
Best Arm Identification in Graphical Bilinear Bandits

Geovani Rizk^{1,2} Albert Thomas² Igor Colin² Rida Laraki^{1,3} Yann Chevaleyre¹

Abstract

We introduce a new graphical bilinear bandit problem where a learner (or a *central entity*) allocates arms to the nodes of a graph and observes for each edge a noisy bilinear reward representing the interaction between the two end nodes. We study the best arm identification problem in which the learner wants to find the graph allocation maximizing the sum of the bilinear rewards. By efficiently exploiting the geometry of this bandit problem, we propose a *decentralized* allocation strategy based on random sampling with theoretical guarantees. In particular, we characterize the influence of the graph structure (e.g. star, complete or circle) on the convergence rate and propose empirical experiments that confirm this dependency.

1. Introduction

In many multi-agent systems the contribution of an agent to a common team objective is impacted by the behavior of the other agents. The agents must coordinate (or be coordinated) to achieve the best team performance. Consider, for instance, the problem of configuring antennas of a wireless cellular network to obtain the best signal quality over the whole network (Siomina et al., 2006). The signal quality of the region covered by a given antenna might be degraded by the behavior of its neighboring antennas due to an increase of interferences or bad user handovers. Another example is the adjustment of the turbine blades of a wind farm where the best adjustment for one turbine may generate turbulence for its neighboring turbines and thus be suboptimal for the global wind farm objective (Bargiacchi et al., 2018).

These real-life problems can be viewed as instances of a *stochastic multi-agent multi-armed bandit* problem (Robbins, 1952; Bargiacchi et al., 2018) where a learner (or a *central entity*) sequentially pulls a joint arm, one arm for each

agent (e.g., all the configuration parameters of the antennas), and receives an associated global noisy reward (e.g., the signal quality over the whole network). The goal of the learner can either be to maximize the accumulated reward, implying a trade-off between exploration and exploitation, or to find the joint arm maximizing the reward, known as *pure exploration* or *best arm identification* (Bubeck et al., 2009; Audibert and Bubeck, 2010).

In this paper we focus on the best arm identification problem in a multi-agent system for which we assume the knowledge of a coordination graph $\mathcal{G} = (V, E)$ representing the agent interactions (Guestrin et al., 2002).

At each round t , a learner

1. chooses for each node $i \in V$ an arm $x_t^{(i)}$ in a finite arm set $\mathcal{X} \subset \mathbb{R}^d$,
2. observes for each edge $(i, j) \in E$ a bilinear reward $r_t^{(i,j)} = x_t^{(i)\top} \mathbf{M}_\star x_t^{(j)} + \eta_t^{(i,j)}$.

Here, we denote by $\mathbf{M}_\star \in \mathbb{R}^{d \times d}$ the unknown parameter matrix, and $\eta_t^{(i,j)}$ a zero-mean σ -sub-Gaussian random variable for all edges $(i, j) \in E$ and round t .

The goal of the central entity is to find, within a minimum number of rounds, the joint arm $(x_\star^{(1)}, \dots, x_\star^{(|V|)})$ such that the expected global reward $\sum_{(i,j) \in E} x_\star^{(i)\top} \mathbf{M}_\star x_\star^{(j)}$ is maximized.

The reward $r_t^{(i,j)}$ reflects the quality of the interaction between the neighboring nodes i and j when pulling respectively the arm $x_t^{(i)}$ and $x_t^{(j)}$ at time t . For instance, when configuring handover parameters of a wireless network, $r_t^{(i,j)}$ can be any criterion assessing the handover quality between antenna i and antenna j , the parameters selected by each antenna both impacting this quantity. The bilinear setting appears as a natural extension of the commonly studied linear setting to model the interaction between two agents. Furthermore, instead of a global reward being a sum of independent linear agent rewards, the global reward is now the result of the interactions between neighboring agents.

As exposed in Jun et al. (2019), the bilinear reward can be

¹PSL - Université Paris Dauphine, CNRS, LAMSADE, Paris, France ²Huawei Noah's Ark Lab ³Liverpool University. Correspondence to: Geovani Rizk <geovani.rizk@dauphine.psl.eu>.

written as a linear reward in a higher dimensional space:

$$r_t^{(i,j)} = \text{vec} \left(x_t^{(i)} x_t^{(j)\top} \right)^\top \text{vec}(\mathbf{M}_\star) + \eta_t^{(i,j)}, \quad (1)$$

where for any matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, $\text{vec}(\mathbf{A})$ denotes the vector in \mathbb{R}^{d^2} which is the concatenation of all the columns of \mathbf{A} .

Since the unknown parameter \mathbf{M}_\star is common to all the edges (i, j) of the graph, the expected global reward at time t can also be written as the scalar product $\left\langle \sum_{(i,j) \in E} \text{vec} \left(x_t^{(i)} x_t^{(j)\top} \right), \text{vec}(\mathbf{M}_\star) \right\rangle$. Hence, solving the best arm identification problem in the described graphical bilinear bandit boils down to solving the same problem in a global linear bandit. Although this trick allows to use classical algorithms in linear bandits, the number of joint arms is growing exponentially with the number of nodes, making such methods impractical.

Another possible way to address this problem based on equation (1) is to consider one linear bandit per edge, with constraints between edges. For more clarity, let us define the arm set $\mathcal{Z} = \{\text{vec}(xx'^\top) \mid (x, x') \in \mathcal{X}^2\}$, and let us refer to any $z \in \mathcal{Z}$ as an *edge-arm* and to any $x \in \mathcal{X}$ as an *node-arm*. At each round t , the learner chooses for each edge (i, j) an edge-arm in \mathcal{Z} with the constraint that for any pair of edges (i, j) and (i, k) , if the edge-arm $\text{vec}(xx'^\top)$ is assigned to the edge (i, j) and the edge-arm $\text{vec}(x''x'''^\top)$ is assigned to the edge (i, k) , then it must be that $x = x''$.

Given this constraint, how do we choose the appropriate sequence of edge-arms in order to build a good estimate of $\text{vec}(\mathbf{M}_\star)$? Moreover, assuming we have built such a good estimator, is there a tractable algorithm to identify the best joint arm, or at least to find a joint arm yielding a high expected reward? In this paper, we answer these questions and provide algorithms and theoretical guarantees.

We show that even with a perfect estimator $\text{vec}(\hat{\mathbf{M}}) = \text{vec}(\mathbf{M}_\star)$, identifying the best joint arm is NP-Hard. To address this issue, we design a polynomial time twofold algorithm. Given $\text{vec}(\hat{\mathbf{M}})$, it first identifies the best edge-arm $z_\star \in \mathcal{Z}$ maximizing $\langle z_\star, \text{vec}(\hat{\mathbf{M}}) \rangle$. Then, it allocates z_\star to a carefully chosen subset of edges. We show that this yields a good approximation ratio in Section 4.

To build our estimator $\text{vec}(\hat{\mathbf{M}})$, we rely on the G-Allocation strategy, as in Soare et al. (2014). We show that there exists a sampling procedure over the node-arms such that the associated edge-arms follow the optimal G-allocation strategy developed in the linear bandit literature. This procedure allows us to avoid the difficulty of having to satisfy the edge-arm constraints explicitly. Furthermore, we analyze the sample complexity of this method. This is detailed in Section 5.

In addition, we highlight the impact of the graph structure

in Section 6 and provide the explicit repercussion on the convergence rate of the algorithm for different types: star, complete, circle and matching graphs. In particular, we show that for favorable graph structures (e.g. circles), our convergence rate matches that of standard linear bandits. Finally, Section 7 evidences the theoretical findings on numerical experiments.

2. Related Work

Best arm identification in linear bandits. There exists a vast literature on the problem of best arm identification in linear bandits (Soare et al., 2014; Xu et al., 2018; Degenne et al., 2020; Kazerouni and Wein, 2019; Zaki et al., 2020; Jedra and Proutiere, 2020), would it be by using greedy strategies (Soare et al., 2014), rounding procedures (Fiez et al., 2019) or random sampling (Tao et al., 2018). Although our problem can be formulated as a linear bandit problem, none of the existing methods would scale-up with the number of agents. Nevertheless, we will be relying on classical techniques, and more specifically those developed in Soare et al. (2014).

Bilinear bandits. Bandits with bilinear rewards have been studied in Jun et al. (2019). The authors derived a no-regret algorithm based on Optimism in the Face of Uncertainty Linear bandit (OFUL) (Abbasi-Yadkori et al., 2011), using the fact that a bilinear reward can be expressed as a linear reward in higher dimension. Our work extends their setting by considering a set of dependent bilinear bandits. Besides, the goal here is to find the best arm rather than minimizing the regret.

Bandits and graphs. Graphs are often used to bring structure to a bandit problem. In Valko et al. (2014) and Mannor and Shamir (2011), the arms are the nodes of a graph and pulling an arm gives information on the rewards of the neighboring arms. The reader can also refer to Valko (2020) for an account on such problems. In Cesa-Bianchi et al. (2013) each node is an instance of a linear bandit and the neighboring nodes are assumed to have similar unknown regression coefficients. The main difference with our setting is that the rewards of the nodes are independent.

Combinatorial and multi-agent bandits. Allocating arms to each node of a graph to then observe a global reward is a combinatorial bandit problem (Cesa-Bianchi and Lugosi, 2012), the number of joint arms scaling exponentially with the number of nodes. This has been extensively studied both in the regret-based (Chen et al., 2013; Perrault et al., 2020) and the pure exploration context (Chen et al., 2014; Cao and Krishnamurthy, 2019; Jourdan et al., 2021; Du et al., 2020). Our problem is closer to the one presented in Amin et al. (2011) and Bargiacchi et al. (2018), where several agents want to maximize a global team reward that

can be decomposed into a sum of observable local rewards as in a *semi-bandit game* (Audibert et al., 2011; Chen et al., 2013). However, we study a more structured context as we assume observable bilinear rewards for each edge of the graph. Furthermore, note that our problem can be solved by the algorithm presented in Du et al. (2020) with a sample complexity increasing in the number of nodes. On the contrary, we propose in this paper an algorithm with a sample complexity decreasing in the number of nodes exploiting the structure of the bilinear reward and the graph. Finally, most of the algorithms developed for combinatorial bandits assume the availability of an oracle to solve the combinatorial optimization problem returning the arm to play or the final best arm recommendation. We make no such assumption.

3. Preliminaries and Notations

Let $\mathcal{G} = (V, E)$ be a directed graph with V the set of nodes, E the set of edges where we assume that if $(i, j) \in E$ then $(j, i) \in E$, and $\mathcal{N}(i)$ the set containing the neighbors of a node $i \in V$. We denote by $n = |V|$ the number of nodes and $m = |E|$ the number of edges. We define the *graphical bilinear bandit* on the graph \mathcal{G} as the setting where a learner sequentially pulls at each round t a joint arm $(x_t^{(1)}, \dots, x_t^{(n)}) \in \mathcal{X}^n$, also called graph allocation or simply allocation when it is clear from the context, and then receives a bilinear reward $r_t^{(i,j)}$ for each edge $(i, j) \in E$. At each round, the joint arm can be constructed simultaneously or sequentially, however all the bilinear rewards are only revealed after the joint arm has been pulled.

We denote $K = |\mathcal{X}|$ the number of node-arms and it is assumed that \mathcal{X} spans \mathbb{R}^d . For each round t of the learning procedure and each node $i \in V$, $x_t^{(i)} \in \mathcal{X}$ represents the node-arm allocated to the node $i \in V$. For each edge $(i, j) \in E$, we denote $z_t^{(i,j)} = \text{vec}(x_t^{(i)} x_t^{(j)\top}) \in \mathcal{Z}$ the associated chosen edge-arm.

The goal is to derive an algorithm that minimizes the number of pulled joint arms required to find the one maximizing the sum of the associated expected bilinear rewards, for a given confidence level. For the sake of simplicity, we assume that the unknown parameter matrix \mathbf{M}_\star in the bilinear reward is symmetric. We provide an analysis of the non-symmetric case in Appendix E.

For any finite set X , $\mathcal{S}_X \triangleq \{\lambda \in [0, 1]^{|X|}, \sum_{x \in X} \lambda_x = 1\}$ denotes the simplex in $\mathbb{R}^{|X|}$. For any vector $x \in \mathbb{R}^d$, $\|x\|$ will denote the ℓ_2 -norm of x . For any square matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, we denote by $\|\mathbf{A}\| \triangleq \sup_{x: \|x\|=1} \|\mathbf{A}x\|$ the spectral norm of \mathbf{A} . Finally, for any vector $x \in \mathbb{R}^d$ and a symmetric positive-definite matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, we define $\|x\|_{\mathbf{A}} \triangleq \sqrt{x^\top \mathbf{A} x}$.

4. An NP-Hard Problem

In this section, we address the problem of finding the best joint arm given \mathbf{M}_\star or a good estimator $\hat{\mathbf{M}}$. If the best edge-arm z_\star is composed of a single node-arm x_\star , that is $z_\star = \text{vec}(x_\star x_\star^\top)$, then finding the best joint arm is trivial and the solution is to assign x_\star to all nodes. Conversely, if z_\star is composed of two distinct node-arms (x_\star, x'_\star) , the problem is harder.

The following theorem states that, even with the knowledge of the true parameter \mathbf{M}_\star , identifying the best joint-arm is NP-Hard with respect to the number of nodes n .

Theorem 4.1. *Consider a given matrix $\mathbf{M}_\star \in \mathbb{R}^{d \times d}$ and a finite arm set $\mathcal{X} \subset \mathbb{R}^d$. Unless $P=NP$, there is no polynomial time algorithm guaranteed to find the optimal solution of*

$$\max_{(x^{(1)}, \dots, x^{(n)}) \in \mathcal{X}^n} \sum_{(i,j) \in E} x^{(i)\top} \mathbf{M}_\star x^{(j)} .$$

The proof of this theorem is in Appendix A and relies on a reduction to the Max-Cut problem. Hence, no matter which estimate $\hat{\mathbf{M}}$ of \mathbf{M}_\star one can build, the learner is not guaranteed to find in polynomial time the joint arm $(x_\star^{(1)}, \dots, x_\star^{(n)})$ maximizing the expected global reward. However, one can notice that, given the matrix \mathbf{M}_\star or even a good enough estimate $\hat{\mathbf{M}}$, identifying the edge-arm $z^\star = \text{vec}(x_\star x_\star^\top) \in \mathcal{Z}$ that maximizes the reward $z_\star^\top \text{vec}(\mathbf{M}_\star)$ requires only K^2 reward estimations (we simply estimate all the linear reward associated to each edge-arm in \mathcal{Z}). Thus, instead of looking for the best joint arm explicitly, we will first identify the best edge-arm z^\star , and then allocate z^\star to the largest number of edges in the graph. We will also show that this approach gives a guarantee on its associated global reward.

Let us consider the graph allocation that places the maximum number of edge-arms z_\star in \mathcal{G} . It is easy to show that the subgraph containing only the edges where z_\star has been pulled is the largest bipartite subgraph included in \mathcal{G} . Recall that a graph $\mathcal{G}' = (V', E')$ is a bipartite if and only if one can partition the node set $V' = (V'_1, V'_2)$ such that

$$(i, j) \in E' \Rightarrow (i, j) \in V'_1 \times V'_2 \text{ or } (j, i) \in V'_1 \times V'_2 .$$

Notice, that if \mathcal{G}' is the largest bipartite subgraph in \mathcal{G} , the number of edges in E' is the maximal number of edge-arms z_\star that can be allocated with a single graph allocation.

Hence, finding the joint arm with the largest number of edge-arms z_\star allocated in the graph is equivalent to finding the largest bipartite subgraph $\mathcal{G}' = (V', E')$ in \mathcal{G} . Once that subgraph is determined, we just need to allocate to all the nodes in V'_1 the node-arm x_\star and to all the nodes of V'_2 the node-arm x'_\star (which is equivalent to allocating to all the edges in E' the edge-arm z_\star).

Furthermore, we know that every m -edge graph contains a bipartite subgraph of at least $m/2$ edges (Erdos, 1975). Therefore, we propose Algorithm 1 which iteratively constructs a bipartite subgraph and allocates the nodes accordingly to create at least $m/2$ edge-arms z_* .

Algorithm 1 Bipartite graph algorithm for Best Arm Identification in Graphical Bilinear Bandits

Input : $\mathcal{G} = (V, E)$, \mathcal{X} , \mathbf{M}
 Find $(x_*, x'_*) \in \arg \max_{(x, x') \in \mathcal{X}^2} x^\top \mathbf{M} x'$
 Set $V_1 = \emptyset$, $V_2 = \emptyset$
for i **in** V **do**
 Set n_1 the number of neighbors of i in V_1
 Set n_2 the number of neighbors of i in V_2
 if $n_1 > n_2$ **then**
 $x^{(i)} = x'_*$
 $V_2 \leftarrow V_2 \cup \{i\}$
 else
 $x^{(i)} = x_*$
 $V_1 \leftarrow V_1 \cup \{i\}$
 end
end
 return $\mathbf{x} = (x^{(1)}, \dots, x^{(n)})$

The following result gives the guarantee on the global reward associated to the joint arm returned by Algorithm 1. We refer the reader to Appendix A for the proof.

Theorem 4.2. *Let us consider the graph $\mathcal{G} = (V, E)$, a finite arm set $\mathcal{X} \subset \mathbb{R}^d$ and the matrix \mathbf{M}_* given as input to Algorithm 1. Then, the expected global reward $r = \sum_{(i,j) \in E} x^{(i)\top} \mathbf{M}_* x^{(j)}$ associated to the returned allocation $\mathbf{x} = (x^{(1)}, \dots, x^{(n)}) \in \mathcal{X}^n$ verifies:*

$$\frac{r - r_{\min}}{r_* - r_{\min}} \geq \frac{1}{2}.$$

where r_* and r_{\min} are respectively the highest and lowest global reward one can obtain with the appropriate joint arm. Finally, the complexity of the algorithm is in $\mathcal{O}(K^2 + n)$.

This type of approximation result is sometimes referred to as *differential approximation* or *z -approximation*, and is often viewed as a more subtle analysis than standard approximation ratio. We emphasize that finding a better ratio than $\frac{1}{2}$ is a very hard task: such a finding would immediately yield an improved differential approximation ratio for the Max-Cut problem, which is an opened problem since 2001 (Hassin and Khuller, 2001).

5. Construction of the Estimate $\hat{\mathbf{M}}$

In the previous section, we designed a polynomial time method that computes a $1/2$ -approximation to the NP-Hard problem of finding the best joint arm given \mathbf{M}_* . Notice that

\mathbf{M}_* is only used to identify the best edge-arm z_* . Thus, using an estimate $\hat{\mathbf{M}}$ of \mathbf{M}_* having the following property:

$$\arg \max_{z \in \mathcal{Z}} z^\top \text{vec}(\hat{\mathbf{M}}) = \arg \max_{z \in \mathcal{Z}} z^\top \text{vec}(\mathbf{M}_*) , \quad (2)$$

would still allow us to identify z_* , and would thus give us the same guarantees.

In this section we tackle the problem of pulling the edge-arms during the learning procedure such that the estimated unknown parameter verifies (2) in as few iterations as possible. To do so, we first formalize the objective related to the linearized version of the problem. Then, we propose an algorithm reaching the given objective with high probability while satisfying the edge-arms constraints.

We denote by $\theta_* = \text{vec}(\mathbf{M}_*)$ the parameter of the linearized problem and $\hat{\theta}_t$ the Ordinary Least Squares (OLS) estimate of θ_* computed with all the data collected up to round t . The empirical gap between two edge-arms z and z' in \mathcal{Z} is denoted $\hat{\Delta}_t(z, z') \triangleq (z - z')^\top \hat{\theta}_t$.

5.1. A Constrained G-Allocation

The goal here is to define the optimal sequence $(z_1, \dots, z_{mt}) \in \mathcal{Z}^{mt}$ that should be pulled in the first t rounds so that (2) is reached as soon as possible. A natural approach is to rely on classical strategies developed for best arm identification in linear bandits. Most of the known strategies (see e.g., Soare et al. (2014); Xu et al. (2018); Fiez et al. (2019)) are based on a bound of the gap error $|(\theta_* - \hat{\theta}_t)^\top (z - z')|$ for all $z, z' \in \mathcal{Z}$. This bound is then used to derive a stopping condition, indicating a sufficient number of rounds t after which the OLS estimate $\hat{\theta}_t$ is precise enough to ensure the identification of the best edge-arm, with high probability.

Let $\delta \in (0, 1)$ and let $\mathbf{A}_t = \sum_{i=1}^{mt} z_i z_i^\top$ be the matrix computed with the mt edge-arms constructed during t rounds. Following the steps of Soare et al. (2014), we can show that if there exists $z \in \mathcal{Z}$ such that for all $z' \in \mathcal{Z}$ the following holds:

$$\|z - z'\|_{\mathbf{A}_t^{-1}} \sqrt{8\sigma^2 \log \left(\frac{6m^2 t^2 K^4}{\delta \pi^2} \right)} \leq \hat{\Delta}_t(z, z') , \quad (3)$$

then with probability at least $1 - \delta$, the OLS estimate $\hat{\theta}_t$ leads to the best edge-arm. Details of the derivation are given in Appendix B.

As mentioned in Soare et al. (2014), by noticing that $\max_{(z, z') \in \mathcal{Z}^2} \|z - z'\|_{\mathbf{A}_t^{-1}} \leq 2 \max_{z \in \mathcal{Z}} \|z\|_{\mathbf{A}_t^{-1}}$, an admissible strategy is to pull edge-arms minimizing $\max_{z \in \mathcal{Z}} \|z\|_{\mathbf{A}_t^{-1}}$ in order to satisfy the stopping condition as soon as possible. More formally, one wants to find the

sequence of edge-arms $\mathbf{z}_{mt}^* = (z_1^*, \dots, z_{mt}^*)$ such that:

$$\mathbf{z}_{mt}^* \in \arg \min_{(z_1, \dots, z_{mt})} \max_{z' \in \mathcal{Z}} z'^{\top} \left(\sum_{i=1}^{mt} z_i z_i^{\top} \right)^{-1} z' . \quad (\text{G-opt-}\mathcal{Z})$$

This is known as *G-allocation* (see e.g., Pukelsheim (2006); Soare et al. (2014)) and is NP-hard to compute (Çivril and Magdon-Ismail, 2009; Welch, 1982). One way to find an approximate solution is to rely on a convex relaxation of the optimization problem (G-opt- \mathcal{Z}) and first compute a real-valued allocation $\lambda^* \in \mathcal{S}_{\mathcal{Z}}$ such that

$$\lambda^* \in \arg \min_{\lambda \in \mathcal{S}_{\mathcal{Z}}} \max_{z' \in \mathcal{Z}} z'^{\top} \left(\sum_{z \in \mathcal{Z}} \lambda_z z z^{\top} \right)^{-1} z' . \quad (\text{G-relaxed-}\mathcal{Z})$$

One could either use random sampling to draw edge-arms as i.i.d. samples from the λ^* distribution or rounding procedures to efficiently convert each λ_z^* into an integer. However, these methods do not take into account the graphical structure of the problem, and at a given round, the m chosen edge-arms may result in two different assignments for the same node. Therefore, random sampling or rounding procedures cannot be straightforwardly used to select edge-arms in \mathcal{Z} . Nevertheless, (G-relaxed- \mathcal{Z}) still gives a valuable information on the number of times, in proportion, each edge-arm $z \in \mathcal{Z}$ must be allocated to the graph. In the next section, we present an algorithm satisfying both the proportion requirements and the graphical constraints.

5.2. Random Allocation over the Nodes

Our algorithm is based on a randomized method directly allocating node-arms to the nodes and thus avoiding the difficult task of choosing edge-arms and trying to allocate them to the graph while ensuring that every node has a unique assignment. The validity of this random allocation is based on Theorem 5.1 below showing that one can draw node-arms in \mathcal{X} and allocate them to the graph such that the associated edge-arms follow the probability distribution λ^* solution of (G-relaxed- \mathcal{Z}).

Theorem 5.1. *Let μ^* be a solution of the following optimization problem:*

$$\min_{\mu \in \mathcal{S}_{\mathcal{X}}} \max_{x' \in \mathcal{X}} x'^{\top} \left(\sum_{x \in \mathcal{X}} \mu_x x x^{\top} \right)^{-1} x' . \quad (\text{G-relaxed-}\mathcal{X})$$

Let $\lambda^* \in \mathcal{S}_{\mathcal{Z}}$ be defined for all $z = \text{vec}(x x'^{\top}) \in \mathcal{Z}$ by $\lambda_z^* = \mu_x^* \mu_{x'}^*$. Then, λ^* is a solution of (G-relaxed- \mathcal{Z}).

Sketch of proof. The objective at the optimum in (G-relaxed- \mathcal{X}) and (G-relaxed- \mathcal{Z}) are respectively equal to

d and d^2 which is the dimension of their respective problem, a result known as the Equivalence Theorem (Kiefer and Wolfowitz, 1960). Thus, by multiplying the optimum value of (G-relaxed- \mathcal{X}) by itself, we can show that for all $z \in \mathcal{Z}$ where $z = \text{vec}(x x'^{\top})$ with $(x, x') \in \mathcal{X}^2$, λ_z^* can be written as the product $\mu_x^* \mu_{x'}^*$. We refer to the Appendix C for the detailed proof.

This theorem implies that, at each round $t > 0$ and each node $i \in V$, if $x_t^{(i)}$ is drawn from μ^* , then for all pairs of neighbors $(i, j) \in E$ the probability distribution of the associated edge-arms $z_t^{(i,j)}$ follows λ^* . Moreover, as μ^* is a distribution over the node-arm set \mathcal{X} , λ^* is a joint (product) probability distribution on \mathcal{X}^2 with marginal μ^* .

We apply the Frank-Wolfe algorithm (Frank et al., 1956) to compute the solution μ^* of (G-relaxed- \mathcal{X}), as it is more suited to optimization tasks on the simplex than projected gradient descent. Although we face a min-max optimization problem, we notice that the function $h(\mu) = \max_{x' \in \mathcal{X}} x'^{\top} \left(\sum_{x \in \mathcal{X}} \mu_x x x^{\top} \right)^{-1} x'$ is convex. We refer the reader to Appendix F and references therein for a proof on the convexity of h and a discussion about using Frank-Wolfe for solving (G-relaxed- \mathcal{X}).

Given the characterization in Theorem 5.1 and our objective to verify the stopping condition in (3), we present our sampling procedure in Algorithm 2. We also note that at each round the sampling of the node-arms can be done in parallel.

Algorithm 2 Randomized G-Allocation strategy for Graphical Bilinear Bandits

Input : graph $\mathcal{G} = (V, E)$, arm set \mathcal{X}

Set $A_0 = I$; $b_0 = 0$; $t = 1$;

Apply the Frank-Wolfe algorithm to find μ^* solution of (G-relaxed- \mathcal{X}).

while stopping condition (3) is not verified **do**

// Sampling the node-arms

Draw $x_t^{(1)}, \dots, x_t^{(n)} \stackrel{\text{iid}}{\sim} \mu^*$ and obtain for all (i, j) in E the rewards $r_t^{(i,j)}$;

// Estimating $\hat{\theta}_t$ with the associated edge-arms

$\mathbf{A}_t = \mathbf{A}_{t-1} + \sum_{(i,j) \in E} z_t^{(i,j)} z_t^{(i,j)\top}$;

$b_t = b_{t-1} + \sum_{(i,j) \in E} z_t^{(i,j)} r_t^{(i,j)}$;

$\hat{\theta}_t = \mathbf{A}_t^{-1} b_t$

$t \leftarrow t + 1$;

end

return $\hat{\theta}_t$

This sampling procedure implies that each edge-arm follows the optimal distribution λ^* . However, if we take the number of times each $z \in \mathcal{Z}$ appears in the m pulled edge-arms of a given round, we might notice that the observed proportion

is not close to λ_z^* , regardless of the size of m . This is due to the fact that the m edge-arms are not independent because of the graph structure (*cf.* Section 6). Conversely, since each group of m edge-arms are independent from one round to another, the proportion of each $z \in \mathcal{Z}$ observed among the mt pulled edge-arms throughout t rounds is close to λ_z^* .

One may wonder if deterministic rounding procedures could be used instead of random sampling on μ^* , as it is done in many standard linear bandit algorithms (Soare et al., 2014; Fiez et al., 2019). Applying rounding procedure on μ^* gives the number of times each node-arm $x \in \mathcal{X}$ should be allocated to the graph. However, it does not provide the actual allocations that the learner must choose over the t rounds to optimally pull the associated edge-arms (*i.e.*, pull edge-arms following λ^*). Thus, although rounding procedures give a more precise number of times each node-arm should be pulled, the problem of allocating them to the graph remains open, whereas by concentration of the measure, randomized sampling methods imply that the associated edge-arms follow the optimal probability distribution λ^* . In this paper, we present a simple and standard randomized G-allocation strategy, but other more elaborated methods could be considered, as long as they include the necessary randomness.

On the choice of the G-allocation problem. We have considered the G-allocation optimization problem (G-opt- \mathcal{Z}), however, one could want to directly minimize $\max_{(z, z') \in \mathcal{Z}^2} \|z - z'\|_{\mathbf{A}_t^{-1}}$, known as the XY-allocation (Soare et al., 2014; Fiez et al., 2019). Hence, one may want to construct edge-arms that follow the distribution λ_{XY}^* solution of the relaxed XY-allocation problem:

$$\min_{\lambda} \max_{z', z''} (z' - z'')^\top \left(\sum_{z \in \mathcal{Z}} \lambda_z z z^\top \right)^{-1} (z' - z'').$$

Although efficient in the linear case, this approach outputs a distribution λ_{XY}^* which is not a joint probability distribution of two independent random variables, and so cannot be decomposed as the product of its marginals. Hence, there is no algorithm that allocates *identically* and *independently* the nodes of the graph to create edge-arms following λ_{XY}^* . Thus, we will rather deal with the upper bound given by the G-allocation as it allows sampling over the nodes.

Static design versus adaptive design. Adaptive designs as proposed for example in Soare et al. (2014) and Fiez et al. (2019) provide a strong improvement over static designs in the case of linear bandits. In our particular setting however, it is crucial to be able to adapt the edge-arms sampling rule to the node-arms, which is possible thanks to Theorem 5.1. This result requires a set of edge-arms \mathcal{Z} expressed as a product of node-arms set \mathcal{X} . Extending the adaptive design of Fiez et al. (2019) to our setting would eliminate edge-arms from \mathcal{Z} at each phase, without trivial guarantees that the newly obtained edge-arms set $\mathcal{Z}' \subset \mathcal{Z}$ could still be

derived from another node-arms set $\mathcal{X}' \subset \mathcal{X}$. An adaptive approach is definitely a natural and promising extension of our method, and is left for future work.

5.3. Convergence Analysis

We now prove the validity of the random sampling procedure detailed in Algorithm 2 by controlling the quality of the approximation $\max_{z \in \mathcal{Z}} z^\top \mathbf{A}_t^{-1} z$ with respect to the optimum of the G-allocation optimization problem $\max_{z' \in \mathcal{Z}} z'^\top \left(\sum_{i=1}^{mt} z_i^* z_i^{*\top} \right)^{-1} z'$ described in (G-opt- \mathcal{Z}). As is usually done in the optimal design literature (see *e.g.*, Pukelsheim (2006); Soare et al. (2014); Sagnol (2010)) we bound the relative error α :

$$\max_{z \in \mathcal{Z}} z^\top \mathbf{A}_t^{-1} z \leq (1 + \alpha) \max_{z' \in \mathcal{Z}} z'^\top \left(\sum_{i=1}^{mt} z_i^* z_i^{*\top} \right)^{-1} z'.$$

Our analysis relies on several results from matrix concentration theory. One may refer for instance to Tropp et al. (2015) and references therein for an extended introduction on that matter. We first introduce a few additional notations.

Let $f_{\mathcal{Z}}$ be the function such that for any non-singular matrix $\mathbf{Q} \in \mathbb{R}^{d^2 \times d^2}$, $f_{\mathcal{Z}}(\mathbf{Q}) = \max_{z \in \mathcal{Z}} z^\top \mathbf{Q}^{-1} z$ and for any distribution $\lambda \in \mathcal{S}_{\mathcal{Z}}$ let $\Sigma_{\mathcal{Z}}(\lambda) \triangleq \sum_{z \in \mathcal{Z}} \lambda_z z z^\top$ be the associated covariance matrix. Finally let $\mathbf{A}_t^* = \sum_{i=1}^{mt} z_i^* z_i^{*\top}$ be the G-optimal design matrix constructed during t rounds.

For $i \in \{1, \dots, n\}$ and $s \in \{1, \dots, t\}$, let $X_s^{(i)}$ be i.i.d. random vectors in \mathcal{X} such that for all $x \in \mathcal{X}$,

$$\mathbb{P}\left(X_1^{(1)} = x\right) = \mu_x^*.$$

Each $X_s^{(i)}$ is to be viewed as the random arm pulled at round s for the node i . Using this notation, the random design matrix \mathbf{A}_t can be defined as

$$\mathbf{A}_t = \sum_{s=1}^t \sum_{(i,j) \in E} \text{vec}\left(X_s^{(i)} X_s^{(j)\top}\right) \text{vec}\left(X_s^{(i)} X_s^{(j)\top}\right)^\top.$$

One can first observe that $f_{\mathcal{Z}}(\mathbf{A}_t)$ can be bounded by the following quantity:

$$\begin{aligned} f_{\mathcal{Z}}(\mathbf{A}_t) &= \max_{z \in \mathcal{Z}} z^\top \left(\mathbf{A}_t^{-1} - (\mathbb{E}\mathbf{A}_t)^{-1} + (\mathbb{E}\mathbf{A}_t)^{-1} \right) z \\ &\leq \max_{z \in \mathcal{Z}} z^\top \left(\mathbf{A}_t^{-1} - (\mathbb{E}\mathbf{A}_t)^{-1} \right) z \\ &\quad + f_{\mathcal{Z}}(mt\Sigma_{\mathcal{Z}}(\lambda^*)) \\ &\leq \max_{z \in \mathcal{Z}} \|z\|^2 \|\mathbf{A}_t^{-1} - (\mathbb{E}\mathbf{A}_t)^{-1}\| \\ &\quad + f_{\mathcal{Z}}(\mathbf{A}_t^*). \end{aligned}$$

Hence, one needs a bound on the maximum eigenvalue of $\mathbf{A}_t^{-1} - (\mathbb{E}\mathbf{A}_t)^{-1}$. Simple linear algebra leads to:

$$\mathbf{A}_t^{-1} - (\mathbb{E}\mathbf{A}_t)^{-1} = \mathbf{A}_t^{-1}(\mathbb{E}\mathbf{A}_t - \mathbf{A}_t)(\mathbb{E}\mathbf{A}_t)^{-1}.$$

Thus, in addition to bounding the maximum eigenvalue of \mathbf{A}_t^{-1} , which is equal to the minimum eigenvalue of \mathbf{A}_t , we need a bound on $\|\mathbf{A}_t - \mathbb{E}\mathbf{A}_t\|$. It may be derived from concentration results on sum of random matrices derived in [Tropp et al. \(2015\)](#). We now state the result controlling the relative error obtained with our randomized sampling allocation. The proof can be found in the Appendix C.

Theorem 5.2. *Let λ^* be a solution of the optimization problem (G-relaxed- \mathcal{Z}). Let $0 \leq \delta \leq 1$ and let t_0 be such that*

$$t_0 = 2Ld^2 \log(2d^2/\delta)/\nu_{\min},$$

where $L = \max_{z \in \mathcal{Z}} \|z\|^2$ and ν_{\min} is the smallest eigenvalue of the covariance matrix $\frac{1}{K^2} \sum_{z \in \mathcal{Z}} zz^\top$. Then, at each round $t \geq t_0$ with probability at least $1 - \delta$, the randomized G-allocation strategy for graphical bilinear bandit in [Algorithm 2](#) produces a matrix \mathbf{A}_t such that:

$$f_{\mathcal{Z}}(\mathbf{A}_t) \leq (1 + \alpha)f_{\mathcal{Z}}(\mathbf{A}_t^*)$$

where

$$\alpha = \frac{Ld^2}{m\nu_{\min}^2} \sqrt{\frac{2v}{t} \log\left(\frac{2d^2}{\delta}\right)} + o\left(\frac{1}{\sqrt{t}}\right),$$

and $v \triangleq \mathbb{E}[(\mathbf{A}_1 - \mathbb{E}\mathbf{A}_1)^2]$.

We have just shown that the approximation value $\max_{z \in \mathcal{Z}} z^\top \mathbf{A}_t^{-1} z$ converges to the optimal value with a rate of $O(\sqrt{v}/(m\sqrt{t}))$. In [Section 6](#), we show that the best case graph implies a $v = O(m)$ matching the convergence rate $O(1/\sqrt{mt})$ of a linear bandit algorithm using randomized sampling to pull mt edge-arms without (graphical) constraints. Moreover, we will see that the worst case graph implies that $v = O(m^2)$.

Since we filled the gap between our constraint objective and the problem of best arm identification in linear bandits, thanks to [Theorem 5.1](#) and [5.2](#), we are able to extend known results for best arm identification in linear bandits on the sample complexity and its associated lower bound.

Corollary 5.3 ([Soare et al. \(2014\)](#), [Theorem 1](#)). *If the G-allocation is implemented with the random strategy of [Algorithm 2](#), resulting in an α -approximation, then with probability at least $1 - \delta$, the best arm obtained with $\hat{\theta}_t$ is z_* and*

$$t \leq \frac{128\sigma^2 d^2 (1 + \alpha) \log\left(\frac{6m^2 t^2 K^4}{\delta\pi}\right)}{m\Delta_{\min}^2},$$

where $\Delta_{\min} = \min_{z \in \mathcal{Z} \setminus \{z_*\}} (z_* - z)^\top \theta_*$.

Moreover, let τ be the number of rounds sufficient for any algorithm to determine the best arm with probability at least $1 - \delta$. A lower bound on the expectation of τ can be

obtained from the one derived for the problem of best arm identification in linear bandits (see *e.g.*, [Theorem 1](#) in [Fiez et al. \(2019\)](#)):

$$\mathbb{E}[\tau] \geq \min_{\lambda \in \mathcal{S}_{\mathcal{Z}}} \max_{z \in \mathcal{Z} \setminus \{z_*\}} \log\left(\frac{1}{2.4\delta}\right) \frac{2\sigma^2 \|z_* - z\|_{\Sigma_{\mathcal{Z}}(\lambda)}^2}{m \left((z_* - z)^\top \theta_*\right)^2}.$$

As observed in [Soare et al. \(2014\)](#) this lower bound can be upper bounded, in the worst case, by $4\sigma^2 d^2 / (m\Delta_{\min}^2)$ which matches our bound up to log terms and the relative error α .

6. Influence of the Graph Structure on v

The convergence bound in [Theorem 5.2](#) depends on $v = \mathbb{E}[(\mathbf{A}_1 - \mathbb{E}\mathbf{A}_1)^2]$. In this section, we characterize the impact of the graph structure on this quantity and, by extension, on the convergence rate. First of all, recall that

$$\mathbf{A}_1 = \sum_{(i,j) \in E} \text{vec}\left(X_1^{(i)} X_1^{(j)\top}\right) \text{vec}\left(X_1^{(i)} X_1^{(j)\top}\right)^\top.$$

Let denote $\mathbf{A}_1^{(i,j)} = \text{vec}\left(X_1^{(i)} X_1^{(j)\top}\right) \text{vec}\left(X_1^{(i)} X_1^{(j)\top}\right)^\top$ such that $\mathbf{A}_1 = \sum_{(i,j) \in E} \mathbf{A}_1^{(i,j)}$ and let define for any random matrices \mathbf{A} and \mathbf{B} the operators $\text{Var}(\mathbf{A}) \triangleq \mathbb{E}[(\mathbf{A} - \mathbb{E}[\mathbf{A}])^2]$ and $\text{Cov}(\mathbf{A}, \mathbf{B}) \triangleq \mathbb{E}[(\mathbf{A} - \mathbb{E}[\mathbf{A}])(\mathbf{B} - \mathbb{E}[\mathbf{B}])]$. We can derive the variance of \mathbf{A}_1 as follows:

$$\begin{aligned} \text{Var}(\mathbf{A}_1) &= \sum_{(i,j) \in E} \text{Var}\left(\mathbf{A}_1^{(i,j)}\right) \\ &+ \sum_{(i,j) \in E} \sum_{\substack{(k,l) \in E \\ (k,l) \neq (i,j)}} \text{Cov}\left(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(k,l)}\right). \end{aligned}$$

One can decompose the sum of the covariances into three groups: a first group where $k \neq i, j$ and $l \neq i, j$ which means that the two edges do not share any node and $\text{Cov}(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(k,l)}) = \mathbf{0}$, and two other groups where the edges share at least one node. For all edges $(i, j) \in E$ we consider either the edges $(i, k) \in E$ where $k \neq j$, yielding $\text{Cov}(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(i,k)})$ or the edges $(j, k) \in E$, yielding $\text{Cov}(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(j,k)})$.

Hence, one has

$$\begin{aligned} \text{Var}(\mathbf{A}_1) &= \sum_{(i,j) \in E} \text{Var}\left(\mathbf{A}_1^{(i,j)}\right) \\ &+ \sum_{i=1}^n \sum_{j \in \mathcal{N}(i)} \sum_{\substack{k \in \mathcal{N}(i) \\ k \neq j}} \text{Cov}\left(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(i,k)}\right) \\ &+ \sum_{i=1}^n \sum_{j \in \mathcal{N}(i)} \sum_{k \in \mathcal{N}(j)} \text{Cov}\left(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(j,k)}\right). \end{aligned}$$

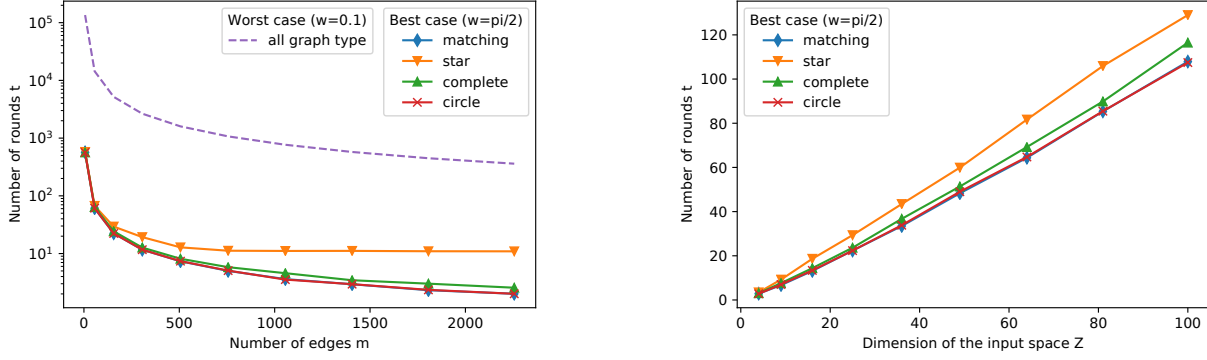


Figure 1. Number of rounds t needed to verify the stopping condition (3) with respect to **left**: the number of edges m where the dimension of the edge-arm space \mathcal{Z} is fixed and equal to 25 and **right**: the dimension of the edge-arm space \mathcal{Z} where the number of edges is fixed and equal to 156. For both experiments we run 100 times and plot the average number of rounds needed to verify the stopping condition.

Let $P \geq 0$ be such that for all $(i, j) \in E$, $\text{Var}(\mathbf{A}_1^{(i,j)}) \preceq P \times \mathbf{I}$ and $M, N \geq 0$ such that for all $(i, j) \in E$:

$$\begin{aligned} \forall k \in \mathcal{N}(i), \text{Cov}(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(i,k)}) &\preceq M \times \mathbf{I} \\ \forall k \in \mathcal{N}(j), \text{Cov}(\mathbf{A}_1^{(i,j)}, \mathbf{A}_1^{(j,k)}) &\preceq N \times \mathbf{I} \end{aligned}$$

We want to compare the quantity $\|\text{Var}(\mathbf{A}_1)\|$ for different types of graphs: star, complete, circle and a matching graph. To have a fair comparison, we want graphs that reveal the same number of rewards at each round of the learning procedure. Hence, we denote respectively n_S , n_{Co} , n_{Ci} and n_M the number of nodes in a star, complete, circle and matching graph of m edges and get:

Star graph:

$$\|\text{Var}(\mathbf{A}_1)\| \leq mP + n_S^2(M + N).$$

Complete graph:

$$\|\text{Var}(\mathbf{A}_1)\| \leq mP + n_{Co}^3(M + N).$$

Circle graph:

$$\|\text{Var}(\mathbf{A}_1)\| \leq mP + n_{Ci}(2M + 4N).$$

Matching graph:

$$\|\text{Var}(\mathbf{A}_1)\| \leq mP + n_M N.$$

We refer the reader to Appendix D for more details on the given upper bounds. Since the star (respectively, complete, circle and matching) graph of m edges has a number of nodes $n_S = m/2 + 1$ (respectively $n_{Co} = (1 + \sqrt{4m + 1})/2$, $n_{Ci} = m/2$ and $n_M = m$), we obtain the bounds stated in Table 1.

Graph	Upper bound on $\ \text{Var}(\mathbf{A}_1)\ $	α
Star	$mP + (M + N)O(m^2)$	$O(1/\sqrt{t})$
Complete	$mP + (M + N)O(m\sqrt{m})$	$O\left(1/\left(m^{\frac{1}{4}}\sqrt{t}\right)\right)$
Circle	$mP + (M + N)O(m)$	$O(1/\sqrt{mt})$
Matching	$mP + mN$	$O(1/\sqrt{mt})$

Table 1. Upper bound on the variance and convergence rate of Algorithm 2 for the star, complete, circle and matching graph with respect to the number of edges m and the number of rounds t .

These four examples evidence the strong dependency of the variance on the structure of the graph. The more independent the edges are (*i.e.*, with no common nodes), the smaller the quantity $\|\text{Var}(\mathbf{A}_1)\|$ is. For a fixed number of edges m , the best case is the matching graph where no edge share the same node and the worst case is the star graph where all the edges share a central node.

7. Experiments

In this section, we consider the modified version of a standard experiment introduced by Soare et al. (2014) and used in most papers on best arm identification in linear bandits (Xu et al., 2018; Tao et al., 2018; Fiez et al., 2019; Zaki et al., 2019) to evaluate the sample complexity of our algorithm on different graphs. We consider $d + 1$ node-arms in $\mathcal{X} \subset \mathbb{R}^d$ where $d \geq 2$. This node-arm set is made of the d vectors $(\mathbf{e}_1, \dots, \mathbf{e}_d)$ forming the canonical basis of \mathbb{R}^d and one additional arm $x_{d+1} = (\cos(\omega), \sin(\omega), 0, \dots, 0)^\top$ with $\omega \in [0, \pi/2]$. Note that by construction, the edge-arm set \mathcal{Z} contains the canonical basis $(\mathbf{e}'_1, \dots, \mathbf{e}'_{d^2})$ of \mathbb{R}^{d^2} . The parameter matrix \mathbf{M}_* has its first coordinate

equal to 2 and the others equal to 0 which makes $\theta_* = \text{vec}(\mathbf{M}_*) = (2, 0, \dots, 0)^\top \in \mathbb{R}^{d^2}$. The best edge-arm is thus $z_* = z^{(1,1)} = \mathbf{e}'_1$. One can note that when ω tends to 0, it is harder to differentiate this arm from $z^{(d+1,d+1)} = \text{vec}(x_{(d+1)}x_{(d+1)}^\top)$ than from the other arms.

We set $\eta_t^{(i,j)} \sim \mathcal{N}(0, 1)$, for all edges (i, j) and round t .

We consider the two cases where $\omega = 0.1$ which makes the edge-arms $z^{(1,1)}$ and $z^{(d+1,d+1)}$ difficult to differentiate, and $\omega = \pi/2$ which makes the edge-arm $z^{(1,1)}$ easily identifiable as the optimal edge-arm. For each of these two cases, we evaluate the influence of the graph structure, the number of edges m and the edge-arm space dimension d^2 on the sampling complexity. Results are shown in Figure 1.

When $\omega = 0.1$, the type of the graph does not impact the number of rounds needed to verify the stopping condition. This is mainly due to the fact that the magnitude of its associated variance is negligible with respect to the number of rounds. Hence, even if we vary the number of edges or the dimension, we get the same performance for any type of graph including the matching graph. This implies that our algorithm performs as well as a linear bandit that draws m edge-arms in parallel at each round. When $\omega = \pi/2$, the number of rounds needed to verify the stopping condition is smaller and the magnitude of the variance is no longer negligible. Indeed, when the number of edges or the dimension increases, we notice that the star graph takes more times to satisfy the stopping condition. Moreover, note that the sample complexities obtained for the circle and the matching graph are similar. This observation is in line with the dependency on the variance shown in Table 1.

8. Conclusion

We introduced a new graphical bilinear bandit setting and studied the best arm identification problem with a fixed confidence. This problem being NP-Hard even with the knowledge of the true parameter matrix \mathbf{M}^* , we first proposed an algorithm that provides a 1/2-approximation. Then, we provided a second algorithm, based on G-allocation strategy, that uses randomized sampling over the nodes to return a good estimate $\hat{\mathbf{M}}$ that can be used instead of \mathbf{M}_* . Finally, we highlighted the impact of the graph structure on the convergence rate of our algorithm and validated our theoretical results with experiments. Promising extensions of the model include considering unknown parameters $\mathbf{M}_*^{(i,j)}$, different for each edge (i, j) of the graph, and investigating XY-allocation strategies.

Acknowledgements

We thank anonymous reviewers, whose comments helped us improve the paper significantly.

References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. (2011). Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320.
- Amin, K., Kearns, M., and Syed, U. (2011). Graphical models for bandit problems. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, page 1–10.
- Audibert, J.-Y. and Bubeck, S. (2010). Best arm identification in multi-armed bandits. In *Proceedings of the 23th Annual Conference on Learning Theory*, pages 41–53.
- Audibert, J.-Y., Bubeck, S., and Lugosi, G. (2011). Minimax policies for combinatorial prediction games. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 107–132.
- Bargiacchi, E., Verstraeten, T., Roijers, D., Nowé, A., and Hasselt, H. (2018). Learning to coordinate with coordination graphs in repeated single-stage multi-agent decision problems. In *International conference on machine learning*, pages 482–490.
- Bubeck, S., Munos, R., and Stoltz, G. (2009). Pure exploration in multi-armed bandits problems. In *International conference on Algorithmic learning theory*, pages 23–37. Springer.
- Cao, T. and Krishnamurthy, A. (2019). Disagreement-based combinatorial pure exploration: Sample complexity bounds and an efficient algorithm. In *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99, pages 558–588.
- Cesa-Bianchi, N., Gentile, C., and Zappella, G. (2013). A gang of bandits. In *Advances in Neural Information Processing Systems*, pages 737–745.
- Cesa-Bianchi, N. and Lugosi, G. (2012). Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404 – 1422.
- Chen, S., Lin, T., King, I., Lyu, M. R., and Chen, W. (2014). Combinatorial pure exploration of multi-armed bandits. In *Advances in Neural Information Processing Systems*, volume 27, pages 379–387.
- Chen, W., Wang, Y., and Yuan, Y. (2013). Combinatorial multi-armed bandit: General framework and applications. In *International Conference on Machine Learning*, pages 151–159.
- Degenne, R., Ménard, P., Shang, X., and Valko, M. (2020). Gamification of pure exploration for linear bandits. *arXiv preprint arXiv:2007.00953*.

- Du, Y., Kuroki, Y., and Chen, W. (2020). Combinatorial pure exploration with full-bandit or partial linear feedback. *arXiv e-prints*, pages arXiv-2006.
- Erdos, P. (1975). Problems and results on finite and infinite graphs. In *Recent advances in graph theory (Proc. Second Czechoslovak Sympos., Prague, 1974)*, pages 183–192.
- Fiez, T., Jain, L., Jamieson, K. G., and Ratliff, L. (2019). Sequential experimental design for transductive linear bandits. In *Advances in Neural Information Processing Systems*, pages 10667–10677.
- Frank, M., Wolfe, P., et al. (1956). An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110.
- Guestrin, C., Lagoudakis, M. G., and Parr, R. (2002). Coordinated reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, page 227–234.
- Hassin, R. and Khuller, S. (2001). z-approximations. *Journal of Algorithms*, 41(2):429–442.
- Jedra, Y. and Proutiere, A. (2020). Optimal best-arm identification in linear bandits. *arXiv preprint arXiv:2006.16073*.
- Jourdan, M., Mutý, M., Kirschner, J., and Krause, A. (2021). Efficient pure exploration for combinatorial bandits with semi-bandit feedback. In *Proceedings of the 31st International Conference on Algorithmic Learning Theory*.
- Jun, K.-S., Willett, R., Wright, S., and Nowak, R. (2019). Bilinear bandits with low-rank structure. In *International Conference on Machine Learning*, pages 3163–3172.
- Kazerouni, A. and Wein, L. M. (2019). Best arm identification in generalized linear bandits. *arXiv preprint arXiv:1905.08224*.
- Kiefer, J. and Wolfowitz, J. (1960). The equivalence of two extremum problems. *Canadian Journal of Mathematics*, 12:363–366.
- Mannor, S. and Shamir, O. (2011). From bandits to experts: On the value of side-observations. In *Advances in Neural Information Processing Systems*, pages 684–692.
- Perrault, P., Boursier, E., Valko, M., and Perchet, V. (2020). Statistical efficiency of thompson sampling for combinatorial semi-bandits. In *Advances in Neural Information Processing Systems*.
- Pukelsheim, F. (2006). *Optimal Design of Experiments*. Society for Industrial and Applied Mathematics.
- Robbins, H. (1952). Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*, 58(5):527–535.
- Sagnol, G. (2010). *Optimal design of experiments with application to the inference of traffic matrices in large networks: second order cone programming and submodularity*. PhD thesis, École Nationale Supérieure des Mines de Paris.
- Siomina, I., Varbrand, P., and Yuan, D. (2006). Automated optimization of service coverage and base station antenna configuration in UMTS networks. *IEEE Wireless Communications*, 13(6):16–25.
- Soare, M., Lazaric, A., and Munos, R. (2014). Best-arm identification in linear bandits. In *Advances in Neural Information Processing Systems*, pages 828–836.
- Tao, C., Blanco, S., and Zhou, Y. (2018). Best arm identification in linear bandits with linear dimension dependency. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 4877–4886.
- Tropp, J. A. et al. (2015). An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230.
- Valko, M. (2020). Bandits on graphs and structures.
- Valko, M., Munos, R., Kveton, B., and Kocák, T. (2014). Spectral bandits for smooth graph functions. In Xing, E. P. and Jebara, T., editors, *International conference on machine learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 46–54.
- Welch, W. (1982). Algorithmic complexity: Three np-hard problems in computational statistics. *Journal of Statistical Computation and Simulation - J STAT COMPUT SIM*, 15:17–25.
- Xu, L., Honda, J., and Sugiyama, M. (2018). A fully adaptive algorithm for pure exploration in linear bandits. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84, pages 843–851.
- Zaki, M., Mohan, A., and Gopalan, A. (2019). Towards optimal and efficient best arm identification in linear bandits. *arXiv preprint arXiv:1911.01695*.
- Zaki, M., Mohan, A., and Gopalan, A. (2020). Explicit best arm identification in linear bandits using no-regret learners. *arXiv preprint arXiv:2006.07562*.
- Çivril, A. and Magdon-Ismail, M. (2009). On selecting a maximum volume sub-matrix of a matrix and related problems. *Theoretical Computer Science*, 410(47):4801–4811.