

MULTI-ARMED BANDITS IN MULTI-AGENT NETWORKS

Shahin Shahrampour*, Alexander Rakhlin†, and Ali Jadbabaie‡

*Department of Electrical Engineering, Harvard University

†Department of Statistics at the Wharton School, University of Pennsylvania

‡Institute for Data, Systems, and Society, Massachusetts Institute of Technology

ABSTRACT

This paper addresses the multi-armed bandit problem in a multi-player framework. Players explore a finite set of arms with stochastic rewards, and the reward distribution of each arm is player-dependent. The goal is to find the best global arm, i.e., the one with the largest expected reward when averaged out among players. To achieve this goal, we develop a distributed variant of the well-known UCB1 algorithm. Confined to a network structure, players exchange information locally to estimate the global rewards, while using a confidence bound relying on the network characteristics. Then, at each round, each player votes for an arm, and the majority vote is played as the network action. The whole network gains the reward of the network action, hoping to maximize the global welfare. The performance of the algorithm is measured via the notion of network regret. We prove that the regret scales logarithmically with respect to time horizon and inversely in the spectral gap of the network. Our algorithm is optimal in the sense that in a complete network it scales down the regret of its single-player counterpart by the network size. We demonstrate numerical experiments to verify our theoretical results.

Index Terms— Sequential decision-making, multi-armed bandits, multi-agent networks, distributed learning.

1. INTRODUCTION

The multi-armed bandit (MAB) problem has been extensively studied in the literature [1–6]. In its classical setting, the problem is defined by a set of *arms* or *actions*, and it captures the exploration-exploitation dilemma for a *learner*. At each time step, the learner chooses an arm and receives its corresponding *payoff* or *reward*. The term *bandit* indicates that only the reward of the chosen arm is revealed to the learner, while the rewards of other arms remain undisclosed at that particular time. The learner then hopes to maximize the total payoff obtained from sequentially selecting the arms. Equivalently, the learner aims to minimize the *regret* by competing with the best single arm in hindsight. The reward model for arms could be stochastic or non-stochastic, and optimal algorithms for both cases are proposed in the literature [6]. While early studies on MAB dates back to nine decades ago, the problem has received considerable attention due to its modern applications. MAB could be an instance of sequential decision-making for ad placement, website optimization, packet routing, cognitive compressive sensing, and etc. [6–9].

In this paper, we depart from the classical setting and address the stochastic MAB in a multi-player network. We consider a scenario where a group of *players* or *agents* collaborate to achieve a team

task. In this framework, each arm may possess different reward distributions for distinct players. The goal is to reach consensus on the arm which best fits the network. We consider this arm to be the arm with the largest expected reward when averaged out among players. As a motivating example, consider a number of brands (arms) selling a product. We have several users (players) in a social network who need to decide on one specific brand for the whole network. The users rate these brands differently and must agree upon using one brand. The decision is made by the majority vote to maximize the global welfare.

Therefore, in our setup, each arm has a true *global* payoff that can be written in terms of an average of *individual* rewards. Agents are not able to identify the best global arm since they do not receive any informative signal about the global reward. Therefore, they need to benefit from side observations gained from local communication. The model has a flavor of distributed detection algorithms where the parameter of interest is not fully observable to an individual agent [10, 11]. However, there is an additional restriction in this work, as we consider a bandit setup where the player only receives the signal of the *chosen* action.

To find the best global arm, players wish to maximize a global objective. As standard in online algorithms, we translate this objective to minimizing a network *regret*. We then propose an algorithm dubbed Distributed Upper Estimated Reward (d-UER) to minimize the network regret. The algorithm estimates the global reward of the arms in a purely *distributed* fashion. To concentrate around the true value of the global rewards with high probability, the algorithm exploits a *confidence bound* that relies on the network topology. Then, players iteratively vote for their favorite arm, and each time the majority vote is played as the network action. The whole network scores the reward of the network action, hoping to maximize the global welfare. This algorithm can be viewed as a distributed version of the well-known UCB1 algorithm.

We prove that the regret of our algorithm scales logarithmically with respect to time horizon. It also scales inversely in the gap between the arms. The inverse dependence to gap is quite natural since similar arms to the best arm are hard to distinguish. Furthermore, the regret relies on the *spectral gap* of the network as in the exploration phase information needs to be propagated throughout the network. Our algorithm is optimal in the sense that in a complete network the regret is scaled down by the network size comparing to its single-player analog. This is an artifact of variance reduction when we distribute samples between agents throughout the network. We finally provide numerical experiments to verify the impact of network size and spectral gap in practice.

Related Work: Several variants of decentralized MAB have been studied in the literature. In [12–14], decentralized MAB has been

The authors gratefully acknowledge the support of ONR BRC Program on Decentralized, Online Optimization, NSF under grants CAREER DMS-0954737 and CCF-1116928, as well as Dean’s Research Fund.

formulated for applications in cognitive radio networks and multi-channel communication systems. Unlike our setup, in these works simultaneous selection of one arm by a few players is not recommended and reduces the reward in some sense. Some other works focused on decentralized MAB with application to advertising systems. In [15], the interaction between users in a social network provides information for an external, centralized decision-maker. In [16] only a single *major* agent in the network has access to its reward sequence, while other agents are aware of the sampling pattern of the major agent. The works of [17, 18] are particularly relevant to our work; however, the arms in those settings are player independent, whereas in our setup they are player dependent. Our work is also related to distributed detection and learning under *full information* setting [11, 19–23]. In these models, the world is governed by a fixed true *state* (arm), aimed to be recovered collaboratively. However, agents receive information about all states per round, whereas our *bandit* setup entails only one-state feedback per round.

Notation

$[n]$	The set $\{1, 2, \dots, n\}$ for any integer n
$\mathbf{1}\{\cdot\}$	The indicator function
$\mathbf{1}$	The vector of all ones
$\mathbb{E}[\cdot]$	The expectation operator
x^\top	Transpose of the vector x
$x(k)$	The k -th element of vector x
$\sigma_i(W)$	The i -th largest singular value of matrix W

2. PROBLEM FORMULATION

We have a multi-agent network with N *players* or *agents* sequentially selecting *arms* or *actions*. The number of arms is a finite number K , a common knowledge in the network. Pulling arm $k \in [K]$ at time $t \in [T]$ rewards the player $i \in [N]$ with a random variable $X_{i,t}(k) \in [0, 1]$. The rewards are independent and identically distributed over time for an individual player. We further assume that they are independent across players and arms. Hence, the expected value $\mu_i = \mathbb{E}[X_{i,t}] \in \mathbb{R}^K$ is fixed over time though the characteristics of each arm is different among players. That is, for arm k and players i, j , the values $\mu_i(k)$ and $\mu_j(k)$ are not equal in general. The *true* reward of arm $k \in [K]$ is the network average as follows

$$\mu(k) := \frac{1}{N} \sum_{i=1}^N \mu_i(k) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[X_{i,t}(k)].$$

Agents aim to maximize the *global* welfare in the sense of finding the arm k^* such that $\mu(k) \leq \mu(k^*)$ for all $k \in [K]$. However, no individual player can make a correct inference by simply collecting individual signals $\{X_{i,t}\}_{t=1}^T$ as these only approximate μ_i over time. Therefore, players must collaborate with each other to estimate the true rewards in a distributed fashion. Based on a distributed protocol discussed in the next section, each agent votes for an arm at each round. We represent by the random variable $I_{i,t}$, the arm that is chosen by player i at time t . The network action I_t at time t is then the majority vote, i.e., the random variable I_t is defined as

$$I_t := \text{the most repeated element of the set } \{I_{i,t}\}_{i=1}^N.$$

In consistent with the bandit setting, after this selection is made, player i only observes the corresponding payoff $X_{i,t}(I_t)$ at period t .

Common to the MAB framework, let us reformulate the problem in terms of *regret*. Without loss of generality, we assume the following order for the true rewards,

$$\mu(1) \geq \mu(2) \geq \dots \geq \mu(K),$$

for the rest of the paper. For any pair $k \leq m$, we define the *gap* as

$$\Delta_{k,m} := \mu(k) - \mu(m) = \frac{1}{N} \sum_{i=1}^N \mu_i(k) - \mu_i(m),$$

to capture the suboptimality of arm m comparing to arm k . We use $n_t(k)$ to denote the number of times that arm k has been chosen by the network as the majority vote until time t . It is easy to observe that maximizing the global welfare is equivalent to minimizing the network *regret* in the following sense,

$$\mathbf{R}_T := T\mu(1) - \sum_{t=1}^T \mathbb{E}[\mu(I_t)] = \sum_{k=2}^K \Delta_{1,k} \mathbb{E}[n_T(k)], \quad (1)$$

where the expectation is taken over the randomness in the choice of arms.

In Section 3, we propose an online, distributed algorithm to solve (1). We prove that the algorithm incurs a strongly sub-linear regret with respect to time. Furthermore, we show that the regret depends on the network characteristics since players communicate with each other to understand the true value of the rewards. In other words, the information dissemination over the network requires a time which appears as a network penalty in the regret. This is in contrast with the classical (one-player) MAB, where the player only investigates the arms in the exploration phase. In our setup (multi-player MAB), an extra information-exchange time is also required in the exploration period to estimate the true rewards.

Network Structure: One individual player is not able to track the best arm in isolation as the signals $\{X_{i,t}\}_{t=1}^T$ are not informative enough for approximation of μ . Therefore, agents communicate with each other iteratively to approximate the true rewards. We use the symmetric and doubly stochastic matrix W to capture the interaction between agents. This matrix has positive diagonals, and a positive $[W]_{ij} > 0$ means that player i assigns a weight $[W]_{ij} = [W]_{ji}$ to observations of player $j \neq i$. When $[W]_{ij} = 0$, agents i and j never communicate with each other directly. Therefore, we have

$$\sum_{j \in \mathcal{N}_i} [W]_{ij} = \sum_{j=1}^N [W]_{ij} = \sum_{i=1}^N [W]_{ij} = \sum_{i \in \mathcal{N}_j} [W]_{ij} = 1,$$

where $\mathcal{N}_i := \{j \in [N] : [W]_{ij} > 0\}$ is the *local* neighborhood of agent i . We assume that the underlying network is *connected*, i.e., there exists a *path* from any player $i \in [N]$ to any player $j \in [N]$. This assumption guarantees the information flow over the network.

3. DISTRIBUTED UPPER ESTIMATED REWARD

We now describe the d-UEER algorithm to minimize regret in (1). The algorithm can be cast as a cooperative version of the well-known UCB1 [3]. As we discussed in Section 2, the feedback setup does not allow an individual player to solely identify the best arm. Therefore, players need to cooperate with each other to collectively explore the arms.

Algorithm 1 Distributed Upper Estimated Reward

Input : Number of agents N and arms K . Parameter $d > 0$.

Initialization : Each action is played once, and the reward of action $k \in [K]$ for player $i \in [N]$ is stored in $\psi_{i,1}(k)$. For each $i \in [N]$ and $k \in [K]$, let $\phi_{i,0}(k) = 0$ and $n_{i,0}(k) = 1$, respectively.

for $t = 1$ to T **do**

for $i = 1$ to N **do**

$$\phi_{i,t} = \sum_{j=1}^N [W]_{ij} \phi_{j,t-1} + \psi_{i,t}. \quad \% \text{ estimation of true rewards}$$

$$C_t(k) = \sqrt{2 \log t \left(\frac{1}{N n_{t-1}(k)} + \frac{2d}{n_{t-1}^2(k)} \right)}. \quad \% \text{ upper confidence bound for each action}$$

$$I_{i,t} = \operatorname{argmax}_{k \in [K]} \left\{ \frac{\phi_{i,t}(k)}{n_{t-1}(k)} + C_t(k) \right\}. \quad \% \text{ agent's action}$$

end for

 Let I_t be the most frequent element of the set $\{I_{i,t}\}_{i=1}^N$. % network's action or majority vote

 For any $k \in [K]$, update the counter as $n_t(k) = n_{t-1}(k) + \mathbf{1}\{k = I_t\}$.

 The network scores $\mu(I_t)$, and player $i \in [N]$ observes $X_{i,t}(I_t)$. % private observation

 Let $\psi_{i,t+1}(k) = X_{i,t}(k) \mathbf{1}\{k = I_t\}$ for any $k \in [K]$.

end for

The d-UER algorithm is summarized in the table above. It provides a completely decentralized method to estimate the true rewards. Players accumulate local observations, and they take into account an upper confidence bound to make individual decisions. The term $\phi_{i,t}$ in the algorithm (normalized by the number of times each arm has been played) is an estimator of the true rewards. The confidence bound C_t allows agents to concentrate around the true value of the rewards with high probability. Then, agent i makes the individual decision $I_{i,t}$ at time t , the majority vote I_t among agents is chosen as the network action, and agents observe the corresponding payoff of the action chosen by the majority vote.

Note that since the setup is multi-player, d-UER exploits a confidence bound relying on the network structure through parameter d . We will see that this parameter must be tuned as an upper bound on a quantity that depends on network size and spectral gap.

The following lemma provides a closed-form solution for $\{\phi_{i,t}\}_{t=1}^T$ and shows the importance of the mixture behavior of the Markov chain W in estimation.

Lemma 1. Any update of the form $\phi_{i,t} = \sum_{j=1}^N [W]_{ij} \phi_{j,t-1} + \psi_{i,t}$ can be expressed as,

$$\phi_{i,t} = \sum_{\tau=1}^t \sum_{j=1}^n [W^{t-\tau}]_{ij} \psi_{j,\tau},$$

whenever the update is initialized at $\phi_{i,0}(k) = 0$, for any $i \in [N]$ and $k \in [K]$. Also, given the connectivity of the network, the doubly stochastic matrix W with positive diagonal satisfies

$$\sum_{\tau=1}^t \sum_{j=1}^N \left| [W^{t-\tau}]_{ij} - \frac{1}{N} \right| \leq dE_1,$$

for any $i \in [N]$, where

$$dE_1 := \frac{2}{1 - \sigma_2(W)} + \frac{\log N}{\log [\sigma_2(W)^{-1}]},$$

and $\sigma_2(W) < 1$ is the second largest singular value of W .

Proof. See the Appendix in [24]. □

The lemma suggests that the update $\phi_{i,t}$ accumulates new information from the environment and averages out the past. It is important to have network connectivity, since it allows $W^t \rightarrow \frac{1}{N} \mathbf{1}\mathbf{1}^\top$ as $t \rightarrow \infty$. Indeed, when the underlying network topology is disconnected, information cannot spread in the whole network. Therefore, players cannot observe some of informative signals dispersed throughout the network. In this case, they are not able to make a correct inference about the true reward of the arms, resulting in the network regret increasing linearly in time.

Finally, the lemma implies that the performance of the algorithm relies on the mixture properties of the Markov chain W . This is proved by the dependence of the quantity dE_1 to $\sigma_2(W)$.

Theorem 2. Let for each $k \in [K]$ and $i \in [N]$, the sequence $\{X_{i,t}(k)\}_{t=1}^T$ be i.i.d. samples from a stationary distribution with $\mu_i(k) = \mathbb{E}[X_{i,t}(k)]$. Let also the independence over arms and agents hold such that

$$\mathbb{E}[X_{i,t}(k)X_{j,t}(k')] = \mathbb{E}[X_{i,t}(k)]\mathbb{E}[X_{j,t}(k')],$$

for $i \neq j$ or $k \neq k'$. Given the connectivity of the network, the regret of d-UER algorithm, defined in (1), satisfies the following bound

$$\begin{aligned} \mathbf{R}_T &\leq \sum_{k=2}^K 4 \max \left\{ \frac{12 \log T}{N \Delta_{1,k}}, Nd \right\} \\ &\quad + K \sum_{k=2}^K \left(2.5 \left(1 + \log \left[\frac{4}{\Delta_{1,k}} \right] \right) dE_1 dE_2 + \frac{2\pi^2}{3} \Delta_{1,k} \right), \end{aligned}$$

whenever $d \geq dE_1$, where

$$dE_1 := \frac{2}{1 - \sigma_2(W)} + \frac{\log N}{\log [\sigma_2(W)^{-1}]},$$

and

$$dE_2 := \frac{\log N}{\log [\sigma_2(W)^{-1}]}.$$

Proof. See the Appendix in [24]. □

Theorem 2 indicates that the regret depends on the network size as well as the second largest singular value of W . Defining the spectral gap of the network as

$$\gamma(W) := 1 - \sigma_2(W),$$

the theorem further shows that the regret bound scales inversely in the spectral gap of the network. Note that regret serves as a *non-asymptotic* performance metric, reiterating the importance of the spectral gap in the finite-time analysis.

The local feedback does not provide each player with adequate information, yielding a delay in proper decision-making. For instance, in cycle and path networks where the diameter is $\mathcal{O}(N)$ the incurred penalty $\mathbf{dE}_1 = \mathcal{O}(N^2 \log N)$ is large, whereas in a complete network $W = \frac{1}{N} \mathbf{1}\mathbf{1}^\top$ the Markov chain is mixed from the outset, and there is no network penalty. The latter can be seen as N copies of a single-player MAB where $\sigma_2(W) = 0$. In this case, the network errors become $\mathbf{dE}_1 = 2$ and $\mathbf{dE}_2 = 0$, and the well-known result of [3] for UCB1 algorithm is recovered (scaled down by a factor of N). The $\frac{1}{N}$ factor is an advantage gained through reducing the variance of samples by distributing N samples among N individuals.

4. NUMERICAL EXPERIMENTS

Our theoretical results indicate that the size and spectral gap of the network are decisive in the performance of d-UER. In this section, we study the impact of these two factors via numerical experiments. In the first scenario, we consider networks of the same size differing in the spectral gap, while in the second scenario we consider complete networks of different sizes. For both cases, we plot the regret versus time horizon to investigate the performance.

For all of our experiments, we deal with $K = 4$ arms. We divide the agents into two groups, for whom the expected value of the rewards are as follows,

	$\mu_i(1)$	$\mu_i(2)$	$\mu_i(3)$	$\mu_i(4)$
$i \in \text{Group 1}$	0.0149	0.7161	0.7944	0.6749
$i \in \text{Group 2}$	0.5144	0.5108	0.9955	0.4778

generated randomly in the unit interval. This gives rise to the following global rewards

$\mu(1)$	$\mu(2)$	$\mu(3)$	$\mu(4)$
0.2646	0.6135	0.8950	0.5764

and the best global arm would be arm 3.

4.1. The Impact of Spectral Gap

For the first scenario, we fix the network size to $N = 50$. We would like to evaluate the performance of d-UER algorithm in three networks: complete, cycle and 4-regular (all with self-loops). For the reward values given in the tables, we run 100 experiments and average out regret over these runs. We then plot Fig. 1 which shows regret for these networks with respect to time ($T = 3000$).

As verified in theoretical results, the regret bound scales inversely with the spectral gap $\gamma(W) = 1 - \sigma_2(W)$. We can observe the impact in Fig. 1 where the networks are sorted correctly with respect to this metric. The complete network (largest spectral gap) has the best performance, while the 4-regular outperforms the cycle (due to its larger spectral gap). The result is also consistent with the intuition that the networks should be sorted according to their connectivity. The more the connectivity, the less the regret.

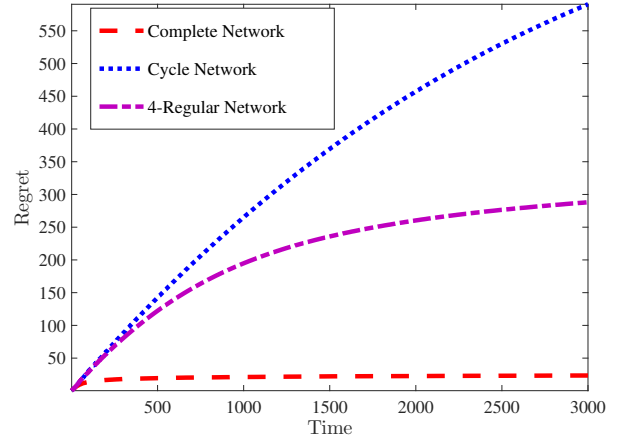


Fig. 1. Performance of d-UER in complete, cycle, and 4-regular networks. The regret scales inversely in the spectral gap. Hence, as the spectral gap grows larger, the regret becomes smaller.

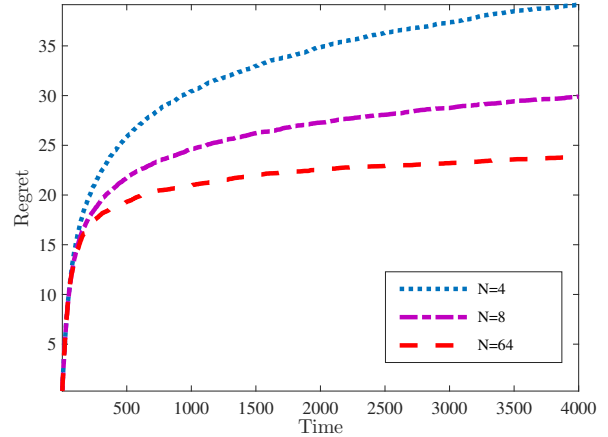


Fig. 2. Performance of d-UER in complete networks of size $N \in \{4, 8, 64\}$. As network grows larger, the estimation variance for agents decreases, and the regret becomes smaller.

4.2. The Impact of Network Size

Theorem 2 indicates that the regret scales logarithmically with time in the long run, and the dominant term has a $1/N$ factor. Therefore, the network size proves to be crucial in the performance. We choose the complete networks such that $W = \frac{1}{N} \mathbf{1}\mathbf{1}^\top$, for different values of $N \in \{4, 8, 64\}$ to investigate the performance of d-UER algorithm. We study complete networks to remove the effect of spectral gap and focus on size. In this case, the spectral gap always remains to be one and does not change with network size.

For the reward values given in the tables, we run 100 experiments and average out regret over these runs. Fig. 2 shows the regret for the three networks with respect to time ($T = 4000$). The simulation certifies that larger network size results in a lower regret.

5. REFERENCES

- [1] Tze Leung Lai and Herbert Robbins, “Asymptotically efficient adaptive allocation rules,” *Advances in applied mathematics*, vol. 6, no. 1, pp. 4–22, 1985.
- [2] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire, “The nonstochastic multiarmed bandit problem,” *SIAM Journal on Computing*, vol. 32, no. 1, pp. 48–77, 2002.
- [3] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer, “Finite-time analysis of the multiarmed bandit problem,” *Machine learning*, vol. 47, no. 2-3, pp. 235–256, 2002.
- [4] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári, “Exploration–exploitation tradeoff using variance estimates in multi-armed bandits,” *Theoretical Computer Science*, vol. 410, no. 19, pp. 1876–1902, 2009.
- [5] Kehao Wang and Lin Chen, “On optimality of myopic policy for restless multi-armed bandit problem: An axiomatic approach,” *IEEE Transactions on Signal Processing*, vol. 60, no. 1, pp. 300–309, 2012.
- [6] Sébastien Bubeck and Nicolo Cesa-Bianchi, “Regret analysis of stochastic and nonstochastic multi-armed bandit problems,” *Machine Learning*, vol. 5, no. 1, pp. 1–122, 2012.
- [7] Aditya Mahajan and Demosthenis Teneketzis, “Multi-armed bandit problems,” in *Foundations and Applications of Sensor Management*, pp. 121–151. Springer, 2008.
- [8] Sattar Vakili, Keqin Liu, and Qing Zhao, “Deterministic sequencing of exploration and exploitation for multi-armed bandit problems,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 7, no. 5, pp. 759–767, 2013.
- [9] Saeed Bagheri and Anna Scaglione, “The restless multi-armed bandit formulation of the cognitive compressive sensing problem,” *IEEE Transactions on Signal Processing*, vol. 63, no. 5, pp. 1183–1198, 2015.
- [10] Dušan Jakovetić, José MF Moura, and João Xavier, “Distributed detection over noisy networks: Large deviations analysis,” *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 4306–4320, 2012.
- [11] Shahin Shahrampour, Alexander Rakhlin, and Ali Jadbabaie, “Distributed detection : Finite-time analysis and impact of network topology,” *IEEE Transactions on Automatic Control*, vol. 61, 2016.
- [12] Keqin Liu and Qing Zhao, “Distributed learning in multi-armed bandit with multiple players,” *IEEE Transactions on Signal Processing*, vol. 58, no. 11, pp. 5667–5681, 2010.
- [13] Cem Tekin and Mingyan Liu, “Performance and convergence of multi-user online learning,” in *Game Theory for Networks*, pp. 321–336. Springer, 2012.
- [14] Dileep Kalathil, Naumaan Nayyar, and Rahul Jain, “Decentralized learning for multiplayer multiarmed bandits,” *IEEE Transactions on Information Theory*, vol. 60, no. 4, pp. 2331–2345, 2014.
- [15] Swapna Buccapatnam, Atila Eryilmaz, and Ness B Shroff, “Multi-armed bandits in the presence of side observations in social networks,” in *IEEE Conference on Decision and Control (CDC)*, 2013, pp. 7309–7314.
- [16] Soumya Kar, H Vincent Poor, and Shuguang Cui, “Bandit problems in networks: Asymptotically efficient distributed allocation rules,” in *IEEE Conference on Decision and Control and European Control Conference*, 2011, pp. 1771–1778.
- [17] Peter Landgren, Vaibhav Srivastava, and Naomi Ehrich Leonard, “On distributed cooperative decision-making in multi-armed bandits,” *arXiv preprint arXiv:1512.06888*, 2015.
- [18] Peter Landgren, Vaibhav Srivastava, and Naomi Ehrich Leonard, “Distributed cooperative decision-making in multi-armed bandits: Frequentist and bayesian algorithms,” *arXiv preprint arXiv:1606.00911*, 2016.
- [19] Shahin Shahrampour and Ali Jadbabaie, “Exponentially fast parameter estimation in networks using distributed dual averaging,” in *IEEE Conference on Decision and Control (CDC)*, 2013, pp. 6196–6201.
- [20] Anusha Lalitha, Anand Sarwate, and Tara Javidi, “Social learning and distributed hypothesis testing,” in *International Symposium on Information Theory (ISIT)*, 2014, pp. 551–555.
- [21] A. Nedic, A. Olshevsky, and C.A. Uribe, “Nonasymptotic convergence rates for cooperative learning over time-varying directed graphs,” in *American Control Conference (ACC)*, July 2015, pp. 5884–5889.
- [22] Liu Qipeng, Zhao Jiuhua, and Wang Xiaofan, “Distributed detection via bayesian updates and consensus,” in *34th Chinese Control Conference (CCC)*. IEEE, 2015, pp. 6992–6997.
- [23] Lili Su and Nitin H Vaidya, “Defending non-bayesian learning against adversarial attacks,” *arXiv preprint arXiv:1606.08883*, 2016.
- [24] Shahin Shahrampour, Alexander Rakhlin, and Ali Jadbabaie, “Multi-armed bandits in multi-agent networks,” <https://sites.google.com/site/shahinshahrampour/publications>, 2016.
- [25] Pascal Massart, *Concentration inequalities and model selection*, vol. 6, Springer, 2007.
- [26] Jeffrey S Rosenthal, “Convergence rates for markov chains,” *SIAM Review*, vol. 37, no. 3, pp. 387–405, 1995.

6. APPENDIX: PROOFS

In this section, we provide the proofs of our technical results. Before proceeding, we would like to reiterate that:

- For the proofs we sometimes assume that $N > 8$ to simplify the bounds. This assumption is made with no loss of generality, and only avoids notational clutter.
- whenever an equation is longer than the text column, we use “ \dots ” to break the equation and start the rest of it from the next line.

- We make use of the McDiarmid's inequality in our proofs. The inequality is standard and its proof can be found in the literature (see e.g. [25]).

Lemma 3. (McDiarmid's Inequality) *Let $X_1, \dots, X_N \in \chi$ be independent random variables and consider the mapping $H : \chi^N \mapsto \mathbb{R}$. If for $i \in \{1, \dots, N\}$, and every sample $x_1, \dots, x_N, x'_i \in \chi$, the function H satisfies*

$$|H(x_1, \dots, x_i, \dots, x_N) - H(x_1, \dots, x'_i, \dots, x_N)| \leq c_i,$$

then for all $\varepsilon > 0$,

$$\mathbb{P}\left\{H(x_1, \dots, x_N) - \mathbb{E}[H(X_1, \dots, X_N)] \geq \varepsilon\right\} \leq \exp\left\{\frac{-2\varepsilon^2}{\sum_{i=1}^N c_i^2}\right\}.$$

Proof of Lemma 1

The proof of the first part is standard (see e.g. Lemma 1 in [11]). For the second part, we follow the lines in the proof of Lemma 2 in [11]. Let \mathbf{e}_i be the i -th unit vector in the standard basis of \mathbb{R}^N . The Markov chain W is irreducible and aperiodic, so by standard properties of stochastic matrices (see e.g. [26]), we have

$$\left\|\mathbf{e}_i^\top W^t - \frac{1}{N} \mathbf{1}^\top\right\|_1 \leq \sqrt{N} \sigma_2(W)^t, \quad (2)$$

for any $i \in [N]$, as $\frac{1}{N} \mathbf{1}^\top$ is the stationary distribution of the transition kernel W . Hence,

$$\sqrt{N} \sigma_2(W)^{t-\tau} \leq 2 \text{ for } t - \tau \geq \tilde{t} := \frac{\log\left[\frac{\sqrt{N}}{2}\right]}{\log[\sigma_2(W)^{-1}]}, \quad (3)$$

and recall that the inequality $\left\|\mathbf{e}_i^\top W^{t-\tau} - \frac{1}{N} \mathbf{1}^\top\right\|_1 \leq 2$ always holds since any power of W is doubly stochastic. With that in mind, we use (2) to break the following sum into two parts to get

$$\begin{aligned} \sum_{\tau=1}^t \sum_{j=1}^N \left| [W^{t-\tau}]_{ij} - \frac{1}{N} \right| &= \sum_{\tau=1}^t \left\| \mathbf{e}_i^\top W^{t-\tau} - \frac{1}{N} \mathbf{1}^\top \right\|_1 \\ &= \sum_{\tau=1}^{t-\tilde{t}} \left\| \mathbf{e}_i^\top W^{t-\tau} - \frac{1}{N} \mathbf{1}^\top \right\|_1 \\ &\quad + \sum_{\tau=t-\tilde{t}+1}^t \left\| \mathbf{e}_i^\top W^{t-\tau} - \frac{1}{N} \mathbf{1}^\top \right\|_1 \\ &\leq \sum_{\tau=1}^{t-\tilde{t}} \sqrt{N} \sigma_2(W)^{t-\tau} + 2\tilde{t} \\ &\leq \frac{\sqrt{N} \sigma_2(W)^{\tilde{t}}}{1 - \sigma_2(W)} + 2\tilde{t}. \end{aligned}$$

for any $i \in [N]$. Substituting \tilde{t} from (3) into above, we get

$$\begin{aligned} \sum_{\tau=1}^t \sum_{j=1}^N \left| [W^{t-\tau}]_{ij} - \frac{1}{N} \right| &\leq \frac{2}{1 - \sigma_2(W)} + \frac{\log\left[\frac{\sqrt{N}}{2}\right]}{\log[\sigma_2(W)^{-1}]} \\ &\leq \frac{2}{1 - \sigma_2(W)} + \frac{\log N}{\log[\sigma_2(W)^{-1}]}, \end{aligned}$$

for any $i \in [N]$. \square

Proof of Theorem 2

Step 1 : Preliminaries

Recall the definition of dE_1 in the statement of the theorem. Throughout the proof we frequently refer to the following quantities

$$\begin{aligned} \ell &:= \max\left\{\frac{48 \log T}{N \Delta_{1,k}^2}, \frac{4Nd}{\Delta_{1,k}}\right\} & \ell' &:= \frac{4dE_1}{\Delta_{1,k}} \\ c_{t,s} &:= \sqrt{2 \log t \left(\frac{1}{Ns} + \frac{2d}{s^2}\right)} & \hat{t} &:= \frac{5 \log\left[\frac{4\sqrt{N}}{\Delta_{1,k}}\right]}{4 \log[\sigma_2^{-1}(W)]}, \end{aligned} \quad (4)$$

listed here for reader's convenience. We suppress the dependence of ℓ and ℓ' to k to avoid clutter in our derivation.

To bound the regret (1), we need to bound the expected number of times that suboptimal arms are played during the entire process. For any $\ell, \ell' > 0$ (and in particular for the choice of ℓ and ℓ' given in (4)), we have

$$\begin{aligned} n_T(k) &= 1 + \sum_{t=1}^T \mathbf{1}\{I_t = k\} \\ &\leq \ell + \sum_{t=1}^T \mathbf{1}\{I_t = k, n_{t-1}(k) \geq \ell\} \\ &\leq \ell + \sum_{t=1}^T \left\lceil \frac{N}{K} \right\rceil^{-1} \sum_{i=1}^N \mathbf{1}\{I_{i,t} = k, n_{t-1}(k) \geq \ell\}, \end{aligned}$$

where the last step follows from the fact that the network's action is at least selected by $\lceil \frac{N}{K} \rceil$ agents, simply because it is the majority vote. We can proceed as

$$n_T(k) \leq \ell + \left\lceil \frac{N}{K} \right\rceil^{-1} \sum_{i=1}^N P_{i,T}(k) + Q_{i,T}(k), \quad (5)$$

where

$$\begin{aligned} P_{i,T}(k) &:= \sum_{t=1}^T \mathbf{1}\{I_{i,t} = k, n_{t-1}(k) \geq \ell, n_{t-1}(1) > \ell'\} \\ Q_{i,T}(k) &:= \sum_{t=1}^T \mathbf{1}\{I_{i,t} = k, n_{t-1}(k) \geq \ell, n_{t-1}(1) \leq \ell'\}. \end{aligned}$$

We now need to bound $P_{i,T}(k)$ and $Q_{i,T}(k)$ to complete the proof. Each step is included separately in the proof as follows.

Step 2 : Bounding $P_{i,T}(k)$

For any $k > 1$ (which refers to a suboptimal arm), we have

$$\begin{aligned} P_{i,T}(k) &= \sum_{t=1}^T \mathbf{1}\{I_{i,t} = k, n_{t-1}(k) \geq \ell, n_{t-1}(1) > \ell'\} \\ &\leq \sum_{t=1}^T \sum_{s_k \geq \ell} \sum_{s_1 > \ell'}^t \mathbf{1}\{I_{i,t} = k, \dots \\ &\quad n_{t-1}(k) = s_k, n_{t-1}(1) = s_1\} \\ &\leq \sum_{t=1}^T \sum_{s_k \geq \ell} \sum_{s_1 > \ell'}^t \mathbf{1}\left\{\frac{\phi_{i,t}(k)}{s_k} + c_{t,s_k} \geq \frac{\phi_{i,t}(1)}{s_1} \dots \right. \\ &\quad \left. + c_{t,s_1}, n_{t-1}(k) = s_k, n_{t-1}(1) = s_1\right\}, \end{aligned} \quad (6)$$

$$+ c_{t,s_1}, n_{t-1}(k) = s_k, n_{t-1}(1) = s_1\}, \quad (7)$$

where we recall the definition of $c_{t,s}$ from (4). Let

$$\mathcal{S}_{k,t} := \{\tau \in [t] : I_{\tau-1} = k\}, \quad (8)$$

be the instances before which player i has selected action $k \in [K]$. Further notice the explicit form of $\phi_{i,t}$ given in Lemma 1, and recall that $\psi_{i,t}(k) = X_{i,t-1}(k) \mathbf{1}\{k = I_{t-1}\}$ for any $k \in [K]$ as described in the d-UEP algorithm. Then, the indicator (7) implies that at least one of the statements (9), (10), or (11) must hold

$$\frac{1}{s_1} \sum_{\tau \in \mathcal{S}_{1,t}} \sum_{j=1}^N [W^{t-\tau}]_{ij} (X_{j,\tau-1}(1) - \mu_j(1)) \leq -c_{t,s_1} \quad (9)$$

$$\frac{1}{s_k} \sum_{\tau \in \mathcal{S}_{k,t}} \sum_{j=1}^N [W^{t-\tau}]_{ij} (X_{j,\tau-1}(k) - \mu_j(k)) \geq c_{t,s_k} \quad (10)$$

$$\begin{aligned} & \frac{1}{s_1} \sum_{\tau \in \mathcal{S}_{1,t}} \sum_{j=1}^N [W^{t-\tau}]_{ij} \mu_j(1) \cdots \\ & - \frac{1}{s_k} \sum_{\tau \in \mathcal{S}_{k,t}} \sum_{j=1}^N [W^{t-\tau}]_{ij} \mu_j(k) < 2c_{t,s_k}, \end{aligned} \quad (11)$$

simply because failing all of them together contradicts the indicator (7). We now want to show all these events occur with small probability. Starting with (11), we have

$$\begin{aligned} \text{LHS of (11)} &= \frac{1}{s_1} \sum_{\tau \in \mathcal{S}_{1,t}} \sum_{j=1}^N \left([W^{t-\tau}]_{ij} - \frac{1}{N} \right) \mu_j(1) \\ & - \frac{1}{s_k} \sum_{\tau \in \mathcal{S}_{k,t}} \sum_{j=1}^N \left([W^{t-\tau}]_{ij} - \frac{1}{N} \right) \mu_j(k) + \Delta_{1,k} \\ & \geq -dE_1 \left(\frac{1}{s_1} + \frac{1}{s_k} \right) + \Delta_{1,k}, \end{aligned}$$

using the second part of Lemma 1 to bound the sums. Noting that $s_k \geq \ell$ and $s_1 > \ell'$ where ℓ and ℓ' are defined in (4) and $d \geq dE_1$ (by assumption), we can simplify above to get

$$\begin{aligned} \text{LHS of (11)} &\geq -dE_1 \left(\frac{1}{s_1} + \frac{1}{s_k} \right) + \Delta_{1,k} \\ &\geq -dE_1 \left(\frac{\Delta_{1,k}}{4dE_1} + \frac{\Delta_{1,k}}{4Nd} \right) + \Delta_{1,k} \\ &\geq -dE_1 \left(\frac{\Delta_{1,k}}{4dE_1} + \frac{\Delta_{1,k}}{4dE_1} \right) + \Delta_{1,k} \\ &= \frac{\Delta_{1,k}}{2}. \end{aligned} \quad (12)$$

On the other hand, we have

$$\text{RHS of (11)} = 2c_{t,s_k} \leq 2c_{T,s_k} \leq \frac{\Delta_{1,k}}{2}, \quad \forall s_k \geq \ell, \quad (13)$$

since by the definition of ℓ in (4) we have

$$\begin{aligned} 4c_{T,s_k}^2 &= \frac{8 \log T}{s_k} \left(\frac{1}{N} + \frac{2d}{s_k} \right) \leq \frac{N\Delta_{1,k}^2}{6} \left(\frac{1}{N} + \frac{2d}{s_k} \right) \\ &\leq \frac{N\Delta_{1,k}^2}{6} \left(\frac{1}{N} + \frac{2\Delta_{1,k}d}{4Nd} \right) \\ &\leq \frac{\Delta_{1,k}^2}{4}. \end{aligned}$$

Combining (11), (12) and (13), we get

$$\text{RHS of (11)} \leq \frac{\Delta_{1,k}}{2} \leq \text{LHS of (11)} < \text{RHS of (11)},$$

which results in a contradiction, and implies (11) never holds for $s_k \geq \ell$ and $s_1 > \ell'$. Therefore, we need to investigate the possibility of conditions (9) and (10).

To study (9) we appeal to McDiarmid's inequality in Lemma 3. First, note that the random variables involved in the sum are $\{X_{j,\tau-1}(1)\}_{j,\tau}$, and they are independent by assumption. Next, when sequences $\{X_{j,\tau-1}(1)\}_{j,\tau}$ and $\{X'_{j,\tau-1}(1)\}_{j,\tau}$ are equal but for one fixed sample (τ', j') , the difference of the following sum is bounded as

$$\begin{aligned} & \left| \frac{1}{s_1} \sum_{\tau \in \mathcal{S}_{1,t}} \sum_{j=1}^N [W^{t-\tau}]_{ij} (X_{j,\tau-1}(1) - X'_{j,\tau-1}(1)) \right| \cdots \\ & \leq \frac{[W^{t-\tau'}]_{ij'}}{s_1}, \end{aligned}$$

which guarantees the bounded difference condition in McDiarmid's inequality. Noting the cardinality of the set $|\mathcal{S}_{1,t}| = s_1$, we now need to compute the sum of the squares of this bound as,

$$\begin{aligned} & \frac{1}{s_1^2} \sum_{\tau' \in \mathcal{S}_{1,t}} \sum_{j'=1}^N [W^{t-\tau'}]_{ij'}^2 \\ &= \frac{1}{Ns_1} + \frac{1}{s_1^2} \sum_{\tau' \in \mathcal{S}_{1,t}} \sum_{j'=1}^N \left([W^{t-\tau'}]_{ij'}^2 - \frac{1}{N^2} \right) \\ &\leq \frac{1}{Ns_1} + \frac{2}{s_1^2} \sum_{\tau' \in \mathcal{S}_{1,t}} \sum_{j'=1}^N \left| [W^{t-\tau'}]_{ij'} - \frac{1}{N} \right| \\ &\leq \frac{1}{Ns_1} + \frac{2dE_1}{s_1^2}, \end{aligned}$$

where the last line is due to the second part of Lemma 1. Using the assumption $d \geq dE_1$, we can simplify above to get

$$\frac{1}{s_1^2} \sum_{\tau' \in \mathcal{S}_{1,t}} \sum_{j'=1}^N [W^{t-\tau'}]_{ij'}^2 \leq \frac{1}{Ns_1} + \frac{2d}{s_1^2}.$$

Therefore, recalling the definition of c_{t,s_1} from (4), we apply McDiarmid's inequality to equation (9) to derive

$$\mathbb{P}\{\text{Eq. (9) holds}\} \leq \exp\{-\log(t^4)\} = \frac{1}{t^4}. \quad (14)$$

A similar statement holds for (10), and combining with (7) we conclude that

$$\begin{aligned} \mathbb{E}[P_{i,T}(k)] &\leq \sum_{t=1}^T \sum_{s_k \geq \ell} \sum_{s_1 > \ell'}^t \mathbb{P}\{\text{Eq. (9) holds}\} \\ &+ \sum_{t=1}^T \sum_{s_k \geq \ell} \sum_{s_1 > \ell'}^t \mathbb{P}\{\text{Eq. (10) holds}\} \\ &\leq \sum_{t=1}^T \sum_{s_k \geq \ell} \sum_{s_1 > \ell'}^t \frac{2}{t^4} \\ &\leq \sum_{t=1}^{\infty} \frac{2}{t^2} = \frac{\pi^2}{3}. \end{aligned} \quad (15)$$

Step 3 : Bounding $Q_{i,T}(k)$

Now it is turn to bound $Q_{i,T}(k)$ to complete the proof. First, note that

$$\begin{aligned} Q_{i,T}(k) &= \sum_{t=1}^T \mathbf{1} \{I_{i,t} = k, n_{t-1}(k) \geq \ell, n_{t-1}(1) \leq \ell'\} \\ &\leq \sum_{s_1=1}^{\ell'} \sum_{t=1}^T \mathbf{1} \{I_{i,t} = k, n_{t-1}(k) \geq \ell, n_{t-1}(1) = s_1\}. \end{aligned}$$

Let us for each $s_1 \in [\ell']$ denote by t_{s_1} the first time that the indicator holds for the particular value of s_1 . Fixing any $\hat{t} > 0$, we have

$$\begin{aligned} Q_{i,T}(k) &\leq \sum_{s_1=1}^{\ell'} \sum_{t=t_{s_1}}^{t_{s_1+1}-1} \mathbf{1} \{I_{i,t} = k, n_{t-1}(k) \geq \ell, n_{t-1}(1) = s_1\} \\ &\leq \ell' \hat{t} + \sum_{s_1=1}^{\ell'} \sum_{t=t_{s_1}+\hat{t}}^{t_{s_1+1}-1} \mathbf{1} \{I_{i,t} = k, \dots \\ &\quad n_{t-1}(k) \geq \ell, n_{t-1}(1) = s_1\} \\ &\leq \ell' \hat{t} + \sum_{s_1=1}^{\ell'} \sum_{t=t_{s_1}+\hat{t}}^{t_{s_1+1}-1} \sum_{s_k=\ell}^t \mathbf{1} \{I_{i,t} = k, \dots \\ &\quad n_{t-1}(k) = s_k, n_{t-1}(1) = s_1\}, \end{aligned} \quad (16)$$

where the last sum is similar to (6) with different indices on the sums. Hence, to satisfy the indicator, at least one of the statements (9), (10) and (11) must hold (for the new indices). Since $s_k \geq \ell$ the analysis of RHS of (11) given in (13) is still valid. Observe that by standard properties of irreducible and aperiodic Markov chains we have [26],

$$\sum_{j=1}^N \left| [W^t]_{ij} - \frac{1}{N} \right| \leq \sqrt{N} \sigma_2^t(W) < \frac{\Delta_{1,k}}{4}, \quad (17)$$

for all

$$t > \frac{\log \left[\frac{4\sqrt{N}}{\Delta_{1,k}} \right]}{\log [\sigma_2^{-1}(W)]}.$$

To analyze the LHS of (11) for new indices, let \hat{t} be defined as in (4) and recall (8). Then, for any $s_1 \in [\ell']$ and $t \in [t_{s_1} + \hat{t}, t_{s_1+1} - 1]$ we have $\mathcal{S}_{1,t} = \mathcal{S}_{1,t_{s_1}}$ by definition of t_{s_1} . Hence, we modify the expression in (12) as

$$\begin{aligned} \text{LHS of (11)} &= \frac{1}{s_1} \sum_{\tau \in \mathcal{S}_{1,t_{s_1}}} \sum_{j=1}^N \left([W^{t-\tau}]_{ij} - \frac{1}{N} \right) \mu_j(1) \\ &\quad + \Delta_{1,k} - \frac{1}{s_k} \sum_{\tau \in \mathcal{S}_{k,t}} \sum_{j=1}^N \left([W^{t-\tau}]_{ij} - \frac{1}{N} \right) \mu_j(k) \\ &\geq -\sqrt{N} \sigma_2^{\hat{t}}(W) - \frac{dE_1}{s_k} + \Delta_{1,k}, \end{aligned}$$

where in the last line we used Lemma 1 and equation (17). We now use the fact that $s_k \geq \ell$ and recall the definition of ℓ and \hat{t} in (4) to get

$$\begin{aligned} \text{LHS of (11)} &\geq -\sqrt{N} \sigma_2^{\hat{t}}(W) - \frac{dE_1}{s_k} + \Delta_{1,k} \\ &\geq \frac{\Delta_{1,k}}{2}. \end{aligned}$$

Combining above with (13) implies that (11) never holds. Notice that our argument about the probability of events (9) and (10) holds for any $s_1, s_k, t > 0$, and therefore, the tail bound (14) holds true again. Employing these facts and returning to (16) we get

$$\mathbb{E}[Q_{i,T}(k)] \leq \ell' \hat{t} + \frac{\pi^2}{3}. \quad (18)$$

Step 4 : Finishing the Proof

Noting $\lceil \frac{N}{K} \rceil \geq \frac{N}{K}$, we substitute (15) and (18) into (5) to obtain the bound

$$\begin{aligned} \mathbb{E}[n_T(k)] &\leq \ell + \frac{K}{N} \sum_{i=1}^N P_{i,T}(k) + Q_{i,T}(k) \\ &\leq \ell + K \ell' \hat{t} + \frac{2K\pi^2}{3}. \end{aligned} \quad (19)$$

Recall the definition of dE_1 and dE_2 from the statement of the theorem and the simplifying assumption that N is large enough ($N > 8$). Then, plugging the above into (1) using quantities defined in (4) concludes the proof. \square