm, the concentration bound,
739 and triangle inequality, $\|\theta_i - \theta_l\|_2 = \|\theta^{j(i)} - \theta^{j(l)}\|_2 \geq \|\hat{\theta}_{i,t} - \hat{\theta}_{l,t}\|_2 - \|\theta^{j(l)} - \theta_{l,t}\|_2 -$

740 $\|\theta^{j(i)} - \theta_{i,t}\|_2 > 0$. Second, we show that if $\|\theta_i - \theta_l\| \geq \gamma_1 > 2\epsilon_*\sqrt{\frac{2}{\lambda_x}}$, the RCLUMB
 741 algorithm will delete the edge (i, l) . This is because if $\|\theta_i - \theta_l\| \geq \gamma_1$, then by the trian-
 742 gle inequality, and $\|\hat{\theta}_{i,t} - \theta^{j(i)}\|_2 < \frac{\gamma_1}{4}$, $\|\hat{\theta}_{l,t} - \theta^{j(l)}\|_2 < \frac{\gamma_1}{4}$, $\theta_i = \theta^{j(i)}$, $\theta_l = \theta^{j(l)}$, we
 743 have $\|\hat{\theta}_{i,t} - \hat{\theta}_{l,t}\|_2 \geq \|\theta_i - \theta_l\| - \|\hat{\theta}_{i,t} - \theta^{j(i)}\|_2 - \|\hat{\theta}_{l,t} - \theta^{j(l)}\|_2 > \gamma_1 - \frac{\gamma_1}{4} - \frac{\gamma_1}{4} = \frac{\gamma_1}{2} >$
 744 $\frac{\sqrt{\lambda} + \sqrt{2\log(\frac{y}{\delta}) + d\log(1 + \frac{T_{i,t}}{\lambda d})}}{\sqrt{\lambda + 2\lambda_x T_{i,t}}} + \epsilon_*\sqrt{\frac{1}{2\lambda_x}} + \frac{\sqrt{\lambda} + \sqrt{2\log(\frac{y}{\delta}) + d\log(1 + \frac{T_{l,t}}{\lambda d})}}{\sqrt{\lambda + 2\lambda_x T_{l,t}}} + \epsilon_*\sqrt{\frac{1}{2\lambda_x}}$, which will trigger
 745 the deletion condition Line 10 in Algo.1.

746 From the above reasoning, we can get that at round t , any user within \bar{V}_t is ζ -close to i_t , and all the
 747 users belonging to $V_{j(i)}$ are contained in \bar{V}_t , which means the algorithm has done a “good partition”
 748 at t by Definition 5.5. \square

749 I Proof of Lemma 5.6

750 We prove the result in two situations: when $\bar{V}_t \in \mathcal{V}$ and when $\bar{V}_t \notin \mathcal{V}$.

751 (1) Situation 1: for any $t \geq T_0$ and $\bar{V}_t \in \mathcal{V}$, which means that the current user i_t is clustered
 752 completely correctly, i.e., $\bar{V}_t = V_{j(i_t)}$, therefore $\theta_l = \theta_{i_t}, \forall l \in \bar{V}_t$, then we have:

$$\begin{aligned} \hat{\theta}_{\bar{V}_t, t-1} - \theta_{i_t} &= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} r_s \right) - \theta_{i_t} \\ &= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} (\mathbf{x}_{a_s}^\top \theta_{i_t} + \epsilon_{a_s}^{i_s, s} + \eta_s) \right) - \theta_{i_t} \\ &= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} (\mathbf{x}_{a_s}^\top \theta_{i_t} + \epsilon_{a_s}^{i_s, s} + \eta_s) \right) - \theta_{i_t} \\ &= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left[\left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right) \theta_{i_t} - \lambda \theta_{i_t} + \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s \right] - \theta_{i_t} \\ &= -\lambda \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \theta_{i_t} + \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \eta_s. \end{aligned}$$

753 Therefore we have

$$\left| \mathbf{x}_a^\top (\hat{\theta}_{\bar{V}_t, t-1} - \theta_{i_t}) \right| \leq \lambda \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \theta_{i_t} \right| + \left| \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} \right| + \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s \right|. \quad (60)$$

754 Next, we bound the three terms in Eq.(60). For the first term:

$$\lambda \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \theta_{i_t} \right| \leq \lambda \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}} \sqrt{\lambda_{\max}(\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1})} \|\theta_{i_t}\|_2 \leq \sqrt{\lambda} \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}}, \quad (61)$$

755 where we use the inequality of matrix norm, the Cauchy–Schwarz inequality, $\|\theta_{i_t}\|_2 \leq 1$, and the
 756 fact that $\lambda_{\max}(\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}) = \frac{1}{\lambda_{\min}(\bar{\mathbf{M}}_{\bar{V}_t, t-1})} \leq \frac{1}{\lambda}$.

757 For the second term in Eq.(60):

$$\begin{aligned}
\left| \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \boldsymbol{\epsilon}_{a_s}^{i_s, s} \right| &\leq \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \boldsymbol{\epsilon}_{a_s}^{i_s, s} \right| \\
&\leq \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \|\boldsymbol{\epsilon}_{a_s}^{i_s, s}\|_\infty \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \right| \\
&\leq \epsilon_* \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \right|, \tag{62}
\end{aligned}$$

758 where in the second inequality we use the Holder's inequality.

759 For the last term, with probability at least $1 - \delta$:

$$\left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s \right| \leq \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}} \left\| \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s \right\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}} \tag{63}$$

$$\begin{aligned}
&\leq \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}} \sqrt{2 \log\left(\frac{1}{\delta}\right) + \log\left(\frac{\det(\bar{\mathbf{M}}_{\bar{V}_t, t-1})}{\det(\lambda \mathbf{I})}\right)} \\
&\leq \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}} \sqrt{2 \log\left(\frac{1}{\delta}\right) + d \log\left(1 + \frac{T}{\lambda d}\right)}, \tag{64}
\end{aligned}$$

760 where the second inequality follows by Theorem 1 in [1], Eq.(64) is because $\det(\bar{\mathbf{M}}_{\bar{V}_t, t-1}) \leq$

$$\left(\frac{\text{trace}(\lambda \mathbf{I} + \sum_{\substack{s \in [t] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top)}{d} \right)^d \leq \left(\frac{\lambda d + T_{\bar{V}_t, t}}{d} \right)^d \leq \left(\frac{\lambda d + T}{d} \right)^d, \text{ and } \det(\lambda \mathbf{I}) = \lambda^d. \tag{65}$$

762 Plugging Eq.(61), Eq.(62) and Eq.(64) into Eq.(60), we can prove Lemma 5.6 in situation 1, i.e., for
763 any $t \geq T_0$ and $\bar{V}_t \in V$, with probability at least $1 - \delta$:

$$\left| \mathbf{x}_a^\top (\hat{\boldsymbol{\theta}}_{\bar{V}_t, t-1} - \boldsymbol{\theta}_{i_t}) \right| \leq \epsilon_* \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \right| + \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}} \left(\sqrt{\lambda} + \sqrt{2 \log\left(\frac{1}{\delta}\right) + d \log\left(1 + \frac{T}{\lambda d}\right)} \right). \tag{65}$$

764 (2) Situation 2: for any $t \geq T_0$ and $\bar{V}_t \notin \mathcal{V}$, which means that the current user is *misclustered* by
765 the algorithm, i.e., $\bar{V}_t \neq V_{j(i_t)}$, but with Lemma H.1, with probability at least $1 - 3\delta$, the current

partition is a “good partition”, i.e., $\|\theta_l - \theta_{i_t}\|_2 \leq 2\epsilon_* \sqrt{\frac{2}{\lambda_x}}, \forall l \in \bar{V}_t$, we have:

$$\begin{aligned}
\hat{\theta}_{\bar{V}_t, t-1} - \theta_{i_t} &= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} r_s \right) - \theta_{i_t} \\
&= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} (\mathbf{x}_{a_s}^\top \theta_{i_s} + \epsilon_{a_s}^{i_s, s} + \eta_s) \right) - \theta_{i_t} \\
&= \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s + \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top \theta_{i_s} - \theta_{i_t} \\
&= \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s + \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\theta_{i_s} - \theta_{i_t}) \\
&\quad + \bar{M}_{\bar{V}_t, t-1}^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right) \theta_{i_t} - \lambda \bar{M}_{\bar{V}_t, t-1}^{-1} \theta_{i_t} - \theta_{i_t} \\
&= \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s + \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\theta_{i_s} - \theta_{i_t}) \\
&\quad - \lambda \bar{M}_{\bar{V}_t, t-1}^{-1} \theta_{i_t}.
\end{aligned}$$

Thus, with Lemma 5.7 and with the previous reasoning, with probability at least $1 - 5\delta$, we have:

$$\begin{aligned}
\left| \mathbf{x}_a^\top (\hat{\theta}_{\bar{V}_t, t-1} - \theta_{i_t}) \right| &\leq \lambda \left| \mathbf{x}_a^\top \bar{M}_{\bar{V}_t, t-1}^{-1} \theta_{i_t} \right| + \left| \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_a^\top \bar{M}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} \right| + \left| \mathbf{x}_a^\top \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s \right| \\
&\quad + \left| \mathbf{x}_a^\top \bar{M}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\theta_{i_s} - \theta_{i_t}) \right| \\
&\leq \epsilon_* \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \left| \mathbf{x}_a^\top \bar{M}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \right| + \|\mathbf{x}_a\|_{\bar{M}_{\bar{V}_t, t-1}^{-1}} \left(\sqrt{\lambda} + \sqrt{2 \log\left(\frac{1}{\delta}\right) + d \log\left(1 + \frac{T}{\lambda d}\right)} \right) \\
&\quad + \frac{\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}}.
\end{aligned}$$

Therefore, combining situation 1 and situation 2, the result of Lemma 5.6 then follows.

J Technical Lemmas and Their Proofs

We first prove the following technical lemma which is used to prove Lemma H.1.

Lemma J.1. *Under Assumption 3.3, at any time t , for any fixed unit vector $\theta \in \mathbb{R}^d$*

$$\mathbb{E}_t[(\theta^\top \mathbf{x}_{a_t})^2 | \mathcal{A}_t] \geq \tilde{\lambda}_x \triangleq \int_0^{\lambda_x} (1 - e^{-\frac{(\lambda_x - x)^2}{2\sigma^2}})^C dx. \quad (66)$$

Proof. The proof of this lemma mainly follows the proof of Claim 1 in [11], but with more careful analysis, since their assumption is more stringent than ours.

Denote the feasible arms at round t by $\mathcal{A}_t = \{\mathbf{x}_{t,1}, \mathbf{x}_{t,2}, \dots, \mathbf{x}_{t,|\mathcal{A}_t|}\}$. Consider the corresponding i.i.d. random variables $\theta_i = (\theta^\top \mathbf{x}_{t,i})^2 - \mathbb{E}_t[(\theta^\top \mathbf{x}_{t,i})^2 | \mathcal{A}_t]$, $i = 1, 2, \dots, |\mathcal{A}_t|$. By Assumption 3.3, θ_i ’s are sub-Gaussian random variables with variance bounded by σ^2 . Therefore, we have that

777 for any $\alpha > 0$ and any $i \in [|\mathcal{A}_t|]$:

$$\mathbb{P}_t(\theta_i < -\alpha | \mathcal{A}_t) \leq e^{-\frac{\alpha^2}{2\sigma^2}},$$

778 where $\mathbb{P}_t(\cdot)$ is the shorthand for the conditional probability
779 $\mathbb{P}(\cdot | (i_1, \mathcal{A}_1, r_1), \dots, (i_{t-1}, \mathcal{A}_{t-1}, r_{t-1}), i_t)$.

780 We also have that $\mathbb{E}_t[(\boldsymbol{\theta}^\top \mathbf{x}_{t,i})^2 | \mathcal{A}_t] = \mathbb{E}_t[\boldsymbol{\theta}^\top \mathbf{x}_{t,i} \mathbf{x}_{t,i}^\top \boldsymbol{\theta} | \mathcal{A}_t] \geq \lambda_{\min}(\mathbb{E}_{\mathbf{x} \sim \rho}[\mathbf{x} \mathbf{x}^\top]) \geq \lambda_x$ by As-
781 sumption 3.3. With the above inequalities, we can get

$$\mathbb{P}_t\left(\min_{i=1, \dots, |\mathcal{A}_t|} (\boldsymbol{\theta}^\top \mathbf{x}_{t,i})^2 \geq \lambda_x - \alpha | \mathcal{A}_t\right) \geq (1 - e^{-\frac{\alpha^2}{2\sigma^2}})^C,$$

782 where C is the upper bound of $|\mathcal{A}_t|$.

783 Therefore, we have

$$\begin{aligned} \mathbb{E}_t[(\boldsymbol{\theta}^\top \mathbf{x}_{a_t})^2 | \mathcal{A}_t] &\geq \mathbb{E}_t\left[\min_{i=1, \dots, |\mathcal{A}_t|} (\boldsymbol{\theta}^\top \mathbf{x}_{t,i})^2 | \mathcal{A}_t\right] \\ &\geq \int_0^\infty \mathbb{P}_t\left(\min_{i=1, \dots, |\mathcal{A}_t|} (\boldsymbol{\theta}^\top \mathbf{x}_{t,i})^2 \geq x | \mathcal{A}_t\right) dx \\ &\geq \int_0^{\lambda_x} (1 - e^{-\frac{(\lambda_x - x)^2}{2\sigma^2}})^C dx \triangleq \tilde{\lambda}_x \end{aligned}$$

784

□

785 Finally, we prove the following lemma which is used in the proof of Theorem 5.3.

Lemma J.2.

$$\sum_{t=T_0+1}^T \min\{\mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{V_j, t-1}^{-1}}^2, 1\} \leq 2d \log(1 + \frac{T}{\lambda d}), \forall j \in [m]. \quad (67)$$

Proof.

$$\begin{aligned} \det(\overline{\mathbf{M}}_{V_j, T}) &= \det\left(\overline{\mathbf{M}}_{V_j, T-1} + \mathbb{I}\{i_T \in V_j\} \mathbf{x}_{a_T} \mathbf{x}_{a_T}^\top\right) \\ &= \det(\overline{\mathbf{M}}_{V_j, T-1}) \det\left(\mathbf{I} + \mathbb{I}\{i_T \in V_j\} \overline{\mathbf{M}}_{V_j, T-1}^{-\frac{1}{2}} \mathbf{x}_{a_T} \mathbf{x}_{a_T}^\top \overline{\mathbf{M}}_{V_j, T-1}^{-\frac{1}{2}}\right) \\ &= \det(\overline{\mathbf{M}}_{V_j, T-1}) \left(1 + \mathbb{I}\{i_T \in V_j\} \|\mathbf{x}_{a_T}\|_{\overline{\mathbf{M}}_{V_j, T-1}^{-1}}^2\right) \\ &= \det(\overline{\mathbf{M}}_{V_j, T_0}) \prod_{t=T_0+1}^T \left(1 + \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{V_j, t-1}^{-1}}^2\right) \\ &\geq \det(\lambda \mathbf{I}) \prod_{t=T_0+1}^T \left(1 + \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{V_j, t-1}^{-1}}^2\right). \end{aligned} \quad (68)$$

786 $\forall x \in [0, 1]$, we have $x \leq 2 \log(1 + x)$. Therefore

$$\begin{aligned} \sum_{t=T_0+1}^T \min\{\mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{V_j, t-1}^{-1}}^2, 1\} &\leq 2 \sum_{t=T_0+1}^T \log\left(1 + \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{V_j, t-1}^{-1}}^2\right) \\ &= 2 \log\left(\prod_{t=T_0+1}^T (1 + \mathbb{I}\{i_t \in V_j\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{V_j, t-1}^{-1}}^2)\right) \\ &\leq 2[\log(\det(\overline{\mathbf{M}}_{V_j, T})) - \log(\det(\lambda \mathbf{I}))] \\ &\leq 2 \log\left(\frac{\text{trace}(\lambda \mathbf{I} + \sum_{t=1}^T \mathbb{I}\{i_t \in V_j\} \mathbf{x}_{a_t} \mathbf{x}_{a_t}^\top)}{\lambda d}\right)^d \\ &\leq 2d \log(1 + \frac{T}{\lambda d}). \end{aligned} \quad (69)$$

787

□

788 K Algorithms of RSCLUMB

789 This section introduces the Robust Set-based Clustering of Misspecified Bandits Algorithm
 790 (RSCLUMB). Unlike RCLUMB, which maintains a graph-based clustering structure, RSCLUMB
 791 maintains a set-based clustering structure. Besides, RCLUMB only splits clusters during the learning
 792 process, while RSCLUMB allows both split and merge operations. A brief illustration is that the
 793 agent will split a user out of its current set(cluster) if it finds an inconsistency between the user and its
 794 set, and if there are two clusters whose estimated preferences are close enough, the agent will merge
 795 them. A detailed discussion of the connection between the graph structure and the set structure can
 796 be found in [25].

797 Now we introduce the details of RSCLUMB. The algorithm first initializes a single set S_1 containing
 798 all users and updates it during the learning process. The whole learning process consists of phases
 799 (Algo. 2 Line 3), where the s -th phase contains 2^s rounds. At the beginning of each phase, the
 800 agent marks all users as "unchecked", and if a user comes later, it will be marked as "checked". If all
 801 users in a cluster are checked, then this cluster will be marked as "checked" meaning it is an accurate
 802 cluster in the current phase. With this mechanism, every phase can maintain an accuracy level, and
 803 the agent can put the accurate clusters aside and focus on exploring the inaccurate ones. For each
 804 cluster V_j , the algorithm maintains two estimated vectors $\hat{\theta}_{V_j}$ and $\tilde{\theta}_{V_j}$, where the $\hat{\theta}_{V_j}$ is similar to
 805 the $\hat{\theta}_{\bar{V}_j}$ in RCLUMB and is used for the recommendation, while the $\tilde{\theta}_{V_j}$ is the average of all the
 806 estimated user preference vectors in this cluster and is used for the split and merge operations.

807 At time t in phase s , the user i_τ comes with the item set \mathcal{D}_τ , where τ represents the index of total
 808 time steps. Then the algorithm determines the cluster and makes a cluster-based recommendation.
 809 This process is similar to RCLUMB. After updating the information (Algo. 2 Line12), the agent
 810 checks if a split or a merge is possible (Algo. 2 Line13-17).

811 By our assumption, users in the same cluster have the same vectors. So a cluster can be regarded
 812 as a good cluster only when all the estimated user vectors are close to the estimated cluster vector.
 813 We call a user is consistent with the cluster if their estimated vectors are close enough. If a user is
 814 inconsistent with its current cluster, the agent will split it out. Two clusters are consistent when their
 815 estimated vectors are close, and the agent will merge them.

816 RSCLUMB maintains two sets of estimated cluster vectors: (i) cluster-level estimation with integrated
 817 user information, which is for recommendations (Line 12 and Line 10 in Algo.2); (ii) the average of
 818 estimated user vectors, which is used for robust clustering (Line 3 in Algo.3 and Line 2 in Algo.4).
 819 The previous set-based CB work [25] only uses (i) for both recommendations and clustering, which
 820 would lead to erroneous clustering under misspecifications, and cannot get any non-vacuous regret
 821 bound in CBMUM.

822 L Main Theorem and Lemmas of RSCLUMB

823 **Theorem L.1** (main result on regret bound for RSCLUMB). *With the same assumptions in Theorem*
 824 *5.3, the expected regret of the RSCLUMB algorithm for T rounds satisfies:*

$$\begin{aligned} R(T) &\leq O\left(u \left(\frac{d}{\tilde{\lambda}_x(\gamma_1 - \zeta_1)^2} + \frac{1}{\tilde{\lambda}_x^2} \right) \log T + \frac{\epsilon_* \sqrt{dT}}{\tilde{\lambda}_x^{1.5}} + \epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T + \epsilon_* \sqrt{\frac{1}{\tilde{\lambda}_x}} T\right) \\ &\leq O(\epsilon_* T \sqrt{md \log T} + d \sqrt{mT} \log T) \end{aligned} \quad (70)$$

825 **Lemma L.2.** *For RSCLUMB, we use T_1 to represent the corresponding T_0 of RCLUMB. Then :*

$$\begin{aligned} T_1 &\triangleq 16u \log\left(\frac{u}{\delta}\right) + 4u \max\left\{\frac{16}{\tilde{\lambda}_x^2} \log\left(\frac{8d}{\tilde{\lambda}_x^2 \delta}\right), \frac{8d}{\tilde{\lambda}_x(\frac{\gamma_1}{6} - \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}})^2} \log\left(\frac{u}{\delta}\right)\right\} \\ &= O\left(u \left(\frac{d}{\tilde{\lambda}_x(\gamma_1 - \zeta_1)^2} + \frac{1}{\tilde{\lambda}_x^2} \right) \log \frac{1}{\delta}\right) \end{aligned}$$

Algorithm 2 Robust Set-based Clustering of Misspecified Bandits Algorithm (RSCLUMB)

- 1: **Input:** Deletion parameter $\alpha_1, \alpha_2 > 0$, $f(T) = \sqrt{\frac{1+\ln(1+T)}{1+T}}$, $\lambda, \beta, \epsilon_* > 0$.
 - 2: **Initialization:**
 - $\mathbf{M}_{i,0} = 0_{d \times d}$, $\mathbf{b}_{i,0} = 0_{d \times 1}$, $T_{i,0} = 0$, $\forall i \in \mathcal{U}$;
 - Initialize the set of cluster indexes by $J = \{1\}$ and the single cluster \mathcal{S}_1 by $\mathbf{M}_1 = 0_{d \times d}$, $\mathbf{b}_1 = 0_{d \times 1}$, $T_1 = 0$, $C_1 = \mathcal{U}$, $j(i) = 1$, $\forall i$.
 - 3: **for all** $s = 1, 2, \dots$ **do**
 - 4: Mark every user unchecked for each cluster.
 - 5: For each cluster V_j , compute $\tilde{T}_{V_j} = T_{V_j}$, $\hat{\boldsymbol{\theta}}_{V_j} = (\lambda \mathbf{I} + \mathbf{M}_{V_j})^{-1} \mathbf{b}_{V_j}$, $\tilde{\boldsymbol{\theta}}_{V_j} = \frac{\sum_{i \in V_j} \hat{\boldsymbol{\theta}}_i}{|V_j|}$
 - 6: **for all** $t = 1, 2, \dots, T$ **do**
 - 7: Compute $\tau = 2^s - 2 + t$
 - 8: Receive the user i_τ and the decision set \mathcal{D}_τ
 - 9: Determine the cluster index $j = j(i_\tau)$
 - 10: Recommend item a_τ with the largest UCB index as shown in Eq. (5)
 - 11: Received the feedback r_τ .
 - 12: Update the information:
$$\begin{aligned} \mathbf{M}_{i_\tau, \tau} &= \mathbf{M}_{i_\tau, \tau-1} + \mathbf{x}_{a_\tau} \mathbf{x}_{a_\tau}^\top, \mathbf{b}_{i_\tau, \tau} = \mathbf{b}_{i_\tau, \tau-1} + r_\tau \mathbf{x}_{a_\tau}, \\ T_{i_\tau, \tau} &= T_{i_\tau, \tau-1} + 1, \hat{\boldsymbol{\theta}}_{i_\tau, \tau} = (\lambda \mathbf{I} + \mathbf{M}_{i_\tau, \tau})^{-1} \mathbf{b}_{i_\tau, \tau} \\ \mathbf{M}_{V_j, \tau} &= \mathbf{M}_{V_j, \tau-1} + \mathbf{x}_{a_\tau} \mathbf{x}_{a_\tau}^\top, \mathbf{b}_{V_j, \tau} = \mathbf{b}_{V_j, \tau-1} + r_\tau \mathbf{x}_{a_\tau}, \\ T_{V_j, \tau} &= T_{V_j, \tau-1} + 1, \hat{\boldsymbol{\theta}}_{V_j, \tau} = (\lambda \mathbf{I} + \mathbf{M}_{V_j, \tau})^{-1} \mathbf{b}_{V_j, \tau}, \\ \tilde{\boldsymbol{\theta}}_{V_j, \tau} &= \frac{\sum_{i \in V_j} \hat{\boldsymbol{\theta}}_{i, \tau}}{|V_j|} \end{aligned}$$
 - 13: **if** i_τ is unchecked **then**
 - 14: Run **Split**
 - 15: Mark user i_τ has been checked
 - 16: Run **Merge**
-

Algorithm 3 Split

- 1: Define $F(T) = \sqrt{\frac{1+\ln(1+T)}{1+T}}$
- 2: **if** $\|\hat{\boldsymbol{\theta}}_{i_\tau, \tau} - \tilde{\boldsymbol{\theta}}_{V_j, \tau}\| > \alpha_1(F(T_{i_\tau, \tau}) + F(T_{V_j, \tau})) + \alpha_2 \epsilon_*$ **then**
- 3: Split user i_τ from cluster V_j and form a new cluster V'_j of user i_τ

$$\begin{aligned} \mathbf{M}_{V_j, \tau} &= \mathbf{M}_{V_j, \tau} - \mathbf{M}_{i_\tau, \tau}, \mathbf{b}_{V_j} = \mathbf{b}_{V_j} - \mathbf{b}_{i_\tau, \tau}, \\ T_{V_j, \tau} &= T_{V_j, \tau} - T_{i_\tau, \tau}, C_{j, \tau} = C_{j, \tau} - \{i_\tau\}, \\ \mathbf{M}_{V'_j, \tau} &= \mathbf{M}_{i_\tau, \tau}, \mathbf{b}_{V'_j, \tau} = \mathbf{b}_{i_\tau, \tau}, \\ T_{V'_j, \tau} &= T_{i_\tau, \tau}, C_{j', \tau} = \{i_\tau\} \end{aligned}$$

826 **Lemma L.3.** For RSCLUMB, after $2T_1 + 1$ rounds: in each phase, after the first u rounds, with
827 probability at least $1 - 5\delta$:

$$\begin{aligned} \left| \mathbf{x}_a^\top (\boldsymbol{\theta}_{i_t} - \hat{\boldsymbol{\theta}}_{\bar{V}_t, t-1}) \right| &\leq \left(\frac{3\epsilon_* \sqrt{2d}}{2\tilde{\lambda}_x^{\frac{3}{2}}} + 6\epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} \right) \mathbb{I}\{\bar{V}_t \notin V\} + \beta \|\mathbf{x}_a\|_{\bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1}} + \epsilon_* \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \left| \mathbf{x}_a^\top \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \mathbf{x}_{a_s} \right| \\ &\triangleq \left(\frac{3\epsilon_* \sqrt{2d}}{2\tilde{\lambda}_x^{\frac{3}{2}}} + 6\epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} \right) \mathbb{I}\{\bar{V}_t \notin V\} + C_{a,t} \end{aligned}$$

Algorithm 4 Merge

1: **for** any two checked clusters V_{j_1}, V_{j_2} satisfying

$$\|\tilde{\theta}_{j_1} - \tilde{\theta}_{j_2}\| < \frac{\alpha_1}{2}(F(T_{V_{j_1}}) + F(T_{V_{j_2}})) + \frac{\alpha_2}{2}\epsilon_*$$

do

2: Merge them:

$$\begin{aligned} \mathbf{M}_{V_{j_1}} &= \mathbf{M}_{j_1} + \mathbf{M}_{j_2}, \mathbf{b}_{V_{j_1}} = \mathbf{b}_{V_{j_1}} + \mathbf{b}_{V_{j_2}}, \\ T_{V_{j_1}} &= T_{V_{j_1}} + T_{V_{j_2}}, C_{V_{j_1}} = C_{V_{j_1}} \cup C_{V_{j_2}} \end{aligned}$$

3: Set $j(i) = j_1, \forall i \in j_2$, delete V_{j_2}

828 M Proof of Lemma L.3

$$\begin{aligned} |\mathbf{x}_a^\top(\theta_i - \hat{\theta}_{\bar{V}_t, t})| &= |\mathbf{x}_a^\top(\theta_i - \theta_{V_t})| + |\mathbf{x}_a^\top(\hat{\theta}_{\bar{V}_t, t} - \theta_{V_t})| \\ &\leq \|\mathbf{x}_a^\top\| \|\theta_i - \theta_{V_t}\| + |\mathbf{x}_a^\top(\hat{\theta}_{\bar{V}_t, t} - \theta_{V_t})| \\ &\leq 6\epsilon_* \sqrt{\frac{1}{2\lambda_x}} + |\mathbf{x}_a^\top(\hat{\theta}_{\bar{V}_t, t} - \theta_{V_t})| \end{aligned} \quad (71)$$

829 where the last inequality holds due to the fact $\|\mathbf{x}_a\| \leq 1$ and the condition of "split" and "merge".

830 For $|\mathbf{x}_a^\top(\hat{\theta}_{\bar{V}_t, t} - \theta_{V_t})|$:

$$\begin{aligned} \hat{\theta}_{\bar{V}_t, t-1} - \theta_{V_t} &= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} r_s \right) - \theta_{V_t} \\ &= \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right)^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} (\mathbf{x}_{a_s}^\top \theta_{i_s} + \epsilon_{a_s}^{i_s, s} + \eta_s) \right) - \theta_{V_t} \\ &= \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s + \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top \theta_{i_s} - \theta_{V_t} \\ &= \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s + \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\theta_{i_s} - \theta_{V_t}) \\ &\quad + \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \left(\sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top + \lambda \mathbf{I} \right) \theta_{V_t} - \lambda \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \theta_{V_t} - \theta_{V_t} \\ &= \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \epsilon_{a_s}^{i_s, s} + \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \eta_s + \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \sum_{\substack{s \in [t-1] \\ i_s \in \bar{V}_t}} \mathbf{x}_{a_s} \mathbf{x}_{a_s}^\top (\theta_{i_s} - \theta_{V_t}) \\ &\quad - \lambda \bar{\mathbf{M}}_{\bar{V}_t, t-1}^{-1} \theta_{V_t}. \end{aligned}$$

831 Thus, with the same method in Lemma 5.7 but replace $\zeta = 4\epsilon_* \sqrt{\frac{1}{2\lambda_x}}$ with $\zeta_1 = 6\epsilon_* \sqrt{\frac{1}{2\lambda_x}}$, and with
832 the previous reasoning, with probability at least $1 - 5\delta$, we have:

$$|\mathbf{x}_a^\top(\hat{\theta}_{\bar{V}_t, t} - \theta_{V_t})| \leq C_{a_t} + \frac{3\epsilon_* \sqrt{2d}}{2\tilde{\lambda}_x^{\frac{3}{2}}} \quad (72)$$

833 The lemma can be concluded.

834 N Proof of Lemma L.2

835 With the analysis in the proof of Lemma H.1, with probability at least $1 - \delta$:

$$\left\| \hat{\theta}_{i,t} - \theta^{j(i)} \right\|_2 \leq \frac{\beta(T_{i,t}, \frac{\delta}{u}) + \epsilon_* \sqrt{T_{i,t}}}{\sqrt{\lambda + \lambda_{\min}(\mathbf{M}_{i,t})}}, \forall i \in \mathcal{U}, \quad (73)$$

836 and the estimated error of the current cluster $\left\| \tilde{\theta}^{j(i)} - \theta^{j(i)} \right\|$ also satisfies this inequality. For
837 set-based clustering structure, to ensure for each user there is only one ζ -close cluster, we let:

$$\frac{\beta(T_{i,t}, \frac{\delta}{u}) + \epsilon_* \sqrt{T_{i,t}}}{\sqrt{\lambda + \lambda_{\min}(\mathbf{M}_{i,t})}} \leq \frac{\gamma_1}{6} \quad (74)$$

838 By assuming $\lambda < 2 \log(\frac{u}{\delta}) + d \log(1 + \frac{T_{i,t}}{\lambda d})$, we can simplify it to

$$\frac{2 \log(\frac{u}{\delta}) + d \log(1 + \frac{T_{i,t}}{\lambda d})}{2\tilde{\lambda}_x T_{i,t}} < \frac{1}{4} \left(\frac{\gamma_1}{6} - \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} \right)^2 \quad (75)$$

839 which can be proved by $\frac{2 \log(\frac{u}{\delta})}{2\tilde{\lambda}_x T_{i,t}} \leq \frac{1}{8} \left(\frac{\gamma_1}{6} - \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} \right)^2$ and $\frac{d \log(1 + \frac{T_{i,t}}{\lambda d})}{2\tilde{\lambda}_x T_{i,t}} \leq \frac{1}{8} \left(\frac{\gamma_1}{6} - \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} \right)^2$. It's
840 obvious that the former one can be satisfied by $T_{i,t} \geq \frac{8 \log(u/\delta)}{\tilde{\lambda}_x (\frac{\gamma_1}{6} - \epsilon_* \sqrt{1/2\tilde{\lambda}_x})^2}$. As for the latter one, by

841 [23] Lemma 9, we can get $T_{i,t} \geq \frac{8d \log(\frac{16}{\tilde{\lambda}_x \lambda (\frac{\gamma_1}{6} - \epsilon_* \sqrt{1/2\tilde{\lambda}_x})^2})}{4\tilde{\lambda}_x (\frac{\gamma_1}{6} - \epsilon_* \sqrt{1/2\tilde{\lambda}_x})^2}$. By assuming $\frac{u}{\delta} \geq \frac{16}{4\tilde{\lambda}_x \lambda (\frac{\gamma_1}{6} - \epsilon_* \sqrt{2/4\tilde{\lambda}_x})^2}$,
842 the lemma is proved.

843 O Proof of Theorem L.1

844 After $2T_1$ rounds, in each phase, at most u times split operations will happen, we use $u \log(T)$ to
845 bound the regret generated in these rounds. Then in the remained rounds the cluster num will be no
846 more than m .

847 For the instantaneous regret R_t at round t , with probability at least $1 - 2\delta$ for some $\delta \in (0, \frac{1}{2})$:

$$\begin{aligned} R_t &= (\mathbf{x}_{a_t^*}^\top \theta_{i_t} + \epsilon_{a_t^*}^{i_t, t}) - (\mathbf{x}_{a_t}^\top \theta_{i_t} + \epsilon_{a_t}^{i_t, t}) \\ &= \mathbf{x}_{a_t^*}^\top (\theta_{i_t} - \hat{\theta}_{\bar{V}_t, t-1}) + (\mathbf{x}_{a_t^*}^\top \hat{\theta}_{\bar{V}_t, t-1} + C_{a_t^*, t}) - (\mathbf{x}_{a_t}^\top \hat{\theta}_{\bar{V}_t, t-1} + C_{a_t, t}) \\ &\quad + \mathbf{x}_{a_t}^\top (\hat{\theta}_{\bar{V}_t, t-1} - \theta_{i_t}) + C_{a_t, t} - C_{a_t^*, t} + (\epsilon_{a_t^*}^{i_t, t} - \epsilon_{a_t}^{i_t, t}) \\ &\leq 2C_{a_t} + 2\epsilon_* + (12\epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} + \frac{3\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}}) \mathbb{I}(\bar{V}_t \notin V) \end{aligned} \quad (76)$$

848 where the last inequality holds due to the UCB arm selection strategy, the concentration bound given
849 in Lemma L.3 and the fact that $\|\epsilon^{i, t}\|_\infty \leq \epsilon_*$.

850 Define such events. Let:

$$\mathcal{E}_2 = \{\text{All clusters } \bar{V}_t \text{ only contain users who satisfy } \|\tilde{\theta}_i - \tilde{\theta}_{\bar{V}_t}\| \leq \alpha_1 \left(\sqrt{\frac{1 + \log(1 + T_{i,t})}{1 + T_{i,t}}} + \sqrt{\frac{1 + \log(1 + T_{\bar{V}_t, t})}{1 + T_{\bar{V}_t, t}}} \right) + \alpha_2 \epsilon_*\}$$

851

$$\mathcal{E}_3 = \{r_t \leq 2C_{a_t} + 2\epsilon_* + 12\epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} + \frac{3\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}}\}$$

852

$$\mathcal{E}' = \mathcal{E}_2 \cap \mathcal{E}_3$$

853 From previous analysis, we can know that $\mathbb{P}(\mathcal{E}_2) \geq 1 - 3\delta$ and $\mathbb{P}(\mathcal{E}_3) \geq 1 - 2\delta$, thus $\mathbb{P}(\mathcal{E}' \geq 1 - 5\delta)$.

854 By taking $\delta = \frac{1}{T}$, we can get:

$$\begin{aligned} E(R_t) &= P(\mathcal{E}) \mathbb{I}\{\mathcal{E}\} R_t + P(\bar{\mathcal{E}}) \mathbb{I}\{\bar{\mathcal{E}}\} R_t \\ &\leq \mathbb{I}\{\mathcal{E}\} R_t + 5 \\ &\leq 2T_1 + 2\epsilon_* T + (12\epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} + \frac{3\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}}) T + 2 \sum_{2T_1}^T C_{a_t} + 5 \end{aligned} \quad (77)$$

855 Now we need to bound $2 \sum_{t=1}^T C_{a_t}$. We already know that after $2T_1$ rounds, in each phase k after
 856 the first u rounds, there will be at most m clusters
 857 Consider phase k , for simplicity, ignore the first u rounds. For the first term in C_{a_t} :

$$\begin{aligned}
 \sum_{t=T_{k-1}}^{T_k} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{\overline{V}_{t,t-1}}}^{-1} &= \sum_{t=T_{k-1}}^{T_k} \sum_{j=1}^{m_t} \mathbb{I}\{i \in \overline{V}_{t,j}\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{\overline{V}_{t,j}}}^{-1} \\
 &\leq \sum_{j=1}^{m_t} \sqrt{\sum_{t=T_{k-1}}^{T_k} \mathbb{I}\{i \in \overline{V}_{t,j}\} \sum_{t=T_{k-1}}^{T_k} \mathbb{I}\{i \in \overline{V}_{t,j}\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{\overline{V}_{t,j}}}^2} \\
 &\leq \sum_{j=1}^{m_t} \sqrt{2T_{k,j} d \log(1 + \frac{T}{\lambda d})} \\
 &\leq \sqrt{2m(T_k - T_{k-1}) d \log(1 + \frac{T}{\lambda d})}
 \end{aligned} \tag{78}$$

858 For all phases:

$$\begin{aligned}
 \sum_{k=1}^s \sqrt{2m(T_{k+1} - T_k) d \log(1 + \frac{T}{\lambda d})} &\leq \sqrt{2 \sum_{k=1}^s 1 \sum_{k=1}^s (T_{k+1} - T_k) m d \log(1 + \frac{T}{\lambda d})} \\
 &\leq \sqrt{2mdT \log(T) \log(1 + \frac{T}{\lambda d})}
 \end{aligned} \tag{79}$$

859 Similarly, for the second term in C_{a_t} :

$$\begin{aligned}
 \sum_{t=T_{k-1}}^{T_k} \sum_{\substack{s \in [t-1] \\ i_s \in \overline{V}_t}} \epsilon_* |\mathbf{x}_{a_t}^T \overline{\mathbf{M}}_{\overline{V}_{t,t-1}}^{-1} \mathbf{x}_{a_s}| &= \sum_{t=T_{k-1}}^{T_k} \sum_{j=1}^{m_t} \mathbb{I}\{i \in \overline{V}_{t,j}\} \sum_{\substack{s \in [t-1] \\ i_s \in \overline{V}_{t,j}}} \epsilon_* |\mathbf{x}_{a_t}^T \overline{\mathbf{M}}_{\overline{V}_{t,j}}^{-1} \mathbf{x}_{a_s}| \\
 &\leq \epsilon_* \sum_{t=T_{k-1}}^{T_k} \sum_{j=1}^{m_t} \mathbb{I}\{i \in \overline{V}_{t,j}\} \sqrt{\sum_{\substack{s \in [t-1] \\ i_s \in \overline{V}_{t,j}}} 1 \sum_{\substack{s \in [t-1] \\ i_s \in \overline{V}_{t,j}}} |\mathbf{x}_{a_t}^T \overline{\mathbf{M}}_{\overline{V}_{t,j}}^{-1} \mathbf{x}_{a_s}|^2} \\
 &\leq \epsilon_* \sum_{t=T_{k-1}}^{T_k} \sum_{j=1}^{m_t} \mathbb{I}\{i \in \overline{V}_{t,j}\} \sqrt{T_{k,j} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{\overline{V}_{t,j}}}^2} \\
 &\leq \epsilon_* \sum_{t=T_{k-1}}^{T_k} \sqrt{\sum_{j=1}^{m_t} \mathbb{I}\{i \in \overline{V}_{t,j}\} \sum_{j=1}^{m_t} \mathbb{I}\{i \in \overline{V}_{t,j}\} T_{k,j} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{\overline{V}_{t,j}}}^2} \\
 &\leq \epsilon_* \sqrt{(T_k - T_{k-1})} \sum_{t=T_{k-1}}^{T_k} \sqrt{\sum_{j=1}^{m_t} \mathbb{I}\{i \in \overline{V}_{t,j}\} \|\mathbf{x}_{a_t}\|_{\overline{\mathbf{M}}_{\overline{V}_{t,j}}}^2} \\
 &\leq \epsilon_* (T_k - T_{k-1}) \sqrt{2md \log(1 + \frac{T}{\lambda d})}
 \end{aligned} \tag{80}$$

860 Then for all phases this term can be bounded by $\epsilon_* T \sqrt{2md \log(1 + \frac{T}{\lambda d})}$.

861 Thus the total regret can be bounded by:

$$\begin{aligned}
 R_t &\leq 2\sqrt{2mTd \log(T) \log(1 + \frac{T}{\lambda d})} (\sqrt{2 \log(T) + d \log(1 + \frac{T}{\lambda d})} + 2\sqrt{\lambda}) \\
 &\quad + 2\epsilon_* T \sqrt{2md \log(1 + \frac{T}{\lambda d})} + 2\epsilon_* T + 12\epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}} T + \frac{3\epsilon_* \sqrt{2d}}{\tilde{\lambda}_x^{\frac{3}{2}}} T + 2T_1 + u \log(T) + 5
 \end{aligned}$$

862 where $T_1 = 16u \log(\frac{u}{\delta}) + 4u \max\{\frac{16}{\tilde{\lambda}_x^2} \log(\frac{8d}{\tilde{\lambda}_x^2 \delta}), \frac{8d}{\tilde{\lambda}_x (\frac{\gamma_1}{6} - \epsilon_* \sqrt{\frac{1}{2\tilde{\lambda}_x}})^2} \log(\frac{u}{\delta})\}$

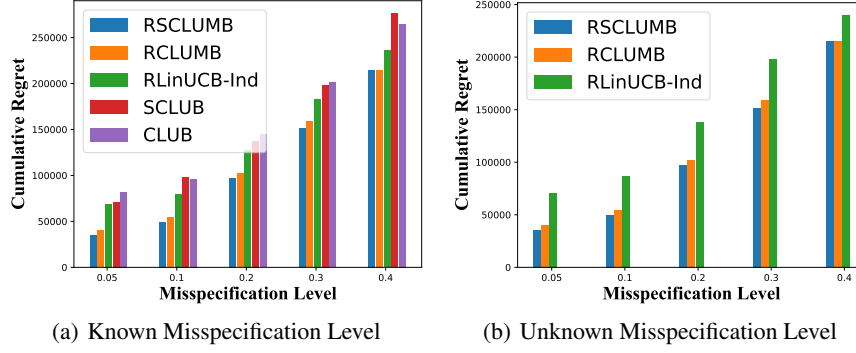


Figure 2: The cumulative regret of the algorithms under different scales of misspecification level.

P More Experiments

For ablation study, we test our algorithms' performance under different scales of deviation. We test RCLUMB and RSCLUMB when $\epsilon^* = 0.05, 0.1, 0.2, 0.3$ and 0.4 in both misspecification level known and unknown cases. In the known case, we set ϵ^* according to the real misspecification level, and we compare our algorithms' performance to the baselines except LinUCB and CW-OFUL which perform worst; in the unknown case, we keep $\epsilon^* = 0.2$, and we compare our algorithms to RLinUCB-Ind as only it has the pre-specified parameter ϵ^* among the baselines. The results are shown in Fig.2. We plot each algorithm's final cumulative regret under different misspecification levels. All the algorithms' performance get worse when the deviation gets larger, and our two algorithms always perform better than the baselines. Besides, the regrets in the unknown case are only slightly larger than the known case. These results can match our theoretical results and again show our algorithms' effectiveness, as well as verify that our algorithm can handle the unknown misspecification level.