



HAL
open science

Low-rank Bandits with Latent Mixtures

Aditya Gopalan, Odalric-Ambrym Maillard, Mohammadi Zaki

► **To cite this version:**

Aditya Gopalan, Odalric-Ambrym Maillard, Mohammadi Zaki. Low-rank Bandits with Latent Mixtures. 2016. hal-01400318

HAL Id: hal-01400318

<https://hal.science/hal-01400318v1>

Preprint submitted on 21 Nov 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Low-rank Bandits with Latent Mixtures

Aditya Gopalan¹

Odalric-Ambrym Maillard²

Mohammadi Zaki¹

ADITYA@ECE.IISC.ERNET.IN

ODALRICAMBRYM.MAILLARD@INRIA.FR

ZAKI@ECE.IISC.ERNET.IN

1. *Electrical Communication Engineering,
Indian Institute of Science,
Bangalore 560012, India*

2. *Inria Saclay – Île de France
Laboratoire de Recherche en Informatique
660 Claude Shannon
91190 Gif-sur-Yvette, France*

Abstract

We study the task of maximizing rewards from recommending items (actions) to users sequentially interacting with a recommender system. Users are modeled as latent mixtures of C many representative user classes, where each class specifies a mean reward profile across actions. Both the user features (mixture distribution over classes) and the item features (mean reward vector per class) are unknown a priori. The user identity is the only contextual information available to the learner while interacting. This induces a low-rank structure on the matrix of expected rewards $r_{a,b}$ from recommending item a to user b . The problem reduces to the well-known linear bandit when either user- or item-side features are perfectly known. In the setting where each user, with its stochastically sampled taste profile, interacts only for a small number of sessions, we develop a bandit algorithm for the two-sided uncertainty. It combines the Robust Tensor Power Method of [Anandkumar et al. \(2014b\)](#) with the **OFUL** linear bandit algorithm of [Abbasi-Yadkori et al. \(2011\)](#). We provide the first rigorous regret analysis of this combination, showing that its regret after T user interactions is $\tilde{O}(C\sqrt{BT})$, with B the number of users. An ingredient towards this result is a novel robustness property of **OFUL**, of independent interest.

Keywords: Multi-armed bandits, online learning, low-rank matrices, recommender systems, reinforcement learning.

1. Introduction

Recommender systems aim to provide targeted, personalized content recommendations to users by learning their responses over time. The underlying goal is to be able to predict which items a user might prefer based on preferences expressed by other related users and items, also known as the principle of collaborative filtering.

A popular approach to model preferences expressed by users in recommender systems is via probabilistic mixture models or *latent class* models ([Hofmann and Puzicha, 1999](#); [Kleinberg and Sandler, 2004](#)). In such a mixture model, we have a set of A items (content) that can be recommended to B users (consumers). Whenever item a is recommended to user b , the system gains an expected reward of $r_{a,b}$. The key structural assumption that captures the relationship between users' preferences is that there exists a set of latent set of C *representative* user types or typical taste profiles.

Formally, each taste profile c is a unique vector $\mathbf{u}_c \equiv (u_{a,c})_a$ of the expected rewards that every item a elicits under the taste profile. Each user b is assumed to sample one of the typical profiles randomly using an individual probability distribution $\mathbf{v}_b \equiv (v_{b,c})_c$; its reward distribution across the items subsequently becomes that induced by the assumed profile.

Our focus is to address the sequential optimization of net reward gained by the recommender, *without any prior knowledge of either the latent user classes or users’ mixture distributions*. Assuming that users arrive to the system repeatedly following an unknown stochastic process and re-sample their profiles over time, according to their respective *unknown* mixtures across latent classes, we seek online learning strategies that can achieve *low regret relative to the best single item that can be recommended to each user*. Note that this is qualitatively different than the task of estimating latent classes or user mixtures in a batch fashion, well-studied by now (Sutskever et al., 2009; Anandkumar et al., 2014a,b); the task of simultaneously optimizing net utility in a bandit fashion in complex expression models like these has received little or no analytical treatment. Our work takes a step towards filling this void.

An especially challenging aspect of online learning in recommender systems is the relatively meager number of available interactions with a same user, which is offset to an extent by the assumption that users can only have a limited number of taste profiles (classes). Indeed, if one can identify the class to which a certain user belongs and aggregate information from all other users in that class, then one can recommend to the user the best item for the class. In practice, classes are latent and not necessarily known in advance, and several works (Gentile et al., 2014; Lazaric et al., 2013; Maillard and Mannor, 2014) study the restricted situation when each user always belongs to one specific class (i.e., when all mixture distributions have support size 1). We go two steps further, since in many situations (a) users cannot be assumed to belong to one class only, such as when a user account is shared by several individuals (e.g. a smart-TV), and (b) the duration of a user-session, that is the number of consecutive recommendations to the same individual connected to a user-account, cannot be assumed to be long¹.

The key challenges that this work addresses are (1) the lack of knowledge of “features” on *both* the user-side and item-side in a linear bandit problem (in this case, both the user mixture weights and the item class reward profiles) and (2) provable regret minimization with very few i.e. $O(1)$ interactions with every user b having a specific taste profile, as opposed to a large number of interactions such as in transfer learning (Lazaric et al., 2013).

Contributions and overview of results. We consider a setting when users are assumed to come from arbitrary mixtures across classes (they are not assumed to fall perfectly in one class as was the assumption in works by Gentile et al. (2014); Maillard and Mannor (2014)). We develop a novel bandit algorithm (Algorithm 3) that combines (a) the *Optimization in the Face of Uncertainty Linear bandit OFUL* algorithm (Abbasi-Yadkori et al., 2011) for bandits with known action features, and (b) a variant of the *Robust Tensor Power (RTP)* algorithm (Anandkumar et al., 2014b) that uses only bandit (partial) estimates of latent user classes with observations coming from a mixture model. More specifically, we introduce a subroutine (Algorithm 1) that makes use of the RTP method to extract item-side attributes (U) and, contributing to its theoretical analysis, show a recovery property (Theorem 1). Note that the RTP method ideally requires (unbiased) estimates of the 2nd and 3rd order moments of actions’ rewards, but with bandit information the learner can access only partial reward information, i.e., a single reward sample from an action. To overcome this, we devise an importance sampling scheme across 3 successive time instants to build the 2nd and 3rd order moment tensor estimates that RTP uses. For the task of issuing recommendations, we develop an algorithm (section 4), essentially based on **OFUL**, instantiated per user, using for each a the *estimated* latent class vectors $\{u_{a,c}\}_c$ (obtained via the RTP subroutine) as arm features, and uncertain parameter vector to be learned \mathbf{v}_b .

1. It is also unlikely to be very short, say, less than 3.

We carry out a rigorous analysis of the algorithm and show that it achieves regret $\tilde{O}(\ell C \sqrt{BT})$ in T rounds of interaction (Theorem 4), provided each arriving user interacts with the system for $\ell \geq 3$ rounds with the same profile. In comparison, the regret of the strategy that completely disregards the latent mixture structure of rewards and employs a standard bandit strategy (e.g. UCB (Auer et al., 2002)) per user, scales as $O(B\sqrt{TA/B}) = O(\sqrt{ABT})$ after T rounds², which is considerably suboptimal in the practical case with a very large number of items but very few representative user classes ($C \ll A$). It is also worth noting that the regret bound we achieve, order-wise, is what would result from applying the **OFUL** or any optimal linear bandit algorithm assuming *a priori knowledge* of all latent user classes $\{u_{a,c}\}_{a,c}$, that is $\tilde{O}(\ell C \sqrt{BT})$. In this sense, our result shows that one can *simultaneously estimate features on both sides of a bilinear reward model and achieve regret performance equivalent to that of a one-sided linear model*, which is the first result of its kind to the best of our knowledge³. Our results are presented for finite time horizons with explicit details of the constants arising from the error analysis of RTP, which at this point are large but possibly improvable.

En route to deriving the regret for our algorithm, we also make a novel contribution that advances the theoretical understanding of **OFUL**, and which is of independent interest. We show that in the standard linear bandit setting, where the expected reward of an arm linearly depends on d features, **OFUL** yields (sub-linear) $\tilde{O}(\rho d \sqrt{T})$ regret even when it makes decisions based on perturbed or inexact feature vectors (Theorem 3), where ρ quantifies the distortion. This property holds whenever the perturbation error is small enough, and we explicitly give both (a) a sufficient condition on the size of the perturbation in terms of the set of actual features, and (b) a bound on the (multiplicative) distortion ρ in the regret due to the perturbation (note that $\rho = 1$ in the ideal linear case).

2. Setup and notation

For any positive integer n , $[n]$ denotes the set $\{1, 2, \dots, n\}$.

At each $n \in \mathbb{N}$, nature selects a user $b_n \in [B]$ according to the probability distribution β over $[B]$, independent of the past, and b_n is revealed to the learner. A user class c_n is subsequently sampled from the probability distribution \mathbf{v}_{b_n} over $[C]$, and c_n (the assumed class of user b_n) interacts with the learner for the next $\ell \geq 3$ consecutive steps. Such an interaction will often be termed a *mini-session*.

In each step $l \in [\ell]$ of a mini-session, the learner plays an action (issues a recommendation) $a_{n,l} \in [A]$ and subsequently receives reward $Y_{n,l} = u_{a_{n,l},c_n} + \eta_{n,l}$, where $\eta_{n,l}$ is a (centered) R -sub-Gaussian i.i.d. random variable independent from $a_{n,l}, c_n$, representing the noise in the reward. We let $\mathbf{u}_a \in \mathbb{R}^C$ represent the vector $(u_{a,c})_{c \in [C]}$ of the mean rewards from action a in each class. Note that $\mathbb{E}[u_{a_{n,l},c_n} | a_{n,l}] = \mathbb{E}[\mathbf{u}_{a_{n,l}}^\top \mathbf{v}_{b_n} | a_{n,l}]$. For convenience, we use the index notation $t \equiv (n, l)$ and introduce $T = N\ell$, where N is the total number of mini-sessions, and T the total number of interactions of the learner with the system. We denote likewise Y_t, a_t, c_t, η_t for $Y_{n,l}, a_{n,l}, c_n, \eta_{n,l}$, and let $u_{\max} \stackrel{\text{def}}{=} \max_{a \in [A], c \in [C]} |u_{a,c}|$.

We are interested in designing an online recommendation strategy, i.e., one that plays actions depending on past observations, achieving low (*cumulative*) regret after $T \equiv (N, \ell)$ mini-sessions, defined as $\mathcal{R}_T \stackrel{\text{def}}{=} \sum_{n \in [N], l \in [\ell]} r_{n,l}$, where $r_{n,l} \stackrel{\text{def}}{=} \max_{a \in [A]} \mathbf{u}_a^\top \mathbf{v}_{b_n} - \mathbf{u}_{a_{n,l}}^\top \mathbf{v}_{b_n}$. In other words, we wish to compete against a strategy that plays for every user an action yielding the highest reward in expectation under its mixture distribution over user classes.

2. Roughly, each UCB per-user plays from a pool of A actions for about T/B rounds, thus suffering regret $O(\sqrt{A(T/B)})$.

3. An earlier result of Djolonga et al. (2013) gets $O(T^{4/5})$ regret while moreover assuming a perfect control of the sampling process (we can't assume this due to the user arrivals).

3. Recovering latent user classes: The EstimateFeatures subroutine

In this section, we provide an estimation algorithm for the matrix U , using the RTP method.⁴

Estimation of tensors. We assume that in mini-session n , when interacting with user b_n , the triplet $\{a_{n,l}\}_{l \leq \ell}$ is chosen from a distribution $p_n(a, a', a'' | b_n)$. Letting $X_{a_{n,l}, b_n, n, l} \stackrel{\text{def}}{=} Y_{n,l} = u_{a_{n,l}, c_n} + \eta_{n,l}$ to explicitly indicate the active user and action chosen at (n, l) , we form the importance-weighted estimates

$$\begin{aligned}\tilde{r}_{a,a',n} &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \frac{X_{a_{i,1}, b_i, i, 1} X_{a_{i,2}, b_i, i, 2}}{p_i(a, a' | b_i)} \mathbb{I}\{a_{i,1} = a, a_{i,2} = a'\}, \\ \tilde{r}_{a,a',a'',n} &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \frac{X_{a_{i,1}, b_i, i, 1} X_{a_{i,2}, b_i, i, 2} X_{a_{i,3}, b_i, i, 3}}{p_i(a, a', a'' | b_i)} \mathbb{I}\{a_{i,1} = a, a_{i,2} = a', a_{i,3} = a''\}.\end{aligned}$$

for the second and third-order tensors⁵.

We introduce the matrices $\widehat{M}_{n,2} \equiv (\tilde{r}_{a,a',n})_{a,a' \in [A]}$ and $M_2 \equiv (m_{a,a'})_{a,a' \in [A]}$ with $m_{a,a'} \stackrel{\text{def}}{=} \mathbb{E}[\tilde{r}_{a,a',n}]$, and the tensors $\widehat{M}_{n,3} \equiv (\tilde{r}_{a,a',a'',n})_{a,a',a'' \in [A]}$ and $M_3 \equiv (m_{a,a',a''})_{a,a',a'' \in [A]}$ with $m_{a,a',a''} \stackrel{\text{def}}{=} \mathbb{E}[\tilde{r}_{a,a',a'',n}]$. The following result decomposes the matrix M_2 and tensor M_3 as weighted sums of outer products.

Lemma 1 *When the user arrivals are i.i.d. according to the law β , i.e., $b_i \stackrel{i.i.d.}{\sim} \beta \forall i \in [n]$, it holds that*

$$\begin{aligned}m_{a,a',n} &= \sum_{c \in [C]} v_{\beta,c} u_{a,c} u_{a',c}, \quad \text{and} \\ m_{a,a',a'',n} &= \sum_{c \in [C]} v_{\beta,c} u_{a,c} u_{a',c} u_{a'',c}.\end{aligned}$$

Having shown the unbiasedness of the empirical 2nd and 3rd moment tensors $\widehat{M}_{n,2}$ and $\widehat{M}_{n,3}$, we next turn to showing concentration to their respective means.

Lemma 2 *Assuming that $p_i(a, a' | b_i) \geq q_{2,i}$ and $p_i(a, a', a'' | b_i) \geq q_{3,i}$ for deterministic $q_{2,i}, q_{3,i}$, for all $i \in \mathbb{N}$, $a, a', a'' \in [A]$, then for all $n \leq N$, with probability higher than $1 - \delta$, it holds simultaneously for all a, a', a'' that*

$$\begin{aligned}|\tilde{r}_{a,a',n} - m_{a,a',n}| &\leq \sqrt{\sum_{i=1}^n q_{2,i}^{-2} \frac{\log(4A^2/\delta)}{2n^2}}, \\ |\tilde{r}_{a,a',a'',n} - m_{a,a',a'',n}| &\leq \sqrt{\sum_{i=1}^n q_{3,i}^{-2} \frac{\log(4A^3/\delta)}{2n^2}}.\end{aligned}$$

An immediate corollary is the following one:

Corollary 1 *Provided that $q_{2,i} = \gamma_i/A^2$ and $q_{3,i} = \gamma_i/A^3$ for some $\gamma_i > 0$, then on an event of probability higher than $1 - \delta$, the following hold simultaneously:*

$$\begin{aligned}e_n^{(2)} &\stackrel{\text{def}}{=} \|\widehat{M}_{n,2} - M_2\| \leq A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^2/\delta)}{2n^2}}, \\ e_n^{(3)} &\stackrel{\text{def}}{=} \|\widehat{M}_{n,3} - M_3\| \leq A^{9/2} \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^3/\delta)}{2n^2}}.\end{aligned}$$

4. We consider $\ell = 3$ to describe the algorithm; $\ell > 3$ is easily handled by repeating the 3-wise sampling $p(a, a', a'')$ for $\lceil \ell/3 \rceil$ times and discarding the remaining (< 3) steps in the mini-session during exploration (leading to a negligible regret overhead).

5. An alternative is the *implicit exploration* method due to [Kocák et al. \(2014\)](#).

Algorithm 1 EstimateFeatures

- 1: **Input:** #sessions n ; #mini-sessions ℓ ; (user, action, reward) tuples $(b_i, a_{i,l}, X_{a_{i,l}, b_i, i, l})_{1 \leq i \leq n, 1 \leq l \leq \ell}$.
 - 2: Compute the $A \times A$ matrix $\widehat{M}_{n,2} = (\tilde{r}_{a,a',n})_{a,a' \in [A]}$ and the $A \times A \times A$ tensor $\widehat{M}_{n,3} = (\widehat{r}_{a,a',a'})_{a,a',a' \in [A]}$.
 - 3: Compute a $A \times C$ whitening matrix \widehat{W}_n of $\widehat{M}_{n,2}$
 {Take $\widehat{W}_n = \widehat{U}_n \widehat{D}_n^{-1/2}$ where \widehat{D}_n is the $C \times C$ diagonal matrix with the top C eigenvalues of $\widehat{M}_{n,2}$, and \widehat{U}_n the $A \times C$ matrix of corresponding eigenvectors.}
 - 4: Form the $C \times C \times C$ tensor $\widehat{T}_n = \widehat{M}_{n,3}(\widehat{W}_n, \widehat{W}_n, \widehat{W}_n)$.
 - 5: Apply the RTP algorithm (Anandkumar et al., 2014b) to \widehat{T}_n , and compute its robust eigenvalues $(\lambda_{n,c})_{c \in [C]}$ with eigenvectors $(\widehat{\varphi}_{n,c})_{c \in [C]}$.
 {The paper of Anandkumar et al. (2014b, Sec. 4) defines eigenvalues/eigenvectors of tensors.}
 - 6: Compute for each $c \in [C]$, $\bar{u}_{n,c} = \lambda_{n,c}(\widehat{W}_n^\top)^\dagger \widehat{\varphi}_{n,c}$ and $\bar{v}_{n,c} = \lambda_{n,c}^{-2}$.
 - 7: **Output: Estimate of latent classes U :** The $A \times C$ matrix \bar{U}_n obtained by stacking the vectors $\bar{u}_{n,c} \in \mathbb{R}^A$ side by side.
-

Reconstruction algorithm. The EstimateFeatures algorithm (Algorithm 1) employs a whitening matrix \widehat{W}_n , of the empirical estimate of the matrix M_2 , to build the empirical tensor \widehat{T}_n . This tensor is then used to recover the columns of the matrix $U = (u_{a,c})_{a \in [A], c \in [C]}$ via the RTP algorithm. For the sake of completeness, we also introduce W , a whitening matrix of M_2 (i.e., $W^\top M_2 W = I$), the corresponding tensor $T = M_3(W, W, W)$, and finally the estimation error $e_n \stackrel{\text{def}}{=} \|\widehat{T}_n - T\|$.

Reconstruction guarantee. Our next result makes use of the following proposition from Anandkumar et al. (2014b, Theorem 5.1), restated here for completeness.

Proposition 1 (Theorem 5.1 of Anandkumar et al. (2014b)) *Let $\widehat{T} = T + E \in \mathbb{R}^{C \times C \times C}$, where T is a symmetric tensor with orthogonal decomposition $T = \sum_{c=1}^C \lambda_c \varphi_c^{\otimes 3}$, where each $\lambda_c > 0$, $\{\varphi_c\}_{c \in [C]}$ is an orthonormal basis, and E is a symmetric tensor with operator norm $\|E\| \leq \varepsilon$. Let $\lambda_{\min} = \min\{\lambda_c : c \in [C]\}$, $\lambda_{\max} = \max\{\lambda_c : c \in [C]\}$. Run the RTP algorithm with input \widehat{T} for C iterations. Let $\{(\widehat{\lambda}_c, \widehat{\varphi}_c)\}_{c \in [C]}$ be the corresponding sequence of estimated eigenvalue/eigenvector pairs returned. Then, there exist universal constants $C_1, C_2 > 0$ for which the following is true. Fix $\eta \in (0, 1)$ and run RTP with parameters (i.e., number of iterations) L, N with $L = \text{poly}(C) \log(1/\eta)$, and $N \geq C_2 \left(\log(C) + \log \log \left(\frac{\lambda_{\max}}{\varepsilon} \right) \right)$. If $\varepsilon \leq C_1 \frac{\lambda_{\min}}{C}$, then with probability at least $1 - \eta$, there exists a permutation $\pi \in \mathbb{S}_C$ such that*

$$\forall c \in [C] : \quad |\lambda_c - \widehat{\lambda}_{\pi(c)}| \leq 5\varepsilon, \quad \|\varphi_c - \widehat{\varphi}_{\pi(c)}\| \leq 8\varepsilon/\lambda_c,$$

$$\text{and} \quad \left\| T - \sum_{c=1}^C \widehat{\lambda}_c \widehat{\varphi}_c^{\otimes 3} \right\| \leq 55\varepsilon.$$

Lemma 1 gives a decomposition of the (symmetric) tensor M_3 , but it may be not orthogonal; standard transformation (Anandkumar et al., 2014b, Sec. 4.3) gives an orthogonal decomposition for the tensor⁶ $M_3(W, W, W)$, with W a matrix that whitens M_2 . We can thus use Proposition 1

6. For a 3rd order tensor $A \in \mathbb{R}^{a \times a \times a}$ and 2nd order tensor or matrix $B \in \mathbb{R}^{a \times b}$, $A(B, B, B) \in \mathbb{R}^{b \times b \times b}$ is the 3rd order tensor defined by $[A(B, B, B)]_{i_1, i_2, i_3} \stackrel{\text{def}}{=} \sum_{j_1, j_2, j_3 \in [a]} A_{j_1, j_2, j_3} B_{j_1, i_1} B_{j_2, i_2} B_{j_3, i_3}$. See Anandkumar et al. (2014b) for more details on notation and results.

with $T = M_3(W, W, W)$, $\widehat{T} = \widehat{T}_n$, $\varepsilon = e_n$ and $\eta = \delta$ in order to prove the following guarantee (Theorem 1) on the recovery error between columns of U and their estimate.

We now introduce mild separability conditions on the mixture weights \mathbf{v}_b and the spectrum of the 2nd moment matrix M_2 needed for the reconstruction guarantee to hold, similar to those assumed for Lazaric et al. (2013, Theorem 2).

Assumption 1 *There exist positive constants v_{\min} , σ_{\min} , σ_{\max} and Γ such that*

$$\begin{aligned} & \min_{b \in [B], c \in [C]} v_{b,c} \geq v_{\min}, \\ & \forall c \in [C], \sigma_c = \sqrt{\lambda_c(M_2)} \in [\sigma_{\min}, \sigma_{\max}] \quad \text{and} \\ & \min_{c \neq c' \in [C] \times [C]} |\sigma_c - \sigma_{c'}| \geq \Gamma, \end{aligned}$$

where $\lambda_c(A)$ denotes the c^{th} top eigenvalue of A .

Theorem 1 (Recovery guarantee for online estimation of user classes U) *Let Assumption 1 hold, and let $\delta \in (0, 1)$. If the number of mini-session satisfies*

$$\frac{n^2}{\sum_{i=1}^n \gamma_i^{-2}} \geq \max \left\{ \frac{2A^6 \log(4A^2/\delta)}{\min\{\Gamma, \sigma_{\min}\}^2}, \frac{A^9 (1 + 10(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}}))(1 + u_{\max}^3)^2 C^5 \log(4A^3/\delta)}{2C_1^2 \sigma_{\min}^3} \right\},$$

then with probability at least $1 - 2\delta$, there exists some permutation $\pi \in \mathbb{S}_C$ such that for all $c \in [C]$, the output \bar{U}_n of the EstimateFeatures algorithm satisfies

$$\|\mathbf{u}_c - \bar{\mathbf{u}}_{n, \pi(c)}\| \leq \diamond A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}}. \quad (1)$$

where $\mathbf{u}_c = (u_{a,c})_{a \in [A]}$. Here, the constant (we use the "diamond" symbol to denote it) is

$$\begin{aligned} \diamond &= \left(\frac{CA}{\sigma_{\min}} \right)^{3/2} \left(13\sqrt{\sigma_{\max}} + 4\sqrt{2 \min\{\Gamma, \sigma_{\min}\}} \right. \\ &+ 5 \left(\frac{\sigma_{\max}}{\Gamma} + \frac{1}{2\sigma_{\max}} \right) \min\{\Gamma, \sigma_{\min}\} \Big) \aleph \\ &+ \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) \frac{1}{v_{\min}^2} \\ &+ 5\sqrt{3/8} \left(\sqrt{\sigma_{\max}} + \sqrt{\min\{\Gamma, \sigma_{\min}\}/2} \right) \left(\frac{2CA}{\sigma_{\min}} \right)^3 \aleph^2 \min\{\Gamma, \sigma_{\min}\}, \end{aligned}$$

with the notation $\aleph = 1 + 10(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}})(1 + u_{\max}^3)$.

The proof strategy follows that of Lazaric et al. (2013, Theorem 2) and is detailed in the appendix for clarity. It consists in relating, on the one hand, the estimation errors $e_n^{(2)}$ of M_2 and $e_n^{(3)}$ of M_3 from Corollary 1 to the condition $\varepsilon \leq C_1 \frac{\lambda_{\min}}{C}$, and, on the other hand, relating the reconstruction error on the columns of U to the control on the terms $|\lambda_c - \widehat{\lambda}_{\pi(c)}|$ and $\|\varphi_c - \widehat{\varphi}_{\pi(c)}\|$ coming from Proposition 1. We note that the bound appearing in the condition on the number of mini-sessions is potentially large (due to the terms A^6 , C^5 , etc.). This is due to the combination of the RTP method with the importance sampling scheme, and it remains unclear if the bound can be significantly improved within this framework.

4. Recovering latent mixture distributions (v_b): robustness of the OFUL algorithm

In order to recover the weights vectors $\mathbf{v}_b \in \mathbb{R}^C$ and thus the matrix V , it would be tempting to use again an instance of the RTP method but this time to aggregate across actions, i.e., by forming a $B \times B$ and $B \times B \times B$ tensor. Unfortunately, aggregation of elements of U fails for two reasons: First, we do not have different views across users b , contrary to what we have for actions a . It is thus hopeless to be able to form an estimate of the 2nd and 3rd moment tensors as before. Second, and rather technically, convex combinations of the $\{u_{a,c}\}_{a \in [A]}$ need not be positive. This prevents the application of the RTP method which requires positive weights to work.

We thus consider a different strategy that uses an algorithm designed for linear bandits. However since the feature matrix U is unknown a priori and can only be estimated, we need to work with perturbed features. A first solution is to propagate the additional error resulting from the error on the features in the standard proof of **OFUL**. However, this leads to a sub-optimal regret that is no longer scaling as $\tilde{O}(\sqrt{T})$ with the time horizon. We overcome this hurdle by showing in Theorem 3 a robustness property of **OFUL** of independent interest, which aids us in controlling the regret of the overall latent class algorithm (Algorithm 3).

Consider **OFUL** run with perturbed (not necessarily linearly realizable) rewards. Formally, consider a finite action set $\mathcal{A} = \{1, 2, \dots, A\}$ and distinct feature vectors $\{\bar{\mathbf{u}}_a \in \mathbb{R}^{C \times 1}\}_{a \in \mathcal{A}}$. Let $\bar{U}^\top := [\bar{\mathbf{u}}_1 \ \bar{\mathbf{u}}_2 \ \dots \ \bar{\mathbf{u}}_A] \in \mathbb{R}^{C \times A}$. The expected reward when playing action $A_t = a$ at time t is denoted by $m_a := \mathbb{E}[Y_t \mid A_t = a]$, with $\mathbf{m} := (m_a)_{a \in \mathcal{A}}$. Let us assume that there exists a unique optimal action for the expected rewards \mathbf{m} , i.e., $\arg \max_{a \in \mathcal{A}} m_a = \{a^*\}$, with the regret at time n being $\mathfrak{R}_n := \sum_{t=1}^n (m_{a^*} - m_{A_t})$. The key point here is that \mathbf{m} need not be linearly realizable w.r.t. the actions' features – we will not require that $\min_{\mathbf{v} \in \mathbb{R}^C} \|\mathbf{m} - \bar{U}\mathbf{v}\|$ be 0.

Algorithm 2 OFUL (Optimism in Face of Uncertainty for Linear bandits) (Abbasi-Yadkori et al., 2011)

Require: Arms' features \bar{U} , regularization parameter λ , norm parameter R_Θ

for all times $t \geq 1$ **do**

1. Form the $C \times (t-1)$ matrix $\bar{\mathbf{U}}_{1:t-1} := [\bar{\mathbf{u}}_{A_1} \ \bar{\mathbf{u}}_{A_2} \ \dots \ \bar{\mathbf{u}}_{A_{t-1}}]$ consisting of all arm features played up to time $t-1$, and $\mathbf{Y}_{1:t-1} := (Y_1, \dots, Y_{t-1})^\top$. Set $V_{t-1} := \lambda I + \sum_{s=1}^{t-1} \bar{\mathbf{u}}_{A_s} \bar{\mathbf{u}}_{A_s}^\top$.

2. Choose the action

$$A_t \in \arg \max_{a \in \mathcal{A}} \max_{\mathbf{v} \in \mathcal{C}_{t-1}} \bar{\mathbf{u}}_a^\top \mathbf{v}, \quad \text{where}$$

$$\mathcal{C}_{t-1} := \{\mathbf{v} \in \mathbb{R}^C : \|\mathbf{v} - \hat{\mathbf{v}}_{t-1}\|_{V_{t-1}} \leq D_{t-1}\},$$

$$D_{t-1} := R \sqrt{2 \log \left(\frac{\det(V_{t-1})^{1/2} \lambda^{-C/2}}{\delta} \right)} + \lambda^{1/2} R_\Theta$$

$$\hat{\mathbf{v}}_{t-1} := V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{Y}_{1:t-1}.$$

end for

OFUL Regret with linearly realizable rewards. The **OFUL** algorithm is stated for the sake of clarity as Algorithm 2. Before studying the linearly non-realizable case, we record the well-known regret bound for it in the unperturbed case, that is when $\forall a \in [A], m_a = \bar{\mathbf{u}}_a^\top \mathbf{v}^*$ for some unknown \mathbf{v}^* .

Theorem 2 (OFUL regret (Abbasi-Yadkori et al., 2011)) Assume that $\|\mathbf{v}^*\|_2 \leq R_\Theta$, and that for all $a \in \mathcal{A}$, $\|\bar{\mathbf{u}}_a\|_2 \leq R_\mathcal{X}$, $|\langle \bar{\mathbf{u}}_a, \mathbf{v}^* \rangle| \leq 1$. Then with probability at least $1 - \delta$, the regret of **OFUL** satisfies: $\forall n \geq 0$,

$$\mathfrak{R}_n \leq 4\sqrt{nC \log(1 + nR_{\mathcal{X}}^2/(\lambda C))} \times \\ \left(\lambda^{1/2} R_{\Theta} + R\sqrt{2 \log(1/\delta) + C \log(1 + nR_{\mathcal{X}}^2/(\lambda C))} \right),$$

provided that the regularization parameter λ is chosen such that $\lambda \geq \max\{1, R_{\mathcal{X}}^2, 1/R_{\Theta}^2\}$.

Regret of OFUL with Perturbed Features. We make a structural definition to present the result. Let $\alpha(\bar{U}) := \max_J \|\mathbf{A}_J^{-1}\|_2$, where $\mathbf{A} = \begin{bmatrix} \bar{U} \\ I_C \end{bmatrix} \in \mathbb{R}^{(A+C) \times C}$, \mathbf{A}_J is the $C \times C$ submatrix of \mathbf{A} formed by picking rows J , and J ranges over all size- C subsets of full-rank rows of \mathbf{A} . We will require for our purposes that $\alpha(\bar{U}^\top)$ is not too large. For intuition regarding α , we refer to [Forsgren \(1996\)](#) (the final 3 paragraphs of p. 770, Corollary 5.4 and section 7). We remark that the condition that $\alpha(\bar{U}^\top)$ be small is analogous to a γ -incoherence type property commonly used in prior work ([Bresler et al., 2014](#), Assumption A2), stating that two distinct feature vectors \mathbf{u}_c and $\mathbf{u}_{c'}$, $c \neq c'$, must have a minimum angle separation.

Let $\mathbf{v}^\circ \in \mathbb{R}^C$ be arbitrary with ℓ^2 norm at most R_{Θ} (it helps to think of $\bar{U}\mathbf{v}^\circ$ as an approximation of \mathbf{m}), $\varepsilon_a := m_a - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ$, $\varepsilon := (\varepsilon_a)_{a \in \mathcal{A}} \in \mathbb{R}^A$. We now state a robustness result for **OFUL** potentially of independent interest.

Theorem 3 (OFUL robustness property) *Suppose $\|\mathbf{v}^\circ\|_2 \leq R_{\Theta}$, $\lambda \geq \max\{1, R_{\mathcal{X}}^2, 1/4R_{\Theta}^2\}$, $\forall a \in \mathcal{A}$, $\|\bar{\mathbf{u}}_a\|_2 \leq R_{\mathcal{X}}$ and $|m_a| \leq 1$. If the deviation from linearity satisfies*

$$\|\varepsilon\|_2 \equiv \|\mathbf{m} - \bar{U}\mathbf{v}^\circ\|_2 < \min_{a \neq a^*} \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2\alpha(\bar{U}^\top) \|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2}, \quad (2)$$

then, with probability at least $1 - \delta$ for all $T \geq 0$,

$$\mathfrak{R}_T \leq 8\rho' \sqrt{TC \log\left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C}\right)} \left(\lambda^{1/2} R_{\Theta} + R\sqrt{2 \log \frac{1}{\delta} + C \log\left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C}\right)} \right),$$

where $\rho' := \max\left\{1, \max_{a \neq a^*} \frac{m_{a^*} - m_a}{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}\right\}$.

Theorem 3 essentially states that when the deviation of the actual mean reward vector from the subspace spanned by the feature vectors is small, the **OFUL** algorithm continues to enjoy a favorable $O(\sqrt{T})$ regret up to a factor $\rho' \geq 1$. The quantity ρ' in the result is a geometric measure of the distortion in the arms' actual rewards \mathbf{m} with respect to the (linear) approximation $\bar{U}\mathbf{v}^\circ$. We control this quantity in the next paragraph. (Note that $\rho' = 1$ in the perfectly linearly realizable case $\varepsilon = 0$, and this gives back the standard **OFUL** regret up to a universal multiplicative constant.)

Applying the Robust analysis of OFUL to the Low-rank Bandit setup. In this paragraph, we translate Theorem 3 to our Low Rank Bandit (LRB) setting in which OFUL uses feature vectors with noisy perturbations (estimated by, say, a Robust Tensor Power (RTP) algorithm). Throughout this section, we fix a user b .

We can now translate Theorem 3 thanks to the correspondence with the perturbed OFUL setting: In our low-rank bandit setting, the matrix $\bar{U} = \bar{U}_n$ depends on the reconstruction algorithm at mini-session n . Moreover, the optimal action $a^* \equiv a_b^*$ now depends on the user b . We denote for a user $b \in [B]$ the minimum gap across suboptimal actions to be $g_b \stackrel{\text{def}}{=} \min_{a \neq a_b^*} (\mathbf{u}_{a_b^*} - \mathbf{u}_a)^\top \mathbf{v}_b$. Likewise, the error vector ε depends on b, n . Its norm $\|\varepsilon\|_2$ appears in the condition (2) and the definition of ρ , and is controlled by the reconstruction error of Theorem 1. It decays with the number of mini-sessions n .

We define $\alpha_n \stackrel{\text{def}}{=} \alpha(\bar{U}_n)$, $\alpha_* \stackrel{\text{def}}{=} \alpha(U)$ and use $\max_b \|\mathbf{v}_b\|$ for R_{Θ} . Using these notations, and adapting the proof of Theorem 3 to handle a variable \bar{U}_n , we can now translate the result of the perturbed **OFUL** to our LRB setting:

OFUL	LRB
\mathbf{v}°	$\mathbf{v}_b \equiv (v_{b,c})_{c \in [C]} \in \mathbb{R}^C$
\bar{U}	$\bar{U}_n \in \mathbb{R}^{A \times C}$
\mathbf{m}	$\mathbf{m}_b \equiv U \mathbf{v}_b \in \mathbb{R}^A$
a^*	$a_b^* := \arg \max_{a \in [A]} \mathbf{u}_a^T \mathbf{v}_b$
ε_a	$(\mathbf{u}_a - \bar{\mathbf{u}}_{n,a})^\top \mathbf{v}_b$
$\varepsilon \equiv (\varepsilon_a)_{a \in \mathcal{A}}$	$(U - \bar{U}_n) \mathbf{v}_b$

Table 1: Correspondences between **OFUL** and Low Rank Bandit (LRB) quantities at time n and for user b

Lemma 3 Let $0 < \delta \leq 1$ and $b \in [B]$. Provided that the number of mini-sessions n_0 satisfies $\frac{n_0^2}{\sum_{i=1}^{n_0} \gamma_i^{-2}} \geq \circlearrowleft_{b,\delta}$, where we introduced the notation

$$\begin{aligned} \circlearrowleft_{b,\delta} = & \max \left\{ \frac{2A^6 \log(4A^2/\delta)}{\min\{\Gamma, \sigma_{\min}\}^2}, \right. \\ & \frac{A^9 (1 + 10(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}}))(1 + u_{\max}^3)^2 C^5 \log(4A^3/\delta)}{2C_1^2 \sigma_{\min}^3} \times \\ & \diamond^2 A^6 C^2 \log(4A^3/\delta) \times \\ & \left. \max \left\{ 2\alpha_*^2, \frac{8A \|\mathbf{v}_b\|_2^2}{g_b^2}, \frac{2^7 \alpha_*^2 C u_{\max}^2 \|\mathbf{v}_b\|_2^2}{g_b^2} + \frac{1}{2} \right\} \right\}, \end{aligned}$$

then with probability at least $1 - 2\delta$, $\|\varepsilon\|_2 = \|(U - \bar{U}_n) \mathbf{v}_b\|_2$ is small enough that for any $n \geq n_0$, condition (2) is satisfied. Consequently, Theorem 3 applies with

$$\begin{aligned} R_\Theta &= \max_b \|\mathbf{v}_b\|_2, \quad R_{\mathcal{X}} = \max_{a \in \mathcal{A}} \|\mathbf{u}_a\|_2 + \frac{\sqrt{A}}{2\alpha_*}, \quad \text{and} \\ \rho' &\equiv \rho'_{n,b} \leq 2. \end{aligned}$$

Thus, provided that the total number of mini-sessions of interaction (not necessarily corresponding to interactions with user b) is large enough, then the **OFUL** algorithm run during interactions with user b will achieve a controlled regret. However, we want to warn that the $\circlearrowleft_{b,\delta}$ resulting from the RTP method, especially the second term of the max, may be potentially large, although being a constant.

5. Putting it together: Online Recommendation algorithm

This section details our main contributions for recommendations in the context of mini-sessions of interactions with unknown mixtures of latent profiles: first Algorithm 3 that combines RTP with **OFUL**, and then a regret analysis in Theorem 4.

The recommendation algorithm we propose (Algorithm 3) uses the RTP method to estimate the matrix U and then applies **OFUL** to determine an optimistic action. Importantly, it finally outputs a distribution that mixes the optimistic action with a uniform exploration. The mixture coefficient goes to 0 with the number of rounds, thus converging to playing **OFUL** only. It ensures that the importance sampling weights are bounded away from 0 in the beginning.

Main analytical result: Regret bound

Theorem 4 (Regret of Algorithm 3) *With Assumption 1 holding, let $\delta \in (0, 1)$, $\circlearrowleft_\delta = \max_{b \in [B]} \circlearrowleft_{b, \delta}$ (from Lemma 3), and let n_0 be the first mini-session at which $\frac{n_0^2}{\sum_{i=1}^{n_0} \gamma_i} \geq \circlearrowleft_\delta$. The regret of Algorithm 3 at time $T = N\ell$ (acting for N mini-sessions of length ℓ) using internal instances of **OFUL** parameterized by $\delta > 0$ satisfies*

$$\begin{aligned} \mathbb{E}[\mathfrak{R}_T] \leq & 16\sqrt{BTC \log\left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C}\right)} \left(\lambda^{1/2} R_\Theta + R \sqrt{2 \log \frac{1}{\delta} + C \log\left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C}\right)} \right) \\ & + \ell(n_0 - 1 + \sum_{n=n_0}^N \gamma_n) + 3\delta T, \end{aligned}$$

provided that $\lambda \geq \min\{1, R_{\mathcal{X}}^2, 1/R_\Theta^2\}$, with $R_\Theta \geq \max_b \|\mathbf{v}_b\|_2$, $R_{\mathcal{X}} \geq \max_{a \in \mathcal{A}} \|\mathbf{u}_a\|_2 + \frac{\sqrt{A}}{2\alpha^*}$. Consequently, choosing $\delta = 1/T$ and $\gamma_n = \sqrt{\log(n+1)/n}$, $n \in \mathbb{N}$, say, yields the order $\mathbb{E}[\mathfrak{R}_T] = O\left(C\sqrt{BT} \log T\right)$.

Discussion. (1) The regret of Algorithm 3 scales with T similar to that of an **OFUL** algorithm run with perfect knowledge of the feature matrix U : $\tilde{O}(C\sqrt{BT})$. This is a non-trivial result as U is not assumed to be known a priori and is estimated by Algorithm 3 using tensor methods.

Algorithm 3 Per-user **OFUL** with exploration

Require: Parameters λ, R_Θ for **OFUL**, exploration rate parameters $\gamma_n, n \geq 1$.

- 1: **for** mini-session $n = 1, \dots, N$ **do**
 - 2: Get user b_n .
 - 3: Let $p_n \sim \text{Bernoulli}(\gamma_n)$
 - 4: **if** $p_n = 0$ **then**
 - 5: {Carry out an **ESTIMATE** mini-session}
 - 6: **for** step $k = 1, 2, \dots, \ell$ **do**
 - 7: Output $a_{n,k} \sim \text{Uniform}([A])$.
 - 8: **end for**
 - 9: Let $\bar{U}_n = \text{EstimateFeatures}$ (Algorithm 1) with input $(b_i, a_{i,l}, X_{a_{i,l}, b_i, i, l})_{1 \leq i \leq n, 1 \leq l \leq \ell, p_n=0}$
 {Update feature estimates using samples from previous **ESTIMATE** mini-sessions}
 - 10: **else**
 - 11: {Carry out an **OFUL** mini-session}
 - 12: **for** step $k = 1, 2, \dots, \ell$ **do**
 - 13: Run one iteration of **OFUL** (Algorithm 2) with features \bar{U}_n , parameters λ and R_Θ , and past actions and rewards $(a_{i,l}, X_{a_{i,l}, b_i, i, l})$, $1 \leq i < n, 1 \leq l \leq \ell$, for which $p_i = 1$ and $b_i = b_n$
 {An instance of **OFUL** for each user using current feature estimates, and observed actions and rewards from previous **OFUL** mini-sessions}
 - 14: Output action $a_{n,k}$ returned by **OFUL**
 - 15: **end for**
 - 16: **end if**
 - 17: **end for**
-

(2) One can also compare the result with the regret of ignoring the mixture (low-rank) structure and simply running an instance of UCB per user, which would scale as $O(\sqrt{ABT})$. This becomes highly suboptimal when the number of actions/items A is much larger than the number of user types C , demonstrating the gain from leveraging the mixed linear structure of the problem. Note

also that we do not need a specific user to interact for a long time but for as few as $\ell \geq 3$ consecutive steps, contrary for instance to the transfer method (Lazaric et al., 2013), where a large number of consecutive interaction steps with the same user is required.

(3) It is worthwhile to contrast the result and approach with that in Djolonga et al. (2013) – the authors there incur an additional regret term due to the error in approximately estimating the low-rank matrix, which requires additional tuning ending up with a regret of $O(T^{4/5})$. On the other hand, we avoid this approximation error by showing and exploiting the robustness property of OFUL, which guarantees \sqrt{T} regret as soon as the estimated features \tilde{U} are within a small radius of the actual ones.

The result (and analysis) does come with a caveat that the model-dependent term \circlearrowleft_δ , although being independent on the time horizon T , is potentially large. With γ_n set as in Theorem 4, it appears as an additive exponential constant term in the regret⁷. This arises from the RTP method, and it is currently unclear if this term can be significantly reduced with the current line of analysis. Numerical evidence, however, indicates that no such large additive constant enters into the regret (Section 5). Also, on the bright side, note that \circlearrowleft_δ does not need to be known by the algorithm.

Numerical Results. The performance of the low-rank bandit strategy (Algorithm 3) is shown in Figure 1, simulated for 20 users arriving uniformly at random, 3 user classes and 200 actions. Both the latent class matrix $U_{200 \times 3}$ the mixture matrix $V_{20 \times 3}$ are random one-shot instantiations. The proposed algorithm (Algorithm 3), with two different exploration rate schedules $\tilde{O}(n^{-1/2})$ and $\tilde{O}(n^{-1/3})$ (‘RTP+OFUL(sqrt)’ and ‘RTP+OFUL(cuberoot)’ in the figure), is compared with (a) basic UCB (‘UCB’ in the figure) ignoring the linear structure of the problem (i.e., UCB per-user with 200 actions), (b) OFUL per-user with complete knowledge of the user classes and $p_n = 1$ always, i.e., no exploration mini-sessions, and (c) An implementation of the Alternating Least Squares estimator (Takács and Tikk, 2012; Mary et al., 2014) for the matrix U along with OFUL per-user. The proposed algorithm, with the theoretically suggested exploration $\tilde{O}(n^{-1/2})$, is observed to exploit the latent structure considerably better than simple UCB, and is not too far from the unrealistic OFUL strategy which enjoys the luxury of latent class information. It is also competitive with performing Alternating Least Squares, which does not come with analytically sound performance guarantees in the bandit learning setting. Also, the large additive constants in the theoretical bounds for Algorithm 3 do not manifest here.

Related work. The popular low-rank matrix completion problem studies the recovery U and V given a small number of entries sampled at random from UV^T with both U and V being tall matrices, see for instance Jain et al. (2013) and citations therein. However, its setting is different than ours for several reasons. It typically deals with batch data arising from a sampling process that is not active but uniform across entries of UV^T . Further, it requires sensing operators having strong properties (such as the RIP property), and most importantly, the performance metric is not regret but reconstruction error (Frobenius or 2-norm).

In the linear bandit literature (Abbasi-Yadkori et al., 2011; Rusmevichientong and Tsitsiklis, 2010; Dani et al., 2008), the key constraining assumption is that either user side (V) or item side (U) features are precisely and completely known a priori. In contrast, the problem of low regret recommendation across users with latent mixtures does not afford us the luxury of knowing either U or V , and so they must be learnt “on the fly”. Another related work in the context of bandit type schemes for latent mixture model recommender systems is that of Bresler et al. (2014), in which, under the very specific uniform mixture model for all users, they exhibit strategies with good regret.

Nguyen et al. (2014) consider an alternating minimization type scheme in linear bandit models with two-sided uncertainty (an alternative model involving latent “factors”). However no rigorous guarantees are given for the bandit schemes they present; moreover, it is not known if alternating minimization finds global minima in general. Another related work is in the transfer learning setting

7. With additional prior knowledge of γ_n , the dependence of the additive term can be made polynomial in \circlearrowleft_δ : choosing $\gamma_n = \min\{1, \sqrt{\circlearrowleft_\delta/n}\}$, it holds that $\ell(n_0 - 1 + \sum_{n=n_0}^N \gamma_n) \leq 2\sqrt{\circlearrowleft_\delta \ell T} + \ell$.

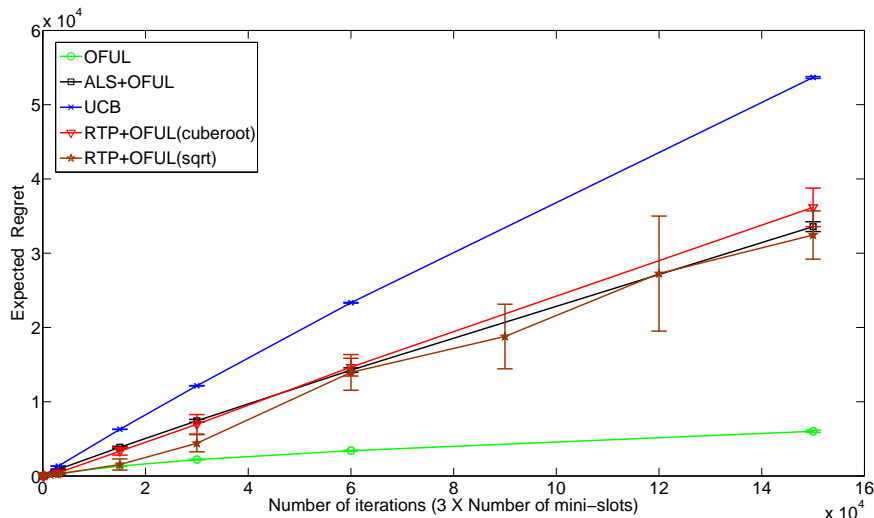


Figure 1: Regret of the proposed algorithm (‘RTP+OFUL’ or Algorithm 3) for two different exploration rate schedules, compared with (a) independent UCB per-user, (b) **OFUL** per-user with perfect knowledge of latent classes U , and (c) Alternating Least Squares estimation for the matrix U , along with **OFUL** per-user. Here, $B = 20$ users, $C = 3$ classes, and $A = 200$, with randomly generated U and V . Plots show the sample mean of cumulative regret with time, with 1 standard deviation-error bars over 10 sample experiments.

from [Lazaric et al. \(2013\)](#): The method combines the RTP method ([Anandkumar et al., 2014b, 2012](#)) essentially with a standard UCB ([Auer et al., 2002](#)), but however works in the setting of a large number interactions with a same user, without assuming access to “user ids”. As a result, the regret bound in this setting scales linearly with the number of rounds. Our result in this paper shows that with additional access to just user identifiers, we can reduce the regret rate to be sublinear in time.

The RTP method has been used as a processing step to the EM algorithm in crowdsourcing ([Zhang et al., 2014](#)), but only convergence properties are considered, which is not enough to provide regret guarantees.

On the theoretical side, our contribution generalizes the setting of *clustered bandits* ([Maillard and Mannor, 2014; Gentile et al., 2014](#)) in which a hard clustering model is assumed (one user is assigned to one class, or equivalently mixture distributions can only have support size 1). In particular, [Maillard and Mannor \(2014\)](#) specifically highlight the benefit of a collaborative gain across users against using a vanilla UCB for each user. However their setting is less general than assuming a soft clustering of users (one user corresponds to a mixture of classes) across various “representative” taste profiles as we study here.

The Alternating Least-Squares (ALS) method ([Takács and Tikk, 2012; Mary et al., 2014](#)) has been shown to yield promising experimental results in similar settings where both U and V are unknown. However, no theoretical guarantees are known for this algorithm that may converge to a local optimum in general.

The work of [Valko et al. \(2014\)](#) studies stochastic bandits with a linear model over a low-rank (graph Laplacian) structure. However, they assume complete knowledge of the graph and hence knowledge of the eigenvectors of the Laplacian, converting it into a bilinear problem with only *one-sided* uncertainty. This is in contrast to our setup where *both* U , V are completely uncertain.

Perhaps the closest work to ours is that of [Djolonga et al. \(2013\)](#) where the authors develop a flexible approach for bandit problems in high dimension but with low-dimensional reward dependence. They use a two-phase algorithm: First a low-rank matrix completion technique (the Dantzig selector) estimates the feature-reward map, then a Gaussian Process-UCB (GP-UCB) bandit algorithm controls the regret, and show that if after n iterations the approximation error between the feature matrix and its estimate is less than η , the final regret is given by the sum of the regret of GP-UCB when given perfect knowledge of the features and of $n + \eta(T - n)$ (due to the learning phase and approximation error). This results in an overall regret scaling with $O(T^{4/5})$. We depart from their results in two fundamental ways: Firstly, they have the possibility of uniformly sampling the entries (a common assumption in low-rank matrix completion techniques). We do not have this luxury in our setting as we do not control the process of user arrivals, that is not constrained to be uniform. Secondly, we prove and exploit a novel robustness property (see [Theorem 8](#)) of the bandit subroutine we use (**OFUL** in our case instead of GP-UCB), which allows us to effectively eliminate the approximation error in their work and obtain a $O(\sqrt{T})$ regret bound (see [Theorem 4](#)).

6. Conclusion & Directions

We consider a full-blown latent class mixture model in which users are described by unknown mixtures across unknown user classes, more general and challenging than when users are assumed to fall perfectly in one class ([Gentile et al., 2014](#); [Maillard and Mannor, 2014](#)).

We provide the first provable sublinear regret guarantees in this setting, when both the canonical classes and user mixture weights are completely unknown, which we believe is striking when compared to existing work in the setting, e.g., alternate minimization typically gets stuck in local minima. We currently use a combination of noisy tensor factorization and linear bandit techniques, and control the uncertainty in the estimates resulting from each one of these techniques. This enables us to effectively recover the latent class structure.

Future directions include reducing the numerical constant (e.g. using an alternative to RTP), and studying how to combine our work with the aggregation of user parameters suggested in [Maillard and Mannor \(2014\)](#).

References

- Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvari. Improved Algorithms for Linear Stochastic Bandits. In *Proc. NIPS*, pages 2312–2320, 2011.
- Animashree Anandkumar, Daniel Hsu, and Sham M Kakade. A method of moments for mixture models and hidden markov models. *arXiv preprint arXiv:1203.0683*, 2012.
- Animashree Anandkumar, Rong Ge, Daniel Hsu, and Sham M Kakade. A tensor approach to learning mixed membership community models. *Journal of Machine Learning Research*, 15(1):2239–2312, 2014a.
- Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *J. Mach. Learn. Res.*, 15(1):2773–2832, January 2014b.
- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2):235–256, 2002.
- Guy Bresler, George H Chen, and Devavrat Shah. A latent source model for online collaborative filtering. In *Proc. NIPS 27*, pages 3347–3355. Curran Associates, Inc., 2014.
- Varsha Dani, Thomas P. Hayes, and Sham M. Kakade. Stochastic Linear Optimization under Bandit Feedback. In *Proc. COLT*, 2008.

- Josip Djolonga, Andreas Krause, and Volkan Cevher. High-dimensional Gaussian process bandits. In *Proc. NIPS*, pages 1025–1033, 2013.
- Anders Forsgren. On linear least-squares problems with diagonally dominant weight matrices. *SIAM Journal on Matrix Analysis and Applications*, 17(4):763–788, 1996.
- Claudio Gentile, Shuai Li, and Giovanni Zappella. Online clustering of bandits. In *Proc. ICML*, pages 757–765, 2014.
- Mohammad Gheshlaghi Azar, Alessandro Lazaric, and Emma Brunskill. Sequential transfer in multi-armed bandit with finite set of models. In *Proc. NIPS*, pages 2220–2228. Curran Associates, Inc., 2013.
- Thomas Hofmann and Jan Puzicha. Latent class models for collaborative filtering. In *IJCAI*, volume 99, pages 688–693, 1999.
- Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *Proc. ACM Symposium on Theory Of computing (STOC)*, pages 665–674. ACM, 2013.
- Jon Kleinberg and Mark Sandler. Using mixture models for collaborative filtering. In *Proc. ACM Symposium on Theory Of Computing (STOC)*, pages 569–578. ACM, 2004.
- Tomáš Kocák, Gergely Neu, Michal Valko, and Rémi Munos. Efficient learning by implicit exploration in bandit problems with side observations. In *Proc. NIPS*, pages 613–621, 2014.
- Alessandro Lazaric, Emma Brunskill, et al. Sequential transfer in multi-armed bandit with finite set of models. In *Proc. NIPS*, pages 2220–2228, 2013.
- Odalric-Ambrym Maillard and Shie Mannor. Latent bandits. In *Proc. ICML*, pages 136–144, 2014.
- Jérémie Mary, Romaric Gaudel, and Preux Philippe. Bandits warm-up cold recommender systems. *arXiv preprint arXiv:1407.2806*, 2014.
- Hai Thanh Nguyen, Jérémie Mary, and Philippe Preux. Cold-start problems in recommendation systems via contextual-bandit algorithms. *arXiv preprint, arXiv:1405.7544*, 2014.
- Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- Gilbert W Stewart, Ji-guang Sun, and Harcourt Brace Jovanovich. *Matrix perturbation theory*, volume 175. Academic press New York, 1990.
- Ilya Sutskever, Joshua B. Tenenbaum, and Ruslan R Salakhutdinov. Modelling relational data using Bayesian clustered tensor factorization. In *Proc. NIPS*, pages 1821–1828. Curran Associates, Inc., 2009.
- Gábor Takács and Domonkos Tikk. Alternating least squares for personalized ranking. In *ACM Conference on Recommender systems*, 2012.
- Michal Valko, Rémi Munos, Branislav Kveton, and Tomáš Kocák. Spectral Bandits for Smooth Graph Functions. In *Proc. ICML*, 2014.
- Yuchen Zhang, Xi Chen, Dengyong Zhou, and Michael I Jordan. Spectral methods meet EM: A provably optimal algorithm for crowdsourcing. In *Proc. NIPS*, pages 1260–1268, 2014.

Appendix A. Proofs of Lemmas 1 and 2

Proof of Lemma 1 This result holds by construction of the estimates $\tilde{r}_{a,a',n}$ and $\tilde{r}_{a,a',a'',n}$. Note that

$$\begin{aligned}
\mathbb{E}[\tilde{r}_{a,a',n}] &= \frac{1}{n} \sum_{i=1}^n \sum_{b \in [B]} \mathbb{E} \left[\frac{X_{a_{i,1},b,i,1} X_{a_{i,2},b,i,2}}{p_i(a,a'|b_i)} \mathbb{I}\{a_{i,1} = a, a_{i,2} = a'\} \middle| b_i = b \right] \beta(b) \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{b \in [B]} \sum_{c \in [C]} \mathbb{E} \left[\frac{X_{a_{i,1},b,i,1} X_{a_{i,2},b,i,2}}{p_i(a,a'|b)} \mathbb{I}\{a_{i,1} = a, a_{i,2} = a'\} \middle| b_i = b, c_i = c \right] v_{b,c} \beta(b) \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{b \in [B]} \sum_{c \in [C]} \mathbb{E} \left[X_{a,b,i,1} X_{a',b,i,2} \middle| b_i = b, c_i = c \right] v_{b,c} \beta(b) \\
&\stackrel{(a)}{=} \frac{1}{n} \sum_{i=1}^n \sum_{b \in [B]} \sum_{c \in [C]} \mathbb{E} \left[X_{a,b,i,1} \middle| b_i = b, c_i = c \right] \mathbb{E} \left[X_{a',b,i,2} \middle| b_i = b, c_i = c \right] v_{b,c} \beta(b) \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{b \in [B]} \sum_{c \in [C]} u_{a,c} u_{a',c} v_{b,c} \beta(b) \\
&= \sum_{c \in [C]} \left(\sum_{b \in [B]} v_{b,c} \beta(b) \right) u_{a,c} u_{a',c} \\
&= \sum_{c \in [C]} v_{\beta,c} u_{a,c} u_{a',c},
\end{aligned}$$

where (a) holds by independence of the sample generated by user b when in the same class c . Note that c_i is the same for all $\ell = 1, 2, 3$ interaction steps, that is $c_i = c_{i,1} = c_{i,2} = c_{i,3}$, where $c_{i,\ell}$ is the class corresponding to sample $X_{a,b,i,\ell}$. This is the reason why we get $u_{a,c} u_{a',c} v_{b,c}$ and not a product $u_{a,c} u_{a',c} v_{b,c}^2$ for instance.

Proof of Lemma 2 Since the rewards generated by each source a, b are i.i.d., the estimate $\tilde{r}_{a,a',n}$ is a sum of i.i.d. random variables bounded in $[0, 1]$, re-weighted by the probability weights $p_i(a, a'|b_i)$, which are measurable functions of the past. Assuming that there exists some deterministic $q_{2,i} > 0$ such that $\forall i \in \mathbb{N}, p_i(a, a'|b_i) \geq q_{2,i}$, we can thus apply a version of Azuma-Hoeffding inequality for bounded martingale difference sequence. Let us recall that by this inequality, for a deterministic time s , and $(Y_m)_{m \leq s} \in [0, 1]$ being a bounded martingale difference sequence, then for all $\delta \in (0, 1)$ it holds

$$\mathbb{P}\left(\left|\frac{1}{s} \sum_{i=1}^s Y_i\right| \geq \sqrt{\frac{\log(2/\delta)}{2s}}\right) \leq \delta.$$

In our case, $Y_i = \frac{X_{a_{i,1},b_i,i,1} X_{a_{i,2},b_i,i,2}}{p_i(a,a'|b_i)} \mathbb{I}\{a_{i,1} = a, a_{i,2} = a'\} - m_{a,a'}$, and we deduce that

$$\mathbb{P}\left(|\tilde{r}_{a,a',n} - m_{a,a'}| \geq \sqrt{\sum_{i=1}^n q_{2,i}^{-2} \frac{\log(2/\delta)}{2n^2}}\right) \leq \delta.$$

Likewise, we get that

$$\mathbb{P}\left(|\tilde{r}_{a,a',a'',n} - m_{a,a',a''}| \geq \sqrt{\sum_{i=1}^n q_{3,i}^{-2} \frac{\log(2/\delta)}{2n^2}}\right) \leq \delta.$$

Taking a union bound over the actions in each case, and then over the two events concludes the proof. \square

Proof of Corollary 1 From Lemma 2, we deduce that on an event of probability higher than $1 - \delta$, it holds simultaneously that

$$e_n^{(2)} \stackrel{\text{def}}{=} \|\widehat{M}_{n,2} - M_2\| \leq A \sqrt{\sum_{m=1}^n q_{2,m}^{-2} \frac{\log(4A^2/\delta)}{2n^2}} \quad \text{and}$$

$$e_n^{(3)} \stackrel{\text{def}}{=} \|\widehat{M}_{n,3} - M_3\| \leq A^{3/2} \sqrt{\sum_{m=1}^n q_{3,m}^{-2} \frac{\log(4A^3/\delta)}{2n^2}}.$$

This indeed holds by relating the norm of the matrix (tensor) with each of the elements. We conclude by replacing the values of $q_{2,i}$ and $q_{3,i}$.

Appendix B. Proof of Theorem 1

We prove in this section a slightly more detailed result, namely, the following:

Theorem 1. Assume that $\{\gamma_i\}_{i \geq 1}$ are chosen such that $n^{-2} \sum_{i=1}^n \gamma_i^{-2} \xrightarrow{n} 0$. Let λ_{\min} be the minimum robust eigenvalue of the tensor $T = M_3(W, W, W)$. Let $\delta \in (0, 1)$. Provided that

$$\frac{n^2}{\sum_{i=1}^n \gamma_i^{-2}} \geq \max \left\{ \frac{2A^6 \log(4A^2/\delta)}{\min\{\Gamma, \sigma_{\min}\}^2}, \frac{A^9 (1 + 10(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}}))(1 + u_{\max}^3)^2 C^5 \log(4A^3/\delta)}{2C_1^2 \lambda_{\min}^2 \sigma_{\min}^3} \right\},$$

with probability higher than $1 - 2\delta$, there exists some permutation $\pi \in \mathbb{S}_C$ such that for all $c \in [C]$,

$$\|u_c - \bar{u}_{n,\pi(c)}\| \leq \Delta A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}} + o(n^{-2} \sum_{i=1}^n \gamma_i^{-2}),$$

where we introduced the problem-dependent constant

$$\Delta = 13\sqrt{\sigma_{\max}} \left(\frac{CA}{\sigma_{\min}} \right)^{3/2} \left(1 + 10\left(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}}\right)(1 + u_{\max}^3) \right) + \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) \frac{1}{v_{\min}^2}.$$

For general $\{\gamma_i\}_{i \geq 1}$ (not necessarily such that $n^{-2} \sum_{i=1}^n \gamma_i^{-2} \xrightarrow{n} 0$), it holds with same probability that

$$\|\mathbf{u}_c - \bar{\mathbf{u}}_{n,\pi(c)}\| \leq \diamond A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}},$$

where, using the notation $\aleph = 1 + 10(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}})(1 + u_{\max}^3)$, we have introduced the constant

$$\begin{aligned} \diamond &= \left(\frac{CA}{\sigma_{\min}} \right)^{3/2} \left(13\sqrt{\sigma_{\max}} + 4\sqrt{2 \min\{\Gamma, \sigma_{\min}\}} + 5\left(\frac{\sigma_{\max}}{\Gamma} + \frac{1}{2\sigma_{\max}}\right) \min\{\Gamma, \sigma_{\min}\} \right) \aleph \\ &+ \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) \frac{1}{v_{\min}^2} + 5\sqrt{3/8} \left(\sqrt{\sigma_{\max}} + \sqrt{\min\{\Gamma, \sigma_{\min}\}/2} \right) \left(\frac{2CA}{\sigma_{\min}} \right)^3 \aleph^2 \min\{\Gamma, \sigma_{\min}\}. \end{aligned}$$

Proof The proof closely follows that of Gheshlaghi Azar et al. (2013). First, note that by property of the rank 1 decomposition ((Anandkumar et al., 2014b, Theorem 4.3)), it holds that $\lambda_c = (\sum_{b \in [B]} v_{b,c} \beta(b))^{-2}$ and thus $v_{\min}^{-2} \geq \lambda_{\max} \geq \lambda_{\min} \geq 1$.

We first decompose the following term to make appear the terms from Proposition 1:

$$\|\mathbf{u}_c - \bar{\mathbf{u}}_{n,\pi(c)}\| \leq \underbrace{|\lambda_c - \hat{\lambda}_{n,\pi(c)}|}_{RTP.1} \underbrace{\|W^{\top\dagger}\|}_b \underbrace{\|\varphi_c\|}_a + \underbrace{|\hat{\lambda}_{\pi(c)}|}_{RTP.3} \underbrace{\|W^{\top\dagger} - \widehat{W}^{\top\dagger}\|}_d \underbrace{\|\varphi_c\|}_1 + \underbrace{|\hat{\lambda}_{\pi(c)}|}_{RTP.3} \underbrace{\|\widehat{W}^{\top\dagger}\|}_c \underbrace{\|\varphi_c - \hat{\varphi}_{n,\pi(c)}\|}_{RTP.2}. \quad (3)$$

Note that φ_c , and $\hat{\varphi}_{n,\pi(c)}$ are both normalized vectors. Thus, (a) is bounded as $\|\varphi_c\| \leq 1$. It holds for (b) that $\|W^{\top\dagger}\| \leq \sqrt{C\sigma_{max}}$, and for (c), on the $1 - \delta$ event Ω from Corollary 1, that

$$\|\widehat{W}^{\top\dagger}\| \leq \sqrt{C\hat{\sigma}_{max}} \leq \sqrt{C}(\sqrt{\sigma_{max}} + \sqrt{e_n^{(2)}}). \quad (4)$$

The term (d) requires a little more work. It holds that

$$\begin{aligned} \|W^{\top\dagger} - \widehat{W}^{\top\dagger}\| &= \|\widehat{U}\widehat{D}^{1/2} - UD^{1/2}\| \\ &\leq \|(\widehat{U} - U)D^{1/2}\| + \|\widehat{U}(\widehat{D}^{1/2} - D^{1/2})\| \\ &\leq \underbrace{\|\widehat{U} - U\|}_{e} \sigma_{max} + \underbrace{\|\widehat{D}^{1/2} - D^{1/2}\|}_{f} \sqrt{C}. \end{aligned}$$

We use the result of Lemma 5 from Gheshlaghi Azar et al. (2013) to control (e) and (f). If $e_n^{(2)} \leq \frac{1}{2}\Gamma$, then it holds

$$\|\widehat{D}^{1/2} - D^{1/2}\| \leq \frac{e_n^{(2)}}{\sigma_{max}} \quad \|\widehat{U} - U\| \leq \frac{2\sqrt{C}e_n^{(2)}}{\Gamma},$$

from which we deduce that

$$\|W^{\top\dagger} - \widehat{W}^{\top\dagger}\| \leq \left(\frac{2\sigma_{max}}{\Gamma} + \frac{1}{\sigma_{max}}\right)\sqrt{C}e_n^{(2)}. \quad (5)$$

At this point, (RTP.1), (RTP.2) and (RTP.3) are controlled by the perturbation method from Anandkumar et al. (2014b), under the condition that $e_n = \|T - \widehat{T}_n\| \leq C_1 \frac{\lambda_{min}}{C}$ (where C_1 is a universal constant). In this case, with probability $1 - \delta$, the RTP algorithm with well-chosen parameters achieves

$$\begin{aligned} |\lambda_c - \hat{\lambda}_{n,\pi(c)}| &\leq 5\|T - \widehat{T}_n\| \\ \|\varphi_c - \hat{\varphi}_{n,\pi(c)}\| &\leq 8 \frac{\|T - \widehat{T}_n\|}{\lambda_c}. \end{aligned}$$

In order to make the condition explicit in our setting, we use the fact that by Lemma 6 from Gheshlaghi Azar et al. (2013), if $e_n^{(2)} \leq \frac{1}{2} \min\{\Gamma, \sigma_{min}\}$ then

$$e_n \leq \left(\frac{C}{\sigma_{min}}\right)^{3/2} \left(e_n^{(3)} + 2(1 + \sqrt{2} + 2)e_n^{(2)}\left(\frac{1}{\Gamma\sigma} + \frac{1}{\sigma_{min}}\right)(e_n^{(3)} + \max_c \|\mathbf{u}_c\|^3)\right). \quad (6)$$

The condition $e_n^{(2)} \leq \frac{1}{2} \min\{\Gamma, \sigma_{min}\}$ holds if the number of sessions n is sufficiently large: Indeed on an event of probability higher than $1 - \delta$, then it is enough that

$$A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^2/\delta)}{2n^2}} \leq \frac{1}{2} \min\{\Gamma, \sigma_{min}\},$$

that is, reordering the terms, that

$$\frac{n^2}{\sum_{i=1}^n \gamma_i^{-2}} \geq \frac{2A^6 \log(4A^2/\delta)}{\min\{\Gamma, \sigma_{\min}\}^2}. \quad (7)$$

Now, in order to satisfy the condition $e_n = \|T - \hat{T}_n\| \leq C_1 \frac{\lambda_{\min}}{C}$, it is enough that

$$\left(\frac{C}{\sigma_{\min}}\right)^{3/2} \left(e_n^{(3)} + 2(1 + \sqrt{2} + 2)e_n^{(2)} \left(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}}\right) (e_n^{(3)} + \max_c \|u_c\|^3) \right) \leq C_1 \frac{\lambda_{\min}}{C}.$$

Let us decompose the left-hand-side term: After some simplifications using $\max_c \|u_c\|^3 \leq A^{3/2} u_{\max}^3$ and $e_n^{(3)} \leq A^{3/2}$, the previous inequality happens when

$$e_n^{(3)} + A^{3/2} \circ e_n^{(2)} \leq C_1 \frac{\lambda_{\min} \sigma_{\min}^{3/2}}{C^{5/2}}.$$

where $\circ = 2(1 + \sqrt{2} + 2) \left(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}}\right) (1 + u_{\max}^3)$. Using the definition of $e_n^{(3)}$ and $e_n^{(2)}$ then we deduce that it is enough that

$$(1 + \circ) A^{9/2} \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^3/\delta)}{2n^2}} \leq C_1 \frac{\lambda_{\min} \sigma_{\min}^{3/2}}{C^{5/2}},$$

that is, reordering the terms that

$$\frac{n^2}{\sum_{i=1}^n \gamma_i^{-2}} \geq \frac{A^9 (1 + \circ)^2 C^5 \log(4A^3/\delta)}{2C_1^2 \lambda_{\min}^2 \sigma_{\min}^3}. \quad (8)$$

Combining the decomposition (3) with (4),(5), and using the fact that $v_{\min}^{-2} \geq \lambda_c \geq 1$, we obtain

$$\begin{aligned} \|\mathbf{u}_c - \bar{\mathbf{u}}_{n,\pi(c)}\| &\leq 5e_n \sqrt{C} \sqrt{\sigma_{\max}} + (\lambda_c + 5e_n) \sqrt{C} \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) e_n^{(2)} \\ &\quad + 8\sqrt{C} (\lambda_c + 5e_n) (\sqrt{\sigma_{\max}} + \sqrt{e_n^{(2)}}) \frac{e_n}{\lambda_c}. \\ &\leq \sqrt{C} \left[13\sqrt{\sigma_{\max}} e_n + \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) \frac{e_n^{(2)}}{v_{\min}^2} + 8\sqrt{e_n^{(2)}} e_n \right. \\ &\quad \left. + 5 \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) e_n^{(2)} e_n + 40(\sqrt{\sigma_{\max}} + \sqrt{e_n^{(2)}}) e_n^2 \right]. \end{aligned}$$

Now, using (6) and unfolding the last inequality, it holds with probability higher than $1 - 2\delta$ that

$$\begin{aligned} &\|\mathbf{u}_c - \bar{\mathbf{u}}_{n,\pi(c)}\| \\ &\leq \sqrt{C} \left[13\sqrt{\sigma_{\max}} \left(\frac{C}{\sigma_{\min}} \right)^{3/2} (e_n^{(3)} + e_n^{(2)} A^{3/2} \circ) + \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) \frac{e_n^{(2)}}{v_{\min}^2} \right. \\ &\quad \left. + 8 \left(\frac{C}{\sigma_{\min}} \right)^{3/2} \sqrt{e_n^{(2)}} (e_n^{(3)} + e_n^{(2)} A^{3/2} \circ) \right. \\ &\quad \left. + 5 \left(\frac{C}{\sigma_{\min}} \right)^{3/2} \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) e_n^{(2)} (e_n^{(3)} + e_n^{(2)} A^{3/2} \circ) + 40(\sqrt{\sigma_{\max}} + \sqrt{e_n^{(2)}}) e_n^2 \right] \\ &\leq \left[13\sqrt{\sigma_{\max}} \left(\frac{CA}{\sigma_{\min}} \right)^{3/2} (1 + \circ) + \left(\frac{2\sigma_{\max}}{\Gamma} + \frac{1}{\sigma_{\max}} \right) \frac{1}{v_{\min}^2} \right] \\ &\quad \times A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}} + o(n^{-2} \sum_{i=1}^n \gamma_i^{-2}), \end{aligned}$$

which, after some cosmetic simplifications, concludes the first part of the proof of Theorem 1.

Alternatively, when $n^{-2} \sum_{i=1}^n \gamma_i^{-2} \not\rightarrow \infty$, we can always resort to the condition that $e_n^{(2)} \leq 1/2 \min\{\Gamma, \sigma_{\min}\}$ in order to simplify the previous derivation. We deduce, similarly, that

$$\begin{aligned}
& \|\mathbf{u}_c - \bar{\mathbf{u}}_{n,\pi(c)}\| \\
& \leq \sqrt{C} \left[13\sqrt{\sigma_{\max}} \left(\frac{C}{\sigma_{\min}}\right)^{3/2} (e_n^{(3)} + e_n^{(2)} A^{3/2} \circ) + \left(\frac{2\sigma_{\max}}{\Gamma_\sigma} + \frac{1}{\sigma_{\max}}\right) \frac{e_n^{(2)}}{v_{\min}^2} \right. \\
& \quad + 8\left(\frac{C}{\sigma_{\min}}\right)^{3/2} \sqrt{e_n^{(2)}} (e_n^{(3)} + e_n^{(2)} A^{3/2} \circ) \\
& \quad \left. + 5\left(\frac{C}{\sigma_{\min}}\right)^{3/2} \left(\frac{2\sigma_{\max}}{\Gamma_\sigma} + \frac{1}{\sigma_{\max}}\right) e_n^{(2)} (e_n^{(3)} + e_n^{(2)} A^{3/2} \circ) + 40(\sqrt{\sigma_{\max}} + \sqrt{e_n^{(2)}}) e_n^{(2)} \right] \\
& \leq \left[\left(13\sqrt{\sigma_{\max}} \left(\frac{CA}{\sigma_{\min}}\right)^{3/2} + 8\left(\frac{CA}{\sigma_{\min}}\right)^{3/2} \sqrt{\min\{\Gamma, \sigma_{\min}\}/2} \right. \right. \\
& \quad \left. \left. + 5\left(\frac{CA}{\sigma_{\min}}\right)^{3/2} \left(\frac{\sigma_{\max}}{\Gamma_\sigma} + \frac{1}{2\sigma_{\max}}\right) \min\{\Gamma, \sigma_{\min}\} \right) (1 + \circ) + \left(\frac{2\sigma_{\max}}{\Gamma_\sigma} + \frac{1}{\sigma_{\max}}\right) \frac{1}{v_{\min}^2} \right. \\
& \quad \left. + 40\left(\sqrt{\sigma_{\max}} + \sqrt{\min\{\Gamma, \sigma_{\min}\}/2}\right) \left(\frac{CA}{\sigma_{\min}}\right)^3 (1 + \circ)^2 \min\{\Gamma, \sigma_{\min}\} \sqrt{3/8} \right] \\
& \quad \times A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}},
\end{aligned}$$

where, in order to control the last term e_n^2 , we used the property that

$$\begin{aligned}
e_n & \leq \left(\frac{CA}{\sigma_{\min}}\right)^{3/2} (1 + \circ) \min\left\{e_n^{(2)} \sqrt{\frac{\log(4A^3/\delta)}{\log(4A^2/\delta)}}, A^{-3/2} e_n^{(3)}\right\} \\
& \leq \left(\frac{CA}{\sigma_{\min}}\right)^{3/2} (1 + \circ) \min\{\sqrt{3/2} e_n^{(2)}, A^{-3/2} e_n^{(3)}\}.
\end{aligned}$$

■

Appendix C. Proof of Theorem 3

Proof Let $\mathbf{M}_{1:t} = (\mathbf{m}_{A_1}, \dots, \mathbf{m}_{A_t})^\top$. The argument used to prove Theorem 2 in Yadkori et al, 2011, can be used to show that

$$\hat{\mathbf{v}}_{t-1} = V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \boldsymbol{\eta}_{1:t-1} + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{M}_{1:t-1}$$

where $\boldsymbol{\eta}_{1:t-1} := (\eta_1, \dots, \eta_{t-1})$ is the observed noise sequence. Let $\mathbf{E}_{1:t-1} := (\varepsilon_{A_1}, \dots, \varepsilon_{A_t})^\top = \mathbf{M}_{1:t-1} - \bar{\mathbf{U}}_{1:t-1} \mathbf{v}^\circ$. We then have

$$\begin{aligned}
\hat{\mathbf{v}}_{t-1} & = V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \boldsymbol{\eta}_{1:t-1} + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{M}_{1:t-1} \\
& = V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \boldsymbol{\eta}_{1:t-1} + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} (\bar{\mathbf{U}}_{1:t-1}^\top \mathbf{v}^\circ + \mathbf{E}_{1:t-1}) \\
& = V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \boldsymbol{\eta}_{1:t-1} + \mathbf{v}^\circ - \lambda V_{t-1}^{-1} \mathbf{v}^\circ + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{E}_{1:t-1}.
\end{aligned}$$

Thus, letting $\mathbf{v}_{t-1}^+ := \mathbf{v}^\circ + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{E}_{1:t-1}$ and using the above with techniques from Yadkori et al together with $\|\mathbf{v}^\circ\|_2 \leq R_\Theta$, we have that

$$\mathbf{v}_{t-1}^+ \in \mathcal{C}_{t-1}$$

with probability at least $1 - \delta$.

Now, let $a_{t-1}^+ \in \arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^\top \mathbf{v}_{t-1}^+$ be an optimal action corresponding to the approximate parameter \mathbf{v}_{t-1}^+ , and define the instantaneous regret at time t with respect to the approximate parameter as

$$r_t^+ := \bar{\mathbf{u}}_{a_{t-1}^+}^\top \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^\top \mathbf{v}_{t-1}^+ \geq 0.$$

We now bound this approximate regret using arguments along the lines of Yadkori et al, 2011. Consider

$$\begin{aligned} r_t^+ &= \bar{\mathbf{u}}_{a_{t-1}^+}^\top \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^\top \mathbf{v}_{t-1}^+ \\ &\leq \bar{\mathbf{u}}_{A_t}^\top \tilde{\mathbf{v}}_t - \bar{\mathbf{u}}_{A_t}^\top \mathbf{v}_{t-1}^+ \quad (\text{since } (A_t, \tilde{\mathbf{v}}_t) \text{ is optimistic}) \\ &= \bar{\mathbf{u}}_{A_t}^\top (\tilde{\mathbf{v}}_t - \mathbf{v}_{t-1}^+) \\ &= \bar{\mathbf{u}}_{A_t}^\top (\tilde{\mathbf{v}}_t - \hat{\mathbf{v}}_{t-1}) + \bar{\mathbf{u}}_{A_t}^\top (\hat{\mathbf{v}}_{t-1} - \mathbf{v}_{t-1}^+) \\ &\leq \|\bar{\mathbf{u}}_{A_t}\|_{V_{t-1}^{-1}} \|\tilde{\mathbf{v}}_t - \hat{\mathbf{v}}_{t-1}\|_{V_{t-1}} + \|\bar{\mathbf{u}}_{A_t}\|_{V_{t-1}^{-1}} \|\hat{\mathbf{v}}_{t-1} - \mathbf{v}_{t-1}^+\|_{V_{t-1}} \quad (\text{Cauchy-Schwarz's inequality}) \\ &\leq 2D_{t-1} \|\bar{\mathbf{u}}_{A_t}\|_{V_{t-1}^{-1}}. \end{aligned} \tag{9}$$

Noting that $m_a \in [-1, 1] \forall a$, the regret can be written as

$$\begin{aligned} R_T &= \sum_{t=1}^T (m_{a^*} - m_{A_t}) = \sum_{t=1}^T \min\{m_{a^*} - m_{A_t}, 2\} \\ &= \rho' \sum_{a \neq a^*} \sum_{t=1}^T \min \left\{ \frac{m_{a^*} - m_a}{\rho'}, \frac{2}{\rho'} \right\} \mathbb{I}\{A_t = a\} \\ &\leq \rho' \sum_{a \neq a^*} \sum_{t=1}^T \min \left\{ \bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ, \frac{2}{\rho'} \right\} \mathbb{I}\{A_t = a\} \quad (\text{using the definition of } \rho') \\ &\stackrel{(a)}{\leq} \rho' \sum_{t=1}^T \min \left\{ 2(\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^\top \mathbf{v}_{t-1}^+), \frac{2}{\rho'} \right\} \stackrel{(b)}{=} 2\rho' \sum_{t=1}^T \min \left\{ \bar{\mathbf{u}}_{a_{t-1}^+}^\top \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^\top \mathbf{v}_{t-1}^+, \frac{1}{\rho'} \right\} \\ &= 2\rho' \sum_{t=1}^T \min \left\{ r_t^+, \frac{1}{\rho'} \right\} = \rho' \sum_{t=1}^T \frac{2}{\rho'} \min \{ \rho' r_t^+, 1 \} \stackrel{(c)}{\leq} \rho' \sum_{t=1}^T \frac{2}{\rho'} \min \left\{ 2\rho' D_{t-1} \|\bar{\mathbf{u}}_{A_t}\|_{V_{t-1}^{-1}}, 1 \right\} \\ &\stackrel{(d)}{\leq} \rho' \sum_{t=1}^T 4D_{t-1} \min \left\{ \|\bar{\mathbf{u}}_{A_t}\|_{V_{t-1}^{-1}}, 1 \right\} \\ &\leq \rho' \sqrt{T \sum_{t=1}^T 16D_T^2 \min \left\{ \|\bar{\mathbf{u}}_{A_t}\|_{V_{t-1}^{-1}}^2, 1 \right\}} \quad (\text{by using Cauchy-Schwarz's inequality}). \end{aligned}$$

In the derivation above,

- Steps (a) and (b) hold because of the following. By Lemma 4 (to follow below), $\|\mathbf{v}_{t-1}^+ - \mathbf{v}^\circ\|_2 = \|\bar{V}_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{E}_{1:t-1}\|_2 \leq \alpha(\bar{U}) \|\varepsilon\|_2$. Since $\arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ$ is uniquely a^* by hypothesis, we have, thanks to Lemma 5 (to follow below), that $\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_a^\top \mathbf{v}_{t-1}^+ > \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2} > 0 \forall a \neq a^*$, establishing (a). This in turn shows that the optimal action for \mathbf{v}_{t-1}^+ is uniquely a^* at all times t , i.e., $a_{t-1}^+ = \arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^\top \mathbf{v}_{t-1}^+ = a^*$, which is equality (b).
- Inequality (c) holds by (9) and (d) holds because $\rho' \geq 1$ by definition, and $D_{t-1} \geq \lambda^{1/2} R_\Theta \geq 1/2$ by hypothesis, implying that $2\rho' D_{t-1} \geq 1$.

The argument from here can be continued in the same way as in [Abbasi-Yadkori et al. \(2011\)](#) to yield

$$R_T \leq 8\rho' \sqrt{TC \log \left(1 + \frac{TR_{\lambda}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{TR_{\lambda}^2}{\lambda C} \right)} \right).$$

This proves the theorem. ■

Lemma 4 (Analysis of the time-varying parameter error $V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{E}_{1:t-1}$) Let $\varepsilon_a = m_a - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ$ be the bias in arm a 's reward due to model error, and let $\varepsilon \equiv (\varepsilon_a)_{a \in \mathcal{A}}$ be the $|\mathcal{A}|$ dimensional vector of arm reward biases. Then,

$$\|V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{E}_{1:t-1}\|_2 \leq \left(\max_J \|\mathbf{A}_J^{-1}\|_2 \right) \|\varepsilon\|_2,$$

where $\mathbf{A}_{(A+C) \times C} = \begin{bmatrix} \bar{\mathbf{U}} \\ I_d \end{bmatrix}$, \mathbf{A}_J is the $C \times C$ submatrix of \mathbf{A} formed by picking rows J , and J ranges over all subsets of full-rank rows of \mathbf{A} .

Proof [Proof of Lemma 4] Let $z_{t-1} := V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1} \mathbf{E}_{1:t-1} = \mathbf{v}_{t-1} - \mathbf{v}^\circ \in \mathbb{R}^C$, with $\|\mathbf{E}_{1:t-1}\|_\infty \leq \|\varepsilon\|_\infty = \|\mathbf{m} - \bar{\mathbf{U}} \mathbf{v}^\circ\|_\infty$. We have

$$\begin{aligned} z_{t-1} &= \left(\sum_{s=1}^{t-1} \bar{\mathbf{u}}_{A_s} \bar{\mathbf{u}}_{A_s}^\top + \lambda I \right)^{-1} \sum_{s=1}^{t-1} \varepsilon_{A_s} \bar{\mathbf{u}}_{A_s} \\ &= \left(\frac{1}{t-1} \sum_{s=1}^{t-1} \bar{\mathbf{u}}_{A_s} \bar{\mathbf{u}}_{A_s}^\top + \frac{\lambda}{t-1} I \right)^{-1} \frac{1}{t-1} \sum_{s=1}^{t-1} \varepsilon_{A_s} \bar{\mathbf{u}}_{A_s} \\ &= \left(\sum_{a \in \mathcal{A}} \bar{\mathbf{u}}_a \bar{\mathbf{u}}_a^\top \frac{\sum_{s=1}^{t-1} \mathbb{I}\{A_s = a\}}{t-1} + \frac{\lambda}{t-1} I \right)^{-1} \sum_{a \in \mathcal{A}} \varepsilon_a \bar{\mathbf{u}}_a \frac{\sum_{s=1}^{t-1} \mathbb{I}\{A_s = a\}}{t-1} \\ &= \left(\sum_{a \in \mathcal{A}} \bar{\mathbf{u}}_a \bar{\mathbf{u}}_a^\top f_a(t-1) + \frac{\lambda}{t-1} I \right)^{-1} \sum_{a \in \mathcal{A}} \varepsilon_a \bar{\mathbf{u}}_a f_a(t-1), \end{aligned}$$

where $f_a(t-1) \equiv f_a$ represents the empirical frequency with which action $a \in \mathcal{A}$ has been played up to and including time $t-1$. This allows us to equivalently interpret z_{t-1} as the solution of a *weighted* ℓ^2 -regularized least squares regression problem with $K = |\mathcal{A}|$ observations (instead of the original interpretation with $t-1$ observations) as follows.

Let $\mathbf{F}^{1/2}$ be the $A \times A$ diagonal matrix with the values $\sqrt{f_1}, \dots, \sqrt{f_A}$ on the diagonal (note: $\sum_{a=1}^A f_a = 1$). With this, we can express z_{t-1} as

$$\begin{aligned} z_{t-1} &= \arg \min_{z \in \mathbb{R}^C} \left\| \mathbf{F}^{1/2} \bar{\mathbf{U}} z - \mathbf{F}^{1/2} \varepsilon \right\|_2^2 + \frac{\lambda}{t-1} \|z\|_2^2 \\ &= \arg \min_{z \in \mathbb{R}^C} \left\| \mathbf{F}^{1/2} (\bar{\mathbf{U}} z - \varepsilon) \right\|_2^2 + \frac{\lambda}{t-1} \|z\|_2^2 \\ &= \arg \min_{z \in \mathbb{R}^C} \left\| \begin{bmatrix} \mathbf{F}^{1/2} & 0 \\ 0 & \sqrt{\frac{\lambda}{t-1}} I_C \end{bmatrix} \left(\begin{bmatrix} \bar{\mathbf{U}} \\ I_C \end{bmatrix} z - \begin{bmatrix} \varepsilon \\ 0 \end{bmatrix} \right) \right\|_2^2 \\ &\equiv \arg \min_{z \in \mathbb{R}^C} \left\| \mathbf{D}^{1/2} (\mathbf{A} z - \mathbf{b}) \right\|_2^2 = (\mathbf{A}^\top \mathbf{D} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{D} \mathbf{b}, \end{aligned}$$

with $\mathbf{D}^{1/2}$ being a $(A + C) \times (A + C)$ diagonal & positive semidefinite matrix, $\mathbf{A}^\top \mathbf{D} \mathbf{A} = \sum_{a \in \mathcal{A}} \bar{\mathbf{u}}_a \bar{\mathbf{u}}_a^\top f_a(t-1) + \frac{\lambda}{t-1} I$ positive definite, and \mathbf{A} having full column rank C . A result of Forsgren (1996, Corollary 2.3) can now be applied to yield

$$\|(\mathbf{A}^\top \mathbf{D} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{D}\|_2 \leq \max_J \|\mathbf{A}_J^{-1}\|_2$$

where J ranges over all subsets of full-rank rows of \mathbf{A} , and \mathbf{A}_J is the $C \times C$ submatrix of \mathbf{A} formed by picking rows J . Thus, $\|z_{t-1}\|_2 \leq (\max_J \|\mathbf{A}_J^{-1}\|_2) \|\varepsilon\|_2$. This proves the lemma. \blacksquare

Lemma 5 (Critical radius) *Let $\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ > \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ \forall a \neq a^*$. Then, the following are equivalent:*

$$\|\mathbf{v} - \mathbf{v}^\circ\|_2 \leq \alpha(\bar{U}) \|\varepsilon\|_2 \Rightarrow \bar{\mathbf{u}}_{a^*}^\top \mathbf{v} - \bar{\mathbf{u}}_a^\top \mathbf{v} > \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2} \quad \forall a \neq a^*, \quad (10)$$

and

$$\|\varepsilon\|_2 < \min_{a \neq a^*} \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2\alpha(\bar{U}) \|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2}. \quad (11)$$

Proof [Proof of Lemma 5] Assuming (11), observe that when \mathbf{v} lies in the interior of an $\alpha(\bar{U}) \|\varepsilon\|_2$ -ball around \mathbf{v}° , we have, for any $a \neq a^*$,

$$\begin{aligned} (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v} &= (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v}^\circ + (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top (\mathbf{v} - \mathbf{v}^\circ) \\ &\geq (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v}^\circ + \min_{\|\psi\|_2 \leq \alpha(\bar{U}) \|\varepsilon\|_2} (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \psi \\ &= (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v}^\circ - \alpha(\bar{U}) \|\varepsilon\|_2 \|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2 \\ &> (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v}^\circ - \alpha(\bar{U}) \|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2 \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2\alpha(\bar{U}) \|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2} \\ &= \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2}, \end{aligned}$$

which proves one direction of the lemma. For the other direction, note that if $\|\varepsilon\|_2 \geq \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2\alpha(\bar{U}) \|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2}$

for some $a \neq a^*$, then by setting $\mathbf{v} = \mathbf{v}^\circ - \frac{(\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ)(\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)}{2\|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2^2}$, we have both

$$\|\mathbf{v} - \mathbf{v}^\circ\|_2 = \left\| \frac{(\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ)(\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)}{2\|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2^2} \right\|_2 = \frac{\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ}{2\|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2} \leq \alpha(\bar{U}) \|\varepsilon\|_2$$

and

$$(\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v} = (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v}^\circ - (\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \frac{(\bar{\mathbf{u}}_{a^*}^\top \mathbf{v}^\circ - \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ)(\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)}{2\|\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a\|_2^2} = \frac{(\bar{\mathbf{u}}_{a^*} - \bar{\mathbf{u}}_a)^\top \mathbf{v}^\circ}{2}$$

which contradicts (10), and we are done. \blacksquare

C.1 Proof of Lemma 3

We begin by establishing some auxiliary technical results, which together imply Lemma 3.

Lemma 6 (Controlling α_n) *If n is large enough so that (1) and*

$$\diamond A^3 C \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^3/\delta)}{2n^2}} \leq \frac{1}{2\alpha_*}, \quad (12)$$

hold, then with probability at least $1 - \delta$,

$$\alpha_n \leq 2\alpha_*. \quad (13)$$

Proof [Proof of Lemma 6] The first step is to estimate the factor α in the analysis of Perturbed OFUL. Towards this, note that the quantity $\alpha \equiv \alpha(\bar{U})$ in our setting becomes

$$\alpha_n \equiv \alpha_n(\bar{U}_n) = \max_J \|(u_n^\diamond)_J^{-1}\|_2,$$

where $u_n^\diamond := \begin{bmatrix} \bar{U}_n \\ I_C \end{bmatrix}$ has rank C , and J ranges over all combinations of its C full-rank rows. For any such subset of C linearly independent rows J , we have, after denoting $u^\diamond := \begin{bmatrix} U \\ I_C \end{bmatrix}$, that

$$\|(u_n^\diamond)_J^{-1}\|_2 \leq \|(u^\diamond)_J^{-1}\|_2 + \|(u_n^\diamond)_J^{-1} - (u^\diamond)_J^{-1}\|_2.$$

The final term above can be bounded using [Anandkumar et al. \(2012, Lemma E.4\)](#) – a version of Theorem 2.5 in [Stewart et al. \(1990\)](#). Assuming $(u^\diamond)_J$ is invertible, and $\|(u^\diamond)_J^{-1} ((u_n^\diamond)_J - (u^\diamond)_J)\|_2 < 1$, then $(u_n^\diamond)_J$ is invertible, and a resulting bound on the norm of its inverse lets us write

$$\|(u^\diamond)_J^{-1}\|_2 + \|(u_n^\diamond)_J^{-1} - (u^\diamond)_J^{-1}\|_2 \leq \|(u^\diamond)_J^{-1}\|_2 + \frac{\|(u_n^\diamond)_J - (u^\diamond)_J\|_2 \|(u^\diamond)_J^{-1}\|_2^2}{1 - \|(u^\diamond)_J^{-1} ((u_n^\diamond)_J - (u^\diamond)_J)\|_2}.$$

Writing $J = J_u \cup J_l$ (u and l stand for “upper” and “lower”) with J_l representing the subset of rows taken from the bottom C rows of u_n^\diamond (i.e., I_C), we have

$$(u_n^\diamond)_J - (u^\diamond)_J = \begin{bmatrix} (\bar{U}_n - U)_{J_u} \\ 0 \end{bmatrix}.$$

Thus, with $\|\cdot\|_F$ denoting the Frobenius norm, and using the dominance of the Frobenius norm over the matrix 2-norm, with probability at least $1 - \delta$,

$$\begin{aligned} \|(u_n^\diamond)_J - (u^\diamond)_J\|_2 &\leq \|(u_n^\diamond)_J - (u^\diamond)_J\|_F = \|(\bar{U}_n - U)_{J_u}\|_F \leq \|\bar{U}_n - U\|_F \\ &= \sqrt{\sum_{c \in [C]} \|\bar{U}_{n,c} - U_c\|_2^2} \\ &\leq \diamond A^3 C \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^3/\delta)}{2n^2}}, \end{aligned} \quad (14)$$

from the RTP error estimate (1).

Now, letting $\alpha \equiv \alpha(U) = \max_J \|(u_J^\diamond)^{-1}\|_2$, the result above implies that for any suitable J ,

$$\begin{aligned} \|(u^\diamond)_J^{-1} ((u_n^\diamond)_J - (u^\diamond)_J)\|_2 &\leq \|(u^\diamond)_J^{-1}\|_2 \|(u_n^\diamond)_J - (u^\diamond)_J\|_2 \\ &\leq \alpha \|(u_n^\diamond)_J - (u^\diamond)_J\|_2 \\ &\leq \alpha \diamond A^3 C \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^3/\delta)}{2n^2}} \\ &< 1/2 \end{aligned}$$

whenever n is large enough to satisfy (12).

When the condition (12) above holds, we get, for any J at time n ,

$$\begin{aligned}
\|(u_n^\diamond)^{-1}\|_2 &\leq \|(u^\diamond)^{-1}\|_2 + \frac{\|(u_n^\diamond)_J - (u^\diamond)_J\|_2 \|(u^\diamond)^{-1}\|_2^2}{1 - \|(u^\diamond)^{-1}((u_n^\diamond)_J - (u^\diamond)_J)\|_2} \\
&\leq \alpha + 2\alpha^2 \|(u_n^\diamond)_J - (u^\diamond)_J\|_2 \\
&\leq \alpha + 2\alpha^2 \diamond A^3 C \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^3/\delta)}{2n^2}} \quad [\text{by (14)}] \\
&\leq \alpha + 2\alpha^2 \frac{1}{2\alpha} = 2\alpha.
\end{aligned}$$

This shows that $\alpha_n = \max_J \|(u_n^\diamond)^{-1}\|_2 \leq 2\alpha$. ■

Lemma 7 (Sufficient condition for (2)) *If n is large enough so that (1), (12) and*

$$\diamond A^3 C \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{\log(4A^3/\delta)}{2n^2}} \leq \min \left\{ \frac{g_b}{4\sqrt{A} \|\mathbf{v}_b\|_2}, \frac{g_b}{16\alpha_* \sqrt{C} u_{\max} \|\mathbf{v}_b\|_2 + g_b} \right\} \quad (15)$$

hold, then (2) is satisfied with probability at least $1 - \delta$.

Proof [Proof of Lemma 7] The term $\|\varepsilon\|_2 = \|(U - \bar{U}_n) \mathbf{v}_b\|_2$ is bounded from above by

$$\begin{aligned}
\|U - \bar{U}_n\|_2 \|\mathbf{v}_b\|_2 &\leq \|U - \bar{U}_n\|_F \|\mathbf{v}_b\|_2 \\
&\leq \sqrt{C} \|\mathbf{v}_b\|_2 \diamond A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}} \quad (\text{by (14)}) \\
&\equiv \sqrt{C} \|\mathbf{v}_b\|_2 \aleph_n, \text{ say.}
\end{aligned} \quad (16)$$

For any $a \neq a^*$,

$$(\bar{\mathbf{u}}_{n,a^*} - \bar{\mathbf{u}}_{n,a})^\top \mathbf{v}_b = (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b + \partial_a \mathbf{v}_b \geq \zeta_a, \quad (17)$$

with $\partial_a^\top := (\bar{\mathbf{u}}_{n,a^*} - \mathbf{u}_{a^*}) - (\bar{\mathbf{u}}_{n,a} - \mathbf{u}_a)$, and $\zeta_a := \inf_{\|\xi\|_2 \leq \|\partial_a\|_2} (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b + \xi^\top \mathbf{v}_b$.

Also, by (14), we have

$$\begin{aligned}
\max_{a \in [A]} \|\bar{\mathbf{u}}_{n,a} - \mathbf{u}_a\|_2 &\leq \sqrt{AC} \max_{c \in [C]} \|\bar{\mathbf{u}}_{n,c} - \mathbf{u}_c\|_2 \\
&\leq \sqrt{AC} \diamond A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}} \\
&=: \aleph_n \sqrt{AC}.
\end{aligned}$$

Thus,

$$\zeta_a \geq \inf_{\|\xi\|_2 \leq 2\aleph_n \sqrt{AC}} (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b + \xi^\top \mathbf{v}_b = (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b - 2\aleph_n \sqrt{AC} \|\mathbf{v}_b\|_2. \quad (18)$$

By (17) and (18), for any $a \neq a^*$,

$$(\bar{\mathbf{u}}_{n,a^*} - \bar{\mathbf{u}}_{n,a})^\top \mathbf{v}_b \geq (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b - 2\aleph_n \sqrt{AC} \|\mathbf{v}_b\|_2. \quad (19)$$

We also have

$$\begin{aligned} \|\bar{\mathbf{u}}_{n,a^*} - \bar{\mathbf{u}}_{n,a}\|_2 &\leq \|\mathbf{u}_{a^*} - \mathbf{u}_a\|_2 + \|\bar{\mathbf{u}}_{n,a^*} - \mathbf{u}_{a^*}\|_2 + \|\bar{\mathbf{u}}_{n,a} - \mathbf{u}_a\|_2 \\ &\leq \|\mathbf{u}_{a^*} - \mathbf{u}_a\|_2 + 2\aleph_n \sqrt{AC} \end{aligned} \quad (20)$$

whenever (12) holds. Putting (16), (19), (20) and the conclusion of Lemma 6 together, we have that condition (2) in our case, i.e.,

$$\|\varepsilon\|_2 \equiv \|(U - \bar{U}_n) \mathbf{v}_b\|_2 \leq \min_{a \neq a^*} \frac{(\bar{\mathbf{u}}_{n,a^*} - \bar{\mathbf{u}}_{n,a})^\top \mathbf{v}_b}{2\alpha_n \|\bar{\mathbf{u}}_{n,a^*} - \bar{\mathbf{u}}_{n,a}\|_2}$$

is satisfied when

$$\sqrt{C} \|\mathbf{v}_b\|_2 \aleph_n \leq \min_{a \neq a^*} \frac{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b - 2\aleph_n \sqrt{AC} \|\mathbf{v}_b\|_2}{4\alpha_* \|\mathbf{u}_{a^*} - \mathbf{u}_a\|_2 + 2\aleph_n \sqrt{AC}}.$$

This, in turn, is satisfied if

$$\begin{aligned} 2\aleph_n \sqrt{AC} \|\mathbf{v}_b\|_2 &\leq \frac{1}{2} \min_{a \neq a^*} (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b = \frac{g_b}{2}, \quad \text{and} \\ \sqrt{C} \|\mathbf{v}_b\|_2 \aleph_n &\leq \frac{g_b/2}{8\alpha_* \sqrt{C} u_{\max} + g_b/(2\|\mathbf{v}_b\|_2)} \\ \Leftrightarrow \aleph_n &\leq \frac{g_b}{16\alpha_* C u_{\max} \|\mathbf{v}_b\|_2 + g_b \sqrt{C}}. \end{aligned}$$

■

Lemma 8 (Control of the distortion ρ due to noisy feature estimates) *If n is large enough so that (1), (12) and (15) hold, then $\rho' \leq 2$ with probability at least $1 - \delta$.*

Proof [Proof of Lemma 8] We begin by considering

$$\max_{a \neq a^*} \frac{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b}{(\bar{\mathbf{u}}_{n,a^*} - \bar{\mathbf{u}}_{n,a})^\top \mathbf{v}_b} \leq \max_{a \neq a^*} \frac{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b}{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b + \partial_a \mathbf{v}_b} \leq \max_{a \neq a^*} \frac{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b}{\zeta_a},$$

with $\partial_a^\top := (\bar{\mathbf{u}}_{n,a^*} - \mathbf{u}_{a^*}) - (\bar{\mathbf{u}}_{n,a} - \mathbf{u}_a)$, and

$$\zeta_a := \inf_{\|\xi\|_2 \leq \|\partial_a\|_2} (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b + \xi^\top \mathbf{v}_b$$

as in the proof of Lemma 7. Also, by (14), we have that with probability at least $1 - \delta$,

$$\begin{aligned} \max_{a \in [A]} \|\bar{\mathbf{u}}_{n,a} - \mathbf{u}_a\|_2 &\leq \sqrt{AC} \max_{c \in [C]} \|\bar{\mathbf{u}}_{n,c} - \mathbf{u}_c\|_2 \\ &\leq \sqrt{AC} \diamond A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}} \\ &=: \aleph_n \sqrt{AC}, \quad \text{say.} \end{aligned}$$

Thus,

$$\begin{aligned}\zeta_a &\geq \inf_{\|\xi\|_2 \leq 2\aleph_n \sqrt{AC}} (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b + \xi^\top \mathbf{v}_b \\ &= (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b - 2\aleph_n \sqrt{AC} \|\mathbf{v}_b\|_2 \\ \Rightarrow \frac{\zeta_a}{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b} &\geq 1 - \frac{2\aleph_n \sqrt{AC} \|\mathbf{v}_b\|_2}{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b} \geq 1 - \frac{2\aleph_n \sqrt{AC} \|\mathbf{v}_b\|_2}{g_b},\end{aligned}$$

where $g_b := \min_{a \neq a^*} (\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b > 0$ is the minimum gap for user b across suboptimal actions.

Provided that (1), (12) and (15) hold, we get that with probability at least $1 - \delta$, $\frac{\zeta_a}{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b} \geq \frac{1}{2}$ for each $a \neq a^*$. Also, by the definition of a^* , the denominator is positive, i.e., $(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b > 0$. Hence,

$$\max_{a \neq a^*} \frac{(\mathbf{u}_{a^*} - \mathbf{u}_a)^\top \mathbf{v}_b}{(\bar{\mathbf{u}}_{n,a^*} - \bar{\mathbf{u}}_{n,a})^\top \mathbf{v}_b} \leq 2,$$

completing the proof of the result. \blacksquare

Lemma 9 (Bounding $R_{\mathcal{X}}$) *If n is large enough so that (1) and (12) hold, then*

$$R_{\mathcal{X}} \leq \frac{\sqrt{A}}{2\alpha_*} + \max_{a \in \mathcal{A}} \|\mathbf{u}_a\|_2,$$

with probability at least $1 - \delta$.

Proof [Proof of Lemma 9] Conditions (1) and (12), together with the estimate (20), imply that for any action a ,

$$\|\bar{\mathbf{u}}_{n,a}\|_2 \leq \|\mathbf{u}_a\|_2 + \|\bar{\mathbf{u}}_{n,a} - \mathbf{u}_a\|_2 \leq \|\mathbf{u}_a\|_2 + \aleph_n \sqrt{AC} \leq \|\mathbf{u}_a\|_2 + \sqrt{A}/(2\alpha_*).$$

with probability at least $1 - \delta$. \blacksquare

In order to conclude the proof of Lemma 3, we gather the conditions from Lemma 6 and Lemma 7. After some simplifications, both conditions are satisfied as soon as

$$\frac{n^2}{\sum_{i=1}^n \gamma_i^{-2}} \geq \diamond^2 A^6 C^2 \log(4A^3/\delta) \max \left\{ 2\alpha_*^2, \frac{8A \|\mathbf{v}_b\|_2^2}{g_b^2}, \frac{2^7 \alpha_*^2 C u_{\max}^2 \|\mathbf{v}_b\|_2^2}{g_b^2} + 1/2 \right\}.$$

Appendix D. Proof of Theorem 4

Proof Let n_0 be the first mini-session such that both conditions in Lemma 3 are satisfied, that is such that

$$\frac{n_0}{\sum_{i=1}^{n_0} \gamma_i^{-2}} \geq \diamond_\delta.$$

The cumulative regret $\mathfrak{R}_T = \sum_{t=1}^T r_t$ of Algorithm 3 satisfies

$$\begin{aligned}\mathfrak{R}_T &= \sum_{n=1}^N \sum_{l=1}^{\ell} r_{n,l} \\ &\leq (n_0 - 1)\ell + \sum_{b \in [B]} \sum_{n=n_0}^N \sum_{l=1}^{\ell} r_{n,l} \mathbb{I}\{b_n = b\}\end{aligned}$$

where $r_t \equiv r_{n,l} \stackrel{\text{def}}{=} \mathbf{u}_{a_{b_n}^*}^\top \mathbf{v}_{b_n} - \mathbf{u}_{a_{n,l}}^\top \mathbf{v}_{b_n}$ is the instantaneous regret of Algorithm 3 at time $t = \ell n + k$ when the current user is $b_n = b$. Using the notations of Algorithm 3, it holds that

$$\begin{aligned} \mathbb{E}[r_t | b_n = b] &= \mathbb{E}[r_t \mathbb{I}\{p_n = 1\} | b_n = b] + \mathbb{E}[r_t \mathbb{I}\{p_n = 0\} | b_n = b] \\ &\leq \mathbb{E}[\mathbf{u}_{a_b^*}^\top \mathbf{v}_b - \mathbf{u}_{\tilde{a}_{n,k}}^\top \mathbf{v}_b] (1 - \gamma_n) + \gamma_n \\ &\leq \mathbb{E}[\mathbf{u}_{a_b^*}^\top \mathbf{v}_b - \mathbf{u}_{\tilde{a}_{n,k}}^\top \mathbf{u}_b] + \gamma_n, \end{aligned}$$

where $\tilde{a}_{n,k}$ is an action output by an instance of **OFUL** for user $b_n = b$. Thus, we have

$$\begin{aligned} &\mathbb{E}[\mathfrak{R}_T | b_1, \dots, b_N] \\ &\leq (n_0 - 1)\ell + \mathbb{E} \left[\sum_{b \in [B]} \sum_{n=n_0}^N \sum_{l=1}^{\ell} \left(\mathbf{u}_{a_b^*}^\top \mathbf{v}_b - \mathbf{u}_{\tilde{a}_{n,l}}^\top \mathbf{v}_b \right) \mathbb{I}\{b_n = b\} \middle| b_1, \dots, b_N \right] + \ell \sum_{n=n_0}^N \gamma_n \\ &= (n_0 - 1)\ell + \underbrace{\sum_{b \in [B]} \mathbb{E} \left[\sum_{\substack{n_0 \leq n \leq N, \\ b_n = b}} \sum_{l=1}^{\ell} \left(\mathbf{u}_{a_b^*}^\top \mathbf{v}_b - \mathbf{u}_{\tilde{a}_{n,l}}^\top \mathbf{v}_b \right) \middle| b_1, \dots, b_N \right]}_{(*)} + \ell \sum_{n=n_0}^N \gamma_n. \end{aligned} \quad (21)$$

For each user $b \in [B]$, the expectation in the right-hand side above corresponds to the cumulative regret of the **OFUL** strategy when interacting with user b in mini-sessions n_0 through N , and when given at each mini-session n the set of perturbed feature vectors $\bar{\mathbf{U}}_n$. Let $N_{b,n_0} = \sum_{n=n_0}^N \mathbb{I}\{b_n = b\}$ count the total number of mini-sessions from n_0 in which user b is present (note that $\sum_{b \in [B]} N_{b,1} = N$ and $\sum_{b \in [B]} \ell N_{b,1} = T$). Let us denote the term $(*)$ in the above explicitly using $\mathfrak{R}_{b, N_{b,n_0}}(\{\bar{\mathbf{U}}_n\}_{n \in [n_0, N], b_n = b})$.

We can now use the **OFUL** robustness guarantee – a natural technical extension⁸ of Theorem 3 along with Lemma 3 – to obtain that, for a given user sequence b_1, \dots, b_N , with probability at least⁹ $1 - 2\delta - \delta = 1 - 3\delta$,

$$\mathfrak{R}_{b, N_{b,n_0}}(\{\bar{\mathbf{U}}_n\}_{n \in [n_0, N], b_n = b}) \leq 16 \sqrt{\ell N_{b,n_0} C \log \left(1 + \frac{\ell N_{b,n_0} R_{\mathcal{X}}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{\ell N_{b,n_0} R_{\mathcal{X}}^2}{\lambda C} \right)} \right).$$

8. Although Theorem 3 holds only for a fixed perturbation ε and feature set $\bar{\mathbf{u}}$, it is not hard to see that a modification of it, with time-varying ε_t , $\bar{\mathbf{u}}_t$ and ρ' being the largest ρ'_t over all times t , yields the same conclusion (regret bound). We provide this extension in Theorem 5 in Appendix E below.

9. Although the time horizons played by each **OFUL** instance per user, N_{b,n_0} , are technically random and unknown to the instance at the start, conditioning on the sequence of users arriving at each time instant lets us use the conclusion of Lemma 3.

This in turn implies that

$$\begin{aligned}
& \sum_{b \in B} \mathbb{E} \left[\mathfrak{R}_{b, N_{b, n_0}}(\{\bar{U}_n\}_{n \in [n_0, N], b_n = b}) \middle| b_1, \dots, b_N \right] \\
& \stackrel{(a)}{\leq} 16 \sum_{b \in B} \sqrt{\ell N_{b, n_0} C \log \left(1 + \frac{\ell N_{b, n_0} R_{\mathcal{X}}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{\ell N_{b, n_0} R_{\mathcal{X}}^2}{\lambda C} \right)} \right) \\
& \quad + \sum_{b \in B} 3\delta \ell N_{b, n_0} \\
& \stackrel{(b)}{\leq} 16 \sum_{b \in B} \sqrt{\ell N_{b, n_0} C \log \left(1 + \frac{\ell N_{b, n_0} R_{\mathcal{X}}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{\ell N_{b, n_0} R_{\mathcal{X}}^2}{\lambda C} \right)} \right) \\
& \quad + 3\delta T.
\end{aligned}$$

The last term on the right-hand side in (a) is due to the fact that with probability at most 3δ , the per-user regret $\mathfrak{R}_{b, N_{b, n_0}}(\{\bar{U}_n\}_{n \in [n_0, N], b_n = b})$ can be as large as $\ell N_{b, n_0}$ (the total number of time slots for which user b interacts with the system). The corresponding term in (b) is by using $\sum_{b \in [B]} \ell N_{b, 1} = T$. Further bounding using the Cauchy-Schwarz inequality $\sum_{b \in B} \sqrt{\ell N_{b, n_0}} \leq \sqrt{BT}$ gives

$$\begin{aligned}
& \sum_{b \in B} \mathbb{E} \left[\mathfrak{R}_{b, N_{b, n_0}}(\{\bar{U}_n\}_{n \in [n_0, N], b_n = b}) \middle| b_1, \dots, b_N \right] \\
& \leq 16 \sum_{b \in B} \sqrt{\ell N_{b, n_0} C \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \right) + 3\delta T \\
& \leq 16 \sqrt{BTC \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \right) + 3\delta T.
\end{aligned}$$

Plugging this estimate into (21), we obtain that

$$\begin{aligned}
\mathbb{E}[R_T | b_1, \dots, b_N] & \leq \ell \left(n_0 - 1 + \sum_{n=n_0}^N \gamma_n \right) \\
& \quad + 16 \sqrt{BTC \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \right) + 3\delta T.
\end{aligned}$$

Expliciting n_0 and tuning γ_n The next step is to control the term $n_0 - 1 + \sum_{n=n_0}^N \gamma_n$. To this end, we explicit n_0 and optimize γ_n . We write $\circ \equiv \circ_{\delta}$ in the sequel for convenience.

If $\gamma_n = \min\{1, \circ^{1/2} n^{-1/2}\}$, then

$$\begin{aligned}
\frac{n^2}{\sum_{m=1}^n \gamma_m^{-2}} & = \frac{n^2}{\lceil \circ \rceil + \frac{1}{\circ} \sum_{m > \lceil \circ \rceil}^n m} \\
& = \frac{2\circ n^2}{\lceil \circ \rceil 2\circ + n(n+1) - \lceil \circ \rceil (\lceil \circ \rceil - 1)} \\
& \geq \frac{2\circ}{1 + 1/n + (\lceil \circ \rceil)/n^2}.
\end{aligned}$$

Thus, this is higher than \circ if $n^2 - n - \lceil \circ \rceil \circ \geq 0$, that is if $n \geq n_0 \stackrel{\text{def}}{=} \lceil 1/2 + \sqrt{\lceil \circ \rceil \circ + 1/4} \rceil$. Since $n_0 \geq \lceil \circ \rceil$, we immediately get

$$\begin{aligned}\sum_{n=n_0}^N \gamma_n &\leq \circ^{1/2} n_0^{-\frac{1}{2}} + 2\circ^{1/2} \left(N^{\frac{1}{2}} - n_0^{\frac{1}{2}} \right) \\ &\leq 1 + 2\circ^{1/2} \left(N^{\frac{1}{2}} - n_0^{\frac{1}{2}} \right).\end{aligned}$$

Thus, we obtain

$$\begin{aligned}n_0 - 1 + \sum_{n=n_0}^N \gamma_n &\leq 2\circ^{1/2} N^{\frac{1}{2}} + n_0 - 2\circ^{1/2} n_0^{1/2} \\ &\leq 2\circ^{1/2} N^{\frac{1}{2}} + n_0 - 2\sqrt{\circ \lceil \circ \rceil}\end{aligned}$$

Using the fact that $\circ > 1$, the bound simplifies to

$$n_0 - 1 + \sum_{n=n_0}^N \gamma_n \leq 2\sqrt{\circ N} + 1.$$

If, on the other hand, a bound on \circ is not readily available beforehand, then choosing $\gamma_n = \sqrt{\log(1+n)/n}$, $n \geq 1$, gives, via a crude bound,

$$\begin{aligned}\sum_{m=1}^n \gamma_m^{-2} &= \sum_{m=1}^n m / \log(1+m) \leq \sum_{m=1}^{\sqrt{n}} m / \log 2 + \sum_{m=\sqrt{n}}^n m / \log(1+\sqrt{n}) \\ &\leq n / \log 2 + n^2 / \log \sqrt{n} \leq 2n^2 / \log \sqrt{n} \\ \Rightarrow \frac{n_0^2}{\sum_{m=1}^{n_0} \gamma_m^{-2}} &\geq \frac{n_0^2}{2n_0^2 / \log \sqrt{n_0}} = \frac{\log n_0}{4}.\end{aligned}$$

The bound above is at least \circ provided $n_0 \geq \exp(4\circ)$. Thus, we finally get that, upon setting $\delta = 1/\sqrt{T}$, the total expected regret satisfies (as an order-wise function of T)

$$\begin{aligned}\mathbb{E}[R_T] &\leq \ell \left(\exp(4\circ) + \sum_{n=1}^N \sqrt{\log(n+1)/n} \right) \\ &\quad + 16\sqrt{BTC \log \left(1 + \frac{TR\chi^2}{\lambda C} \right)} \left(\lambda^{1/2} R_\Theta + R\sqrt{\log T + C \log \left(1 + \frac{TR\chi^2}{\lambda C} \right)} \right) + 3\sqrt{T} \\ &= O \left(C\sqrt{BT} \log T \right).\end{aligned}$$

■

Appendix E. Extension of Theorem 3: Robustness of OFUL's regret with time-varying features

We now control the robust regret for user b $\mathfrak{R}_{b, N_b, n_0}(\{\bar{U}_n\}_{n \in [n_0, N], b_n = b}) = \sum_{\substack{n_0 \leq n \leq N, \\ b_n = b}} \sum_{l=1}^{\ell} \left(\mathbf{u}_{a_b^*}^\top \mathbf{v}_b - \mathbf{u}_{\bar{a}_{n,l}}^\top \mathbf{v}_b \right)$, when OFUL is run with evolving feature matrices $\{\bar{U}_n\}_{n \in [n_0, N], b_n = b}$ with decreasing feature error $\varepsilon_n = (U - \bar{U}_n) \mathbf{v}_b$, instead of a fixed \bar{U} with fixed error $\varepsilon = (U - \bar{U}) \mathbf{v}_b$.

We reindex the $n \in [n_0, N]$, $b_n = b$ as $t = 1, \dots, \dots$ and prove the following result.

Theorem 5 (OFUL robustness result, extension of Theorem 3 for time-varying features) Assume $\|\mathbf{v}^\circ\|_2 \leq R_\Theta$, $\lambda \geq \max\{1, R_{\mathcal{X}}^2, 1/4R_\Theta^2\}$, $\forall a \in \mathcal{A}$, $t \leq T$, $\|\bar{\mathbf{u}}_a^{(t)}\|_2 \leq R_{\mathcal{X}}$ and $|m_a| \leq 1$, and that for all $t \leq T$, $\arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}^\circ = \{a^*\}$ (i.e., the linearly realizable approximation with respect to the current features has a^* as its unique optimal action). If

$$\|\varepsilon^{(t)}\|_2 \equiv \|\mathbf{m} - \bar{U}^{(t)} \mathbf{v}^\circ\|_2 < \min_{a \neq a^*} \frac{\bar{\mathbf{u}}_{a^*}^{(t)\top} \mathbf{v}^\circ - \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}^\circ}{2\alpha(\bar{U}^{(t)\top}) \|\bar{\mathbf{u}}_{a^*}^{(t)} - \bar{\mathbf{u}}_a^{(t)}\|_2}, \quad (22)$$

then with probability at least $1 - \delta$, for all $T \geq 0$,

$$\mathfrak{R}_T \leq 8\rho' \sqrt{TC \log\left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C}\right)} \left(\lambda^{1/2} R_\Theta + R \sqrt{2 \log \frac{1}{\delta} + C \log\left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C}\right)} \right),$$

where $\rho' := \max_t \max\left\{1, \max_{a \neq a^*} \frac{m_{a^*} - m_a}{\bar{\mathbf{u}}_{a^*}^{(t)\top} \mathbf{v}^\circ - \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}^\circ}\right\}$.

Proof Let $\mathbf{M}_{1:t} = (\mathbf{m}_{A_1}, \dots, \mathbf{m}_{A_t})^\top$. The argument used to prove Theorem 2 in Yadkori et al, 2011, shows that

$$\hat{\mathbf{v}}_{t-1} = V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \boldsymbol{\eta}_{1:t-1} + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \mathbf{M}_{1:t-1}$$

where $\boldsymbol{\eta}_{1:t-1} := (\eta_1, \dots, \eta_{t-1})$ is the observed noise sequence, and where $\bar{\mathbf{U}}_{1:t-1}^{(t)}$ is the matrix built from the time varying features at time t and the action sequence thus far. Let $\mathbf{E}_{1:t-1}^{(t)} := (\varepsilon_{A_1}^{(t)}, \dots, \varepsilon_{A_t}^{(t)})^\top = \mathbf{M}_{1:t-1} - \bar{\mathbf{U}}_{1:t-1}^{(t)} \mathbf{v}^\circ$. We then have

$$\begin{aligned} \hat{\mathbf{v}}_{t-1} &= V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \boldsymbol{\eta}_{1:t-1} + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \mathbf{M}_{1:t-1} \\ &= V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \boldsymbol{\eta}_{1:t-1} + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \left(\bar{\mathbf{U}}_{1:t-1}^{(t)\top} \mathbf{v}^\circ + \mathbf{E}_{1:t-1}^{(t)} \right) \\ &= V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \boldsymbol{\eta}_{1:t-1} + \mathbf{v}^\circ - \lambda V_{t-1}^{-1} \mathbf{v}^\circ + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \mathbf{E}_{1:t-1}^{(t)}. \end{aligned}$$

Thus, letting $\mathbf{v}_{t-1}^+ := \mathbf{v}^\circ + V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \mathbf{E}_{1:t-1}^{(t)}$, and using the above with techniques from Yadkori et al together with $\|\mathbf{v}^\circ\|_2 \leq R_\Theta$, we have that

$$\mathbf{v}_{t-1}^+ \in \mathcal{C}_{t-1}$$

with probability at least $1 - \delta$.

Now, let $a_{t-1}^+ \in \arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}_{t-1}^+$ be an optimal action corresponding to the approximate parameter \mathbf{v}_{t-1}^+ and approximate feature $\bar{\mathbf{u}}_{a_{t-1}^+}^{(t)\top}$, and define the instantaneous regret at time t with respect to the approximate parameter as

$$r_t^+ := \bar{\mathbf{u}}_{a_{t-1}^+}^{(t)\top} \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^{(t)\top} \mathbf{v}_{t-1}^+ \geq 0.$$

We now bound this approximate regret using arguments along the lines of Yadkori et al, 2011 as follows. Write

$$\begin{aligned}
r_t^+ &= \bar{\mathbf{u}}_{a_{t-1}^+}^{(t)\top} \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^{(t)\top} \mathbf{v}_{t-1}^+ \\
&\leq \bar{\mathbf{u}}_{A_t}^{(t)\top} \tilde{\mathbf{v}}_t - \bar{\mathbf{u}}_{A_t}^{(t)\top} \mathbf{v}_{t-1}^+ \quad (\text{since } (A_t, \tilde{\mathbf{v}}_t) \text{ is optimistic}) \\
&= \bar{\mathbf{u}}_{A_t}^{(t)\top} (\tilde{\mathbf{v}}_t - \mathbf{v}_{t-1}^+) \\
&= \bar{\mathbf{u}}_{A_t}^{(t)\top} (\tilde{\mathbf{v}}_t - \hat{\mathbf{v}}_{t-1}) + \bar{\mathbf{u}}_{A_t}^{(t)\top} (\hat{\mathbf{v}}_{t-1} - \mathbf{v}_{t-1}^+) \\
&\leq \left\| \bar{\mathbf{u}}_{A_t}^{(t)} \right\|_{V_{t-1}^{-1}} \|\tilde{\mathbf{v}}_t - \hat{\mathbf{v}}_{t-1}\|_{V_{t-1}} + \left\| \bar{\mathbf{u}}_{A_t}^{(t)} \right\|_{V_{t-1}^{-1}} \|\hat{\mathbf{v}}_{t-1} - \mathbf{v}_{t-1}^+\|_{V_{t-1}} \quad (\text{Cauchy-Schwarz's inequality}) \\
&\leq 2D_{t-1} \left\| \bar{\mathbf{u}}_{A_t}^{(t)} \right\|_{V_{t-1}^{-1}}. \tag{23}
\end{aligned}$$

Noting that $m_a \in [-1, 1] \forall a$, the regret can be written as

$$\begin{aligned}
R_T &= \sum_{t=1}^T (m_{a^*} - m_{A_t}) = \sum_{t=1}^T \min\{m_{a^*} - m_{A_t}, 2\} \\
&= \rho' \sum_{a \neq a^*} \sum_{t=1}^T \min \left\{ \frac{m_{a^*} - m_a}{\rho'}, \frac{2}{\rho'} \right\} \mathbb{I}\{A_t = a\} \\
&\leq \rho' \sum_{a \neq a^*} \sum_{t=1}^T \min \left\{ \bar{\mathbf{u}}_{a^*}^{(t)\top} \mathbf{v}^\circ - \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}^\circ, \frac{2}{\rho'} \right\} \mathbb{I}\{A_t = a\} \quad (\text{using the definition of } \rho') \\
&\stackrel{(a)}{\leq} \rho' \sum_{t=1}^T \min \left\{ 2 \left(\bar{\mathbf{u}}_{a^*}^{(t)\top} \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^{(t)\top} \mathbf{v}_{t-1}^+ \right), \frac{2}{\rho'} \right\} \stackrel{(b)}{=} 2\rho' \sum_{t=1}^T \min \left\{ \bar{\mathbf{u}}_{a_{t-1}^+}^{(t)\top} \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_{A_t}^{(t)\top} \mathbf{v}_{t-1}^+, \frac{1}{\rho'} \right\} \\
&= 2\rho' \sum_{t=1}^T \min \left\{ r_t^+, \frac{1}{\rho'} \right\} = \rho' \sum_{t=1}^T \frac{2}{\rho'} \min \{ \rho' r_t^+, 1 \} \stackrel{(c)}{\leq} \rho' \sum_{t=1}^T \frac{2}{\rho'} \min \left\{ 2\rho' D_{t-1} \left\| \bar{\mathbf{u}}_{A_t}^{(t)} \right\|_{V_{t-1}^{-1}}, 1 \right\} \\
&\stackrel{(d)}{\leq} \rho' \sum_{t=1}^T 4D_{t-1} \min \left\{ \left\| \bar{\mathbf{u}}_{A_t}^{(t)} \right\|_{V_{t-1}^{-1}}, 1 \right\} \\
&\leq \rho' \sqrt{T \sum_{t=1}^T 16D_T^2 \min \left\{ \left\| \bar{\mathbf{u}}_{A_t}^{(t)} \right\|_{V_{t-1}^{-1}}^2, 1 \right\}} \quad (\text{by using Cauchy-Schwarz's inequality}).
\end{aligned}$$

In the derivation above,

- Steps (a) and (b) hold because of the following. By Lemma 10 (to follow below), $\|\mathbf{v}_{t-1}^+ - \mathbf{v}^\circ\|_2 = \left\| V_{t-1}^{-1} \bar{\mathbf{U}}_{1:t-1}^{(t)} \mathbf{E}_{1:t-1}^{(t)} \right\|_2 \leq \alpha(\bar{U}_t) \|\varepsilon^{(t)}\|_2$. Since $\arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}^\circ$ is uniquely a^* by hypothesis, we have, thanks to Lemma 5, that $\bar{\mathbf{u}}_{a^*}^{(t)\top} \mathbf{v}_{t-1}^+ - \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}_{t-1}^+ > \frac{\bar{\mathbf{u}}_{a^*}^{(t)\top} \mathbf{v}^\circ - \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}^\circ}{2} > 0 \forall a \neq a^*$, establishing (a). This in turn shows that the optimal action for \mathbf{v}_{t-1}^+ is uniquely a^* at all times t , i.e., $a_{t-1}^+ = \arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}_{t-1}^+ = a^*$, which is precisely equality (b).
- *Remark.* In the above, Lemma 5 is written for generic $\bar{\mathbf{u}}_a$, ε , so in particular applies to each time varying $\bar{\mathbf{u}}_a^{(t)}$, $\varepsilon^{(t)}$. We also used an extended version of Lemma 4 to the case of varying $\bar{\mathbf{u}}_a^{(t)}$, $\varepsilon^{(t)}$, which we state and prove below as Lemma 10.
- Inequality (c) holds by (23) and (d) holds because $\rho' \geq 1$ by definition, and $D_{t-1} \geq \lambda^{1/2} R_\Theta \geq 1/2$ by hypothesis, implying that $2\rho' D_{t-1} \geq 1$.

The argument from here can be continued in the same way as in [Abbasi-Yadkori et al. \(2011, proof of Theorem 3\)](#) to yield

$$R_T \leq 8\rho' \sqrt{TC \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \left(\lambda^{1/2} R_{\Theta} + R \sqrt{2 \log \frac{1}{\delta} + C \log \left(1 + \frac{TR_{\mathcal{X}}^2}{\lambda C} \right)} \right).$$

This proves the theorem. ■

Lemma 10 (Extension of Lemma 4 to time-varying feature sets) *Let $\varepsilon_a^{(t)} = m_a - \bar{\mathbf{u}}_a^{(t)\top} \mathbf{v}^\circ$ be the bias in arm a 's reward due to model error, with respect to the features \bar{U}_t , and let $\varepsilon^{(t)} \equiv \left(\varepsilon_a^{(t)} \right)_{a \in \mathcal{A}}$. Then, we have*

$$\left\| V_{t-1}^{-1} \bar{U}_{1:t-1} \mathbf{E}_{1:t-1}^{(t)} \right\|_2 \leq \left(\max_J \left\| \mathbf{A}_J^{(t)-1} \right\|_2 \right) \left\| \varepsilon^{(t)} \right\|_2,$$

where $\mathbf{A}_{(A+C) \times C}^{(t)} = \begin{bmatrix} \bar{U}^{(t)} \\ I_d \end{bmatrix}$, $\mathbf{A}_J^{(t)}$ is the $C \times C$ submatrix of $\mathbf{A}^{(t)}$ consisting of rows in J , and J ranges over all subsets of full-rank rows of $\mathbf{A}^{(t)}$.

Proof [Proof of Lemma 10] Let $z_{t-1}^{(t)} := V_{t-1}^{-1} \bar{U}_{1:t-1} \mathbf{E}_{1:t-1}^{(t)} = \mathbf{v}_{t-1}^+ - \mathbf{v}^\circ \in \mathbb{R}^C$, thus $\left\| \mathbf{E}_{1:t-1}^{(t)} \right\|_\infty \leq \left\| \varepsilon^{(t)} \right\|_\infty = \left\| \mathbf{m} - \bar{U}^{(t)} \mathbf{v}^\circ \right\|_\infty$. We now write

$$\begin{aligned} z_{t-1}^{(t)} &= \left(\sum_{s=1}^{t-1} \bar{\mathbf{u}}_{A_s}^{(t)} \bar{\mathbf{u}}_{A_s}^{(t)\top} + \lambda I \right)^{-1} \sum_{s=1}^{t-1} \varepsilon_{A_s}^{(t)} \bar{\mathbf{u}}_{A_s}^{(t)} \\ &= \left(\frac{1}{t-1} \sum_{s=1}^{t-1} \bar{\mathbf{u}}_{A_s}^{(t)} \bar{\mathbf{u}}_{A_s}^{(t)\top} + \frac{\lambda}{t-1} I \right)^{-1} \frac{1}{t-1} \sum_{s=1}^{t-1} \varepsilon_{A_s}^{(t)} \bar{\mathbf{u}}_{A_s}^{(t)} \\ &= \left(\sum_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^{(t)} \bar{\mathbf{u}}_a^{(t)\top} \frac{\sum_{s=1}^{t-1} \mathbb{I}\{A_s = a\}}{t-1} + \frac{\lambda}{t-1} I \right)^{-1} \sum_{a \in \mathcal{A}} \varepsilon_a^{(t)} \bar{\mathbf{u}}_a^{(t)} \frac{\sum_{s=1}^{t-1} \mathbb{I}\{A_s = a\}}{t-1} \\ &= \left(\sum_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^{(t)} \bar{\mathbf{u}}_a^{(t)\top} f_a(t-1) + \frac{\lambda}{t-1} I \right)^{-1} \sum_{a \in \mathcal{A}} \varepsilon_a^{(t)} \bar{\mathbf{u}}_a^{(t)} f_a(t-1), \end{aligned}$$

where $f_a(t-1)$ is the empirical frequency with which action $a \in \mathcal{A}$ has been played up to and including time $t-1$. This allows us to equivalently interpret z_{t-1} as the solution of a *weighted* ℓ^2 -regularized least squares regression problem with $K = |\mathcal{A}|$ observations (instead of the original interpretation with $t-1$ observations) as follows (we suppress the dependence of f_a on t as per the context for clarity of notation).

Let $\mathbf{F}^{1/2}$ be the $A \times A$ diagonal matrix with the values $\sqrt{f_1}, \dots, \sqrt{f_A}$ on the diagonal (note: $\sum_{a=1}^A f_a = 1$). With this, we can express z_{t-1} as

$$\begin{aligned} z_{t-1}^{(t)} &= \arg \min_{z \in \mathbb{R}^C} \left\| \mathbf{F}^{1/2} \bar{U}^{(t)} z - \mathbf{F}^{1/2} \varepsilon^{(t)} \right\|_2^2 + \frac{\lambda}{t-1} \|z\|_2^2 \\ &= \arg \min_{z \in \mathbb{R}^C} \left\| \mathbf{F}^{1/2} \left(\bar{U}^{(t)} z - \varepsilon^{(t)} \right) \right\|_2^2 + \frac{\lambda}{t-1} \|z\|_2^2 \\ &= \arg \min_{z \in \mathbb{R}^C} \left\| \begin{bmatrix} \mathbf{F}^{1/2} & 0 \\ 0 & \sqrt{\frac{\lambda}{t-1}} I_C \end{bmatrix} \left(\begin{bmatrix} \bar{U}^{(t)} \\ I_C \end{bmatrix} z - \begin{bmatrix} \varepsilon^{(t)} \\ 0 \end{bmatrix} \right) \right\|_2^2 \\ &\equiv \arg \min_{z \in \mathbb{R}^C} \left\| \mathbf{D}^{1/2} (\mathbf{A}z - \mathbf{b}) \right\|_2^2 = (\mathbf{A}^\top \mathbf{D} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{D} \mathbf{b}, \end{aligned}$$

with $\mathbf{D}^{1/2}$ being a $(A + C) \times (A + C)$ diagonal & positive semidefinite matrix, $\mathbf{A}^\top \mathbf{D} \mathbf{A} = \sum_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^{(t)} \bar{\mathbf{u}}_a^{(t)\top} f_a(t-1) + \frac{\lambda}{t-1} I$ being positive definite, and \mathbf{A} having full column rank C . A result of [Forsgren \(1996, Corollary 2.3\)](#) now gives

$$\|(\mathbf{A}^\top \mathbf{D} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{D}\|_2 \leq \max_J \|\mathbf{A}_J^{-1}\|_2$$

where J ranges over all subsets of full-rank rows of \mathbf{A} , and \mathbf{A}_J is the $C \times C$ submatrix of \mathbf{A} formed by picking rows J . Thus, $\|z_{t-1}^{(t)}\|_2 \leq (\max_J \|\mathbf{A}_J^{-1}\|_2) \|\varepsilon^{(t)}\|_2$. This proves the lemma. \blacksquare

Appendix F. Unregularized Least squares

In our setting where we consider finitely many arms, one way wonder whether it is possible to remove the regularization parameter λ . Following [Rusmevichientong and Tsitsiklis \(2010\)](#), this is indeed possible under the assumption that the minimum eigenvalue of $\sum_{a \in \mathcal{A}} \mathbf{u}_a \mathbf{u}_a^\top$ is away from 0. Then, we first play each arm once (once for all users B , not for each of them) before running [Algorithm 3](#), where **OFUL** is used with $\lambda = 0$ and with D_{t-1} redefined to be $4R^2 \left(A \log(t) + \log(A/\delta) \right)$. This leads essentially to similar bounds, with α^* replaced by $\max_J \|\mathbf{U}_J^{-1}\|_2$, as we show below.

Let $\mathcal{U} \subset \mathbb{R}^C$. We receive at time s , observation $y_s = \mathbf{u}_s^\top \mathbf{v}^* + \eta_s \in \mathbb{R}$ where $\mathbf{v}^* \in \mathbb{R}^C$ and $\mathbf{u}_s \in \mathcal{U}$.

We make the following

Assumption 2 *There exists $R_{\mathcal{X}}, R, \lambda_0 \in \mathbb{R}_*^+$ such that*

1. $\forall s, \|\mathbf{u}_s\| \leq R_{\mathcal{X}}$
2. $\forall \lambda \in \mathbb{R}, \log \mathbb{E} \exp(\lambda \eta_s) \leq \lambda^2 R^2 / 2$.
3. $\lambda_{\min}(\sum_{s=1}^t \mathbf{u}_s \mathbf{u}_s^\top) \geq \lambda_0$.

[Assumption 2.3](#) is satisfied for instance when there are C points $(\mathbf{u}_{0,i})_{i \in [C]}$ in \mathbb{R}^d such that $\lambda_{\min}(\sum_{i=1}^C \mathbf{u}_{0,i} \mathbf{u}_{0,i}^\top) = \lambda_0 > 0$, and $\mathbf{u}_s = \mathbf{u}_{0,s}$ for $s \in [C]$. We consider the least-squares estimate

$$\mathbf{v}_t = \left(\sum_{s=1}^t \mathbf{u}_s \mathbf{u}_s^\top \right)^{-1} \sum_{s=1}^t \mathbf{u}_s y_s,$$

F.1 Preliminary

In case \mathcal{U} is finite, one can get the following result

Theorem 6 *Let us introduce the confidence set*

$$\mathcal{C}_t = \left\{ w \in \mathbb{R}^C : w^\top G_t w \leq D_{t,\delta} \right\}, \text{ where } G_t = \sum_{s=1}^t \mathbf{u}_s \mathbf{u}_s^\top$$

$$\text{and } D_{t,\delta} = 4R^2 \left(|\mathcal{U}| \log(t) + \log(|\mathcal{U}|/\delta) \right).$$

Then, under [Assumption 2](#), it holds

$$\mathbb{P} \left(\mathbf{v}_t - \mathbf{v}^* \in \mathcal{C}_t \right) \geq 1 - \delta.$$

In the general case, it holds

Theorem 7 *Let us introduce the confidence set*

$$\mathcal{C}_t = \left\{ w \in \mathbb{R}^C : w^\top G_t w \leq D_{t,\delta} \right\}, \text{ where } G_t = \sum_{s=1}^t \mathbf{u}_s \mathbf{u}_s^\top$$

$$\text{and } D_{t,\delta} = 16R^2 \left[1 + \log \left(1 + \frac{36R_\chi^2}{\lambda_0} \right) \right] \left[C \log \left(\frac{36R_\chi^2}{\lambda_0} t \right) + \log(1/\delta) \right] \log(t).$$

Then, under Assumption 2, and if $t \geq \frac{\lambda_0}{12R_\chi^2}$ it holds

$$\mathbb{P} \left(\mathbf{v}_t - \mathbf{v}^* \in \mathcal{C}_t \right) \geq 1 - \delta.$$

Proof: Indeed, let $z_t = \sum_{s=1}^t \mathbf{u}_s \eta_s$. Since G_t is invertible, it holds that $\mathbf{v}_t = \mathbf{v}^* + G_t^{-1} z_t$, and thus

$$(\mathbf{v}_t - \mathbf{v}^*)^\top G_t (\mathbf{v}_t - \mathbf{v}^*) = z_t G_t^{-1} z_t$$

In the case when \mathcal{U} is finite, using the Proof of Theorem B.1 in [Rusmevichientong and Tsitsiklis \(2010\)](#) then we further get for all $\varepsilon > 0$,

$$\mathbb{P} \left(z_t G_t^{-1} z_t \geq \varepsilon^2 R^2 \right) \leq |\mathcal{U}| t^{|\mathcal{U}|} e^{-\varepsilon^2/4},$$

Thus, choosing $\varepsilon = 2\sqrt{\log(|\mathcal{U}|t^{|\mathcal{U}|}/\delta)}$, we obtain that

$$\mathbb{P} \left(z_t G_t^{-1} z_t \geq 4R^2 \left(|\mathcal{U}| \log(t) + \log(|\mathcal{U}|/\delta) \right) \right) \leq \delta,$$

which concludes the proof of Theorem 6.

From the Proof of Theorem B.2 in [Rusmevichientong and Tsitsiklis \(2010\)](#), it holds that for all $\varepsilon > 2$,

$$\mathbb{P} \left(z_t G_t^{-1} z_t \geq \varepsilon^2 k_0^2 R^2 \log(t) \right) \leq (36R_\chi^2 t / \lambda_0)^C e^{-\varepsilon^2/4},$$

where $k_0 = 2\sqrt{1 + \log(1 + 36R_\chi^2/\lambda_0)}$, which leads to

$$\mathbb{P} \left(z_t G_t^{-1} z_t \geq 4(1 + \log(1 + 36R_\chi^2/\lambda_0)) R^2 \log(t) \varepsilon^2 \right) \leq (36R_\chi^2 t / \lambda_0)^C e^{-\varepsilon^2/4},$$

Thus, let us use $\varepsilon = 2\sqrt{\log((36R_\chi^2 t / \lambda_0)^C / \delta)}$, which satisfies $\varepsilon > 2$ as soon as $t > \frac{\lambda_0 e^{1/C}}{36R_\chi^2}$, thus in particular if $t \geq \frac{\lambda_0}{12R_\chi^2}$. Now, introducing the constant $c = 36R_\chi^2 / \lambda_0$, we obtain

$$\mathbb{P} \left(z_t G_t^{-1} z_t \geq 16R^2 (1 + \log(1 + c)) \log(t) \left(C \log(ct) + \log(1/\delta) \right) \right) \leq \delta,$$

which concludes the proof of theorem 7. \square

F.2 Application to Low-Rank bandits

In order to apply this result to the low-rank bandit problem, we need to show that G_t is invertible. In our case, this matrix is at mini-session n $\tilde{M}_t = \sum_{s=1}^t \tilde{\mathbf{u}}_{n,a_s} \tilde{\mathbf{u}}_{n,a_s}^\top$.

Let us assume that all actions are sample at least once in the beginning. Thus, in this case $\lambda_{\min}(\tilde{M}_t) \geq \lambda_{\min}(\tilde{A})$, where $\tilde{A} = \sum_{a \in [A]} \tilde{\mathbf{u}}_{n,a} \tilde{\mathbf{u}}_{n,a}^\top$. For convenience, let us also introduce the $C \times C$ matrix $A = \sum_{a \in [A]} \mathbf{u}_a \mathbf{u}_a^\top = U^\top U$.

In order to show that \tilde{M}_t is invertible, it is enough to show that $\lambda_{\min}(\tilde{A}) > 0$.

Now, by the result of reconstruction of the feature matrix M , we know that there exists with high probability a permutation π such that the columns are well estimated:

$$\forall c, \|\mathbf{u}_{\pi(c)} - \tilde{\mathbf{u}}_{n,c}\| \leq \diamond A^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}}.$$

Thus, we study $E = \tilde{A} - A$. Let λ be any eigenvalue of E , then it holds

$$\begin{aligned} \lambda &\leq \text{trace}(E) = \sum_{a \in [A]} \text{trace}\left(\tilde{\mathbf{u}}_{n,a} \tilde{\mathbf{u}}_{n,a}^\top - \mathbf{u}_a \mathbf{u}_a^\top\right) \\ &\leq \sum_{a \in [A]} \|\tilde{\mathbf{u}}_{n,a}\|^2 - \|\mathbf{u}_a\|^2 \\ &\leq \sum_{a \in [A]} \sum_{c \in [C]} \tilde{u}_{n,a,c}^2 - u_{a,c}^2 \\ &\leq \sum_{a \in [A]} \sum_{c \in [C]} (\tilde{u}_{n,a,c} - u_{a,c})^2 + 2u_{a,c}(\tilde{u}_{n,a,c} - u_{a,c}) \\ &\leq \sum_{c \in [C]} \|\tilde{\mathbf{u}}_{n,c} - \mathbf{u}_c\|^2 + 2 \sum_{c \in [C]} \sqrt{\sum_{a \in [A]} u_{a,c}^2} \sqrt{\sum_{a \in [A]} (\tilde{u}_{n,a,c} - u_{a,c})^2} \\ &\leq \sum_{c \in [C]} \|\tilde{\mathbf{u}}_{n,c} - \mathbf{u}_c\|^2 + 2\|\mathbf{u}_c\| \|\tilde{\mathbf{u}}_{n,c} - \mathbf{u}_c\| \\ &\leq (2u_{max} + 1) \sum_{c \in [C]} \|\tilde{\mathbf{u}}_{n,c} - \mathbf{u}_c\|. \end{aligned}$$

Thus, provided that n is large enough that

$$\lambda_{\min}(A) > 2(2u_{max} + 1) \sum_{c \in [C]} \|\tilde{\mathbf{u}}_{n,c} - \mathbf{u}_c\|,$$

we deduce that \tilde{M}_t is invertible. Using the fact that $A = U^\top U$, This translates to the condition

$$\lambda_{\min}(U^\top U) > 2\diamond(2u_{max} + 1)CA^3 \sqrt{\sum_{i=1}^n \gamma_i^{-2} \frac{C \log(4A^3/\delta)}{2n^2}}$$

that is

$$\frac{n^2}{\sum_{i=1}^n \gamma_i^{-2}} > \frac{4\diamond^2(2u_{max} + 1)^2 C^3 A^6 \log(4A^3/\delta)}{\lambda_{\min}^2(U^\top U)}.$$

Thus, assuming that all actions are chosen at least once in the beginning, and that

$$\frac{n^2}{\sum_{m=1}^n \gamma_m^{-2}} > \frac{4\diamond^2(2u_{max} + 1)^2 C^3 A^6 \log(4A^3/\delta)}{\lambda_{\min}^2(U^\top U)},$$

then $\lambda_{\min}(\tilde{M}_t) \geq \lambda_{\min}(U^\top U)/2 = \lambda_0/2 > 0$ and Theorem 6 and Theorem 7 both apply.

In order to control the regret of the unregularized version of **OFUL**, we now use the proof of [Rusmevichientong and Tsitsiklis \(2010, Theorem 4.1\)](#) combined with the fact that $\lambda_{\min}(\tilde{M}_t) \geq \lambda_0/2$ to get

$$\sum_{t=A+1}^n \min\{\|\bar{\mathbf{u}}_{A_t}\|_{\tilde{M}_{t-1}^{-1}}^2, 1\} \leq 2 \max\{1, \frac{2R_{\mathcal{X}}^2}{\lambda_0}\} \left(C \log(\max\{1, \frac{2R_{\mathcal{X}}^2}{\lambda_0}\}) + (C+1) \log(n+1) \right).$$

A straightforward adaptation of the proof of Theorem 3 then gives

$$\begin{aligned} \mathfrak{R}_n &\leq \rho' \sqrt{n(A + 16D_{n,\delta}^2 \sum_{t=A+1}^n \min\{\|\bar{\mathbf{u}}_{A_t}\|_{\tilde{M}_{t-1}^{-1}}^2, 1\})} \\ &\leq 16\rho' R^2 \left(A \log(n) + \log(A/\delta) \right) \sqrt{n \left(2 + \frac{4R_{\mathcal{X}}^2}{\lambda_0} \right) \left(C \log \left(1 + \frac{2R_{\mathcal{X}}^2}{\lambda_0} \right) + (C+1) \log(n+1) \right)} \\ &\quad + \rho' \sqrt{An}. \end{aligned}$$

Following the same steps as for Lemma 3, we finally obtain the result:

Theorem 8 (Unregularized OFUL robustness result) *Assume $\|\mathbf{v}^\circ\|_2 \leq R_\Theta$, for all $a \in \mathcal{A}$, $\|\bar{\mathbf{u}}_a\|_2 \leq R_{\mathcal{X}}$ and $|m_a| \leq 1$, and that $\arg \max_{a \in \mathcal{A}} \bar{\mathbf{u}}_a^\top \mathbf{v}^\circ = \{a^*\}$ (i.e., the linearly realizable approximation has a^* as its unique optimal action). Assume that each action has been played once. Let $0 < \delta \leq 1$. Provided that the number of mini-sessions n_0 is large enough to satisfy*

$$\frac{n_0^2}{\sum_{i=1}^{n_0} \gamma_i^{-2}} \geq \tilde{\mathcal{O}}_{b,\delta}$$

where

$$\begin{aligned} \tilde{\mathcal{O}}_{b,\delta} &= \max \left\{ \frac{2A^6 \log(4A^2/\delta)}{\min\{\Gamma, \sigma_{\min}\}^2}, \frac{A^9 (1 + 10(\frac{1}{\Gamma} + \frac{1}{\sigma_{\min}}))(1 + u_{\max}^3)^2 C^5 \log(4A^3/\delta)}{2C_1^2 \sigma_{\min}^3} \right. \\ &\quad \left. \frac{4\Diamond^2 (2u_{\max} + 1)^2 C^3 A^6 \log(4A^3/\delta)}{\lambda_{\min}^2(U^\top U)}, \right. \\ &\quad \left. \Diamond^2 A^6 C^2 \log(4A^3/\delta) \max \left\{ 2\bar{\alpha}_*^2, \frac{8A \|\mathbf{v}_b\|_2^2}{g_b^2}, \frac{2^7 \bar{\alpha}_*^2 C u_{\max}^2 \|\mathbf{v}_b\|_2^2}{g_b^2} + 1/2 \right\} \right\}, \end{aligned}$$

then with probability at least $1 - \delta$ for all $T \geq 0$, the regret $\mathfrak{R}_{A+1:n}$ of the **OFUL** algorithm from decision $A+1$ to n satisfies

$$\mathfrak{R}_{A+1:n} \leq 32R^2 \left[A \log(n) + \log(A/\delta) \right] \sqrt{n \left(2 + \frac{4\bar{R}_{\mathcal{X}}^2}{\lambda_0} \right) \left(C \log \left(1 + \frac{2\bar{R}_{\mathcal{X}}^2}{\lambda_0} \right) + (C+1) \log(n+1) \right)},$$

where we introduced

$$\bar{R}_{\mathcal{X}} = \max_{a \in \mathcal{A}} \|\mathbf{u}_a\|_2 + \frac{\sqrt{A}}{2\bar{\alpha}_*} \quad \text{and} \quad \bar{\alpha}_* = \min_j \|U_j^{-1}\|.$$

This result enables to get the corresponding variant of Theorem 4 using an unregularized **OFUL**.