
Stochastic Low-Rank Latent Bandits

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To be written.

2 1 Introduction

3 STORY: We address a recommendation problem in the hard setting where no feature is available to
4 the learner. Blah blah: recommendation and bandits, major problem, blah blah.

5 We rely on the assumption that the underlying click-through rate matrix has a latent sructure that we
6 cannot directly observe but that we propose to leverage nonetheless. We formulate a rank- d bandit
7 problem that generalizes previous works on rank-1 and on latent bandits (quote, quote). We propose
8 a meta algorithm that uses two layers of bandit algorithms in order to learn 1/the best set of items
9 overall and 2/ the individual preferences. This is a novel and efficient bandit startegy for the latent
10 bandits and an elegant generalization of the rank-1 setting. We show a regret bound for our algorithm
11 and run experiments on simulated and real data.

12 XXXXXXXXXXXXXXXXXXXXXXX

13 Cla: I haven't changed this section yet, wanted to make sure the story is right before.

14 In this paper, we study the problem of recommending the best items to users who are coming
15 sequentially. The learner has access to very less prior information about the users and it has to adapt
16 quickly to the user preferences and suggest the best item to each user. Furthermore, we consider the
17 setting where users are grouped into clusters and within each cluster the users have the same choice
18 of the best item, even though their quality of preference may be different for the best item. These
19 clusters along with the choice of the best item for each user are unknown to the learner. Also, we
20 assume that each user has a single best item preference.

21 This complex problem can be conceptualized as a low rank stochastic bandit problem where there
22 are K users and L items. The reward matrix, denoted by $\bar{M} \in [0, 1]^{K \times L}$, generating the rewards
23 for user, item pair has a low rank structure. The online learning game proceeds as follows, at every
24 timestep t , nature reveals one user (or row) from \bar{M} where user is denoted by i_t . The learner selects
25 some items (or columns) from \bar{M} , where an item is denoted by $j_t \in [L]$. Then the learner receives
26 one noisy feedback $r_t(i_t, j_t) \sim \mathcal{D}(\bar{M}(i_t, j_t))$, where \mathcal{D} is a distribution over the entries in \bar{M} and
27 $\mathbb{E}[r_t(i_t, j_t)] = \bar{M}(i_t, j_t)$. Then the goal of the learner is to minimize the cumulative regret by quickly
28 identifying the best item j^* for each $i \in [K]$ where $\bar{M}(i, j^*) = \arg \max_{j \in [L]} \{\bar{M}(i, j)\}$.

29 1.1 Notation and Learning Setting

30 Throughout the paper, we denote $[n] = \{1, 2, \dots, n\}$. An instance of the *Low-Rank Bandit* problem is
31 a matrix $R \in [0, 1]^{K \times L}$ representing the expected click-through rates (CTRs) for each user $k \in [K]$
32 on each item $l \in [L]$. If, $J \subset [L]$ is a subset of columns, we denote $R(:, J) \in [0, 1]^{K \times |J|}$ the
33 corresponding submatrix containing the $|J|$ columns of R .

We assume that there exists a latent structure, i.e that $R = UV^T$ where the rows of U and V contain the hidden users' and item's features. It is important to notice that none of those features are observable, meaning that we cannot build on a linear bandit model, and in particular our problem cannot be seen as a *clustering of bandits* problem Gentile et al. (2014). However, the rank of the CTR matrix is assumed to be low, that is $d \ll \min\{L, K\}$. This is the key assumption of our model. It implies, by definition, the following property.

Observation 1. Let $M \in \mathbb{R}^{K \times L}$ be a rank- d matrix. Then,

- There exists a basis J^* of d column such that all the L columns' latent features are linear combinations of the vectors in J^* ;
- There exists a basis I^* of d users such that all the K users' latent features are linear combinations of the vectors in I^* .

Without loss of generality, the above mentioned bases can be chosen of maximal volume such that the corresponding transformation matrix is the least singular possible.

Proof. The existence of the basis on both dimensions comes directly by definition of the low rank assumption. The choice of the spanning vectors is arbitrary and maximising the volume means choosing vectors with larger norm and hence potentially larger payoff. ■

Cl: Here state the result on the existence of a best set of d items, I'm not sure how to state it. It is not an "assumption" though, it is a Lemma or a Fact but not an assumption. It is a consequence of the low rank assumption :)

The interaction at round $t \geq 1$ of the learner with the online recommender system characterized by R goes as follows:

- a user $i_t \in [K]$ shows up – it corresponds to the index of a row of the matrix. It can be seen as an unobserved context generated by the environment;
- the learner chooses a set $J_t \subset [L]$ such that $|J_t| = d$ to be sequentially presented to the user;
- the user browses those d options and send an individual feedback for each of them (semi-bandit setting): $\forall j \in J_t$, the learner observes $Y_{t,j} = R(i_t, j) + \eta_{t,j}$ where $(\eta_{t,j})_{t,j \geq 0}$ is a sequence of i.i.d centered random variables.

Cl: fix your noise model here. Bernoulli ??

For each user $i \in [K]$, there exists one unique best item $j^*(i) \in [L]$

Cl: Define the best item, define the expected regret

The objective of the learning agent is to minimize the expected cumulative regret up to horizon n . We define the cumulative regret, denoted by \mathcal{R}_n as,

1.2 Related Works

In Maillard and Mannor (2014) the authors propose the Latent Bandit model where there are two sets: 1) set of arms denoted by \mathcal{A} and 2) set of types denoted by \mathcal{B} which contains the latent information regarding the arms. The latent information for the arms are modeled such that the set \mathcal{B} is assumed to be partitioned into $|\mathcal{C}|$ clusters, indexed by $\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_C \in \mathcal{C}$ such that the distribution $v_{a,b}, a \in \mathcal{A}, b \in \mathcal{B}_c$ across each cluster is same. Note, that the identity of the cluster is unknown to the learner. At every timestep t , nature selects a type $b_t \in \mathcal{B}_c$ and then the learner selects an arm $a_t \in \mathcal{A}$ and observes a reward $r_t(a, b)$ from the distribution $v_{a,b}$.

Another way to look at this problem is to imagine a matrix of dimension $|\mathcal{A}| \times |\mathcal{B}|$ where again the rows in \mathcal{B} can be partitioned into $|\mathcal{C}|$ clusters, such that the distribution across each of this clusters are same. Now, at every timestep t one of this row is revealed to the learner and it chooses one column such that the $v_{a,b}$ is one of the $\{v_{a,c}\}_{c \in \mathcal{C}}$ and the reward for that arm and the user is revealed to the learner.

77 This is actually a much simpler approach than the setting we considered because note that the
 78 distributions across each of the clusters $\{v_{a,c}\}_{c \in \mathcal{C}}$ are identical and estimating one cluster distribution
 79 will reveal all the information of the users in each cluster.

80 **2 MetaBand: learning efficiently on a rank- d matrix.**

81 Let $\bar{M} = UV^\top$, where U is non-negative and V is hott topics. Let j_1^* and j_2^* be the indices of
 82 hott-topics vectors. Then

$$(j_1^*, j_2^*) = \arg \max_{j_1, j_2 \in [L]} f(\{j_1, j_2\}),$$

83 where $f(S) = \frac{1}{K} \sum_{i \in [K]} \max_{j \in S} R(i, j)$

84 The key observation is that f is monotone and submodular in S . Therefore, the problem of learning
 85 j_1, j_2 online is an online submodular maximization problem.

86 So, when $d = 2$, $|\mathcal{B}_t| = 2$ and there are two EXP3 Column-Bandits.

87 After observing the reward r_1, r_2 for $j_1, j_2 \in \mathcal{B}_t$ we update,

88 $EXP_1, \hat{r}_{1,j_1} = r_1$.

89 $EXP_2, \hat{r}_{2,j_2} = \max\{r_1, r_2\} - r_1$.

Algorithm 1 Low Rank Bandit Strategy

```

1: Input: Time horizon  $n$ ,  $\text{Rank}(\bar{M}) = d$ .
2: for  $t = 1, \dots, n$  do
3:   Nature reveals user  $i_t$ . ▷ Nature chooses user
4:   Column-Bandits suggests  $\mathcal{B}_t \subseteq [L]$  items.  $|\mathcal{B}_t| = d$ 
5:   if Exploration condition satisfied then
6:     User Bandits suggests each item in  $\mathcal{B}_t$ , once to user  $i_t$  and receive feedback.
7:     Update Column-Bandits and User Bandits on feedback received.
8:   else
9:     Suggest best item in  $\mathcal{B}_t$   $d$  times to user  $i_t$  and receive feedback.
```

Algorithm 2 Low Rank Bandit Greedy (LRG)

```

1: Input: Time horizon  $n$ ,  $\text{Rank}(\bar{R}) = d$ .
2: Explore Parameters:  $\epsilon \in (0, 1)$ .
3: for  $t = 1, \dots, n$  do
4:   Nature reveals user  $i_t$ . ▷ Nature chooses user
5:   Column-EXP3 suggests  $\mathcal{B}_t \subseteq [L]$  items.  $|\mathcal{B}_t| = d$ 
6:   With  $\epsilon$  probability do ▷ Exploration
7:     User Bandit suggests each arm  $j \in \mathcal{B}_t$  once to user  $i_t$  and receive feedback.
8:   Or With  $(1 - \epsilon)$  probability do ▷ Exploitation
9:     User Bandit suggests arm  $j \in \arg \max_{j \in \mathcal{B}_t} \{\hat{R}(i_t, j)\}$ ,  $d$  times to user  $i_t$  and receive feedback.
10:  Update Column-Bandits and User Bandit on feedback received.
```

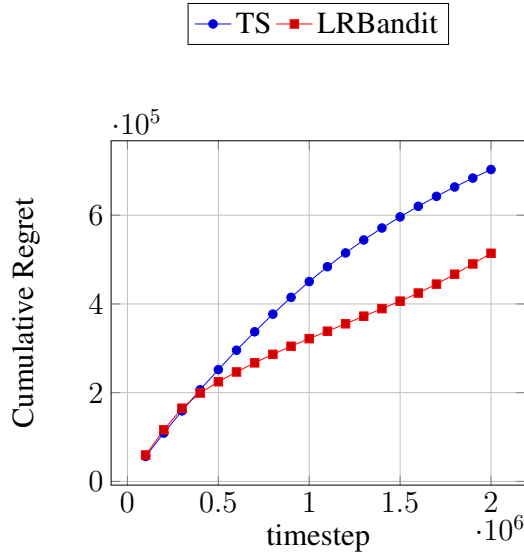
Algorithm 3 Low Rank Bandit UCB (LRUCB)

```

1: Input: Time horizon  $n$ ,  $\text{Rank}(\bar{R}) = d$ .
2: Definition:  $U(i, j) = \sqrt{\frac{2 \log n}{N_{i,j}}}$ .
3: for  $t = 1, \dots, n$  do
4:   Nature reveals user  $i_t$ . ▷ Nature chooses user
5:   Column-EXP3 suggests  $\mathcal{B}_t \subseteq [L]$  items.  $|\mathcal{B}_t| = d$ 
6:   if  $\left( \hat{R}(i_t, j) - U(i_t, j) \leq \hat{R}(i_t, j') + U(i_t, j') \right), \forall j, j' \in \mathcal{B}_t$  then ▷ Confidence interval
       overlap, Exploration
7:     User Bandit suggests each arm  $j \in \mathcal{B}_t$  once to user  $i_t$  and receive feedback.
8:   else ▷ Exploitation
9:     User Bandit suggests arm  $j \in \arg \max_{j \in \mathcal{B}_t} \left\{ \hat{R}(i_t, j) + U(i_t, j) \right\}$ ,  $d$  times to user  $i_t$  and
       receive feedback.
10:  Update Column-Bandits and User Bandits on feedback received.

```

90 3 Experiments



(a) Expt-1: 1024 Users, 128 arms, Round-Robin, Noisy Setting, Rank 2, equal sized clusters

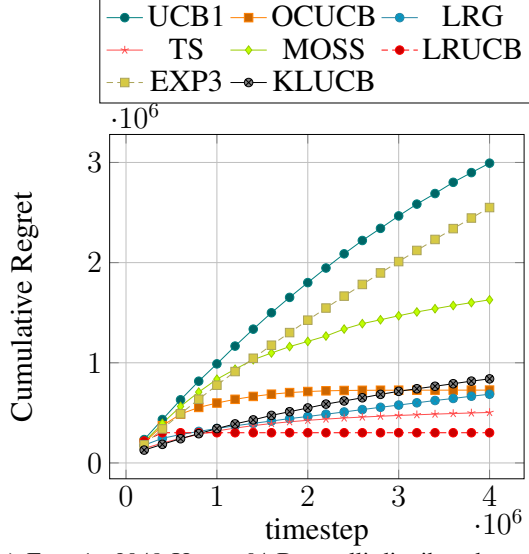
Figure 1: A comparison of the cumulative regret by MRLG and MRLUCB.

91 4 Conclusions and Future Direction

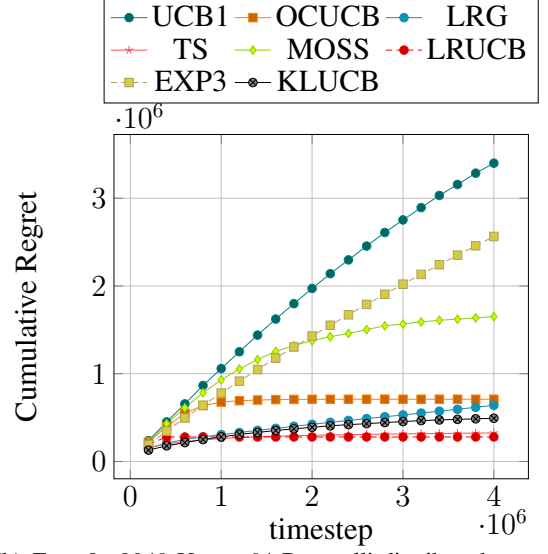
92 To be written.

93 References

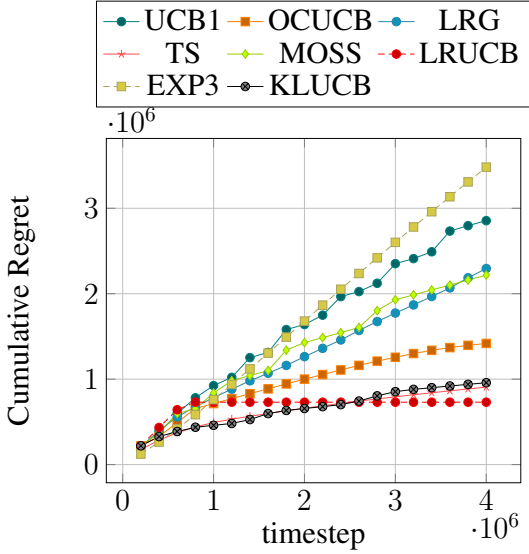
- 94 Gentile, C., Li, S., and Zappella, G. (2014). Online clustering of bandits. In *International Conference*
95 *on Machine Learning*, pages 757–765.
- 96 Maillard, O.-A. and Mannor, S. (2014). Latent bandits. In *International Conference on Machine*
97 *Learning*, pages 136–144.



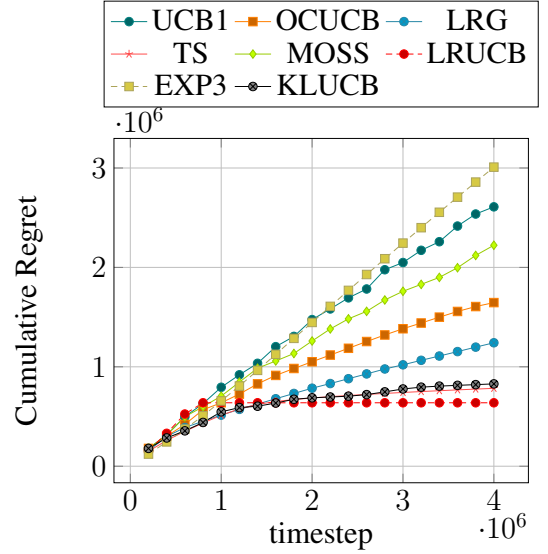
(a) Expt-1: 2048 Users, 64 Bernoulli-distributed arms, Round-Robin, Noisy Setting, Rank 2, equal sized clusters



(b) Expt-2: 2048 Users, 64 Bernoulli-distributed arms, Round-Robin, Noisy Setting, Rank 2, un-equal sized clusters, 70:30 split



(c) Expt-3: 4096 Users, 128 Bernoulli-distributed arms, Round-Robin, Noisy Setting, Rank 2, equal sized clusters



(d) Expt-4: 4096 Users, 128 Bernoulli-distributed arms, Round-Robin, Noisy Setting, Rank 2, un-equal sized clusters, 80:20 split

Figure 2: A comparison of the cumulative regret incurred by the various bandit algorithms.