

Group Learning and Opinion Diffusion in a Broadcast Network

Yang Liu, Mingyan Liu
Electrical Engineering and Computer Science
University of Michigan, Ann Arbor
{youngliu, mingyan}@umich.edu

Abstract—We analyze the following group learning problem in the context of opinion diffusion: Consider a network with M users, each facing N options. In a discrete time setting, at each time step, each user chooses K out of the N options, and receive randomly generated rewards, whose statistics depend on the options chosen as well as the user itself, and are unknown to the users. Each user aims to maximize their expected total rewards over a certain time horizon through an online learning process, i.e., a sequence of exploration (sampling the return of each option) and exploitation (selecting empirically good options) steps. Different from a typical regret learning problem setting (also known as the class of multi-armed bandit problems), the group of users share information regarding their decisions and experiences in a broadcast network. The challenge is that while it may be helpful to observe others’ actions in one’s own learning (i.e., second-hand learning), what is considered desirable option for one user may be undesirable for another (think of restaurant choices), and this difference in preference is in general unknown a priori. Even when two users happen to have the same preference (e.g., they agree one option is better than the other), they may differ in their absolute valuation of each individual option.

Within this context we consider two group learning scenarios, (1) users with uniform preferences and (2) users with diverse preferences, and examine how a user should construct its learning process to best extract information from others’ decisions and experiences so as to maximize its own reward. Performance is measured in *weak regret*, the difference between the user’s total reward and the reward from a user-specific best single-action policy (i.e., always selecting the set of options generating the highest mean rewards for this user). Within each scenario we also consider two cases: (i) when users exchange full information, meaning they share the actual rewards they obtained from their choices, and (ii) when users exchange limited information, e.g., only their choices but not rewards obtained from these choices. We show the gains from group learning compared to individual learning from one’s own choices and experiences.

I. INTRODUCTION

We analyze the following group learning problem in the context of opinion diffusion: Consider a network with M users, each facing N options. In a discrete time setting, at each time step, each user chooses K out of the N options, and receive randomly generated rewards, whose statistics depend on the options chosen as well as the user itself, and are unknown to the users. Each user aims to maximize its

expected total reward over a certain time horizon through an online learning process, i.e., a sequence of exploration (sampling the return of each option) and exploitation (selecting empirically good options) steps. Taken separately, an individual user’s learning process may be mapped into a standard multi-armed bandit (MAB) problem which has been extensively studied, see e.g., [5], [2], [3].

Our interest in this study, however, is on how an individual’s learning process may be affected by “second-hand learning”, i.e., by observing how others in the group act. The challenge is that while it may be helpful to observe others’ actions to speed up one’s own learning, what is considered desirable option for one may be undesirable for another (think of restaurant choices: one Yelp user’s recommendation may or may not be useful for another), and this difference in preference is in general unknown a priori. Moreover, even when two users happen to have the same preference (e.g., they agree one option is better than the other), they may differ in their absolute valuation of each individual option (again think of restaurant choices: two Yelp users may agree restaurant A is better than B, but one user may rate them 5 and 4 stars respectively, while the other may rate them 4 and 3 stars, respectively).

Consequently it seems that if an individual wants to take others’ actions into account in its own learning process, it would also need to figure out whether their preferences are aligned, which may add to the overhead in the learning process. This raises the interesting question of whether learning from group behavior is indeed beneficial to an individual, and if so what type of learning algorithm can effectively utilize the group information in addition to its own direct observations. This is what we aim to address in this paper.

We will assume that users are heterogeneous in general, i.e., when using the same option they obtain rewards driven by different random processes with different mean values. We then consider two scenarios. (1) In the first, users have *uniform preference ordering* of the N options. This means that even though they may value each options differently (i.e., have different reward statistics), they always have the same preference ordering. (2) In the second, users have *diverse preference orderings* of the N options, meaning that one user’s best options are not so for another. Within each scenario we also consider two cases: (i) when users

exchange full information, meaning they disclose the actual rewards they obtained from their choices, and (ii) when users exchange limited information, e.g., only their choices but not rewards obtained from these choices. For each of these cases we examine how a user should construct its learning process to best extract information from others' decisions and experiences so as to maximize its own reward. Performance is measured in *weak regret*, the difference between the user's total reward and the reward from a user-specific best single-action policy (i.e., always selecting the set of options generating the highest mean rewards for this user).

This problem can also be viewed as a learning problem with contextual information (or side information in some literature), see e.g., [4], [7], [6]. However, in these studies statistical information linking a user's own information and the side information is required in the following sense. Denote by X a user's observation and by Y the side information (say shared information from other users), the knowledge of the conditional probability of observing X (i.e., $p(X|Y)$) needs to be given or assumed. In contrast, we do not require such statistical information; instead we examine how a user can estimate and learn from the shared, and possibly imperfect side information.

The paper is organized as follows. Section II gives the system model. Sections III and IV analyze the uniform and diverse preference scenarios, respectively. Numerical results are presented in Section V and Section VI concludes the paper.

II. PROBLEM FORMULATION AND SYSTEM MODEL

Consider a system or network of M users indexed by the set $\mathcal{U} = \{1, 2, \dots, M\}$ and a set of available options denoted by $\Omega = \{1, 2, \dots, N\}$. The system works in discrete time indexed by $t = 1, 2, \dots$. At each time step a user can choose up to K options. For user i an option j generates an IID reward denoted by random variable X_j^i , with a mean reward given by $\mu_j^i := \mathbb{E}[X_j^i]$. We will assume that $\mu_l^i \neq \mu_j^i, l \neq j, \forall i \in \mathcal{U}$, i.e., different options present distinct values to a user. We will denote the set of top K options (in terms of mean rewards) for user i as N_K^i and its complement \overline{N}_K^i . Denote by $a^i(t)$ the set of choices made by user i at time t ; the sequence $\{a^i(t)\}_{t=1,2,\dots}$ constitutes user i 's policy.

Following the classical regret learning literature, we will adopt the *weak regret* as a performance metric, which measures the gap between the total reward (up to some time T) of a given learning algorithm and the total reward of the best single-action policy given a priori average statistics, which in our case is the sum reward generated by the top K options for a user. This is formally given as follows for user i adopting algorithm a :

$$R^{i,a}(T) = T \cdot \sum_{j \in N_K^i} \mu_j^i - \mathbb{E} \left[\sum_{t=1}^T \sum_{j \in a^i(t)} X_j^i \right] \quad (1)$$

The goal of a learning algorithm is to minimize the above regret measure.

As mentioned in the introduction, we consider two scenarios. In the first case, users share the same preference ordering over the N options, i.e., if $\mu_{j_1}^i > \mu_{j_2}^i, j_1, j_2 \in \Omega$, then $\mu_{j_1}^k > \mu_{j_2}^k, \forall k \neq i, k \in \mathcal{U}$. This implies that $N_K^i = N_K^k, \forall i, k \in \mathcal{U}$. This will be referred to as the *uniform preference* scenario.

In the second, the *diverse preference* scenario, users have different preference orderings over the N options. Specifically, in this case we will assume that the M users may be classified into G distinct groups, indexed by the set $\mathcal{G} = \{1, 2, \dots, G\}$, with users within the same group (say group l) having a unique K -preferred set N_K^l , assumed to be public knowledge. Note that even with the same preferred set, users may be further classified based on the actual ordering of these top K options. Our model essentially bundles these sub-classes into the same group, provided their top K choices are the same. This is because as a user is allowed K choices at a time, further distinguishing their preferences within these K options will not add to the performance of an algorithm.

Under each scenario, we further consider two types of information shared/exchanged by the users. Under the first type, users disclose *full information*: they not only announce the decisions they make (the options they choose), but also the observations following the decisions, i.e., the actual rewards received from those options. Such announcements may be made at the end of each time step, or may be made periodically but at a lesser frequency. The second type of exchange is *partial information* where users disclose only part of decisions and/or observations. Specifically, we will assume that the users only share their decision information, i.e., the set of choices they make, at the beginning of each time step, but withhold the actual observation/reward information following the decisions.

III. GROUP LEARNING WITH UNIFORM PREFERENCE

Without loss of generality, we will assume that under the uniform preference ordering we have $\mu_1^i > \mu_2^i > \dots > \mu_N^i, \forall i \in \mathcal{U}$.

A. Uniform preference, full information (U_FULL)

This case will be referred to as U_FULL. Under this model users not only broadcast their decisions within the network, but also release observations of selected options' quality/rewards at the end of each time step. Since users have the same preference ordering, a fact assumed to be known to the users, it would seem straightforward that one user could easily learn from another. The challenge here lies in the fact that the statistics driving the rewards are not identical for all users even when using the same option. So information obtained from another user may need to be treated differently from one's own observations.

In general the reward user i obtains from option j may be modeled as

$$X_j^i = f(X_j, \mathcal{N}_i, \mathcal{L}_i), \quad (2)$$

where $f(\cdot)$ is some arbitrary unknown function, X_j describes certain *intrinsic* or *objective* value of option j that is independent of the specific user (e.g., the bandwidth of a channel, or the rating given to a restaurant by AAA, and so on), \mathcal{N}_i is a noise term, and \mathcal{L}_i captures user-specific features that affect the *perceived* value of this option to user i (e.g., user i 's location information or transceiver specification which may affect its perceived channel quality, or user i 's dietary restriction which may affect its preference for different types of restaurants). For simplicity in this study we will limit our attention to the following special case of user-specific valuations, where the rewards received by two users from the same option are given by a linear relationship:

$$\mu_j^i / \mu_j^k = \delta_j^{i,k}. \quad (3)$$

The scaling factor $\delta_j^{i,k}$ will be referred to as the *distortion* or *distortion factor* between two users.

Under this model it can be seen that a user could recov-

er/convert observations from other users for its own use by estimating the distortion. Consider two users i and k , and option j . Denote by $r_j^i(t)$ the sample mean reward collected by i directly itself from option j up to time t . This quantity is not only available to user i , but also to all other users k due to the full information disclosure, and vice versa. User i then estimates the distortion between itself and user k by $\tilde{\delta}_j^{i,k}(t) = r_j^i(t) / r_j^k(t)$.

With this quantity we then make the following simple modification to the well-known UCB algorithm introduced in [3]. In the original UCB (or rather, a trivial multiple-play extension of it), user i 's decision $a^i(t)$ at time t is entirely based on its own observations. Specifically, denote by $n_j^i(t)$ the number of times user i has selected option j up to time t . The original UCB then selects option j at time t , if its index value given below is among the K highest:

$$\text{UCB index: } r_j^i(t) + \sqrt{\frac{2 \log t}{n_j^i(t)}}. \quad (4)$$

Under the modified algorithm (referred to as the U_FULL algorithm), option j is selected at time t if its index value defined below is among the K highest:

$$\text{U_FULL index: } \frac{\sum_{m=1}^t X_j^i(m) \cdot \mathbf{I}_{j \in a^i(m)} + \sum_{k \neq i} \sum_{m=1}^t \tilde{\delta}_j^{i,k}(m) X_j^k(m) \cdot \mathbf{I}_{j \in a^k(m)}}{\sum_{i \in \mathcal{U}} n_j^i(t)} + \sqrt{\frac{2 \log t}{\sum_{i \in \mathcal{U}} n_j^i(t)}}. \quad (5)$$

We have the following results on algorithm U_FULL.

Theorem 3.1: The weak regret of user i under U_FULL is upper bounded by

$$R_{\text{U_FULL}}^i(t) \leq \sum_{j \in \bar{N}_K^i} \left\lceil \frac{8 \log t}{M \cdot \Delta_j^i} \right\rceil + \text{const.} \quad (6)$$

where $\Delta_j^i = \mu_K^i - \mu_j^i$.

Proof 1: The proof can be found in Appendix-A, in which Lemma A.2 and Theorem 3.1 are proved simultaneously.

Under the original UCB algorithm [3] a single user's weak regret is upper bounded by (the superscript i is suppressed here because the result applies to any single user)

$$R_{\text{UCB}}(t) \leq \sum_{j \in \bar{N}_K} \left\lceil \frac{8 \log t}{\Delta_j} \right\rceil + \text{const.} \quad (7)$$

Therefore we see that there is potential gain in group learning. Note however the improvement is not guaranteed as it appears in an upper bound, which does not necessarily imply better performance. The performance comparison is shown later via simulation.

It can be shown that similar result exists when the full information is broadcast at periodic intervals but not neces-

sarily at the end of each time step.

B. Uniform preference, partial information (U_PART)

We now consider the case where users only share their decisions/actions, but not their direct observations. This case (and the associated algorithm) will be referred to as U_PART. The difficulty in this case comes from the fact that to a user i , even though other users' actions reflect an option's relative value to them (and by positive association to user i itself), the actions do not directly reveal the actual observations.

Denote by $n_j(t)$ the total number of times option j has been selected by the entire group up to time t . Then $\beta_j(t) := \frac{n_j(t)}{\sum_{l \in \Omega} n_l(t)}$ denotes the frequency at which option j is being used by the group up to time t . This will be referred to as the group recommendation or behavior. Several observations immediately follow. Firstly, we have $\sum_{j \in \Omega} \beta_j(t) = 1, \forall t$. Secondly, as time goes on, we would like better options j to be selected with higher frequency, i.e., larger $\beta_j(t)$.

With these observations, we construct the following algorithm U_PART, by biasing toward potentially good options as indicated by the group behavior. Under the U_PART algorithm, option j is selected at time t if its index value

defined below is among the K highest:

$$\begin{aligned} & \text{U_PART index:} \\ r_j^i(t) - \alpha(1 - \beta_j(t))\sqrt{\frac{\log t}{t}} + \sqrt{\frac{2 \log t}{n_j^i(t)}}, \end{aligned} \quad (8)$$

where α is a weighting factor over the group recommendation.

A few remarks are in order. In the above index expression, the middle, bias term serves as a penalty: a larger group frequency $\beta_j(t)$ means a smaller penalty. But its effect diminishes as t increases. This reflects the notion that as time goes on a user becomes increasingly more confident in its own observations and relies less and less on the group recommendation. Lastly, the weight α captures how much the user values the group recommendation compared to its own observations, with a small value indicating a small weight.

We have the following result on the U_PART algorithm.

Theorem 3.2: The weak regret of user i under U_PART is upper bounded by

$$R_{\text{U_PART}}^i(t) \leq \sum_{j \in \bar{N}_K^i} \left\lceil \frac{4(\sqrt{2} - \alpha(1 - \beta_j(t)))^2 \log t}{\Delta_j^i} \right\rceil + \text{const.}, \quad (9)$$

and $\beta_j(t) \leq \mathcal{O}(\frac{\log t}{t})$.

Proof 2: Proof can be found in Appendix-B.

Again, compared to the bounds from the original UCB, we potentially achieve a better performance as the bound constant decreases from 8 to $4(\sqrt{2} - \alpha(1 - \beta_j(t)))^2$ with the group recommendation mechanism, but with the same cautionary note on the upper bound. The performance comparison is shown in simulation results later.

IV. GROUP LEARNING WITH DIVERSE PREFERENCES

A. Diverse preferences, full information (D_FULL)

To capture simultaneously the group difference and the individual difference, we will assume that the mean rewards

$$\text{Assign user } i \text{ to group } g_i(t) \text{ if: } g_i(t) = l^* = \operatorname{argmax}_{l \in G} D^{i,l}(t) = |\tilde{N}_K^i(t) \cap N_K^l|, \quad (11)$$

with ties broken randomly.

We again estimate the pair-wise distortion factor in a manner similar to the uniform preference case:

received by two users from option j are related as follows: Since different groups may have its own identity information, the direct attenuation factor model may fail to catch the observation differences between different groups. To illustrate the idea better we consider a slightly different model compared to the uniform group case. Consider the following model:

$$\mu_j^i = d_{ik} + \delta_j^{i,k} \cdot \mu_j^k, \quad (10)$$

where the additive term d_{ik} is the preference distance between groups: $d_{ik} = 0$ for i, k belonging to the same group, while the multiplicative term $\delta_j^{i,k}$ remains the distortion between two individuals beyond the group difference.

For group identity to be meaningful one would expect the group distance values d to be on a higher order of magnitude than the individual distortion values δ . One way to learn from the shared information is then to identify one's group affiliation with respect to others', which is then used to weigh different observations from different individuals. Here we assume users know the group preference distances d_{ik} . We introduce the following sample frequency based group identity classification mechanism. Each user keeps the same set of statistics $n_j^i(t)$ as before: the number of times user i is seen using option j . From these a user tries to estimate another's preference by ordering the statistics: at time t user i 's preference is estimated to be the set $\tilde{N}_K^i(t)$, which contains elements/options j whose frequency $n_j^i(t)$ is among the K highest of all i 's frequencies. User i is then put in the preference group l with whose (known) preferred set N_K^l its estimated preference $\tilde{N}_K^i(t)$ is the closest in distance, defined as follows:

$\tilde{\delta}_j^{i,k}(t) = r_j^i(t)/(r_j^k(t) - d_{ik})$. The resulting D_FULL algorithm run by user i then selects, at time t , an option j if its index value given below is among the K highest:

$$\text{D_FULL index: } \frac{\sum_{m=1}^t X_j^i(m) \cdot \mathbf{I}_{j \in a^i(m)} + \sum_{k \neq i} \sum_{m=1}^t \tilde{\delta}_j^{i,k}(m) X_j^k(m) \cdot \mathbf{I}_{j \in a^k(m)}}{\sum_{i \in \mathcal{U}} n_j^i(t)} + \sqrt{\frac{2 \log t}{\sum_{i \in \mathcal{U}} n_j^i(t)}}. \quad (12)$$

We have the following result on algorithm D_FULLL.

Theorem 4.1: For each user j belonging to group r we have that the probability of incorrect classification at time t is bounded as

$$P(g_j(t) \neq r) \leq C_1 \cdot \frac{\log t}{t}, \forall (j, r), t. \quad (13)$$

for some positive constant C_1 .

Proof 3: We provide a sketch of the proof here. The main idea is to bound the probability that the number of times sub-optimal arms are played being higher than that of the optimal arms by time t . Consider $j \in \bar{N}_K$ and $* \in N_K$.

$$\begin{aligned} & P\left(\sum_{n=1}^t I_{a(n)=j} \geq \sum_{n=1}^t I_{a(n)=*}\right) \\ &= P\left(\sum_{n=1}^t I_{a(n)=j} \geq \sum_{n=1}^t I_{a(n)=*} \mid \sum_{n=1}^t I_{a(n)=*} \geq \mathcal{O}(t)\right) \\ &\cdot P\left(\sum_{n=1}^t I_{a(n)=*} \geq \mathcal{O}(t)\right) \\ &+ P\left(\sum_{n=1}^t I_{a(n)=j} \geq \sum_{n=1}^t I_{a(n)=*} \mid \sum_{n=1}^t I_{a(n)=*} < \mathcal{O}(t)\right) \\ &\cdot P\left(\sum_{n=1}^t I_{a(n)=*} < \mathcal{O}(t)\right) \end{aligned} \quad (14)$$

Each term can in turn be bounded using Markov inequality to get the desired result. Details are omitted for brevity.

Theorem 4.2: Under algorithm D_FULLL, user i 's weak regret is upper bounded by

$$R_{\text{D_FULL}}^i(t) \leq \sum_{j \in \bar{N}_K^i} \left\lceil \frac{8 \log t}{M \cdot \Delta_j^i} \right\rceil + \text{const.}, \quad (15)$$

Proof 4: The proof of Theorem 4.2 follows similarly the case of uniform preference, and is thus omitted.

B. Diverse preferences, partial information (D_PART)

As in the case of D_FULLL we can track for each user $n_j^i(t)$ and obtain the frequency of choices $\beta_j^i(t)$, and use this information to perform group classification. A user i then assigns a weight to an option j given by the ratio. Our algorithm proceeds in parallel with the uniform group case except for the following difference: each group will assign another group with certain discount for their observations instead of raw statistics. To be specific, user i will assign the following weight to the j th option:

$$\beta_j^i(t) = \frac{\sum_{k \in \mathcal{U}} (n_j^k(t))^{\omega^{i,k}}}{\sum_{m \in \Omega} \sum_{k \in \mathcal{U}} (n_m^k(t))^{\omega^{i,k}}}, \quad (16)$$

where weights $\omega^{i,k} = 1$ if i estimates user k to be in the same group as itself, and $\omega^{i,k} < 1$ otherwise. $\omega^{i,k}$ can also be chosen as a function of the group distance.

The resulting algorithm D_PART is as follows, where user i chooses option j if its index value is among the top K

highest:

$$r_j^i(t) - \alpha(1 - \beta_j^i(t)) \sqrt{\frac{\log t}{t}} + \sqrt{\frac{2 \log t}{n_j^i(t)}}, \quad (17)$$

The upper bound on the weak regret under D_PART is the same as in the case of U_PART.

Theorem 4.3: Under algorithm D_PART, user i 's weak regret is upper bounded by

$$R_{\text{D_PART}}^i(t) \leq \sum_{j \in \bar{N}_K^i} \left\lceil \frac{4(\sqrt{2} - \alpha(1 - \beta_j^i(t)))^2 \log t}{\Delta_j^i} \right\rceil + \text{const.}, \quad (18)$$

where $\beta_j^i(t) \leq \mathcal{O}(\frac{1}{t^{\omega_d}})$ for some $0 < \omega_d < 1$.

Proof 5: The proof follows similarly as in the uniform reference case with partial information exchange. Below we only sketch the main difference which comes from bounding the sample frequency of other users. First we note that

$$\sum_{l \in \Omega} \sum_{j \in \mathcal{U}} n_l^{\omega^{i,g_j}}(t) = \mathcal{O}(t). \quad (19)$$

For user j within the same group as user i we have

$$E[n_m^{\omega^{i,g_j}}(t)] = E[n_m^j(t)] = \mathcal{O}(\log t). \quad (20)$$

For user j from a different group we know

$$E[n_m^{\omega^{i,g_j}}(t)] < E[n_m^j(t)] \leq \mathcal{O}(t) \quad (21)$$

since $\omega^{i,g_j} < 1$. Meanwhile the chance of mis-classifying a user from a different group as one from the same group is upper bounded by $\mathcal{O}(\frac{\log t}{t})$, and the number of such mis-classification is at most

$$\mathcal{O}(t) \cdot \mathcal{O}\left(\frac{\log t}{t}\right) = \mathcal{O}(\log t) \quad (22)$$

which helps establish $\beta_j(t) \leq \mathcal{O}(\frac{1}{t^{\omega_d}})$.

V. NUMERICAL RESULTS

We start with U_FULLL. In our simulation we have three users with five independent options; each user targets the top three options at each time, i.e., $M = K = 3, N = 5$. Furthermore the five options' reward statistics are given by exponentially distributed random variables. The distortion factor at each user for each option is modeled as a Gaussian random variable with certain mean and variance 1.

From Fig. 1 we see with full information exchange the system's performance can be greatly improved compared with individual learning. Moreover, its performance is comparable to a centralized scheme (denoted UCB Centralized in the figure), whereby the M users are centrally controlled and coordinated in their learning using UCB, and allowing simultaneous selection of the same options by multiple users.

Next we show the performance of U_PART. The simulation setting is the same as the one above and is not repeated here. From Fig. 2a we see that U_PART outperforms multiuser UCB with individual learning. We also see from

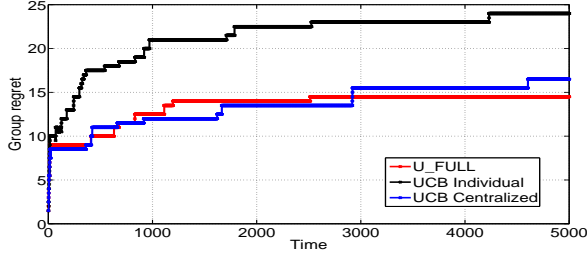
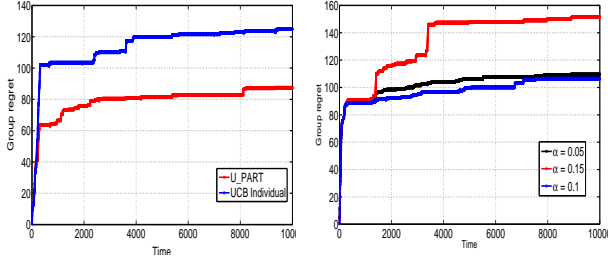


Fig. 1: Compare U_FULL and UCB individual.



(a) Performance comparison

(b) Different α

Fig. 2: Compare U_PART and UCB Individual.

Fig. 2b that though a larger α results in a larger upper bound, the actual performance does not necessarily increase with α .

We end this section by simulating a network with diverse group preferences. As we mentioned in previous sections, the major difference between learning algorithms of diverse preferences and uniform preference is each user estimates other users' group identity before taking actions over observed/reported samples. Therefore instead of presenting similar regrets results as in the previous cases, we present the mis-classification rate of our algorithm, given in Fig. 3.

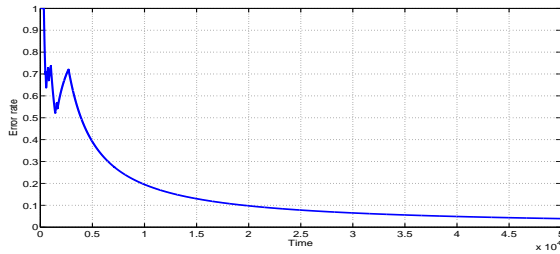


Fig. 3: Performance of error rate

VI. CONCLUSION

In this paper we considered group learning problem in the context of opinion diffusion and analyzed two scenarios: uniform group preference vs. diverse group preferences. For each case we also considered sharing full vs. partial information, and constructed UCB-like index based group learning algorithms and derived their associated upper bounds on weak regret. These upper bounds are in general better than

the original upper bound obtained by UCB when a user learns in isolation. This points to the potential gain by combining first-hand and second hand learning.

REFERENCES

- [1] V. Anantharam, P. Varaiya, and J. Walrand. Asymptotically efficient allocation rules for the multiarmed bandit problem with multiple plays-part i: I.i.d. rewards. *Automatic Control, IEEE Transactions on*, 32(11):968 – 976, nov 1987.
- [2] Venkat Anantharam, Pravin Varaiya, and Jean Walrand. Asymptotically Efficient Allocation Rules for the Multiarmed Bandit Problem with Multiple Plays Part I: I.I.D. Rewards, Part II: Markovian Rewards. Technical Report UCB/ERL M86/62, EECS Department, University of California, Berkeley, 1986.
- [3] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time Analysis of the Multiarmed Bandit Problem. *Mach. Learn.*, 47:235–256, May 2002.
- [4] Chih chun Wang, Student Member, Sanjeev R. Kulkarni, and H. Vincent Poor. Bandit problems with side observations. *IEEE Transactions on Automatic Control*, 50:338–355, 2005.
- [5] T. L. Lai and H. Robbins. Asymptotically Efficient Adaptive Allocation Rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- [6] John Langford and Tong Zhang. The Epoch-Greedy Algorithm for Multi-armed Bandits with Side Information. In *NIPS*, 2007.
- [7] Tyler Lu, Dvid Pl, and Martin Pal. Contextual multi-armed bandits. *Journal of Machine Learning Research*, 9:485–492, 2010.

APPENDIX A

PROOF OF THEOREM 3.1:

We follow the notation used in [3] where UCB is first introduced and analyzed, and bound the number of times sub-optimal arms are played. Consider the total number of times that option j has been used by user i up to time t :

$$\begin{aligned}
 T_j^i(t) &\leq l + \sum_{n=l+1}^t \{I_n^i = j, T_j^i(n-1) \geq l\} \\
 &\leq l + \sum_{n=l+1}^t \left\{ \min_{0 < s^i < n} r_{s^i}^{i,*} + c_{n-1, s^*}^i \leq \max_{l \leq s_j^i < n} r_{j, s_j^i}^i + c_{n-1, s_j}^i \right\} \\
 &\leq l + \sum_{n=l+1}^t \left\{ \min_{0 < s^i < n} r_{s^i}^{i,*} + c_{n-1, s^*}^i \leq \max_{l \leq s_j^i < n} r_{j, s_j^i}^i + c_{n-1, s_j}^i \right\} \\
 &\leq l + \sum_{n=1}^{\infty} \sum_{s^i=1}^{n-1} \sum_{s_j^i=l}^{n-1} \{r_{s^i}^{i,*} + c_{n, s^*}^i \leq r_{j, s_j^i}^i + c_{n, s_j}^i\}. \quad (23)
 \end{aligned}$$

Here we use $r_{j, s_j^i}^i$ to denote the estimated sample mean (as defined in Eqn. 5 with more details) of option j for user i with s_j^i local plays under our group learning algorithm. And $c_{t, j}^i$ is the time bias as in the classical UCB defined as following,

$$c_{t, s_j}^i = \sqrt{\frac{2 \log t}{\sum_{u \in \mathcal{U}} s_j^u}} \quad (24)$$

with s_j^u being number of plays of option j from user u .

Observe that $r_{s^i}^{i,*} + c_{n, s^*}^i \leq r_{j, s_j^i}^i + c_{n, s_j}^i$ implies that at least one of the following must hold:

$$r_{s^i}^{i,*} \leq \mu_j^{i,*} - c_{n, s^*}^i, r_{j, s_j^i}^i \geq \mu_j^i + c_{n, s_j}^i, \mu_j^{i,*} \leq \mu_j^i + 2c_{n, s_j}^i. \quad (25)$$

We now bound each term. We proceed by examining two cases, depending on whether the following assumption is true.

Assumption 1: For any user i and option j , under the proposed index policy

$$E[T_j^i(t)] \geq \mathcal{O}(\log t), j \in \bar{N}_K. \quad (26)$$

$$E[T_j^i(t)] = \mathcal{O}(t), j \in N_K. \quad (27)$$

Remark A.1: It has been shown in the literature that for a standard stochastic multi-armed bandit problem the above must hold for a sub-linear regret policy and that $\mathcal{O}(\log t)$ is the lower bound on weak regret. Thus under our index policy U_FULL we know that for $j \in \bar{N}_K$ we must have $E[T_j^i(t)] \geq \mathcal{O}(\log t)$. Later we will show that indeed under our policy we have $E[T_j^i(t)] = \mathcal{O}(\log t)$ and $E[T_j^i(t)] = \mathcal{O}(t)$ for $j \in N_K$.

Consider now the case that the above assumption is true. And we have the following lemma.

Lemma A.2: $\forall \epsilon > 0$,

$$P(|\tilde{\delta}_j^{i,k}(t) - \delta_j^{i,k}| > \epsilon) \leq 1/t^{d_U} \quad (28)$$

with d_U being some finite positive constant.

Proof 6: For simplicity of presentation, in this proof we omit all sub and super-scripts when there is no confusion; and further denote by $\delta(t)$ the estimate at time t and δ^* the true value.

$$P(|\delta(t) - \delta^*| > \epsilon) = P(\delta(t) > \delta^* + \epsilon) + P(\delta(t) < \delta^* - \epsilon)$$

Consider $P(\delta(t) > \delta^* + \epsilon)$ and denote $c = \delta^* + \epsilon$.

$$P(\delta(t) > c) = P\left(\frac{r_j^i(t)}{r_j^k(t)} > c\right) = P(r_j^i(t) > c \cdot r_j^k(t)).$$

$$\begin{aligned} P(r_j^i(t) > c \cdot r_j^k(t)) &= \int_x P(r_j^i(t) > c \cdot x) \cdot P(r_j^k(t) = x) dx \\ &= \int_{x \leq \mu_j^k - \epsilon} P(r_j^i(t) > c \cdot x) P(r_j^k(t) = x) dx \\ &+ \int_{\mu_j^k - \epsilon < x < \mu_j^k + \epsilon} P(r_j^i(t) > c \cdot x) \cdot P(r_j^k(t) = x) dx \\ &+ \int_{x \geq \mu_j^k + \epsilon} P(r_j^i(t) > c \cdot x) \cdot P(r_j^k(t) = x) dx. \end{aligned} \quad (29)$$

where for simplicity we have used $P()$ to denote the density function.

Note that for $x \notin (\mu_j^k - \epsilon, \mu_j^k + \epsilon)$, we have $P(r_j^k(t) = x) \leq e^{-2\epsilon^2 T_j^k(t)}$. Therefore we have

$$\begin{aligned} &\int_{x \leq \mu_j^k - \epsilon} P(r_j^i(t) > c \cdot x) \cdot P(r_j^k(t) = x) dx \\ &\leq \int_{x \leq \mu_j^k - \epsilon} P(r_j^i(t) > c \cdot x) \cdot e^{-2\epsilon^2 T_j^k(t)} dx \\ &\leq e^{-2\epsilon^2 T_j^k(t)} \int_{x \leq \mu_j^k - \epsilon} 1 dx = \mathcal{O}(e^{-T_j^k(t)}) \end{aligned} \quad (30)$$

$$\begin{aligned} &\int_{\mu_j^k - \epsilon \leq x \leq \mu_j^k + \epsilon} P(r_j^i(t) > c \cdot x) \cdot P(r_j^k(t) = x) dx \\ &\leq \int_{\mu_j^k - \epsilon \leq x \leq \mu_j^k + \epsilon} P(r_j^i(t) > c \cdot x) dx \\ &\leq \int_{\mu_j^k - \epsilon \leq x \leq \mu_j^k + \epsilon} e^{-2\epsilon^2 T_j^i(t)} dx = \mathcal{O}(e^{-T_j^i(t)}) \end{aligned} \quad (31)$$

$$\begin{aligned} &\int_{x \geq \mu_j^k + \epsilon} P(r_j^i(t) > c \cdot x) \cdot P(r_j^k(t) = x) dx \\ &\leq e^{-2\epsilon^2 T_j^i(t)} \int_{x \geq \mu_j^k + \epsilon} P(r_j^k(t) = x) dx \\ &\leq e^{-2\epsilon^2 T_j^i(t)} P(r_j^k(t) \geq \mu_j^k + \epsilon) \\ &\leq e^{-2\epsilon^2 T_j^i(t)} \cdot e^{-2\epsilon^2 T_j^k(t)} = \mathcal{O}(e^{-(T_j^i(t) + T_j^k(t))}) \end{aligned} \quad (32)$$

As $E[T_j^k(t)] \geq \mathcal{O}(\log t)$ we have

$$\mathcal{O}(E[e^{-T_j^k(t)}]) \leq 1/t^d \quad (33)$$

here d is some finite positive number. Other terms can be similarly analyzed, proving the lemma. ■

With the help of the above Lemma we have the following bound (for user k)

$$\sum_{n=1}^{s_j^k} \frac{1}{n^{d_U}} \leq \frac{(s_j^k)^{1-d_U} - 1}{1 - d_U}. \quad (34)$$

And we have the distortion factor in the sample mean bounded by $C \cdot \frac{(s_j^k)^{1-d_U}}{\sum_{u \in \mathcal{U}} s_j^u}$. Next we show

$$\sqrt{\frac{\log n}{\sum_{u \in \mathcal{U}} s_j^u}} > C \cdot \frac{(s_j^k)^{1-d_U}}{\sum_{u \in \mathcal{U}} s_j^u}. \quad (35)$$

For $s_j^i \leq \mathcal{O}(\log n)$, since $d_U > 0$ we have $(\sum_{u \in \mathcal{U}} s_j^u)^{1/2} > C \cdot \frac{(s_j^k)^{1-d_U}}{\sqrt{\log n}}$. If $s_j^i > \mathcal{O}(\log n)$, consider two cases. If $s_j^k = \mathcal{O}(\log n)$, the above holds obviously. If $s_j^k > \mathcal{O}(\log n)$, through the proof of Lemma A.2 we know $d_U \geq 1$ (details omitted) and again we have (35) hold. Therefore we have

$$\begin{aligned} &P\{r_{s_i}^{i,*} \leq \mu^{i,*} - \sqrt{2} \sqrt{\frac{\log n}{\sum_{u \in \mathcal{U}} s_j^u}} \pm \frac{C' \cdot \sum_{k \in \mathcal{U}} (s_j^k)^{1-d_U}}{\sum_{u \in \mathcal{U}} s_j^u}\} \\ &\approx P\{r_{s_i}^{i,*} \leq \mu^{i,*} - \sqrt{2} \sqrt{\frac{\log n}{\sum_{u \in \mathcal{U}} s_j^u}}\} \\ &\leq e^{-4 \log n} = n^{-4} \end{aligned} \quad (36)$$

And similarly

$$P\{r_{j,s_j^i}^i \geq \mu_j^i + c_{n,s_j}^i\} \leq n^{-4}. \quad (37)$$

For the last term $\mu^{i,*} \leq \mu_j^i + 2c_{n,s_j}^i$, let $l = \lceil \frac{8}{M} \log t / (\Delta_j^i)^2 \rceil$

we have

$$\begin{aligned}
& \mu^{i,*} - \mu_j^i - 2c_{n,s_j}^i \\
& \geq \mu^{i,*} - \mu_j^i - 2\sqrt{2} \sqrt{\log t / \sum_{u \in \mathcal{U}} s_j^u} \\
& \geq \mu^{i,*} - \mu_j^i - \Delta_j^i = 0
\end{aligned} \tag{38}$$

The rest of the proof regarding weak performance bound follows [3] and is thus omitted.

Now consider the case where the assumption does not hold. As mentioned earlier we must have $E[T_j(t)] \geq \mathcal{O}(\log t)$ for $j \in \bar{N}_K$. Therefore the only possibility for the assumption to not hold is $E[T_j(t)] < \mathcal{O}(t)$ for $j \in N_K$.

Suppose there is an optimal arm N^* that is sample at order less than $\mathcal{O}(t)$. Then there must be a sub-optimal arm $j \in \bar{N}_K$ that is sampled with order at least $\mathcal{O}(t)$.

According to previous proof under the first case (when the assumption holds), we have $E[T_j^i(t)] \leq \sum_{j \in \bar{N}_K} \lceil \frac{8 \log t}{M \cdot (\Delta_j^i)^2} \rceil$ whenever $d_U > 0$. Therefore we know that for this assumption to not hold we must have $d_U = 0$. However we see from the proof of Lemma A.2 that the number of plays of the corresponding optimal arms must be bounded as constant $\mathcal{O}(1)$. Specifically, recall we have

$$\begin{aligned}
& P(|\tilde{\delta}_{N^*}^{i,k}(t) - \delta_{N^*}^{i,k}| > \epsilon) \\
& \leq \max\{\mathcal{O}(e^{-T_{N^*}^i(t)}), \mathcal{O}(e^{-T_{N^*}^k(t)})\}
\end{aligned} \tag{39}$$

We have $d_U > 0$ (at time t) except for the case

$$\min\{T_{N^*}^i(t), T_{N^*}^k(t)\} = \mathcal{O}(1),$$

i.e., the user plays arm N^* only constant number of times.

To simplify the analysis, we will additionally assume that for each option each user will use no more samples (regarding order) from other users than it has locally. Now we check U_FULL index. For the optimal arm we have $\sqrt{\frac{\log t}{\mathcal{O}(1)}}$ as the bias term in the index while for j it is $\sqrt{\frac{\log t}{\mathcal{O}(t)}}$. From above argument we know with a sufficiently large t , the index of the optimal arm must be larger than j which contradicts the fact that the optimal arm is only sampled a constant number of time, proving that this second case (the assumption does not hold) cannot be true. The theorem is thus proved.

APPENDIX B PROOF OF THEOREM 3.2:

In this proof for simplicity we denote

$$\hat{c}_{t,s_j}^i = \sqrt{\frac{2 \log t}{s_j}}, \tag{40}$$

$$\tilde{c}_{t,s_j}^i = \sqrt{\frac{2 \log t}{s_j}} - \alpha[1 - \beta_j(t)] \sqrt{\frac{\log t}{t}} \tag{41}$$

Following the same method introduced in Auer's work [3], to bound the regret we need to bound the number of times

the sub-optimal arms are played. Suppose $j \in \bar{N}_K$ we have

$$\begin{aligned}
T_j^i(t) & \leq l + \sum_{n=l+1}^t \{I_n^i = j, T_j^i(n-1) \geq l\} \\
& \leq l + \sum_{n=l+1}^t \left\{ \min_{0 < s^i < n} r_{s^i}^{i,*} + \hat{c}_{n-1,s^i}^i \leq \max_{l \leq s_j^i < n} r_{j,s_j^i}^i + \hat{c}_{n-1,s_j^i}^i \right\} \\
& \leq l + \sum_{n=1}^{\infty} \sum_{s^i=1}^{n-1} \sum_{s_j^i=l}^{n-1} \{r_{s^i}^{i,*} + \hat{c}_{n,s^i}^i \leq r_{j,s_j^i}^i + \hat{c}_{n,s_j^i}^i\}
\end{aligned}$$

The following analysis applies to any i and thus we omit the i in the sub and super-scripts. Observe that $r_s^* + \hat{c}_{n,s^*} \leq r_{j,s_j} + c_{n,s_j}$ implies that at least one of the following must hold,

$$r_s^* \leq \mu^* - \hat{c}_{n,s^*}, r_{j,s_j} \geq \mu_j + \hat{c}_{n,s_j}, \mu^* \leq \mu_j + 2\hat{c}_{n,s_j} \tag{42}$$

We bound each term as follows.

$$\begin{aligned}
P\{r_s^* \leq \mu^* - \hat{c}_{n,s^*}\} & \leq P\{r_s^* \leq \mu^* - (\sqrt{2} - \alpha) \sqrt{\frac{\log n}{n}}\} \\
& \leq e^{-2(\sqrt{2}-\alpha)^2 \log n} = n^{-2(\sqrt{2}-\alpha)^2}
\end{aligned} \tag{43}$$

And similarly $P\{r_{j,s_j} \geq \mu_j + \hat{c}_{n,s_j}\} \leq n^{-2(\sqrt{2}-\alpha)^2}$. Let $l = \lceil (2\sqrt{2} - \alpha[1 - \beta_j(t)])^2 \log t / \Delta_j^2 \rceil$

$$\mu^* - \mu_j - 2\hat{c}_{n,s_j} \geq \mu^* - \mu_j - \Delta_j = 0 \tag{44}$$

Therefore

$$\begin{aligned}
E[T_j(t)] & \leq \lceil \frac{(2(\sqrt{2} - \alpha \cdot [1 - \beta_j(t)]))^2 \log t}{\Delta_j^2} \rceil \\
& + \sum_{n=1}^{\infty} \sum_{s=1}^{n-1} \sum_{\substack{s_j = \lceil (2\sqrt{2} - \\ \alpha \cdot (1 - \beta_j(t)))^2 \log t / \Delta_j^2 \rceil}}^{t-1} (P\{r_s^* \leq \mu^* - \hat{c}_{n,s^*}\}) \\
& + P\{r_{j,s_j} \geq \mu_j + \hat{c}_{n,s_j}\} \\
& \leq \lceil \frac{4(\sqrt{2} - \alpha \cdot [1 - \beta_j(t)])^2 \log t}{\Delta_j^2} \rceil \\
& + \sum_{n=1}^{\infty} \sum_{s=1}^n \sum_{s_j=1}^n 2n^{-2(\sqrt{2}-\alpha)^2}
\end{aligned} \tag{45}$$

First consider the second term. If $-2(\sqrt{2} - \alpha)^2 < -3$, i.e., $\sqrt{2} - \sqrt{3}/2 > \alpha$, the sum converges to a constant, i.e.,

$$E[T_j(t)] \leq \lceil \frac{4(\sqrt{2} - \alpha \cdot [1 - \beta_j(t)])^2 \log t}{\Delta_j^2} \rceil + \text{const.} \tag{46}$$

Next we bound $\beta_j(t), \forall j \in \bar{N}_K$. Recall that $\beta_j(t) = \frac{n_j(t)}{\sum_k n_k(t)}$ and we know $n_j(t) \leq \lceil \frac{8 \log t}{\Delta_j^2} \rceil, \forall j \in \bar{N}_K$. As there are $M \cdot K$ observations from the group at each time, we have $\sum_k n_k(t) = M \cdot K \cdot t$ and therefore

$$\beta_j(t) \leq \frac{M \cdot \lceil \frac{8 \log t}{\Delta_j^2} \rceil}{M \cdot K \cdot t} \leq \mathcal{O}\left(\frac{\log t}{t}\right). \tag{47}$$