
Stochastic Low-Rank Latent Bandits

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To be written.

2 1 Introduction

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

10 This complex problem can be conceptualized as a low rank stochastic bandit problem where there
11 are K users and L items. The reward matrix, denoted by $\bar{M} \in [0, 1]^{K \times L}$, generating the rewards
12 for user, item pair has a low rank structure. The online learning game proceeds as follows, at every
13 timestep t , nature reveals one user (or row) from \bar{M} where user is denoted by i_t . The learner selects
14 some items (or columns) from \bar{M} , where an item is denoted by $j_t \in [L]$. Then the learner receives
15 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
16 $\mathbb{E}[r_t(i_t, j_t)] = \bar{M}(i_t, j_t)$. Then the goal of the learner is to minimize the cumulative regret by quickly
17 identifying the best item j^* for each $i \in [K]$ where $\bar{M}(i, j^*) = \arg \max_{j \in [L]} \{\bar{M}(i, j)\}$.

18 1.1 Notations, Problem Formulation and Assumptions

19 We define $[n] = \{1, 2, \dots, n\}$ and for any two sets A and B , A^B denotes the set of all vectors who
20 take values from A and are indexed by B . Let, $M \in [0, 1]^{K \times L}$ denote any matrix, then $M(I, :)$
21 denote any submatrix of k rows such that $I \in [K]^k$ and similarly $R(:, J)$ denote any submatrix of j
22 columns such that $J \in [L]^j$.

23 Let \bar{M} be reward matrix of dimension $K \times L$ where K is the number of user or rows and L is the
24 number of arms or columns. Also, let us assume that this matrix \bar{M} has a low rank structure of rank
25 $d \ll \min\{L, K\}$. Let U and V denote the latent matrices for the users and items, which are not
26 visible to the learner such that,

$$\bar{M} = UV^\top \quad \text{s.t.} \quad U \in [\mathbb{R}^+]^{K \times d}, V \in [0, 1]^{L \times d}$$

27 Furthermore, we put a constraint on V such that, $\forall j \in [L], \|V(j, :)\|_1 \leq 1$.

28 **Assumption 1.** We assume that there exists d -column base factors, denoted by $V(J^*, :)$, such that
29 all rows of V can be written as a convex combination of $V(J^*, :)$ and the zero vector and $J^* = [d]$.
30 We denote the column factors by $V^* = V(J^*, :)$. Therefore, for any $i \in [L]$, it can be represented by

$$V(i, :) = a_i V(J^*, :),$$

31 where $\exists a_i \in [0, 1]^d$ and $\|a_i\|_1 \leq 1$.

32 **Assumption 2.** For each user i_t revealed by the nature at round t , the learner is allowed to select
 33 atmost d -items, where d is the rank of the matrix \bar{R} .

34 The above assumption 2 can be conceptualized in this real-world scenario where the learner has to
 35 suggest movies to users and each movie belongs to a different genre (say thriller, romance, comedy,
 36 etc). So, the learner can suggest d movies belonging to different genres to each user, and the user can
 37 click one, or all, or none of the recommended movies.

38 The main goal of the learning agent is to minimize the cumulative regret until the end of horizon n .
 39 We define the cumulative regret, denoted by \mathcal{R}_n as,

$$\mathcal{R}_n = \sum_{t=1}^n \left\{ \sum_{z=1}^d \left(r_t(i_t, j_{t,z}^*) - r_t(i_t, j_{t,z}) \right) \right\}$$

40 where, $j_{t,z}^* = \arg \max_{j \in [L]} \{\bar{M}(i_t, j)\}$ and $j_{t,z}$ be the suggestion of the learner for the i_t -th user for
 41 $z = 1, 2, \dots, d$. Note that $r_t(i_t, j_{t,z}^*) \sim \mathcal{D}(\bar{M}(i_t, j_{t,z}^*))$ and $r_t(i_t, j_{t,z}) \sim \mathcal{D}(\bar{M}(i_t, j_{t,z}))$. Taking
 42 expectation over both sides, we can show that,

$$\mathbb{E}[\mathcal{R}_n] = \mathbb{E} \left[\sum_{t=1}^n \left\{ \sum_{z=1}^d \left(r_{z,t}(i_t, j_{t,z}^*) - r_{z,t}(i_t, j_{t,z}) \right) \right\} \right] = \mathbb{E} \left[\sum_{t=1}^n \sum_{z=1}^d \left(N_{i_t, j_{t,z}, t} \right) \right] \Delta_{i_t, j_{t,z}, t}$$

43 where, $\Delta_{i_t, j_{t,z}, t} = \bar{M}(i_t, j_{t,z}^*) - \bar{M}(i_t, j_{t,z})$ and $N_{i_t, j_{t,z}, t}$ is the number of times the learner has
 44 observed the $j_{t,z}$ -th item for the i_t -th user for $z = 1, 2, \dots, d$. Let, $\Delta = \min_{i \in [K], j \in [L]} \{\Delta_{i,j}\}$ be the
 45 minimum gap over all the user, item pair in \bar{M} .

46 1.2 Related Works

47 In Maillard and Mannor (2014) the authors propose the Latent Bandit model where there are two
 48 sets: 1) set of arms denoted by \mathcal{A} and 2) set of types denoted by \mathcal{B} which contains the latent
 49 information regarding the arms. The latent information for the arms are modeled such that the set \mathcal{B}
 50 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
 51 $v_{a,b}, a \in \mathcal{A}, b \in \mathcal{B}_c$ across each cluster is same. Note, that the identity of the cluster is unknown to
 52 the learner. At every timestep t , nature selects a type $b_t \in \mathcal{B}_c$ and then the learner selects an arm
 53 $a_t \in \mathcal{A}$ and observes a reward $r_t(a, b)$ from the distribution $v_{a,b}$.

54 Another way to look at this problem is to imagine a matrix of dimension $|\mathcal{A}| \times |\mathcal{B}|$ where again the
 55 rows in \mathcal{B} can be partitioned into $|\mathcal{C}|$ clusters, such that the distribution across each of this clusters are
 56 same. Now, at every timestep t one of this row is revealed to the learner and it chooses one column
 57 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
 58 learner.

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

62 2 Contributions

63 To be written.

64 3 Proposed Algorithms

65 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
 66 hott-topics vectors. Then

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

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

68 The key observation is that f is monotone and submodular in S . Therefore, the problem of learning

69 j_1, j_2 online is an online submodular maximization problem.

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

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

72 $EXP_1, \hat{r}_{1,j_1} = r_1$.

73 $EXP_2, \hat{r}_{2,j_2} = \max\{r_1, r_2\} - r_1$.

Algorithm 1 Low Rank Bandit Strategy

1: **Input:** Time horizon n , $Rank(\bar{M}) = d$.

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 , $Rank(\bar{R}) = d$.

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** ϵ 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 , $Rank(\bar{R}) = d$.

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. $|\mathcal{B}_t| = d$

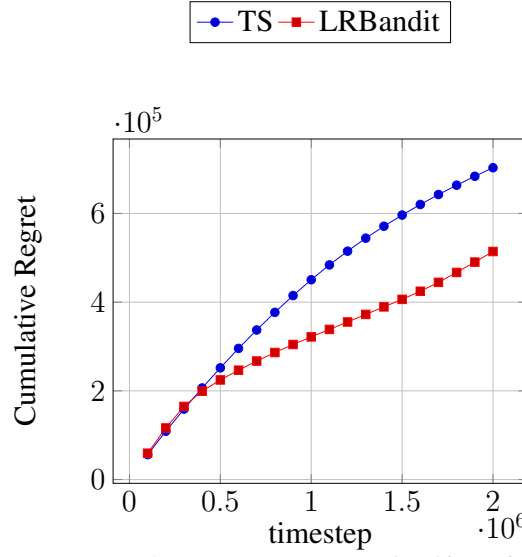
6: **if** $(\hat{R}(i_t, j) - U(i_t, j) \leq \hat{R}(i_t, j') + U(i_t, j'))$, $\forall j, j' \in \mathcal{B}_t$ **then** ▷ Confidence interval overlap, Exploration

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 $j \in \arg \max_{j \in \mathcal{B}_t} \{\hat{R}(i_t, j) + U(i_t, j)\}$, d times to user i_t and receive feedback.

10: Update Column-Bandits and User Bandits on feedback received.



(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.

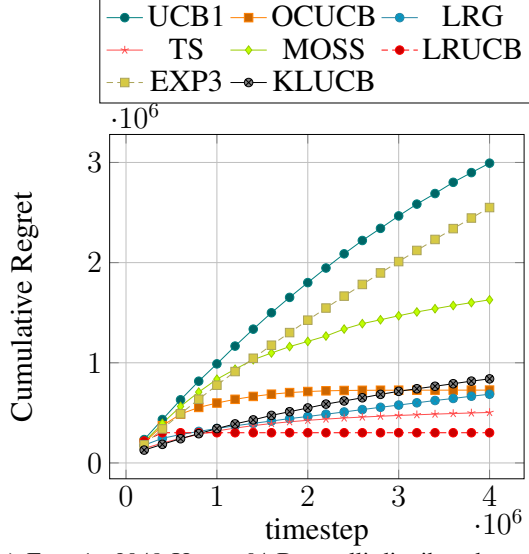
74 4 Experiments

75 5 Conclusions and Future Direction

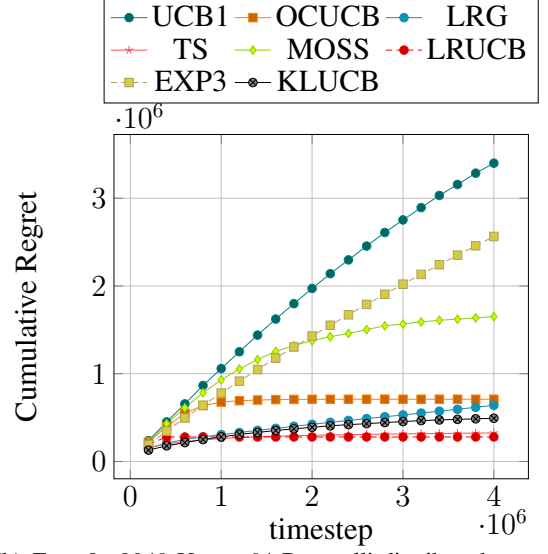
76 To be written.

77 References

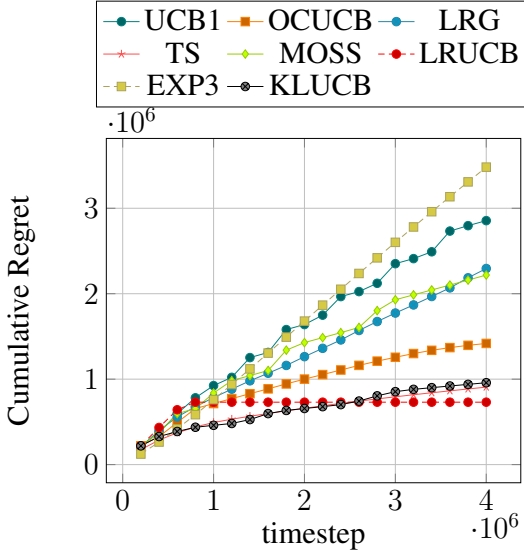
78 Maillard, O.-A. and Mannor, S. (2014). Latent bandits.



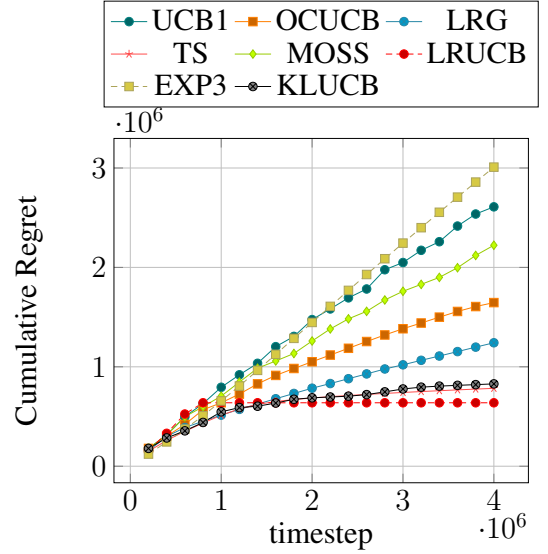
(a) Expt-1: 2048 Users, 64 Bernoulli-distributed arms, Round-Robin, Noisy Setting, Rank 2, equal sized clusters



(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.