# Relational Weight Optimization for Enhancing Team Performance in Multi-Agent Multi-Armed Bandits

**Monish Reddy Kotturu** [*] **Saniya Vahedian Movahed** [**]
**Paul Robinette** [***] **Kshitij Jerath** [****] **Amanda Redlich** [†]
**Reza Azadeh** [*]

[*] *Miner School of Computer & Information Sciences, University of Massachusetts Lowell, Lowell, MA 01854*
[**] *University of Texas San Antonio, San Antonio, TX 78249*
[***] *Electrical & Computer Engineering, University of Massachusetts Lowell, Lowell MA 01854*
[****] *Mechanical and Industrial Engineering, University of Massachusetts Lowell, Lowell MA 01854*
[†] *Mathematics & Statistics, University of Massachusetts Lowell, Lowell MA 01854*

**Abstract:**
We introduce an approach to improve team performance in a Multi-Agent Multi-Armed Bandit (MAMAB) framework using Fastest Mixing Markov Chain (FMMC) and Fastest Distributed Linear Averaging (FDLA) optimization algorithms. The multi-agent team is represented using a fixed relational network and simulated using the Coop-UCB2 algorithm. The edge weights of the communication network directly impact the time taken to reach distributed consensus. Our goal is to shrink the timescale on which the convergence of the consensus occurs to achieve optimal team performance and maximize reward. Through our experiments, we show that the convergence to team consensus occurs slightly faster in large constrained networks.

## 1. INTRODUCTION

Multi-Armed Bandits (MABs) are a class of reinforcement learning problems where an agent is presented with a set of arms (i.e., actions), with each arm giving a reward drawn from a probability distribution unknown to the agent (Lattimore and Szepesvári, 2020). The goal of the agent is to maximize its total reward which requires balancing exploration and exploitation. MABs offer a simple model to simulate decision-making under uncertainty. Practical applications of MAB algorithms include news recommendations (Yang and Toni, 2018), online ad placement (Aramayo et al., 2022), dynamic pricing (Babaioff et al., 2015), and adaptive experimental design (Rafferty et al., 2019). In contrast to single-agent cases, in certain applications such as search and rescue, a team of agents should cooperate with each other to accomplish goals by maximizing team performance. Such problems are solved using Multi-Agent Multi-Armed Bandit (MAMAB) algorithms (Xu et al., 2020). Most existing algorithms rely on the presence of multiple agents and try to solve the problem using the shared information among them ignoring the relationship between team members (Sankararaman et al., 2019; Rangi and Franceschetti, 2018).

Few works, on the other hand, use graph representation (Shahrampour et al., 2017; Madhushani and Leonard, 2019) and grouping (Sankararaman et al., 2019) to estab-lish a specific team structure. Using a graph to represent the team behavior ensures that the relationship between the agents are held. However, existing works either do not consider the weight of each relationship (graph edges) (Madhushani and Leonard, 2020; Agarwal et al., 2021) or expect the user to manually set those weights (Moradipari et al., 2022).

In this paper, we propose a new approach that combines graph optimization and MAMAB algorithms to enhance team performance by expediting the convergence to consensus of arm means. Our proposed approach:

- improves team performance by optimizing the edge weights in the graph representing the team structure in large constrained teams,
- does not require manual tuning of the graph weights,
- is independent of the MAMAB algorithm and only depends on the consensus formula, and
- formulates the problem as a convex optimization, which is computationally efficient for large teams.

We evaluate Coop-UCB2 combined with six different graph optimization approaches in various team structures, measured by time taken to converge to consensus of the best arm mean. Our results show that the proposed method outperforms existing graph-based MAMAB algorithms with manual tuning or other adjustment heuristics in large, constrained teams, but does not have a significant effect in small networks.

## 2. RELATED WORK

Algorithms similar to Coop-UCB2 have been previously devised to simulate distributed learning in multi-agent networks. A distributed version of the Upper Confidence Bound (UCB) 1 algorithm (Shahrampour et al., 2017; Auer et al., 2002) called d-UER was proposed that aims to minimize network regret while also using a network structure. It maintains the same principle as UCB in that it uses the upper confidence bound that relies on the network topology. Players vote for their estimated best arm iteratively and the network takes an action based on the majority vote. However, it focused on the case where the arms are dependent on the players. In Sankararaman et al. (2019), the authors developed an algorithm in which agents choose to communicate with other agents uniformly and independently at random and only communicate the arm IDs. They show a significant reduction in the network regret in both full communication and no communication scenarios, mimicking different levels of collaboration. This and other works such as Chawla et al. (2020); Zhu et al. (2021); Martínez-Rubio et al. (2019) use gossip style algorithms to solve the MAMAB problem. However, the above algorithms do not optimize the consensus step of the communication process for faster convergence.

Literature on convex optimization methods involving distributed consensus models includes *least-mean-square consensus* (LMSC) (Xiao et al., 2007), which takes fast linear iterations (Xiao and Boyd, 2003) further by including the total mean-square deviation of each variable that converges to a steady-state value and minimizing this deviation. Carli et al. (2007) compared the behavior and performance of linear averaging algorithms for consensus problems in cases where transmission noise existed and where communication is quantized. They show that the algorithms fail to converge to a consensus when there is added white noise (with zero mean and bounded covariance), and they provide some solutions that involve discarding information from agents who are too close to the source of noise to prevent drift from initial average. Zhou et al. (2013) also examine the discrete average consensus problem with bounded noise and show that an increase upper bound for noise decreases the convergence accuracy.

## 3. BACKGROUND

### 3.1 Stochastic Multi-Armed Bandits

A stochastic Multi-Armed Bandit (MAB) is a collection of probability distributions $\nu = (P_a : a \in \mathcal{A})$ over all the available arms (i.e., actions) $\mathcal{A}$. In each round, $t \in \{1, \ldots, T\}$, the agent interacts with the MAB (i.e., environment) by selecting an arm $a(t) \in \mathcal{A}$, and receives a reward $r(t) \in \mathbb{R}$ sampled from $P_{a(t)}$. Throughout this paper, we consider MABs from the Sub-Gaussian environment class $\mathcal{E}_{SG}^k(\sigma^2) = \{P_a : \sigma - \text{Sub-Gaussian}\}$. For a $\sigma-$Sub-Gaussian random variable $X$, $\mathbb{P}(X \geq \epsilon) \leq \exp(\frac{-\epsilon^2}{2\sigma^2})$ holds for any $\epsilon \geq 0$. Additionally, the tail of $X$ decays approximately as fast as that of $\mathcal{N}(0, \sigma^2)$ (Lattimore and Szepesvári, 2020). The Gaussian, Bernoulli, and uniform distributions are examples of Sub-Gaussian distributions.

In single-agent MAB problems, the goal of the agent is to find a policy $\pi$ that maximizes the total reward $\sum_{t=1}^{T} r(t)$, where $T \in \mathbb{N}$ is the horizon. This is equivalent to minimizing the regret which is defined by $R_T(\pi, \nu) = T\mu^*(\nu) - \mathbb{E}[\sum_{t=1}^{T} r(t)]$, where $\mu^* = \arg\max_{a \in \mathcal{A}} \mu_a(\nu)$ indicates the mean of the best arm and $\mu_a(\nu)$ is the mean of arm $a$ in the bandit $\nu$. The regret can be represented as $R_n(\pi, \nu) = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}[n_a(T)]$, where $\Delta_a(\nu) = \mu^*(\nu) - \mu_a(\nu)$ is the immediate regret (i.e., the suboptimality gap) in the bandit $\nu$, and $n_a(\tau) = \sum_{t=1}^{\tau} \mathbb{I}\{a(t) = a\}$ is the total number of times action $a$ was selected after the end of round $\tau$.

### 3.2 Multi-Agent Multi-Armed Bandits

The goal of a team playing a Multi-Agent Multi-Armed Bandit is to maximize individual rewards over the horizon $T$, while playing the same bandit, cooperatively. In this setup, agents communicate with each other over a network, modeled by an undirected graph $\mathcal{G} = (V, E)$, where $V$ is the set of vertices representing $M$ agents and $E$ is the set of edges representing the connections between agents. Each agent, $k \in \{1, \ldots, M\}$, uses its current estimation to select an arm $A_t^k$ from the finite set of arms $\mathcal{A} = \{A_i\}$ for $i \in \{1, \ldots, N\}$, where $N \in \mathbb{N}$ and $|\mathcal{A}| = N$. The agent then receives a real value reward, updates its estimates, and shares the reward information of the selected arm with connected neighbors through $\mathcal{G}$. The objective of the cooperative multi-agent MAB problem is to maximize the expected total group reward that is equivalent to minimizing the expected group regret defined by $R_n(\pi, \nu) = \sum_{i=1}^{N} \sum_{k=1}^{M} \Delta_i \mathbb{E}[n_i^k(T)]$.

In such a cooperative setting, agents update their estimations by combining their observations (i.e., realized rewards) and their neighbors' observations usually through a consensus process. The simplest and most common form of consensus, which is used in forming opinions in social learning, is the distributed averaging algorithm presented as $x(t + 1) = Px(t)$, where the vector $x(t)$ represents agents' opinions in round $t$ and $P$ is the transition matrix in which $P_{ij}$ defines the weight that agent $i$ gives to agent $j$'s opinion. In such consensus process with no disturbances or new updates, all agents' opinions converge asymptotically to $x(0)$. The process of averaging consensus along its variants have been greatly studied in literature (Zhou et al., 2013; Boyd et al., 2004b; Xiao et al., 2007; Xiao and Boyd, 2003; Carli et al., 2007). In the presence of an external signal, the distributed averaging consensus problem can be shown as $x(t + 1) = P(x(t) + z(t))$, where $z(t)$ represents new updates, which in a multi-agent MAB problem is the vector of rewards realized by agents in round $t$.

## 4. METHODOLOGY

### 4.1 Problem Description

The transition matrix, or Perron matrix, $P$ plays a significant role in the convergence behavior of the above-mentioned distributed averaging consensus process. In a cooperative multi-agent MAB algorithm, agents share the knowledge of observed reward information with their

neighbors. The consensus process then updates the current estimate of the arms (i.e., actions) using the transition matrix $P$. Therefore, varying the values of the $P$ directly affects the performance of each agents and the team. In this paper, we aim at enhancing the multi-agent team performance by expediting the convergence to consensus. To achieve this goal, we propose an approach that combines graph optimization and multi-agent MAB algorithms. Given a relational network (i.e., the graph $\mathcal{G}$), our approach finds the optimal relational weights (i.e., elements of the transition matrix $P$) and uses those for the consensus process in the multi-agent MAB algorithm.

### 4.2 Multi-Agent MAB Algorithm

Our proposed method is independent of the choice of the multi-agent MAB algorithm and only relies on the use of distributed averaging consensus in the estimation of the arms' mean values through shared realized reward between agents. Among existing algorithms, we focus on Coop-UCB2 (Landgren et al., 2021) which is a multi-agent extension of the Upper Confidence Bound (UCB) algorithm (Auer et al., 2002) in which a group of agents perform cooperative decision making to solve a multi-armed bandit problem. Coop-UCB2 uses distributed consensus in a graph network structure for cooperative decision making. The agents communicate with each other through a network represented by a fixed undirected graph $\mathcal{G} = (V, E)$. The algorithm starts by each agent sampling each arm once. In the next rounds, each agent selects the arm with maximum estimated mean according to the following estimation:

$$Q_i^k = \frac{\hat{s}_i^k}{\hat{n}_i^k} + \sigma_g \sqrt{\frac{2\gamma}{G(\eta)} \cdot \frac{\hat{n}_i^k + f(t-1)}{M\hat{n}_i^k} \cdot \frac{\ln(t-1)}{\hat{n}_i^k}}, \quad (1)$$

where $Q_i^k$ is the estimation of the mean for arm $i$ generated by agent $k$, $\hat{s}_i^k$ is the estimation of the total reward given by arm $i$ until time-step $t$ for agent $k$, $\hat{n}_i^k$ is the number of times arm $i$ was selected by agent $k$ until time $t$, $f(t)$ is an increasing sub-logarithmic function of $t$ (e.g., $\sqrt{\log(t)}$), $\sigma_g \in \mathbb{R}_+$, $\gamma > 1$, $\eta \in (0,4)$, and $G(\eta) = 1 - \eta^2/16$ (Landgren et al., 2021).

Agents update their estimates cooperatively through a distributed consensus process represented using the following equations:

$$\hat{n}_i(t) = P(\hat{n}_i(t-1) + \xi_i(t)), \quad (2)$$
$$\hat{s}_i(t) = P(\hat{s}_i(t-1) + r_i(t)), \quad (3)$$

where $\xi_i(t)$ is a vector representing the number of times the action $i$ was selected and $r_i(t)$ is the reward vector at time-step $t$. In Coop-UCB2, $P$ is defined as a row stochastic matrix (also known as Perron matrix) obtained from the following equation:

$$P = I_M - \frac{\kappa}{d_{max}} L, \quad (4)$$

where $I_M$ is the identity matrix of order $M$ (the number of agents), $\kappa \in (0,1]$ is a step size parameter, and $d_{max} = max\{deg(i)|i \in \{1,...,M\}\}$. $L$ is the Laplacian matrix calculated from the adjacency matrix $A = [A_{ij}]$ for graph $\mathcal{G}$ as

$$L_{ij} = \begin{cases} \sum_{k=1, k\neq i}^{M} A_{ik}, & j = i \\ -A_{ij}, & j \neq i \end{cases}. \quad (5)$$

However, this method of calculating the edge weights of $P$ in (4) is independent of the relational network's structure and complexity. Consequently, using this methods does not result in the optimal team performance in many situations. We consider a team performance optimal if the network of agents converges to the optimal consensus as fast as possible.

### 4.3 Optimization using Fastest Mixing Markov Chain

With the goal of discovering an optimal $P$ matrix, we consider formulating the problem in the context of Markov chain on the given graph. We can define a Markov chain over the undirected graph $\mathcal{G} = (V, E)$ representing our multi-agent team. Our goal is to optimize the weight $P_{ij}$ of each edge $(i,j) \in E$, which is the probability of transitioning from node $i$ to node $j$. The assigned probabilities must be non-negative and for each node, the sum of the probabilities of links connected to the node must equal to one. In other words, $P_{ij} = p(x(t+1) = i|x(t) = j)$ for $i, j = 1,...,M$, where $x(t) \in \{1,...,M\}$, for $t \in \mathbb{Z}_+$ represents states. The transition matrix must satisfy $P_{ij} \geq 0$, $\mathbf{1}^\top P = \mathbf{1}^\top$, $P = P^\top$, and $P_{ij} = 0$ for $(i,j) \notin E$, where $\mathbf{1}$ is the vector of all ones. The second largest eigenvalue of $P$ is called the mixing rate and is defined as $max\{\lambda_2, -\lambda_n\}$, where $1 = \lambda_1 \geq \lambda_2 \geq ... \geq \lambda_n$ in which $\lambda_i$s are the eigenvalues of $P$. The smaller mixing rate results in faster mixing in the Markov chain towards the equilibrium distribution $(1/M)\mathbf{1}$. The Fastest Mixing Markov Chain (FMMC) is the optimization problem to find the minimum mixing rate. By applying the idea of FMMC to our consensus process, we can obtain optimal relational network weights that expedite the convergence.

The optimization problem can be represented as the following convex form:

$$\begin{aligned} &\text{minimize} \quad \|P - (1/M)\mathbf{1}\mathbf{1}^\top\|_2 \\ &\text{subject to} \quad P \geq 0, P\mathbf{1} = \mathbf{1}, P = P^\top, \\ &\qquad\qquad P_{ij} = 0, (i,j) \notin E, \end{aligned} \quad (6)$$

where $\|.\|_2$ represents the spectral norm.

The formulation in (6) can be expressed as a Semi Definite Program (SDP) (Alizadeh, 1995; Boyd et al., 2004a) by introducing a scalar slack variable $s \in \mathbb{R}$ as follows:

$$\begin{aligned} &\text{minimize} \quad s \\ &\text{subject to} \quad -sI \preceq P - (1/M)\mathbf{1}\mathbf{1}^\top \preceq sI, \\ &\qquad\qquad P \geq 0, P\mathbf{1} = \mathbf{1}, P = P^\top, \\ &\qquad\qquad P_{ij} = 0, (i,j) \notin E, \end{aligned} \quad (7)$$

where the symbol $\preceq$ denotes matrix inequality. The optimization problem in (7) can be solved using interior-point algorithms for SDP (Alizadeh, 1995).

### 4.4 Optimization using Fastest Distributed Linear Averaging

Alternatively, a distributed consensus over a network of agents can be optimized by finding a linear iteration that is able to calculate the average weights of nodes and edges

in the given graph (Xiao and Boyd, 2003). Fastest Linear Distributed Linear Averaging (FDLA) achieves this with a set of constraints on $P$ to solve the problem faster. The FDLA optimization problem can be formulated as

$$\begin{aligned} \text{minimize} \quad & \rho(P - (1/M)\mathbf{1}\mathbf{1}^T), \\ \text{subject to} \quad & P \in \mathscr{S}, P = P^T, P\mathbf{1} = \mathbf{1}, \end{aligned} \tag{8}$$

where $\rho$ is the spectral radius of $P$ and $\mathscr{S} = \left\{ P \in \mathbb{R}^{M \times M} | P_{ij} = 0 \text{ if } i,j \notin E \text{ and } i \neq j \right\}$.

The formulation in (8) can be expressed as an SDP by introducing a scalar slack variable $s \in \mathbb{R}$ as follows:

$$\begin{aligned} \text{minimize} \quad & s \\ \text{subject to} \quad & -sI \preceq P - (1/M)\mathbf{1}\mathbf{1}^\top \preceq sI, \\ & P \in \mathscr{S}, P\mathbf{1} = \mathbf{1}, P = P^\top. \end{aligned} \tag{9}$$

Similar to FMMC, the optimization problem in (9) can be solved using interior-point algorithms for SDP (Alizadeh, 1995). Unlike FMMC, FDLA, however, allows for negative values in $P$. It has to be noted that although negative values are not obvious choices for weights that represent the relationship between the agents, they can result in faster convergence in certain situations.

### 4.5 On the Convergence of the Running Consensus

Our running average consensus works according to $x(t+1) = P(x(t) + z(t))$, where $\mathbf{1}^\top P = \mathbf{1}^\top$, $P\mathbf{1} = \mathbf{1}$, i.e., $P$ is doubly stochastic, and $z(t) \sim \mathcal{N}(\mu, \sigma_2)$ is external observation (i.e. rewards) obtained by playing an arm with the true mean $\mu$, and variance $\sigma^2$. If $z(t) = 0$, our formula reduces to the discrete average consensus and converges to $\bar{x} = \frac{1}{M}\sum_{i=1}^{M} x(0) = \frac{\mathbf{1}\mathbf{1}^\top}{M} x(0)$. When $z(t) \sim \mathcal{N}(0, \sigma^2)$ (i.e., white noise), it has been shown that the consensus system is stable (Carli et al., 2007; Hatano et al., 2005). Here, however, we assume that $z(t)$ is the observation of the agent received after taking an action. So, it can be written:

$$x(t) = P^t x(0) + \sum_{k=1}^{t} z(t-k+1)P^k \tag{10}$$

$$= P^t x(0) + \sum_{k=1}^{t} \mu P^k + \sum_{k=1}^{t} \varepsilon(t-k+1)P^k$$

where $\varepsilon$ represents the noise. We know that $\lim_{t \to \infty} P^t x(0) = \bar{x}\mathbf{1}$. Since $P$ is doubly stochastic with the maximum eigenvalue equal to one, it can be shown that $\lim_{t \to \infty} \sum_{k=1}^{t} z(t-k+1)P^k = C\mathbf{1}$, where $C$ is a constant. Consequently, we conclude that the running consensus converges to a constant $\lim_{t \to \infty} x(t) = \bar{x}\mathbf{1} + C\mathbf{1} = C'\mathbf{1}$.

## 5. EXPERIMENTS

### 5.1 Experimental Setups

Our experiments consist of 100-armed bandits where for each arm $a$, the true mean $\mu_a$ is sampled from $\mathcal{N}(0,1)$. For each arm $a$, the rewards are sampled from $\mathcal{N}(\mu_a, \sigma^2)$, where $\sigma_a = 1$ unless mentioned otherwise. Each experiment involve 10,000 runs each of which includes playing a different bandit for 1,000 time steps. The networks we use in our experiments are shown in Fig. 1. They
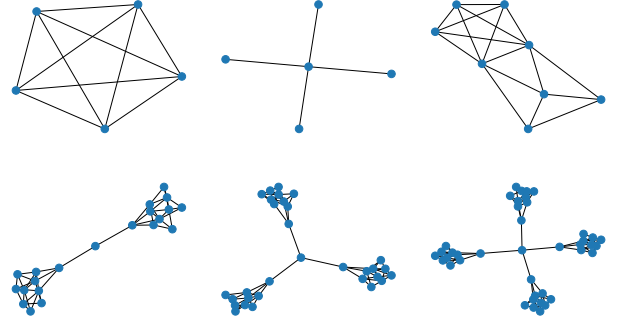


Fig. 1. Graphical representation of the networks

include 5-agent networks adopted from Landgren et al. (2021), namely, all-to-all and star networks, and an 8-agent network adopted from Xiao and Boyd (2003). We also introduce our own set of networks that are created by joining small clusters by a common "parent" node. This is to increase the complexity of the network by decreasing the closeness centrality between nodes from different clusters and increasing the betweenness centrality of the parent node and thus, introducing a constraint on how fast the information is propagated throughout the team. In each experiment, we performed the FMMC and FDLA optimization on the given graph as presented in Section 4 and used the Coop-UCB2 algorithm in (1) with the distributed consensus in (2) and (3).

### 5.2 Heuristics

Here, we introduce a few heuristics that could be used for computing an optimal $P$ that guarantees convergence of distributed linear averaging iteration.

**Constant-edge weights:** A heuristic that achieves this goal is *constant-edge weights* that simply sets $P = I_M - \alpha L$. In this method, all the self weights add up to 1 ($P\mathbf{1} = \mathbf{1}$) and the edge weights are equal to $\alpha$.

**Maximum-degree weights:** another heuristic that follows the same principle, but uses $\alpha = 1/d_{max}$, where $d_{max}$ is the maximum degree in the graph. It is proven that distributed linear averaging converges using the constant-edge weights method where $\alpha \in (0, 1/d_{max})$ (Xiao and Boyd, 2003).

**Local-degree weights:** Another heuristic that generates $P$ by assigning the weights based on the maximum of the two neighboring nodes as $P_{ij} = \frac{1}{max\{d_i, d_j\}}$, where $d_i$ and $d_j$ are the degrees of nodes $i$ and $j$ in the graph, respectively.

### 5.3 Evaluation Metrics

We evaluate the performance of our optimized weight matrix $P$ using two metrics: the spectral radius and the convergence time.

**Spectral radius ($\rho$):** represents the asymptotic convergence factor defined as

$$\rho = \sup_{x(0) \neq \bar{x}} \lim_{t \to \infty} \left( \frac{\|x(t) - \bar{x}\|_2}{\|x(0) - \bar{x}\|_2} \right)^{1/t}. \tag{11}$$
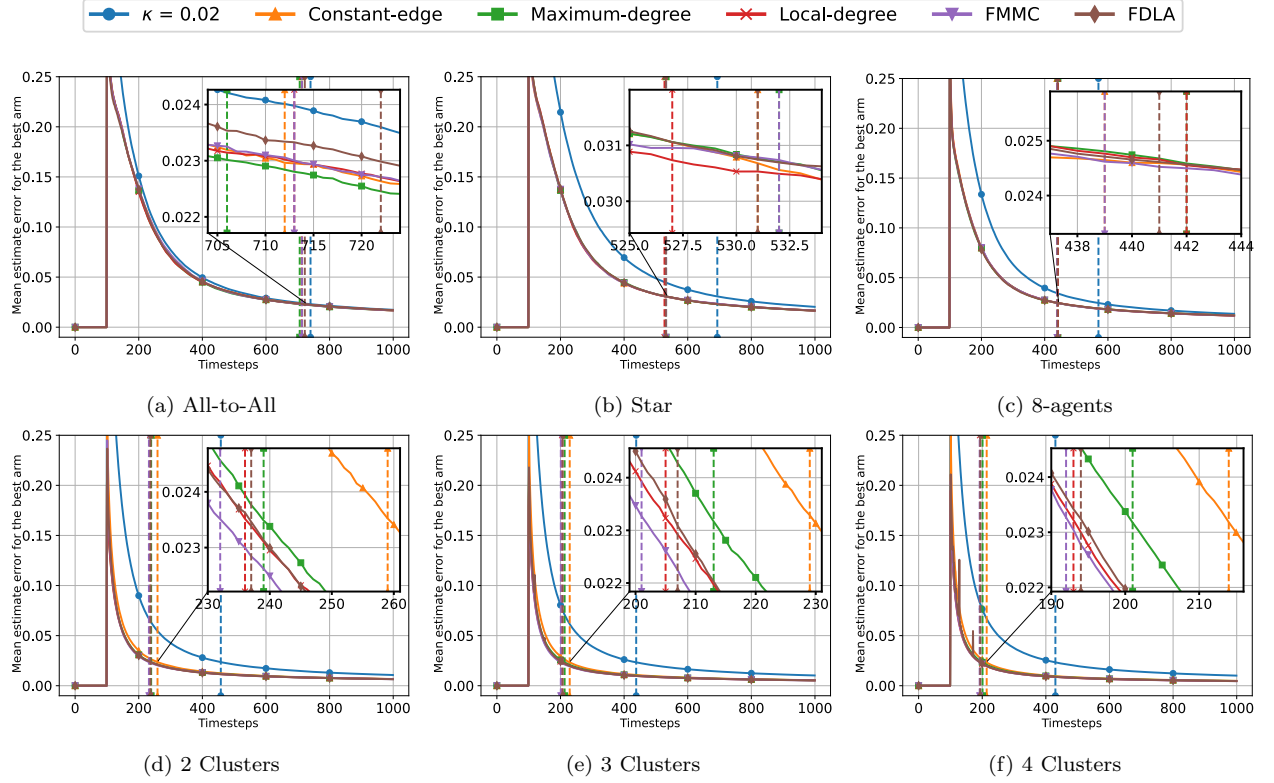
Fig. 2. Comparison of team average of the errors between the estimated and true means of the best arm in different networks. The vertical dashed lines represent the time taken by the network to reach 5% of the final value of the largest error among all algorithms.

Table 1. Convergence factors in networks with weights optimized using various methods.

| Edge weight generation method | All-to-All | | Star | | 8-agents | | 2-Cluster | | 3-Cluster | | 4-Cluster | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\rho$ | $\tau$ | $\rho$ | $\tau$ | $\rho$ | $\tau$ | $\rho$ | $\tau$ | $\rho$ | $\tau$ | $\rho$ | $\tau$ |
| $\kappa = 0.02$ | 0.975 | 39.498 | 0.995 | 199.5 | 0.995 | 196.5 | 0.999 | 3838.7 | 0.999 | 3950.2 | 0.999 | 3952.5 |
| Constant-edge | **3.3e-16** | **0.028** | **0.667** | **2.466** | 0.655 | 2.363 | 0.981 | 52.011 | 0.982 | 53.899 | 0.982 | 54.155 |
| Maximum-degree | 0.250 | 0.721 | 0.750 | 3.476 | 0.746 | 3.416 | 0.987 | 76.283 | 0.987 | 78.513 | 0.987 | 78.559 |
| Local-degree | 0.250 | 0.721 | 0.750 | 3.476 | 0.743 | 3.369 | 0.984 | 60.277 | 0.983 | 57.623 | 0.983 | 57.668 |
| FMMC | 5.2e-08 | 0.060 | 0.750 | 3.476 | 0.667 | 2.466 | 0.974 | 37.718 | 0.977 | 43.724 | 0.981 | 52.279 |
| FDLA | 2.9e-08 | 0.058 | **0.667** | **2.466** | **0.600** | **1.958** | **0.969** | **31.983** | **0.974** | **37.653** | **0.977** | **42.667** |

**Convergence time ($\tau$):** defined as

$$\tau = \frac{1}{log(1/\rho)}. \tag{12}$$

Assuming $*$ indicates the best arm in a given MAB, we compare the average estimated mean for the best arm among all above-mentioned methods defined as $\frac{\hat{s}^*}{\hat{n}^*}$, where $\hat{s}^*$ and $\hat{n}^*$ are estimation of the total reward given by the best arm and the number of times the best arm was selected, respectively. We define the team average of the errors between the estimated and true means as $\delta$:

$$\delta = \frac{1}{M} \cdot \sum_{k=1}^{M} \left( \frac{\hat{s}_k^*}{\hat{n}_k^*} - \mu^* \right). \tag{13}$$

### 5.4 Results and Research Insights

All teams optimized using FMMC and FDLA theoretically perform better in terms of $\rho_{\text{asym}}$ and $\tau_{\text{asym}}$ as reported in Table 1. However, in terms of $\delta$, the smaller 5 and 8-agent teams perform inconsistently between the heuristics, FMMC, and FDLA. The heuristics perform better in every

case as represented by the vertical dashed lines, which signify the time taken to reach 5% of the largest error among all heuristic methods and algorithms.

As the network complexity increases when transitioning to the clustered networks, we see a trend where FMMC does consistently better than most of the heuristics, usually closely followed by local-degree or FDLA. Even if $\rho$ and $\tau$ show that FDLA should be the most optimal, we hypothesize that it does not translate to an actual performance improvement due to negative weights produced by FDLA. Although not shown in the figures, it should be noted that a network with the same number of nodes and edges as the clustered networks but without the constraint of having one parent node do not have the same improvement when optimized using FMMC and FDLA. This maybe due to this large randomly generated graph still having similar properties to the smaller networks that affect communication, such as high closeness centralities among all nodes. Therefore, we think that the problem and team size play a significant role in the effectiveness of optimization.

While networks optimized with both FMMC and FDLA performed better than others in the larger networks, it is important to consider the differences between FMMC and FDLA. FMMC ensures that the resulting transition matrix $P$ has positive weights as its optimization problem imposes the constraint that $P \geq 0$. On the other hand, FDLA does not impose the same constraint. So, the resulting matrix may contain negative weights, which may be important to consider depending on the application of this network optimization approach.

## 6. CONCLUSIONS

In this paper, we propose a method that uses consensus optimization methods to optimize edge weights of a communication graph and improve the team performance of a multi-agent network playing a MAB. The results of our experiments show that this optimization step can improve the team performance in terms of the time taken to reach a consensus on the best arm in large constrained networks, which consist of several small clusters connected by a single parent node. In our future work, we would like to consider real-world networks with different properties that may also benefit from optimized edge weights.

## REFERENCES

Agarwal, M., Aggarwal, V., and Azizzadenesheli, K. (2021). Multi-agent multi-armed bandits with limited communication. *arXiv preprint arXiv:2102.08462*.

Alizadeh, F. (1995). Interior point methods in semidefinite programming with applications to combinatorial optimization. *SIAM journal on Optimization*, 5(1), 13–51.

Aramayo, N., Schiappacasse, M., and Goic, M. (2022). A multiarmed bandit approach for house ads recommendations. *Marketing Science*.

Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2), 235–256.

Babaioff, M., Dughmi, S., Kleinberg, R., and Slivkins, A. (2015). Dynamic pricing with limited supply.

Boyd, S., Boyd, S.P., and Vandenberghe, L. (2004a). *Convex optimization*. Cambridge university press.

Boyd, S., Diaconis, P., and Xiao, L. (2004b). Fastest mixing markov chain on a graph. *SIAM Review*, 46(4), 667–689. doi:10.1137/S0036144503423264. URL https://doi.org/10.1137/S0036144503423264.

Carli, R., Fagnani, F., Frasca, P., Taylor, T., and Zampieri, S. (2007). Average consensus on networks with transmission noise or quantization. In *2007 European Control Conference (ECC)*, 1852–1857. doi: 10.23919/ECC.2007.7068829.

Chawla, R., Sankararaman, A., Ganesh, A., and Shakkottai, S. (2020). The gossiping insert-eliminate algorithm for multi-agent bandits. In *International conference on artificial intelligence and statistics*, 3471–3481. PMLR.

Hatano, Y., Das, A.K., and Mesbahi, M. (2005). Agreement in presence of noise: pseudogradients on random geometric networks. In *Proceedings of the 44th IEEE Conference on Decision and Control*, 6382–6387. IEEE.

Landgren, P., Srivastava, V., and Leonard, N. (2021). Distributed cooperative decision making in multi-agent multi-armed bandits. *Automatica*, 125, 109445. doi: 10.1016/j.automatica.2020.109445.

Lattimore, T. and Szepesvári, C. (2020). *Bandit algorithms*. Cambridge University Press.

Madhushani, U. and Leonard, N.E. (2019). Heterogeneous stochastic interactions for multiple agents in a multi-armed bandit problem. In *2019 18th European Control Conference (ECC)*, 3502–3507. IEEE.

Madhushani, U. and Leonard, N.E. (2020). A dynamic observation strategy for multi-agent multi-armed problem. In *2020 European Control Conference (ECC)*, 1677–1682. IEEE.

Martínez-Rubio, D., Kanade, V., and Rebeschini, P. (2019). Decentralized cooperative stochastic bandits. *Advances in Neural Information Processing Systems*, 32.

Moradipari, A., Ghavamzadeh, M., and Mahnoosh, A. (2022). Collaborative multi-agent stochastic linear bandits. In *2022 American Control Conference (ACC)*, 2761–2766. IEEE.

Rafferty, A., Ying, H., Williams, J., et al. (2019). Statistical consequences of using multi-armed bandits to conduct adaptive educational experiments. *Journal of Educational Data Mining*, 11(1), 47–79.

Rangi, A. and Franceschetti, M. (2018). Multi-armed bandit algorithms for crowdsourcing systems with online estimation of workers' ability. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, 1345–1352.

Sankararaman, A., Ganesh, A., and Shakkottai, S. (2019). Social learning in multi agent multi armed bandits. 3(3). doi:10.1145/3366701. URL https://doi.org/10.1145/3366701.

Shahrampour, S., Rakhlin, A., and Jadbabaie, A. (2017). Multi-armed bandits in multi-agent networks. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2786–2790. doi: 10.1109/ICASSP.2017.7952664.

Xiao, L. and Boyd, S. (2003). Fast linear iterations for distributed averaging. In *42nd IEEE International Conference on Decision and Control (IEEE Cat. No.03CH37475)*, volume 5, 4997–5002 Vol.5. doi: 10.1109/CDC.2003.1272421.

Xiao, L., Boyd, S.P., and Kim, S.J. (2007). Distributed average consensus with least-mean-square deviation. *J. Parallel Distributed Comput.*, 67, 33–46.

Xu, X., Tao, M., and Shen, C. (2020). Collaborative multi-agent multi-armed bandit learning for small-cell caching. *IEEE Transactions on Wireless Communications*, 19(4), 2570–2585.

Yang, K. and Toni, L. (2018). Graph-based recommendation system. In *2018 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, 798–802. IEEE.

Zhou, M., He, J., Cheng, P., and Chen, J. (2013). Discrete average consensus with bounded noise. In *52nd IEEE Conference on Decision and Control*, 5270–5275. doi: 10.1109/CDC.2013.6760718.

Zhu, Z., Zhu, J., Liu, J., and Liu, Y. (2021). Federated bandit: A gossiping approach. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 5(1), 1–29.