# Maximizing and Satisficing in Multi-armed Bandits with Graph Information\*

Parth K. Thaker Arizona State University pkthaker@asu.edu

Nikhil Rao Microsoft nikhilrao860gmail.com Mohit Malu Arizona State University mmalu@asu.edu

Gautam Dasarathy Arizona State University gautamd@asu.edu

# Abstract

Pure exploration in multi-armed bandits has emerged as an important framework for modeling decision making and search under uncertainty. In modern applications however, one is often faced with a tremendously large number of options and even obtaining one observation per option may be too costly rendering traditional pure exploration algorithms ineffective. Fortunately, one often has access to similarity relationships amongst the options that can be leveraged. In this paper, we consider the pure exploration problem in stochastic multi-armed bandits where the similarities between the arms is captured by a graph and the rewards may be represented as a smooth signal on this graph. In particular, we consider the problem of finding the arm with the maximum reward (i.e., the maximizing problem) or one that has sufficiently high reward (i.e., the satisficing problem) under this model. We propose novel algorithms **GRUB** (GRaph based UcB) and  $\zeta$ -**GRUB** for these problems and provide theoretical characterization of their performance which specifically elicits the benefit of the graph side information. We also prove a lower bound on the data requirement that shows a large class of problems where these algorithms are near-optimal. We complement our theory with experimental results that show the benefit of capitalizing on such side information.

# 1 Introduction

The multi-armed bandit has emerged as an important paradigm for modeling sequential decision making and learning under uncertainty. Practical applications include design policies for sequential experiments [44], combinatorial online leaning tasks [9], collaborative learning on social media networks [30, 4], latency reduction in cloud systems [23] and many others [8, 59, 50, 24]. In the traditional multi-armed bandit problem, the goal of the agent is to sequentially choose among a set of actions or arms to maximize a desired performance criterion or reward. This objective demands a delicate tradeoff between exploration (of new arms) and exploitation (of promising arms). An important variant of the reward maximization problem is the identification of arms with the highest (or near-highest) expected reward. This *best arm identification* [41, 13] problem, which is one of pure exploration, has a wide range of important applications like identifying and testing drugs to treat infectious diseases like COVID-19, finding relevant users to run targeted ad campaigns, hyperparameter optimization in neural networks and recommendation systems. The broad range of

<sup>\*</sup>This work was supported in part by the National Science Foundation through the awards CCF-2048223, CNS-2003111, CCF-2029044, and OAC-1934766. This work was also supported partly by the ASU SenSIP Center

<sup>36</sup>th Conference on Neural Information Processing Systems (NeurIPS 2022).

applications of this paradigm is unsurprising given its ability to essentially model any optimization problem of black-box functions on discrete (or discretizable) domains with noisy observations.

While pure exploration problems in bandits show considerable promise, there are significant hurdles to their practical usage. In modern applications, one is often faced with a tremendously large number of options (sometimes in the order of hundreds of millions) that need to be considered for decision making. In such cases, playing (i.e., obtaining a random sample from) each bandit arm even once could be intractable. This renders traditional approaches to pure exploration ineffective. Fortunately, in several applications, the arms and their rewards are related to each other and information about the reward of one arm may be deduced from plays of similar arms. In this paper, we consider the pure exploration problem in stochastic multi-armed bandits where the similarities between arms is captured by a graph and the rewards may be represented as a smooth signal on this graph. Such graph side information is available in a wide range of applications: search and recommendation systems have graphs that capture similarities between items [17, 43, 53, 11]; drugs, molecules and their interactions can be represented on a graph [19]; targeted advertising considers users connected to each other in a social network [20], and hyperparameters for training neural network are often inter-related [57]. It is worth noting that such graphs are sometimes intrinsic to the problem (e.g., spatial coordinates or social/computer networks), or may be inferred based on a similarity metrics defined on arm features; a recent line of work considers constructing such graphs to enable more effective learning [see e.g., 58, 31].

**Our Contributions:** We consider the pure exploration in multi-arm bandits problem when a graph that captures similarities between the arms is available. In particular, we consider the problem of finding the arm with the maximum reward (i.e., the maximizing problem) or one that has sufficiently high reward (i.e., the satisficing problem<sup>2</sup>) under the assumption that arm rewards are smooth with respect to a known graph. Our main contributions may be summarized as follows:

(a) We devise a novel algorithm GRUB for the best arm identification problem (i.e., the maximizing problem) that specifically exploits the *homophily* (strong connections imply similar average rewards) on the graph (Section 3).

(b) We provide a theoretical characterization of the performance of GRUB. To this end, we define a novel measure  $\Im$  that we dub the "*influence factor*" which depends on the resistance distance of the underlying graph. This measure captures the benefit of the graph side information and plays a central role in the analysis of GRUB. In the traditional (graph-free) best arm identification problem, the sample complexity is know to scale as  $\sum_{i=1}^{n} \frac{1}{\Delta_i^2}$ , where  $\Delta_i$  is the gap between the expected rewards of the best arm and arm *i*. On the other hand, we show that GRUB roughly has a complexity that scales like  $\sum_{i \in \mathcal{H}} \frac{1}{\Delta_i^2}$  samples where the set  $\mathcal{H}$  is a set dependent on the influence factor, which contains arms which are hard to distinguish from optimal arm. For a broad range of problems  $|\mathcal{H}| \ll n$ , yielding significant improvement over traditional best arm identification algorithms (Section 4).

(c) In Section 5, we provide lower bounds on the minimum number of samples required for identification of the optimal arm when a graph encoding arm similarities is available. This shows the near-optimality of GRUB for an important class of representative problems.

(d) In many real world scenarios, the aim of finding the absolute best arm can often be too costly or even intractable. In these situations, it may be more appropriate to solve the *satisficing* problem, where the algorithm returns an arm that is good enough. We propose a variant of GRUB, dubbed  $\zeta$ -GRUB for this important setting in Section 6

(e) Finally, in Section 7, we complement our theoretical results with an empirical evaluation of our algorithms. We further provide algorithmic improvements to GRUB and discuss novel sampling policies for best arm identification in the presence of graph information.

#### 1.1 Related Work

The textbook [32] is an excellent resource for the general problem of multi-armed bandits. The pure exploration variant of the bandit problem is more recent, and has also received considerable attention in the literature [6, 7, 15, 14, 4, 21]. These lines of work treat the bandit arms or actions as independent entities and playing a particular arm yields no information about any other arm. This leads to great difficulty in scaling such methods, since in the problem setups with large number of

<sup>&</sup>lt;sup>2</sup>named after Herbert Simon's celebrated alternative model of decision making [48]

arms, attempting to play *all* arms is not practical. We resolve this precise roadblock by introducing a convenient way of of appending graph side information into the mix which provably accelerates the process of sub-optimal arm elimination (potentially without playing it even once!)

A recent line of work [35, 32, 18, 56, 16, 39] has proposed the leveraging of structural sideinformation for the multi-armed bandit problem for regret minimization. Such topology-based bandit methods work under the assumption that pulling an arm reveals information about other, correlated arms [18, 47], which help in developing better regret methods. Similarly, spectral bandits [29, 56, 51] assume user features are modelled as signals defined on an underlying graph, and use this to assist in learning. The works [3] and [54] consider similar graph information models, albiet at a degraded level. The authors in [33] use the graphs to improve the regret bounds in a thresholding bandit setting. Work revolving around spectral bandits utilize the *spectrum* of the graph laplacian. In contrast, we focus on the *combinatorial properties* of the graphs to devise algorithms and analyse them. Another line of work [12, 52, 36, 37] considers search problems on graphs under a different model and there is an opportunity for future work to combine these techniques.

Most of the aforementioned works focus on regret minimization in the presence of graph information. The problem of pure exploration with similarity graphs has received far less attention. The authors in [29] were the first to attempt at filling this gap for the spectral bandit setting. They provide an information-theoretic lower bound and a gradient-based algorithm to estimate this lower bound to sample the arms. The authors provide performance guarantees for the algorithm, but these results only indirectly capture the benefit brought by the graph; our results on the other hand are based on a novel complexity measure that explicitly elicits the benefit of having the graph side information.

Note that, similarity graph information considered in this work is fundamentally different from linear rewards assumption in contextual/linear bandits. In the linear bandits problem, the reward behavior is assumed to be low dimensional and this is crucial for the improved regret bounds and sample complexity guarantees [32, 49]. In the current work we do not make any assumptions on low dimensionality of the rewards but still show improvements in sample complexity provided a good arm-similarity graph is available. We show a toy example in Appendix H where a low dimensional linear bandit cannot be competitive with the corresponding graph-bandit setting.

# 2 Problem Setup and Notation

We consider an *n*-armed bandit problem with the set of arms given by  $[n] \triangleq \{1, 2, 3, \ldots, n\}$ . Each arm  $i \in [n]$  is associated with a  $\sigma$ -sub-Gaussian distribution  $\nu_i$ . That is,  $\mathbb{E}_{X \sim \nu_i} [\exp(s(X - \mu_i))] \leq \exp\left(\frac{\sigma^2 s^2}{2}\right) \forall s \in \mathbb{R}$ , where  $\mu_i = \mathbb{E}_{\nu_i} [X]$  is said to be the (expected or mean) reward associated to arm *i*. We will let  $\boldsymbol{\mu} \in \mathbb{R}^n$  denote the vector of all the arm rewards. A "play" of an arm *i* is simply an observation of an independent sample from  $\nu_i$ ; this can be thought of as a noisy observation of the corresponding mean  $\mu_i$ . The goal of the best-arm identification problem is to identify, from such noisy samples, the arm  $a^* \triangleq \arg \max_{i \in [n]} \mu_i$  that has the maximum expected reward, denoted by  $\mu^*$ . For each arm  $i \in [n]$ , we will let  $\Delta_i \triangleq \mu^* - \mu_i$  denote the sub-optimality of the arm.

As discussed in Section 1, our goal is to consider the best-arm identification where one has additional access to information about the similarity of the arms under consideration. In particular, we model this side information as a weighted undirected graph  $G = (V_G, E_G, A_G)$  where the vertex set,  $V_G = [n]$ , is identified with the set of arms, the edge set  $E_G \subseteq {\binom{[n]}{2}}$ , and adjacency matrix  $A_G \in \mathbb{R}^{n \times n}$  describes the weights of the edges E between the arms which captures the similarity in means of connected arms; the higher the weight, the more similar the rewards from the corresponding arms. We will let  $L_G = D_G - A_G$  denote the combinatorial Laplacian<sup>3</sup> of the graph [10], where  $D_G = \text{diag}(A_G \times \mathbb{1}_n)$  is a diagonal matrix containing the weighted degrees of the vertices. We will suppress the dependence on G when the context is clear. Subsequently, we show that if one has access to this graph and the vector of rewards  $\mu$  is *smooth* with respect to the graph (that is, highly similar arms have highly similar rewards), then one can solve the pure exploration problem extremely efficiently. We will capture the degree of smoothness of  $\mu$  with respect to the graph using

<sup>&</sup>lt;sup>3</sup>All our results continue to hold if this is replaced with the normalized, random walk, or generalized Laplacian.

the following seminorm<sup>4</sup>:

$$\|\boldsymbol{\mu}\|_{G}^{2} \triangleq \langle \boldsymbol{\mu}, L_{G}\boldsymbol{\mu} \rangle = \sum_{\{i,j\} \in E_{G}} A_{ij} (\mu_{i} - \mu_{j})^{2}.$$

$$\tag{1}$$

The second equality above can be verified by a straightforward calculation. Also, notice that  $\|\boldsymbol{\mu}\|_G$  being small implies  $\mu_i \approx \mu_j$  for  $(i, j) \in E$ . In such scenario we say that the mean vector  $\boldsymbol{\mu}$  is smooth over graph G. This observation has inspired the use of the Laplacian in several lines of work to enforce smoothness on the vertex-valued functions [2, 51, 60, 33]. For  $\epsilon > 0$ , we say that arms (rewards) are  $\epsilon$ -smooth with respect to a graph G if  $\|\boldsymbol{\mu}\|_G \leq \epsilon$ .

Let  $\mathcal{C}(G) \subset 2^{[n]}$  denote the set of all connected components and let  $k(G) \triangleq |\mathcal{C}(G)|$  denote the number of connected components of the graph G. For a vertex  $i \in [n]$ , we will let  $C_i(G) \in \mathcal{C}(G)$  denote the connected component that contains i. When the context is clear we sometimes let  $C_i(G)$  also refer all the nodes in the connected component. We say a graph G = ([n], E) has k-isolated cliques if it can be divided into fully connected sub-graphs  $G_i = (V_i, E_i)$  such that  $V_i \subseteq [n], E_i = {V_i \choose 2}$  for all  $i \in [k], V_i \cap V_j = \emptyset, E_i \cap E_j = \emptyset$  for all  $i, j \in [k]$ , and  $\bigcup_{i=1}^k V_i = [n], \bigcup_{i=1}^k E_i = E$ . Notice that we only have one clique if G is fully connected.

To solve the best-arm identification problem, we need a sampling policy to sequentially and interactively select the next arm to play, and a stopping criterion. For any time  $t \in \mathbb{N}$ , the sampling policy  $\pi_t = {\pi_s}_{s \leq t}$  is a function that maps t to an arm in [n] given the history of observations up to time t - 1. With slight abuse of notation, we will let  $\pi_t$  denote the arm chosen by an agent at time t. Let  $r_{t,\pi_t}$  denote the random reward observed at time t from arm  $\pi_t$ . We use  $t_i(\pi_t)$  (referred as  $t_i$  for simplicity) to denote the number of times arm i is played under the sampling policy  $\pi_t$ . In this paper we tackle the following problems:

**P1 (Best arm identification):** Given n arms and an arbitrary graph G capturing similarity between the arms, can we design a policy  $\pi_T$  that exploits the similarity to find the best arm efficiently?

**P2** ( $\zeta$ -best arm identification): Under the setting in **P1**, can we design a similarity exploiting policy  $\pi_T$  so as to find an arm belonging to the set  $B(\zeta) \triangleq \{i \in [n] : |\mu_i - \mu_{a^*}| \le \zeta\}$  efficiently?

# **3** The GRUB Algorithm

We now introduce GRUB (GRaph based Upper Confidence Bound), a novel but natural algorithm for best arm identification in the presence of graph side information. We begin with an intuitive description of how GRUB incorporates the graph side information into an *upper confidence bound* (UCB) strategy. Most UCB algorithms [32, 51] compute the estimates of mean and variance, and use these to eliminate arms that have been deduced to be sub-optima. The key idea behind GRUB is that the arm similarity information allows us to create high-quality estimates of mean rewards and confidence intervals for arms that have not been (sufficiently) sampled yet. In what follows, we describe the building blocks of GRUB.

#### 3.1 Leveraging Graph Side Information

We introduce two key ideas that lie at the heart of the GRUB algorithm. First, at each step, GRUB computes a regularized estimate of the means of *all the arms*; the regularization based on the graph Laplacian essentially promotes the smoothness of the mean vector on the given graph. This allows the algorithm to estimate means of arms it has *never sampled*. To do this, at any given time step T, the algorithm solves the following Laplacian-regularized least-squares optimization program:

$$\hat{\boldsymbol{\mu}}_T = \operatorname*{arg\,min}_{\boldsymbol{\mu} \in \mathbb{R}^n} \left\{ \left[ \sum_{t=1}^T (r_{t,\pi_t} - \mu_{\pi_t})^2 \right] + \rho \langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle \right\},\tag{2}$$

 $<sup>{}^{4}</sup>L_{G}$  is not positive definite, and can be verified to have as many zero eigenvalues as the number of connected components in G

where  $\rho > 0$  is a tunable parameter. Equation (2) admits a closed form solution of the form

$$\hat{\boldsymbol{\mu}}_T = \left(\sum_{t=1}^T \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^\top + \rho L_G\right)^{-1} \left(\sum_{t=1}^T \mathbf{e}_{\pi_t} r_{t,\pi_t}\right),$$

provided the matrix  $V_T \triangleq \sum_{t=1}^{T} \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^\top + \rho L_G$  is invertible;  $\mathbf{e}_i$  denotes the *i*-th standard basis vector for the Euclidean space  $\mathbb{R}^n$ . In Appendix **B** we show that invertibility holds if and only if the sampling policy yields at least one sample per connected component of *G*. This is a rather mild condition that we arrange for explicitly in our algorithm, given that we know the graph *G*. In what follows we assume that every connected component of graph *G* is sampled at least once. This regularized mean estimation procedure yields an estimate of the mean that is both in agreement with observations and smooth on the graph – thereby allowing information sharing among similar arms.

The second key idea of our algorithm is the utilization of the graph G in tracking the confidence bounds of *all the arms simultaneously*. Intuitively, for identifying the best arm, we must be reasonably certain about the sub-optimality of the other arms. This in turn would require the algorithm to track a high-probability confidence bound on the means of all the arms. In the traditional (graph-free) best arm identification problem, the confidence interval of an arm's mean estimate depends on the number of times the arm has been played. Requiring multiple plays of all suboptimal arms for obtaining high confidence bounds is potentially disastrous when the number of arms is very large. In our setup, we show that the knowledge of the similarity graph greatly improves this situation. In particular, we show that a play of any arm not only tightens its own confidence interval, but also has an impact on the confidence intervals of *all connected arms*. To quantify the benefit of graph information for the confidence bounds, we will define a novel quantity for each arm – the effective number of plays.

**Definition 3.1** (Effective Number of Plays). Let  $\rho > 0$  and  $\{t_i\}_{i=1}^n$  denote the number of plays of each of the *n* arms when a sampling policy  $\pi_T$  is employed for *T* time steps. Suppose that for each connected component  $C \in C(G)$ , there is at least one arm  $i_C \in C$  such that  $t_{i_C} > 0$ . Then the effective number of plays for each arm  $i \in [n]$  is defined as  $t_{\text{eff},i} \triangleq \left[ (N_T + \rho L_G)^{-1} \right]_{i_i}^{-1}$ , where  $N_T$  is a diagonal matrix of  $\{t_i\}_{i=1}^n$ , and  $L_G$  denotes the Laplacian of the given graph *G*.

Effective number of plays  $t_{\text{eff},i}$  for any arm i is influenced by two factors: (a) the number of samples of arm i itself, and (b) the number of samples of any arm in the connected component  $j \in C(i), j \neq i$ . It can be shown that for any arm  $i, t_{\text{eff},i}$  depends on the number of connections of node i in graph G and its value increases as the connectivity of the node increases. The choice of the terminology for this quantity is justified by the following lemma, which provides a high confidence bound for the mean estimate of each arm .

**Lemma 3.2** (Concentration inequality). For any T > k(G), the following holds with probability at least  $1 - \delta$ :

$$|\hat{\mu}_T^i - \mu_i| \le \sqrt{\frac{1}{t_{eff,i}}} \left( 2\sigma \sqrt{14 \log\left(\frac{2w_i(\boldsymbol{\pi}_T)}{\delta}\right)} + \rho \|\boldsymbol{\mu}\|_G \right), \quad \forall i \in [n]$$
(3)

where  $w_i(\boldsymbol{\pi}_T) = a_0 n t_{eff,i}^2$  for any constant  $a_0 > 0$ ,  $\hat{\mu}_T^i$  is the *i*-th coordinate of the estimate from (2)

Notice that the *effective number of plays* has a similar role as the number of plays in traditional pure exploration algorithms [13]. Indeed, in the absence of graph information,  $t_{\text{eff},i}$  reduces to  $t_i$ , the total number of plays of individual arms. Lemma 3.2 recovers high confidence bounds for standard best-arm identification problem [13]. It should be noted that while our work is the first to identify this interpretable quantity explicitly, the result of Lemma 3.2 in other forms has appeared before in the literature [1, 51, 56].

We introduce our algorithm GRUB for best arm identification when the arms can be approximately cast as nodes on a graph. GRUB uses insights from graph-based mean estimation (2) and upper confidence bound estimation (3) for its elimination policies to search for the optimal arm.

GRUB accepts as input a graph G on n arms (and its Laplacian  $L_G$ ), a regularization parameter  $\rho > 0$ , a smoothness parameter  $\epsilon > 0$ , and an error tolerance parameter  $\delta \in (0, 1)$ . It is composed of the following major blocks.

Initialization: First, GRUB identifies the clusters in the G using a Cluster-Identification

routine. Any algorithm that can efficiently partition a graph can be used here, e.g METIS [25]. GRUB then samples one arm from each cluster. This ensures  $V_T \succ 0$ , which enables GRUB to estimate  $\hat{\mu}_T$ using the closed form solution of eq. (2). A great advantage of GRUB is that the initialization phase only requires steps equal to the number of disconnected components in the graph. This is in direct contrast with traditional best arm identification algorithms, which require atleast one sample from every arm initially.

**Sampling policy:** At each round, GRUB obtains a sample from the arm returned by the routine Sampling-Policy, which cyclically samples arms from different clusters while ensuring that no arm is resampled before all arms in consideration have the same number of samples. This is distinct from standard cyclic sampling policies that is traditionally used for best arm identification [13], but any of them may be modified readily to provide a cluster-aware sampling policy for GRUB. In our experiments, we show that replacing cyclic sampling with more statistics- and structure-aware sampling greatly improves performance; a theoretical analysis of these is a promising avenue for future work. One of the major advantage of GRUB is the lite nature of the computation. Every loop just requires a rank-1 inverse update which can be performed very efficiently and it does not need any subroutines, unlike [29]

**Bad arm elimination :** At any time t, let A be the set of all arms in consideration for being optimal. Using the uncertainty bound from (3), GRUB uses the following criteria for sub-optimal arm elimi-

nation. At each iteration, GRUB identifies an arm  $a_{\max} \in A$ ,  $a_{\max} = \underset{i \in A}{\operatorname{argmax}} \left[ \hat{\mu}_t^i - \beta_i(t) \sqrt{t_{\text{eff},i}^{-1}} \right]$ ,

where  $\beta_i(t) = \left(2\sigma\sqrt{14\log\left(\frac{2na_0t_{\text{eff},i}^2}{\delta}\right)} + \rho\epsilon\right)$ , with the *highest lower bound* on its mean estimate. Following this, GRUB removes arms from the set A according to the following elimination policy,

$$A \leftarrow \left\{ \mathbf{a} \in A \mid \hat{\mu}_t^{a_{\max}} - \hat{\mu}_t^a \le \beta_a(t) \sqrt{t_{\text{eff},a}^{-1}} + \beta_{a_{\max}}(t) \sqrt{t_{\text{eff},a_{\max}}^{-1}} \right\}.$$
 (4)

Note that GRUB does not require any optimization innerloop as in [29]. This potentially provides GRUB with a significant computation advantage, especially when the dimensionality of the problem is very large. The pseudocode for GRUB can be found in Appendix E.

Next, we derive performance guarantees on the sample complexity for GRUB to return the best arm with high probability.

#### 4 Theoretical Analysis of GRUB

In this section we provide a formal statement of the sample complexity of GRUB. To do this, we first introduce a novel quantity we call *influence factor*. The influence factor of an arm is derived from resistance distance, a classical graph theoretic concept. This adds to the interpretability and understanding of the instances where using graph side information might be of tremendous use to the application. The usage of graph through the influence factor allows us to identify arms that can be eliminated quickly from consideration.

#### 4.1 Resistance Distance and Influence Factor

We first recall the definition of resistance distance in a graph.

**Definition 4.1** (Resistance Distance). [5] For any graph G with n nodes, given a constant  $\delta > 0$ , the **resistance distance**  $r_{\delta,G}(i,j)$  between two nodes i, j is defined as,

$$r_{\delta,G}(i,j) = R_{ii} + R_{jj} - R_{ij} - R_{ji}, \tag{5}$$

where  $R \triangleq (L_G + \delta \mathbb{1}\mathbb{1}^T)^{\dagger}$ ;  $\dagger$  denotes the Moore-Penrose inverse,  $L_G$  is the Laplacian of graph G, and  $\mathbb{1} \in \mathbb{R}^n$  is the vector of all 1's.

When the context is clear we denote the resistance distance simply as  $r_G(\cdot, \cdot)$ . The terminology comes from circuit theory: Suppose that an graph G = ([n], E) is thought of as a resistor network on the nodes [n] where each edge  $\{i, j\}$  has a unit resistance. Then, the effective resistance between two nodes i and j is precisely the resistance distance r(i, j). It can be shown in general that nodes

that are close by or connected by several paths have a small resistance distance. Given its ability to capture closeness of nodes in graph, the resistance distance has found a broad range of applications and has been the subject of much study; see e.g., [28, 5, 55].

Using the notion of resistance distance, we define the influence factor  $\Im(\cdot, G)$  of a vertex below. This novel measure quantifies the impact of the graph on the parameter estimation of arm j, and in particular, allows us to use the combinatorial properties of the graph and the arm means to classify arms into two sets: competitive and non-competitive; the definition of these sets follows right after. As our theory will show, the competitive arms are sampled as though we were in the traditional graph-free setting; on the other hand, non-competitive arms are eliminated rapidly, often with zero plays! Indeed, the smoother the reward vector is with respect to the graph, the fewer competitive arms there are – it is this phenomenon that is captured using the influence factor.

**Definition 4.2** (Influence Factor). Let G be a graph on the vertex set [n]. For each  $j \in [n]$ , define influence factor  $\Im(j, G)$  as:

$$\Im(j,G) = \begin{cases} \min_{i \in C_j(G), i \neq j} \{ r_G(i,j)^{-1} \}, & \text{if } |C_j(G)| > 1, \\ 0, & \text{otherwise }. \end{cases}$$
(6)

Here,  $r_G(i, j)$  is the resistance distance between arm i and j in G as in Definition 4.1.

**Definition 4.3** (Competitive and Non-Competitive Arms). Fix  $\mu \in \mathbb{R}^n$ , graph D, regularization parameter  $\rho$ , confidence parameter  $\delta$ , and smoothness parameter  $\epsilon$ . We define  $\mathcal{H}_D$  to be the set of competitive arms and  $\mathcal{N}_D$  to be the set of non-competitive arms as follows:

$$\mathcal{H}_D = \left\{ j \in [n] \middle| \Delta_i \le 2\sqrt{\frac{2}{\rho \Im(i, D)}} \left( 2\sigma \sqrt{14 \log\left(\frac{2a_0 n \rho^2 \Im(i, D)^2}{\delta}\right)} + \rho \epsilon \right) \right\}$$
(7)

and  $\mathcal{N}_D \triangleq [n] \setminus \mathcal{H}_D$ .

As the name suggests, the arms in  $\mathcal{H}$  are close to the optimal arm  $a^*$  in mean (competitive compared to the optimal arm  $a^*$ ) and requires several plays before they can be discarded, as shown in the theorem below. Note from the above definition that an arm is more likely to be part of this set if its mean is high (i.e.,  $\Delta_i$  is low) and its influence factor is low. Similarly, the non-competitive set is composed of arms whose means are not competitive with the optimal arm.

Armed with these definitions, we are now ready to state our main theorem that characterizes the performance of GRUB.

#### 4.2 Sampling policy performance

Cyclic sampling policies have been traditionally used in multi-armed bandit problems for best-arm identification [13]. The sample complexity bound for GRUB with cyclic sampling is as follows:

**Theorem 4.4** (GRUB Sample Complexity). Consider *n*-armed bandit problem with mean vector  $\mu \in \mathbb{R}^n$ . Let G = (V, E) be the similarity graph with the vertex set V = [n] and edge set E, let  $\mathcal{G}$  be the set of subgraphs of G, and further suppose that  $\mu$  is  $\epsilon$ -smooth i.e.,  $\|\mu\|_G \leq \epsilon$ . Define

$$T_{sufficient} \triangleq \underset{D \in \mathcal{G}}{\arg\min} \sum_{C \in \mathcal{C}_D} \left[ \sum_{\substack{i \in C \cap \mathcal{H}_D \\ i \neq 1}} \frac{1}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] + \max_{i \in C \cap \mathcal{N}_D} \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right],$$

where  $\Delta_i = \mu^* - \mu_i$  for all suboptimal arms,  $\mathcal{H}_D$  and  $\mathcal{N}_D$  are as in Definition 4.3,  $\mathcal{C}_D$  is the set of connected components of a given graph D and  $c_1, c_2$  are constants independent of system parameters. Then, with probability at least  $1 - \delta$ , GRUB: (a) terminates in no more than  $T_{sufficient}$  rounds, and (b) returns the best arm  $a^* = \arg \max_i \mu_i$ .

*Remark* 4.5. The required number of samples for successful elimination of suboptimal arms, and therefore the successful identification of the best arm, can be split into two categories based on the sets defined in Definition 4.3. Each sub-optimal *highly competitive arm*  $j \in \mathcal{H}$  requires  $\mathcal{O}(1/\Delta_j^2)$  samples, which is comparable to the classical (graph-free) best-arm identification problem. Additionally, the non-competitive arms  $\mathcal{N}$  can be eliminated without being played, depending on the influence factor:

one round of the cyclic sampling suffices to eliminate these arms (even if they are never played!). We refer the reader to Appendix E for a more detailed discussion. Indeed, the smaller  $|\mathcal{H}|$  is, the more the graph side information benefits GRUBand vice-versa.

*Remark* 4.6. Note that  $T_{\text{sufficient}}$  in Theorem 4.4 involves the minimum over all subgraphs. As we show in Lemma I.8 in the appendix,  $\mathfrak{I}$  can actually increase if one restricts their attention to certain subgraphs of G; this in turn increases the size of  $\mathcal{N}$  and decreases the size of  $\mathcal{H}$ , hence, giving a tighter upper bound on the performance of the algorithm. GRUB *automatically adapts to the best subgraph* to maximize the influence factor  $\mathfrak{I}(\cdot, \cdot)$  to obtain the best possible sample complexity and this is reflected in the statement of Theorem 4.4.

The complete proof of Theorem 4.4 can be found in Appendix E, where we also provide more insights on the behavior of the confidence bound as a function of the number of samples acquired. These results may be of independent interest to the reader.

## 5 Lower Bounds

Let us consider an *n*-armed bandit setup with arm indices [1, ..., n]. Let  $\mu^*$  indicate the mean of the optimal arm and  $\mu_i$  indicate the mean values of all other arms such that  $\mu_i < \mu^*$ . For the rest of this section, without loss of generality, let the index of optimal arm be 1.

**Theorem 5.1.** Given an n-armed bandit model with associated mean vector  $\boldsymbol{\mu} \in \mathbb{R}^n$  and similarity graph G smooth on  $\boldsymbol{\mu}$ , i.e.  $\langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle \leq \epsilon$ , for any  $0 < \epsilon < \epsilon_0$ . Let G = ([n], E) be the graph with only isolated cliques and w.l.o.g let arm 1 be the optimal arm. Then define

$$T_{necessary} = \sum_{C \in \mathcal{C}_G/C^*} \min_{j \in C} \left\{ \frac{4\sigma^2 \log 5}{(\Delta_j - \sqrt{\epsilon})^2} \right\} + \sum_{j \in C^*/1} \frac{4\sigma^2 \log 5}{\Delta_j^2},\tag{8}$$

where  $C^*$  is the clique with the optimal arm and  $\epsilon_0 := \min_{i \in [n]/1, j \in C(i)} \left[ \Delta_j \left[ 1 - \frac{\Delta_i}{\sqrt{\Delta_i^2 + \Delta_j^2}} \right] \right]^2$ . Then any  $\delta$ -PAC algorithm will need at-least  $T_{necessary}$  steps to terminate, provided  $\delta \leq 0.1$ .

Using Theorem 5.1, we can show that GRUB is minimax optimal for a *n*-armed bandit problems for certain class of similarity graph G. The following result shows that the upperbound on the sample complexity provided in Theorem 4.4 matches the lower bound established in Theorem 5.1 in  $\Delta_i$  up to a constant factor.

**Corollary 5.2** (Isolated clusters). Consider the setup as in Theorem 5.1 with the further restriction that graph G be such that the optimal node is isolated and  $\epsilon < \min_{j \in [n]} \frac{\Delta_j^2}{2}$ . Define,

$$T_{necessary} \ge \sum_{C \in \mathcal{C}_G / \{1\}} \max_{j \in C} \left\{ \frac{8\sigma^2 \log 5}{\Delta_j^2} \right\}.$$
(9)

Then any algorithm that takes fewer than  $T_{necessary}$  samples will have a probability of error at least 0.1.

As can be seen in Corollary 5.2, the lower bound expression can scale as standard *n*-armed bandit (implying no added advantage of having graph side-information) or can behave as a  $|C_G|$ -armed bandit problem (scales as the number of clusters in graph *G* rather than number of nodes *n*) purely by changing the similarity graph *G*. The difference between  $C_F$  (connected components in the subgraph constructed by making optimal arm isolated) and  $C_G$  (connected components in the given similarity graph) can lead to more interesting behaviour in terms of lower bound expressions on sample complexity.

# 6 $\zeta$ -best-arm identification

It can be observed from Theorem 4.4 that the fact that the means are  $\epsilon$ -smooth implies that distinguishing arm j from  $a^*$  would require at least  $O(\epsilon^{-2})$  samples. A tighter upper bound on the violation  $\epsilon$  and an edge between j and  $a^*$  would make the suboptimal arm j harder to eliminate. However, it stands to reason that in such situations, it might be more practical to not demand for the absolute best

arm, but rather an arm that is nearly optimal. Indeed, in several modern applications we discuss in Section 1, finding an approximate best arm is tantamount to solving the problem. In such cases, a simple modification of GRUB can be used to quickly eliminate definitely suboptimal arms, and then output an arm that is guaranteed to be nearly optimal. To formalize this, we consider the  $\zeta$ -best arm identification problem as follows.

**Definition 6.1.** For a given  $\zeta > 0$ , arm *i* is called  $\zeta$ -best arm if  $\mu_i \ge \mu_{a^*} - \zeta$ , where  $a^* = \arg \max_i \mu_i$ 

The goal of the  $\zeta$ -best arm identification problem is to return an arm  $\tilde{a}$  that is  $\zeta$ -optimal. We achieve this by a simple modification to GRUB, which we dub  $\zeta$ -GRUB, which ensures that all the remaining arms *i* satisfy  $4\beta(t_i)\sqrt{t_{\text{eff},i}^{-1}} \leq \zeta$ . It then outputs the best arm amongst those that are remaining. The following theorem characterizes the sample complexity for  $\zeta$ -GRUB:

**Theorem 6.2.** Consider *n*-armed bandit problem with mean vector  $\boldsymbol{\mu} \in \mathbb{R}^n$ . Let *G* be the given similarity graph on vertex set [n], and further suppose that  $\boldsymbol{\mu}$  is  $\epsilon$ -smooth. Let *C* be the set of connected components of *G*. Define,

$$T_{sufficient} \triangleq \underset{D \in \mathcal{G}}{\arg\min} \sum_{C \in \mathcal{C}_D} \left| \sum_{i \in C \cap \mathcal{H}_D} \frac{1}{(\Delta_i \vee \zeta)^2} \left[ c_1 \log \frac{c_2}{\delta(\Delta_i \vee \zeta)} + \frac{\rho \epsilon}{2} \right] + \underset{i \in C \cap \mathcal{N}_D}{\max} \left\{ \frac{2}{(\Delta_i \vee \zeta)^2} \left[ c_1 \log \frac{c_2}{\delta(\Delta_i \vee \zeta)} + \frac{\rho \epsilon}{2} \right] \right\} \right],$$
(10)

where  $\Delta_i = \mu^* - \mu_i$  for all suboptimal arms,  $\mathcal{H}_D$  and  $\mathcal{N}_D$  are as in Definition 4.3,  $\mathcal{C}_D$  is the set of connected components of a given graph Dand  $\Delta_i \lor \zeta = \max\{\zeta, \Delta_i\}$  and  $c_1, c_2$  are constants independent of system parameters. Then, with probability at least  $1 - \delta$ ,  $\zeta$ -GRUB: (a) terminates in no more than  $T_{\text{sufficient}}$  rounds, and (b) returns a  $\zeta$ -best arm.

The pseudocode for the  $\zeta$ -GRUB can be found in Appendix G.

# 7 Experiments

For all our experiments, we use Intel® Core<sup>TM</sup> i7-10875H CPU @ 2.30GHz × 16 with 32 GB memory. We set  $\delta = 1e - 3$ ,  $\rho = 2.0$ ,  $\sigma = 2.0$ . We evaluate GRUB with different sampling strategies from section J and compare its performance to standard UCB algorithm on both synthetic and real datasets.

**Better sampling strategies:** Theorem 4.4 established a baseline w.r.t. sampling protocol by solving  $T_{\text{sufficient}}$  for naive cyclic sampling policy (a sampling policy which does not exploit the graph properties). Note that, even if the sampling policy does not utilize any graph properties, the similarity graph is still being utilized in computing the mean estimate and the confidence widths. For the safe elimination of suboptimal arms, the ultimate goal of GRUB is to shrink the confidence bounds  $\beta_i \sqrt{(t_{\text{eff},i})^{-1}}$  as quickly as possible. For the complete description of all the alternatives please refer to Appendix J

**Synthetic Data:** We consider an *n*-armed bandit setup with the aim of finding the best arm. The number of arms scale from n = 50 to 200 in steps of 50. We consider 2 cases: *G* is a Stochastic Block model(SBM) with parameters  $(p, q) = (0.9, 1e^{-4})$  and *G* is a Barabási–Albert(BA) graph with parameter m = 2, both containing 10 clusters. We run every setup for 20 runs and record the stopping time for all runs. In Figure 1, we compare the baseline cyclic algorithm (Nograph-UCB) with GRUB and its variants (GRUB-MVM, JVM-O, JVM-N), more details on this in Appendix J.

As can be seen in Figure 1, all graph-aware algorithms (GRUB with different sampling policies) outperform the standard UCB based best-arm identification algorithm. Within the different GRUB, different sampling policies exploit the graph infromation in different ways, leading to variations in their performance. GRUB (cyclic sampling based) is outperformed by all other sampling based GRUB methods.

We show additional experiments with different graph parameters for Stochastic block model and Barabási-Albert graphs and different cluster sizes as well as real data in Appendix K. The full code used for conducting experiments can be found at the following Github repository. Discussion about limitations, future works and broader impact are provided in Appendix A.



Figure 1: (Best seen in color) Performance of GRUB and its variant sampling protocols for SBM  $((p,q) = (0.95, 1e^{-4}))$  [Left] and BA (m = 2) [Right]. GRUB outperforms the standard cyclic UCB method

### References

- Yasin Abbasi-yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- [2] Rie Kubota Ando and Tong Zhang. Learning on graph with laplacian regularization. Advances in neural information processing systems, 19:25, 2007.
- [3] Alexia Atsidakou, Orestis Papadigenopoulos, Constantine Caramanis, Sujay Sanghavi, and Sanjay Shakkottai. Asymptotically-optimal Gaussian bandits with side observations. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 1057–1077. PMLR, 17–23 Jul 2022.
- [4] Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-armed bandits. In *COLT*, pages 41–53, 2010.
- [5] Ravindra B Bapat and Somit Gupta. Resistance distance in wheels and fans. *Indian Journal of Pure and Applied Mathematics*, 41(1):1–13, 2010.
- [6] Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in multi-armed bandits problems. In *International conference on Algorithmic learning theory*, pages 23–37. Springer, 2009.
- [7] Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412(19):1832–1852, 2011.
- [8] Wei Cao, Jian Li, Yufei Tao, and Zhize Li. On top-k selection in multi-armed bandits and hidden bipartite graphs. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [9] Shouyuan Chen, Tian Lin, Irwin King, Michael R Lyu, and Wei Chen. Combinatorial pure exploration of multi-armed bandits. In *NIPS*, pages 379–387, 2014.
- [10] Fan RK Chung and Fan Chung Graham. Spectral graph theory. Number 92. American Mathematical Soc., 1997.
- [11] G Dasarathy, N Rao, and R Baraniuk. On computational and statistical tradeoffs in matrix completion with graph information. In *Signal Processing with Adaptive Sparse Structured Representations Workshop SPARS*, 2017.
- [12] Gautam Dasarathy, Robert Nowak, and Xiaojin Zhu. S2: An efficient graph based active learning algorithm with application to nonparametric classification. In *Conference on Learning Theory*, pages 503–522. PMLR, 2015.

- [13] Eyal Even-Dar, Shie Mannor, Yishay Mansour, and Sridhar Mahadevan. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal* of machine learning research, 7(6), 2006.
- [14] Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In NIPS-Twenty-Sixth Annual Conference on Neural Information Processing Systems, 2012.
- [15] Aurélien Garivier and Emilie Kaufmann. Non-asymptotic sequential tests for overlapping hypotheses and application to near optimal arm identification in bandit models, 2019.
- [16] Claudio Gentile, Shuai Li, and Giovanni Zappella. Online clustering of bandits, 2014.
- [17] Jiafeng Guo, Xueqi Cheng, Gu Xu, and Huawei Shen. A structured approach to query recommendation with social annotation data. In *Proceedings of the 19th ACM international conference* on Information and knowledge management, pages 619–628, 2010.
- [18] Samarth Gupta, Shreyas Chaudhari, Gauri Joshi, and Osman Yağan. Multi-armed bandits with correlated arms, 2020.
- [19] Vassilis N. Ioannidis, Xiang Song, Saurav Manchanda, Mufei Li, Xiaoqin Pan, Da Zheng, Xia Ning, Xiangxiang Zeng, and George Karypis. Drkg - drug repurposing knowledge graph for covid-19. https://github.com/gnn4dr/DRKG/, 2020.
- [20] Mohsen Jamali and Martin Ester. Trustwalker: a random walk model for combining trust-based and item-based recommendation. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 397–406, 2009.
- [21] Kevin Jamieson and Robert Nowak. Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. In 2014 48th Annual Conference on Information Sciences and Systems (CISS), pages 1–6. IEEE, 2014.
- [22] Ming Ji and Jiawei Han. A variance minimization criterion to active learning on graphs. In Neil D. Lawrence and Mark Girolami, editors, *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*, volume 22 of *Proceedings of Machine Learning Research*, pages 556–564, La Palma, Canary Islands, 21–23 Apr 2012. PMLR.
- [23] Gauri Joshi, Emina Soljanin, and Gregory Wornell. Efficient redundancy techniques for latency reduction in cloud systems, 2017.
- [24] Kirthevasan Kandasamy, Gautam Dasarathy, Barnabas Poczos, and Jeff Schneider. The multifidelity multi-armed bandit. In Advances in Neural Information Processing Systems, pages 1777–1785, 2016.
- [25] George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. SIAM JOURNAL ON SCIENTIFIC COMPUTING, 20(1):359–392, 1998.
- [26] Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of a/b testing, 2015.
- [27] Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best arm identification in multi-armed bandit models, 2016.
- [28] Douglas J Klein and Milan Randić. Resistance distance. *Journal of mathematical chemistry*, 12(1):81–95, 1993.
- [29] Tomáš Kocák and Aurélien Garivier. Best arm identification in spectral bandits. *arXiv preprint arXiv:2005.09841*, 2020.
- [30] Ravi Kumar Kolla, Krishna Jagannathan, and Aditya Gopalan. Collaborative learning of stochastic bandits over a social network. *IEEE/ACM Transactions on Networking*, 26(4):1782– 1795, 2018.
- [31] Dan Kushnir and Luca Venturi. Diffusion-based deep active learning. *CoRR*, abs/2003.10339, 2020.

- [32] Tor Lattimore and Csaba Szepesvári. Bandit algorithms. Cambridge University Press, 2020.
- [33] Daniel LeJeune, Gautam Dasarathy, and Richard Baraniuk. Thresholding graph bandits with grapl. In *International Conference on Artificial Intelligence and Statistics*, pages 2476–2485. PMLR, 2020.
- [34] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. http://snap.stanford.edu/data, June 2014.
- [35] Shuai Li, Alexandros Karatzoglou, and Claudio Gentile. Collaborative filtering bandits, 2016.
- [36] John Lipor and Gautam Dasarathy. Quantile search with time-varying search parameter. In 2018 52nd Asilomar Conference on Signals, Systems, and Computers, pages 1016–1018. IEEE, 2018.
- [37] Yifei Ma, Roman Garnett, and Jeff Schneider. σ -optimality for active learning on gaussian random fields. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, Advances in Neural Information Processing Systems, volume 26. Curran Associates, Inc., 2013.
- [38] Yifei Ma, Roman Garnett, and Jeff G Schneider.  $\sigma$ -optimality for active learning on gaussian random fields. In *NIPS*, pages 2751–2759, 2013.
- [39] Yifei Ma, Tzu-Kuo Huang, and Jeff Schneider. Active search and bandits on graphs using sigma-optimality. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pages 542–551, 2015.
- [40] Yifei Ma, Tzu-Kuo Huang, and Jeff Schneider. Active search and bandits on graphs using sigma-optimality. UAI'15, page 542–551, Arlington, Virginia, USA, 2015. AUAI Press.
- [41] Shie Mannor and John N. Tsitsiklis. The sample complexity of exploration in the multi-armed bandit problem. 5:623–648, December 2004.
- [42] Kenneth Nordström. Convexity of the inverse and moore–penrose inverse. *Linear Algebra and its Applications*, 434(6):1489–1512, 2011.
- [43] Nikhil Rao, Hsiang-Fu Yu, Pradeep Ravikumar, and Inderjit S Dhillon. Collaborative filtering with graph information: Consistency and scalable methods. In *NIPS*, volume 2, page 7. Citeseer, 2015.
- [44] Herbert Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*, 58(5):527–535, 1952.
- [45] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding, 2019.
- [46] Benedek Rozemberczki and Rik Sarkar. Characteristic Functions on Graphs: Birds of a Feather, from Statistical Descriptors to Parametric Models. In *Proceedings of the 29th ACM International Conference on Information and Knowledge Management (CIKM '20)*, page 1325–1334. ACM, 2020.
- [47] Ohad Shamir. A variant of azuma's inequality for martingales with subgaussian tails, 2011.
- [48] Herbert A Simon. A behavioral model of rational choice. *The quarterly journal of economics*, 69(1):99–118, 1955.
- [49] Marta Soare, Alessandro Lazaric, and Rémi Munos. Best-arm identification in linear bandits, 2014.
- [50] Cem Tekin and Eralp Turğay. Multi-objective contextual multi-armed bandit with a dominant objective. *IEEE Transactions on Signal Processing*, 66(14):3799–3813, 2018.
- [51] Michal Valko, Rémi Munos, Branislav Kveton, and Tomáš Kocák. Spectral bandits for smooth graph functions. In *International Conference on Machine Learning*, pages 46–54. PMLR, 2014.

- [52] Dingyu Wang, John Lipor, and Gautam Dasarathy. Distance-penalized active learning via markov decision processes. In 2019 IEEE Data Science Workshop (DSW), pages 155–159. IEEE, 2019.
- [53] Liwei Wu, Hsiang-Fu Yu, Nikhil Rao, James Sharpnack, and Cho-Jui Hsieh. Graph dna: Deep neighborhood aware graph encoding for collaborative filtering. In *International Conference on Artificial Intelligence and Statistics*, pages 776–787. PMLR, 2020.
- [54] Yifan Wu, András György, and Csaba Szepesvári. Online learning with gaussian payoffs and side observations, 2015.
- [55] Wenjun Xiao and Ivan Gutman. Resistance distance and laplacian spectrum. *Theoretical chemistry accounts*, 110(4):284–289, 2003.
- [56] Kaige Yang, Xiaowen Dong, and Laura Toni. Laplacian-regularized graph bandits: Algorithms and theoretical analysis, 2020.
- [57] M Todd Young, Jacob Hinkle, Arvind Ramanathan, and Ramakrishnan Kannan. Hyperspace: Distributed bayesian hyperparameter optimization. In 2018 30th International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD), pages 339–347. IEEE, 2018.
- [58] Jifan Zhang, Julian Katz-Samuels, and Robert D. Nowak. GALAXY: graph-based active learning at the extreme. *CoRR*, abs/2202.01402, 2022.
- [59] Yuan Zhou, Xi Chen, and Jian Li. Optimal pac multiple arm identification with applications to crowdsourcing. In *International Conference on Machine Learning*, pages 217–225. PMLR, 2014.
- [60] Xiaojin Jerry Zhu. Semi-supervised learning literature survey. 2005.

# Checklist

- 1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes] Please refer to Appendix A for a full discussion.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] Please refer to Appendix A.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
  - (b) Did you include complete proofs of all theoretical results? [Yes] we include them in the supplementary material.
- 3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We include the code in supplementary material.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] We specify them in the experiment section 7
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes]
  - (b) Did you mention the license of the assets? [No] We did not.
  - (c) Did you include any new assets either in the supplemental material or as a URL? [No] We did not.
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] Data is public.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No]
- 5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

# Appendix

The appendix is organized as follows. In Appendix A we provide a discussion of the results of the paper, future work, and broader impacts. Appendices B-D and Appendix I provide various supporting results and insights into our main theoretical results. Appendix E and Appendix G provide sample complexity guarantees for GRUB and  $\zeta$ -GRUB respectively. Appendix F states and proves necessary conditions on the sample complexity, and Appendix H presents a discussion on the incomparability of our graph bandits problem with that of linear bandits. Finally, Appendix J and K contain better sampling strategies and additional experiments respectively.

An anonymized repository containing the code that supports the algorithmic and experimental results of this paper may be found here (see also Appendix L):

### A Discussion and Broader Impacts

In this work, we consider the problem of best arm identification (and approximate best arm identification) when one has access to information about the similarity between the arms in the form of a graph. We propose a novel algorithm GRUB for this important family of problems and establish sample complexity guarantees for the same. In particular, our theory explicitly demonstrated that benefit of this side information (in terms of the properties of the graph) in quickly locating the best or approximate best arms. We support these theoretical findings with experimental results in both simulated and real settings.

**Future Work and Limitations.** We outline several sampling policies inspired by our theory in Section 7; an extension of our theoretical results to account for these improved sampling policies is a natural candidate for further exploration. The algorithms and theory of this paper assume knowledge of (an upper bound) on the smoothness of the reward vector with respect to the graph. While this is where one uses domain expertise, this could be hard to estimate in certain real world problems. A generalization of the algorithmic and theoretical framework proposed here that is *adaptive* to the unknown graph-smoothness is an exciting avenue for future work<sup>5,6</sup>. The sub-Gaussianity assumption of this work can also be generalized to other tail behaviors in follow up work. Another limitation of this work is that the statistical benefit of the graph-based quadratic penalization comes at a computational cost – each mean estimation step involves the inversion of an  $n \times n$  matrix which has a complexity of  $O(n^2 \log(n))$ . However, an exciting recent line of work suggests that this matrix inversion can be made significantly faster when coupled with a spectral sparsification of the graph  $G^{7,8}$  while controlling the statistical impact of such a modification. In the context of this problem, this suggests a compelling avenue for future work that studies the statistics-vs-computation tradeoffs in using graph side information.

For this work, we demonstrated the advantages of this side information in pure exploration problems, given knowledge of such an  $\epsilon$ . Extensions that consider goodness-of-fit and misspecification with respect to the graph G and smoothness parameters  $\epsilon$  are interesting avenues for follow up work. Finally, we focus on the ridge-type regularizer of the form  $\langle \mu, L_G \mu \rangle$ . For future work, it may be productive to expand to a much broader class of regularizers such as those of the form of  $||A\mu||_q^p$ , where A represents a information/ structural constraint matrix and p, q are some positive numbers.

**Potential Negative Social Impacts.** Our methods can be used for various applications such as drug discovery, advertising, and recommendation systems. In scientifically and medically critical applications, the design of the reward function becomes vital as this can have a significant impact on the output of the algorithm. One must take appropriate measures to ensure a fair and transparent outcome for various downstream stakeholders. With respect to applications in recommendation and targeted advertising systems, it is becoming increasingly evident that such systems may exacerbate

<sup>&</sup>lt;sup>5</sup>T. Tony Cai, Ming Yuan "Adaptive covariance matrix estimation through block thresholding," The Annals of Statistics, Ann. Statist. 40(4), 2014-2042, (August 2012)

<sup>&</sup>lt;sup>6</sup>Banerjee, T., Mukherjee, G., & Sun, W. (2020). Adaptive sparse estimation with side information. Journal of the American Statistical Association, 115(532), 2053-2067.

<sup>&</sup>lt;sup>7</sup>Spielman, D. A., & Teng, S. H. (2011). Spectral sparsification of graphs. SIAM Journal on Computing, 40(4), 981-1025.

<sup>&</sup>lt;sup>8</sup>Vishnoi, Nisheeth K. "Lx= b." Foundations and Trends in Theoretical Computer Science 8.1–2 (2013): 1-141.

polarization and the creation of filter-bubbles. Especially the techniques proposed in this paper could reinforce emerging polarization (which would correspond to more clustered graphs and therefore better recommendation performance) when used in such contexts. It will of course be of significant interest to mitigate such adverse outcomes by well-designed interventions or by considering multiple similarity graphs that capture various dimensions of similarity. This is a compelling avenue for future work.

# **B** Parameter estimation

At any time T, GRUB, along with the graph-side information, uses data gathered to estimate the mean  $\hat{\mu}_T$  in order to decide the sampling and elimination protocols. The following lemma gives the estimation routine used for GRUB.

**Lemma B.1.** The closed form expression of  $\hat{\mu}_T$  is given by,

$$\hat{\boldsymbol{\mu}}_T = \left(\sum_{t=1}^T \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^T + \rho L_G\right)^{-1} \left(\sum_{t=1}^T \mathbf{e}_{\pi_t} r_t^{\pi_t}\right)$$
(11)

*Proof.* Using the reward data  $\{r_{t,\pi_t}\}_{t=1}^T$  gathered up-to time T and the sampling policy  $\pi_T$ , the mean vector estimate  $\hat{\mu}_T$  is computed by solving the following laplacian-regularized least-square optimization schedule:

$$\hat{\boldsymbol{\mu}}_{T} = \underset{\boldsymbol{\mu} \in \mathbb{R}^{n}}{\operatorname{arg\,min}} \sum_{t=1}^{T} \left( \mu_{\pi_{t}} - r_{t,\pi_{t}} \right)^{2} + \rho \langle \boldsymbol{\mu}, L_{G} \boldsymbol{\mu} \rangle$$
(12)

where  $\rho > 0$  is a tunable penalty parameter. The above optimization problem can be equivalently written in the following quadratic form:

$$\hat{\boldsymbol{\mu}}_T = \operatorname*{arg\,min}_{\boldsymbol{\mu} \in \mathbb{R}^n} \left( \langle \boldsymbol{\mu}, V(\boldsymbol{\pi}_T, G) \boldsymbol{\mu} \rangle - 2 \left\langle \boldsymbol{\mu}, \left( \sum_{t=1}^T \mathbf{e}_{\pi_t} r_{t, \pi_t} \right) \right\rangle + \sum_{t=1}^T r_{t, \pi_t}^2 \right)$$

where  $V(\boldsymbol{\pi}_T, G)$  denotes,

$$V(\boldsymbol{\pi}_T, G) = \sum_{t=1}^T \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^T + \rho L_G$$
(13)

In order to obtain  $\hat{\mu}_T$ , we compute vanishing point of the gradient as follows,

$$\left( \langle \boldsymbol{\mu}, V(\boldsymbol{\pi}_T, G) \boldsymbol{\mu} \rangle - 2 \left\langle \boldsymbol{\mu}, \left( \sum_{t=1}^T \mathbf{e}_{\pi_t} r_{t, \pi_t} \right) \right\rangle + \sum_{t=1}^T r_{t, \pi_t}^2 \right) |_{\boldsymbol{\mu} = \hat{\boldsymbol{\mu}}_T} = 0$$
  

$$\Rightarrow \quad \hat{\boldsymbol{\mu}}_T = V(\boldsymbol{\pi}_T, G)^{-1} \left( \sum_{t=1}^T \mathbf{e}_{\pi_t} r_t^{\pi_t} \right)$$
(14)

The sampling policy in GRUB uses the mean estimates and their high probability confidence bounds to eliminate suboptimal arm. In the following lemma we compute the high probability confidence bounds on the estimates of the mean and introduces the idea of effective samples of each arm given the graph side information.

**Lemma B.2.** For any T > k(G) and  $i \in [n]$ , the following holds with probability no less than  $1 - \frac{\delta}{w_i(\boldsymbol{\pi}_T)}$ :

$$|\hat{\mu}_T^i - \mu_i| \le \sqrt{\frac{1}{t_{eff,i}}} \left( 2\sigma \sqrt{14 \log\left(\frac{2w_i(\boldsymbol{\pi}_T)}{\delta}\right)} + \rho \|\boldsymbol{\mu}\|_G \right)$$
(15)

where  $w_i(\boldsymbol{\pi}_T) = a_0 n t_{eff,i}^2$  for some constant  $a_0 > 0$ ,  $\hat{\mu}_T^i$  is the *i*-th coordinate of the estimate from *B.1* and,

$$t_{eff,i} = \frac{1}{\left[\left(\sum_{t=1}^{T} \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^{\top} + \rho L_G\right)^{-1}\right]_{ii}}$$

*Proof.* Let the sequence of bounded variance noise and data gathered up-to time T be denoted by  $\{\eta_t, r_{\pi_t, t}\}_{t=1}^T$ . Let  $S_T = \sum_{t=1}^T \eta_t \mathbf{e}_{\pi_t}$  and  $N_T = \sum_{t=1}^T \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^T$ . Using the closed form expression of  $\hat{\boldsymbol{\mu}}_T$  from eq. B.1, the difference between the estimate and true value  $\hat{\mu}_T^i - \mu_i$  can be obtained as follows:

$$\hat{\mu}_T^i - \mu_i = \langle \mathbf{e}_i, \hat{\boldsymbol{\mu}}_T - \boldsymbol{\mu} \rangle = \langle \mathbf{e}_i, V_T^{-1} S_T - \rho V_T^{-1} L_G \boldsymbol{\mu} \rangle$$

The deviation  $\hat{\mu}_T^i - \mu_i$  can be upper-bounded as follows:

$$|\langle \mathbf{e}_i, \hat{\boldsymbol{\mu}}_T - \boldsymbol{\mu} \rangle| \le |\langle \mathbf{e}_i, V_T^{-1} S_T \rangle| + |\langle \mathbf{e}_i, \rho V_T^{-1} L_G \boldsymbol{\mu} \rangle|$$

Further, in order to obtain the variance of the estimate  $\hat{\boldsymbol{\mu}}_T$ , we bound the deviation  $|\mu_T^i - \mu_i|$  by separately bounding  $|\langle \mathbf{e}_i, V_T^{-1}S_T \rangle|$  and  $|\langle \mathbf{e}_i \rho V_T^{-1}L_G \boldsymbol{\mu} \rangle|$ .

With regards to the first term  $\langle \mathbf{e}_i, V_T^{-1} S_T \rangle$ , note that

Using a variant of Azuma's inequality [47, 51], for any  $\kappa > 0$  the following inequality holds,

$$\mathbb{P}\left(|\langle \mathbf{e}_{i}, V_{T}^{-1}S_{T}\rangle|^{2} \leq \kappa^{2}\right) \geq 1 - 2\exp\left\{-\frac{\kappa^{2}}{56\sigma^{2}\sum_{t=1}^{T}\left(\langle \mathbf{e}_{i}, V_{T}^{-1}\mathbf{e}_{\pi_{t}}\rangle\right)^{2}}\right\}$$
(16)

Using the fact that  $V_T \succ \left(\sum_{t=1}^T \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^T\right)$ , we can further simplify the above bound using the following computation,

$$\sum_{t=1}^{T} \left( \left\langle \mathbf{e}_{i}, V_{T}^{-1} \mathbf{e}_{\pi_{t}} \right\rangle \right)^{2} = \left\langle V_{T}^{-1} \mathbf{e}_{i}, \left( \sum_{t=1}^{T} \mathbf{e}_{\pi_{t}} \mathbf{e}_{\pi_{t}}^{T} \right) V_{T}^{-1} \mathbf{e}_{i} \right\rangle$$
$$\leq \left\langle \mathbf{e}_{i}, V_{T}^{-1} \mathbf{e}_{i} \right\rangle = [V_{T}^{-1}]_{ii} \tag{17}$$

Substituting  $\delta' = 2 \exp\left\{-\frac{\kappa^2}{56\sigma^2 \sum_{t=1}^T \left(\langle \mathbf{e}_i, V_T^{-1} \mathbf{e}_{\pi_t} \rangle \right)^2}\right\}$ , we can finally conclude that given the historical data  $\mathcal{F}_{T-1}$  till time T-1, following is true with probability  $1-\delta'$ ,

$$|\langle \mathbf{e}_i, V_T^{-1} S_T \rangle|^2 \le 56\sigma^2 [V_T^{-1}]_{ii} \log\left(\frac{2}{\delta'}\right) \tag{18}$$

Second term  $\langle \mathbf{e}_i, \rho V_T^{-1} L_G \boldsymbol{\mu} \rangle$  can be upperbounded using cauchy-schwartz inequality,

$$\begin{aligned} |\langle \mathbf{e}_{i}, \rho V_{T}^{-1} L_{G} \boldsymbol{\mu} \rangle| &= \rho \langle \mathbf{e}_{i}, L_{G} \boldsymbol{\mu} \rangle_{V_{T}^{-1}} \\ &\leq \rho \sqrt{\langle \mathbf{e}_{i}, V_{T}^{-1} \mathbf{e}_{i} \rangle} \sqrt{\langle L_{G} \boldsymbol{\mu}, V_{T}^{-1} L_{G} \boldsymbol{\mu} \rangle} \\ &\leq \rho \sqrt{[V_{T}^{-1}]_{ii}} \| \boldsymbol{\mu} \|_{G} \end{aligned}$$
(19)

Combining the upperbound (19), (18) and substituting  $\delta' = \frac{\delta}{w(\pi_T)}$  we get Lemma 3.2. Hence proved.

# **C** Influence Factor

A key component in our characterization of the performance of GRUB is the *influence factor* for each arm; recall that for a given graph D,  $C_i(D)$  denotes the connected component that contains i. The influence factor for each arm is defined as,

**Definition C.1.** Let *D* be a graph on the vertex set [n]. For each  $j \in [n]$ , define **influence factor**  $\Im(j, D)$  as:

$$\Im(j,D) = \begin{cases} \min_{i \in C_j(D), i \neq j} \{r_D(i,j)^{-1}\} & \text{if } |C_j(D)| > 1\\ 0 & \text{otherwise} \end{cases}$$
(20)

where,  $r_D(i, j)$  is the resistance distance between arm i and j on graph D as in Definition 4.1.

Note that we refer the resistance distance without the parameter  $\delta$ , as the value of resistance distance is independent of the value of  $\delta$ . This happens due to the cancellation of  $\delta$  factor in  $R_{ii}+R_{jj}-R_{ji}-R_{ij}$ . The influence factor can also be thought of as the minimum influence any arm i in the connected component of arm j has over the arm j

# **D** Effective Samples

**Theorem D.1.** Let  $\pi_T$  indicate the sampling policy until time T. Let G be the given graph,  $\mathfrak{I}(.,G)$  indicates the minimum influence factor for arms. Then effective samples can be lower bounded by,

$$t_{eff,i} \ge t_i + \frac{1}{2} \lfloor \min\{\rho \mathfrak{I}(i,G), \sum_{j \in C(i)} t_j\} \rfloor$$
(21)

where  $t_i$  indicates the no. of samples of arm i and  $|\cdot|$  indicates the floor.

*Proof.* Using Lemma 1.5, we have the following bound on  $[V_T^{-1}]_{ii}$ ,

$$[V(\boldsymbol{\pi}_T, G)^{-1}]_{ii} \le \max\left\{\frac{1}{t_i + \frac{\rho\Im(i, G)}{2}}, \frac{1}{t_i + \frac{t_C - t_i}{2}}\right\}$$
(22)

where T is the total number of samples and  $t_C$  is all the samples from the connected component C(i) apart from arm *i*. Thus rewriting the equation for  $t_{\text{eff},i}$ , we get,

$$t_{\text{eff},i} \ge t_i + \frac{1}{2} \min\{\rho \mathfrak{I}(i,G), \sum_{j \in C(i)} t_j\}$$

$$(23)$$

Hence proved.

### **E GRUB** Sample complexity

In order to compute the sample complexity for GRUB, we classify the arms into two categories: competitive and non-competitive. The split of arms into these two categories is not required for the algorithm, but provides tighter complexity bounds as will be observed in this appendix. The division of the arms is contingent on its suboptimality and the structure of the provided graph side information. A modified version of the Definition (4.3) of competitive set and non-competitive set is as follows:

**Definition E.1.** Fix  $\mu \in \mathbb{R}^n$ , graph *D*, regularization parameter  $\rho$ , confidence parameter  $\delta$ , and smoothness parameter  $\epsilon$  and noise variance  $\sigma$ . We define  $\mathcal{H}$  to be the set of competitive arms and  $\mathcal{N}$  to be the set of non-competitive arms as follows:

$$\mathcal{H}(D,\boldsymbol{\mu},\boldsymbol{\delta},\boldsymbol{\rho},\boldsymbol{\epsilon}) = \left\{ j \in [n] \middle| \Delta_i \leq 2\sqrt{\frac{2}{\boldsymbol{\rho}\mathfrak{I}(i)}} \left( 2\sigma \sqrt{14 \log\left(\frac{2a_0 n \boldsymbol{\rho}^2 \mathfrak{I}(i)^2}{\boldsymbol{\delta}}\right) + \boldsymbol{\rho}\boldsymbol{\epsilon}} \right) \right\},\$$
$$\mathcal{N}(D,\boldsymbol{\mu},\boldsymbol{\delta},\boldsymbol{\rho},\boldsymbol{\epsilon}) \triangleq [n] \setminus \mathcal{H}(D,\boldsymbol{\mu},\boldsymbol{\delta},\boldsymbol{\rho},\boldsymbol{\epsilon})$$

When the context is clear, we will use suppress the dependence on the parameters in Definition E.1.

Further, we derive an expression for the worst-case sample complexity by analysing the number of samples required to eliminate arms with different difficulty levels, i.e. arms in competitive set and non-competitive set. We first derive the sample complexity results for the case when graph G is connected and then extend it to disconnected graphs.

**Lemma E.2.** Consider *n*-armed bandit problem with mean vector  $\mu \in \mathbb{R}^n$ . Let G be a given connected similarity graph on the vertex set [n], and further suppose that  $\mu$  is  $\epsilon$ -smooth. Define

$$T_{sufficient} \triangleq \sum_{i \in \mathcal{H}} \frac{1}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] + \max_{i \in \mathcal{N}} \left\{ \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right\}$$
(24)

Then, with probability at least  $1 - \delta$ , GRUB: (a) terminates in no more than  $T_{sufficient}$  rounds, and (b) returns the best arm  $a^* = \arg \max_i \mu_i$ .

*Proof.* With out loss of generality, assume that  $a^* = 1$ . Let  $\{t_i\}_{i=1}^n$  denote the number of plays of each arm up to time T. By Lemma 3.2, we can state that,

$$\mathbb{P}\left(|\hat{\mu}_{T}^{i} - \mu_{i}| \ge \gamma_{i}(\boldsymbol{\pi}_{T})\right) \le \frac{2\delta}{a_{0}nt_{\text{eff},i}^{2}}$$

$$(25)$$

where,  $\gamma_i(\boldsymbol{\pi}_T) = \beta_i(\boldsymbol{\pi}_T) \sqrt{t_{\text{eff},i}^{-1}}$  and  $\beta_i(\boldsymbol{\pi}_T) = \left(2\sigma \sqrt{14 \log\left(\frac{2a_0 n t_{\text{eff},i}^2}{\delta}\right)} + \rho \|\boldsymbol{\mu}\|_G\right)$ .

As is reflected in the elimination policy (4), at any time t, arm 1 can be mistakenly eliminated in GRUB only if  $\hat{\mu}_t^i > \hat{\mu}_t^1 + \gamma_i(\boldsymbol{\pi}_t) + \gamma_1(\boldsymbol{\pi}_t)$ . Let  $T_s$  be the stopping time of GRUB, then the total failure probability for GRUB can be upper-bounded as,

$$\mathbb{P}(\text{Failure}) \leq \sum_{t=2}^{T_s} \sum_{i=2}^n \mathbb{P}\left(\hat{\mu}_t^i \geq \hat{\mu}_t^1 + \gamma_i(\boldsymbol{\pi}_t) + \gamma_1(\boldsymbol{\pi}_t)\right)$$

Note that  $\mathbb{P}\left(\hat{\mu}_t^i \geq \hat{\mu}_t^1 + \gamma_i(\boldsymbol{\pi}_t) + \gamma_1(\boldsymbol{\pi}_t)\right) \leq \left[\mathbb{P}\left(\hat{\mu}_t^i \geq \mu^i + \gamma_i(\boldsymbol{\pi}_t)\right) + \mathbb{P}\left(\hat{\mu}_t^1 \leq \mu^1 - \gamma_1(\boldsymbol{\pi}_t)\right)\right]$ , provided that  $\gamma_i(\boldsymbol{\pi}_t), \gamma_1(\boldsymbol{\pi}_t) \leq \frac{\Delta_i}{2}$ . Hence the failure probability can be upperbounded as,

$$\mathbb{P}(\text{Failure}) \leq \sum_{i=2}^{n} \sum_{t=2}^{T_s} \left[ \mathbb{P}\left( \hat{\mu}_t^i \geq \mu^i + \gamma_i(\boldsymbol{\pi}_t) \right) + \mathbb{P}\left( \hat{\mu}_t^1 \leq \mu^1 - \gamma_1(\boldsymbol{\pi}_t) \right) \right]$$
(26)

conditioned on  $\gamma_i(\boldsymbol{\pi}_T), \gamma_1(\boldsymbol{\pi}_T) \leq \frac{\Delta_i}{2}$ .

Let  $a_0 \ge 4 \sum_{t=1}^{\infty} t_{\text{eff},i}^{-2}$ , then from Lemma 3.2,

$$\mathbb{P}(\text{Failure}) \leq \sum_{i=2}^{n} \sum_{t=2}^{T_s} \frac{2\delta}{a_0 n t_{\text{eff},i}^2} \leq \delta$$
(27)

۸

The finiteness of the infinite sum of  $t_{\text{eff},i}^{-2}$  can be found in Lemma I.13.

Thus, in order to keep  $\mathbb{P}(\text{Failure}) \leq \delta$ , it is sufficient if, at the time of elimination of arm *i*, we have enough samples to ensure,

$$\gamma_i(\boldsymbol{\pi}_T) \leq \frac{\Delta_i}{2}$$

$$\sqrt{\frac{1}{t_{\text{eff},i}}} \left( 2\sigma \sqrt{14 \log\left(\frac{2a_0 n t_{\text{eff},i}^2}{\delta}\right)} + \rho \epsilon \right) \leq \frac{\Delta_i}{2}$$
(28)

In the absence of graph information, equation (28) devolves to the same sufficiency condition for number of samples required for suboptimal arm elimination as [13], upto constant factor. Rewriting the above equation,

$$\frac{\log(a_i)}{a_i} \le \sqrt{\frac{\delta}{d_1}} \frac{\Delta_i^2}{d_0}$$
(29)

where  $d_0 = 64 \times 14\sigma^2$ ,  $d_1 = 2na_0 e^{\frac{\rho^2 \epsilon^2}{4 \times 14\sigma^2}}$  and  $a_i = \sqrt{\frac{d_1}{\delta}} t_{\text{eff},i}$ . The following bound on  $a_i$  is sufficient to satisfy eq. (29),

$$a_i \geq 2\sqrt{\frac{d_1}{\delta}} \frac{d_0}{\Delta_i^2} \log\left(\sqrt{\frac{d_1}{\delta}} \frac{d_0}{\Delta_i^2}\right)$$

Resubstituting  $t_{\text{eff},i}$ , we obtain the sufficient number of plays required to eliminate arm i as,

$$t_{\text{eff},i} \ge \frac{c_1}{\Delta_i^2} \left[ \log\left(\frac{c_2}{\delta^{\frac{1}{2}} \Delta_i^2}\right) + c_3 \right]$$
(30)

where  $c_1 = 2 \times 64 \times 14\sigma^2$ ,  $c_2 = 64 \times 14\sigma^2 \sqrt{2na_0}$  and  $c_3 = \frac{\rho^2 \epsilon^2}{8 \times 14\sigma^2}$ . In the further text we are suppressing the powers of  $\delta$ ,  $\Delta_i$  within the log factor as it adds only a constant multiple to the lower bound.

The further part of the proof we use the following bound on  $t_{\text{eff}}$ , from Theorem D.1 as follows:

$$t_{\text{eff},i} \ge t_i + \frac{1}{2}\min\left\{\rho\mathfrak{I}(i), T - t_i\right\} \quad \forall i \in [n]$$
(31)

Hence a sufficiency condition for the GRUB to produce the best-arm with probability  $1 - \delta$  is given when both the following conditions are satisfied,

$$t_i + \frac{\rho \mathfrak{I}(i)}{2} \ge \frac{1}{\Delta_i^2} \left[ c_1 \log \left( \frac{c_2}{\delta \Delta_i} \right) + \frac{\rho \epsilon}{2} \right]$$
(32)

and,

$$T + t_i \ge T \ge \frac{2}{\Delta_i^2} \left[ c_1 \log\left(\frac{c_2}{\delta \Delta_i}\right) + \frac{\rho \epsilon}{2} \right]$$
(33)

From the Definition E.1 of competitive arms  $\mathcal{H}$  and non-competitive arms  $\mathcal{N}$ , we have,

$$\mathcal{H} = \left\{ j \in [n] \middle| \Delta_i \le 2\sqrt{\frac{2}{\rho \mathfrak{I}(i)}} \left( 2\sigma \sqrt{14 \log\left(\frac{2a_0 n \rho^2 \mathfrak{I}(i)^2}{\delta}\right)} + \rho \epsilon \right) \right\}$$
(34)

After the first  $\max_{i \in \mathcal{N}} \left\{ \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right\}$  samples, all arms in  $\mathcal{N}$  are eliminated. Further, let  $k_1$  be the index of the first arm to be eliminated (in  $\mathcal{H}$ ) and  $t_{k_1}^*$  be the number of samples of arm  $k_1$  before getting eliminated then the total number of additional time steps played until the arm  $k_1$  is eliminated is at most  $|\mathcal{H}|t_{k_1}^*$ . Let  $k_2$  be the index of the next arm in  $\mathcal{H}$  to be eliminated. The number of additional plays until the next arm is eliminated is given by  $(|\mathcal{H}| - 1)[t_{k_2}^* - t_{k_1}^*]$  and so on.

Summing up all the samples required to converge to the optimal arm is given by, (let  $t_{k_0}^* = 0$ )

$$\sum_{h=1}^{|\mathcal{H}|} (|\mathcal{H}| - h))[t_{k_h}^* - t_{k_{h-1}}^*] = \sum_{h=1}^{|\mathcal{H}| - 1} t_{k_h}^* = \sum_{i \in \mathcal{H}/1} t_i^*$$
(35)

Hence the final sample complexity can be computed as follows:

• Number of plays required for arms in  $\mathcal{H}$ :

$$\sum_{i \in \mathcal{H}/1} t_i^* \ge \sum_{i \in \mathcal{H}/1} \frac{1}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right]$$
(36)

• Number of plays required for all the arms in  $\mathcal{N} := [n]/\mathcal{H}$  to be eliminated:

$$T \ge \max_{i \in \mathcal{N}} \left\{ \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right\}$$
(37)

Hence the final sample complexity can be given by,

$$T_{\text{sufficient}} \triangleq \max_{i \in \mathcal{N}} \left\{ \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right\} + \sum_{i \in \mathcal{H}/1} \frac{1}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right]$$
(38)

Hence proved.

We extend Lemma E.2 to the case when graph G has disconnected clusters.

Note: The following theorem stated in the main paper has a typographical error in the equation for  $T_{\text{sufficient}}$  in place of  $\arg \min it$  is supposed to be min.

**Theorem E.3.** Consider *n*-armed bandit problem with mean vector  $\boldsymbol{\mu} \in \mathbb{R}^n$ . Let  $\mathcal{G}$  be the set of subgraphs of given similarity graph G on the vertex set [n], and further suppose that  $\boldsymbol{\mu}$  is  $\epsilon$ -smooth. Define

$$T_{sufficient} \triangleq \min_{D \in \mathcal{G}} \sum_{C \in \mathcal{C}_D} \left| \sum_{i \in C \cap \mathcal{H}_D} \frac{1}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] + \max_{i \in C \cap \mathcal{N}_D} \left\{ \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right\} \right|$$
(39)

where  $\Delta_i = \mu^* - \mu_i$  for all suboptimal arms,  $\mathcal{H}_D$  and  $\mathcal{N}_D$  are as in Definition E.1,  $\mathcal{C}_D$  is the set of connected components of a subgraph  $D \in \mathcal{G}$  and  $c_1, c_2$  are constants independent of system parameters. Then, with probability at least  $1 - \delta$ , GRUB: (a) terminates in no more than  $T_{sufficient}$ rounds, and (b) returns the best arm  $a^* = \arg \max_i \mu_i$ .

*Proof.* Let  $C_G$  denote the connected components of graph G. From Lemma E.2, the number of samples for each connected component  $C \in C_G$  can be given as,

$$T_{\text{sufficient}} = \left[ \sum_{i \in C \cap \mathcal{H}} \frac{1}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] + \max_{i \in C \cap \mathcal{N}} \left\{ \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right\} \right]$$
(40)

We can obtain the sample complexity for obtaining the best arm by summing it over all the components  $C \in C$ , gives us the sample complexity for GRUB while considering graph G.

$$T_{\text{sufficient}} = \sum_{C \in \mathcal{C}_G} \left[ \sum_{i \in C \cap \mathcal{H}} \frac{1}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] + \max_{i \in C \cap \mathcal{N}} \left\{ \frac{2}{\Delta_i^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i} + \frac{\rho \epsilon}{2} \right] \right\} \right]$$
(41)

Any subgraph D of graph G satisfies,

$$\langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle \le \epsilon \Rightarrow \langle \boldsymbol{\mu}, L_D \boldsymbol{\mu} \rangle \le \epsilon$$
 (42)

As seen in Definition E.1, the influence factor is instrumental in deciding the competitive and non-competitive sets, which further dictates the sample complexity bounds. Further, notice from Lemma I.8 that the influence factor  $\Im(i, D)$  is not monotonic when considering subgraph D of graph G. Hence considering a subgraph of G could potentially increase the number of non-competitive arms and provide us with a tighter bound on the performance for GRUB.

Hence  $T_{\text{sufficient}}$  in (40) can be made tighter by considering the minimum value over the entire set of subgraphs  $\mathcal{G}$ .

We next derive sample complexity upper bounds for GRUB in certain illuminating special cases.

**Corollary E.4** (Isolated clusters). Consider the setup as in Theorem 4.4 with the further restriction that G consists of a subgraph F such that optimal node is isolated and arms [2, ..., n] are split in k

clusters and 
$$\Delta_i \ge 2\sqrt{\frac{2}{\rho\Im(i,F)}} \left(2\sigma\sqrt{14\log\left(\frac{2a_0n\rho^2\Im(i,F)^2}{\delta}\right)} + \rho\epsilon\right), \forall i \in [2,\dots,n].$$
 Define  

$$T_{sufficient} \triangleq \sum_{C \in \mathcal{C}_F/1} \max_{j \in C} \frac{2}{\Delta_j^2} \left[c_1\log\left(\frac{c_2}{\delta\Delta_i}\right) + \frac{\rho\epsilon}{2}\right]$$
(43)

Then, with probability at least  $1 - \delta$ , GRUB: (a) terminates in no more than  $T_{sufficient}$  rounds, and (b) returns the best arm  $a^* = \arg \max_i \mu_i$ .

#### Algorithm 1 GRUB

**Input:** Regularization parameter  $\rho$ , Smoothness parameter  $\epsilon$ , Error bound  $\delta$ , Total arms n, Laplacian  $L_G$ , Sub-gaussianity parameter  $\sigma$  $t \leftarrow 0$  $A = \{1, 2, \dots, n\}$ t = 0 $V_0 \leftarrow \rho L_G$  $\mathcal{C}(G) \leftarrow \texttt{Cluster-Identification}(L_G)$ for  $C \in \mathcal{C}(G)$  do  $t \leftarrow t + 1$ Pick random arm  $k \in C$  to observe reward  $r_{t,k}$  $V_t \leftarrow V_{t-1} + \mathbf{e}_k \mathbf{e}_k^T$ , and  $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + r_{t,k} \mathbf{e}_k$ end for while |A| > 1 do  $t \leftarrow t + 1$ for  $i \in A$  do  $t_{\text{eff},i} \leftarrow ([V_t^{-1}]_{ii})^{-1}$  $\beta_i(t) \leftarrow 2\sigma \sqrt{14 \log\left(\frac{2nt_{\text{eff},i}^2}{\delta}\right)} + \rho \epsilon$ end for  $k \leftarrow \text{Sampling-Policy}(t, V_t, A, \mathcal{C}(G))$ Sample arm k to observe reward  $r_{t,k}$  $V_t \leftarrow V_{t-1} + \mathbf{e}_k \mathbf{e}_k^T$  $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + r_{t,k} \mathbf{e}_k$  $\hat{\boldsymbol{\mu}}_t \leftarrow V_t^{-1} \mathbf{x}_t$  $\begin{aligned} \mu_t &\leftarrow \operatorname{arg\,max}_{i \in A} \left[ \hat{\mu}_t^i - \beta(t_i) \sqrt{t_{\mathrm{eff},i}^{-1}} \right] \\ A &\leftarrow \left\{ \mathbf{a} \in A \mid \hat{\mu}_t^{a_{\max}} - \hat{\mu}_t^a \leq \beta_a(t) \sqrt{t_{\mathrm{eff},a}^{-1}} \right. \\ &\left. + \beta_{a_{\max}}(t) \sqrt{t_{\mathrm{eff},a_{\max}}^{-1}} \right\} \end{aligned}$ end while return A

Corollary E.4 shows that in scenarios where the arms are well clustered, the sample complexity of GRUB can scale with the number of clusters, a quantity that is typically significantly smaller than the total number of nodes in the graph.

**Corollary E.5** (Star graph). Consider the setup as in Theorem 4.4 with the further restriction that G consists of a star subgraph with the central node as the optimal arm and  $\Delta_i \leq$ 

$$2\sqrt{\frac{2}{\rho\Im(i,F)}} \left(2\sigma\sqrt{14\log\left(\frac{2a_0n\rho^2\Im(i,F)^2}{\delta}\right)} + \rho\epsilon\right), \forall i \in [2,\dots,n]. \text{ Define}$$
$$T_{sufficient} \triangleq \sum_{i=2}^n \frac{1}{\Delta_i^2} \left[c_1\log\left(\frac{c_2}{\delta\Delta_i}\right) + \frac{\rho\epsilon}{2}\right]$$
(44)

Then, with probability at least  $1 - \delta$ , GRUB: (a) terminates in no more than  $T_{sufficient}$  rounds, and (b) returns the best arm  $a^* = \arg \max_i \mu_i$ .

In Corollary E.5,  $T_{\text{sufficient}}$  is the same sample complexity as vanilla best arm identification, upto constant factors which is due to the fact that pulling one of the spoke arms does not yield much information about the other spoke arms, and this is the exact situation in the standard pure exploration setting.

# F Lower bounds

In this section we give a lower bound on the sample complexity for any  $\delta$ -PAC to return the best arm for a n armed bandit problem along with graph side information.

**Theorem F.1.** Given an n-armed bandit model with associated mean vector  $\boldsymbol{\mu} \in \mathbb{R}^n$  and similarity graph G smooth on  $\boldsymbol{\mu}$ , i.e.  $\langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle \leq \epsilon$ , for any  $0 < \epsilon < \epsilon_0$ . Let G = ([n], E) be the graph with only k isolated cliques and w.l.o.g let arm 1 be the optimal arm. Then define

$$T_{necessary} = \sum_{C \in \mathcal{C}_G/C^*} \min_{j \in C} \left\{ \frac{4\sigma^2 \log 5}{(\Delta_j - \sqrt{\epsilon})^2} \right\} + \sum_{j \in C^*/1} \frac{4\sigma^2 \log 5}{\Delta_j^2}$$
(45)

where  $C^*$  is the clique with the optimal arm and  $\epsilon_0 := \min_{i \in [n]/1, j \in C(i)} \left[ \Delta_j \left[ 1 - \frac{\Delta_i}{\sqrt{\Delta_i^2 + \Delta_j^2}} \right] \right]^2$ . Then any  $\delta$ -PAC algorithm will need at-least  $T_{necessary}$  steps to terminate, provided  $\delta \leq 0.1$ .

*Proof.* We prove the theorem in two steps: Firstly, we construct the sample complexity lower bound for the similarity graph with the isolated optimal arm and a clique of rest of the sub-optimal arms, followed by step 2 the sample complexity lower bound for a graph with single cluster

#### Step 1:

Consider a n + 1 armed bandit problem with mean vector  $\boldsymbol{\mu} \in \mathbb{R}^{n+1}$  and similarity graph M with an isolated optimal arm (arm 1) and *n*-clique cluster of suboptimal arms, satisfying the condition for smoothness of rewards over the graph, i.e.,  $\langle \boldsymbol{\mu}, L_M \boldsymbol{\mu} \rangle \leq \epsilon$ . Then the following holds

$$\max_{i \neq 1} \mu_i \le \min_{j \neq 1} \left\{ \mu_j + \sqrt{\epsilon} \right\} \tag{46}$$

Assume that ordering of mean in *n*-clique of suboptimal arms is known. From [26], there exists a  $\delta$ -PAC algorithm, for  $\delta \leq 0.1$ , which can successful identify the best arm for the subproblem with just the optimal arm and arm with the maximum mean in the *n*-clique cluster, i.e.  $j' = \arg \max \mu_j$ 

with the total number of samples given by,

$$T \ge \frac{4\log 5\sigma^2}{\Delta_{i'}^2} \tag{47}$$

Now consider the case where the ordering of the mean in *n*-clique is unknown. In order to remove all the suboptimal arms provided  $\epsilon \leq \min_{j \neq 1} \Delta_j^2$  and (46) holds, it is suffices to be able to distinguish between the optimal arm and a hypothetical suboptimal arm with mean  $\mu_j + \sqrt{\epsilon}$  where *j* is any arm from suboptimal *n*-clique, and the minimum number of samples required by any  $\delta$ -PAC algorithm to successfully identify the best arm with  $\delta \leq 0.1$  is given by,

$$T \ge \frac{4\log 5\sigma^2}{(\Delta_j - \sqrt{\epsilon})^2} \tag{48}$$

The best performance in terms of sample complexity out of all the random choice of arm from the suboptimal *n*-clique cluster is,

$$T \ge \min_{j \ne 1} \left\{ \frac{4\log 5\sigma^2}{(\Delta_j - \sqrt{\epsilon})^2} \right\}$$
(49)

Given  $\epsilon_0 := \min_{i \in [n]/1, j \in C(i)} \left[ \Delta_j \left[ 1 - \frac{\Delta_i}{\sqrt{\Delta_i^2 + \Delta_j^2}} \right] \right]^2$  and  $\epsilon < \epsilon_0$ , it can be verified that for any arm  $i, j \neq 1$ ,

$$\min_{j \neq 1} \frac{4\log 5\sigma^2}{(\Delta_j - \sqrt{\epsilon})^2} < \frac{4\log 5\sigma^2}{\Delta_i^2} + \frac{4\log 5\sigma^2}{\Delta_j^2}$$
(50)

where the left hand side corresponds to the sample complexity lower bound of removing the suboptimal arms i, j with the graph side information and the right hand side corresponds to the same without graph side information.

Hence it can be inferred that it is inefficient to remove the arms individually (disregarding the graph information).

### Step 2 :

Consider a n + 1 armed bandit problem with mean vector  $\boldsymbol{\mu} \in \mathbb{R}^{n+1}$  with a given similarity graph N such that  $\langle \boldsymbol{\mu}, L_N \boldsymbol{\mu} \rangle \leq \epsilon$ . Let all the suboptimal arms be connected to the optimal arm.

Here we show by an adversarial example that it is not possible to have a lower bound on the sample complexity which scales better than,

$$T \ge \sum_{j \ne 1} \frac{4 \log 5\sigma^2}{\Delta_j^2} \tag{51}$$

There exists a  $\delta$ -PAC algorithm which can determine that arms  $j = 3, \ldots, n$  are suboptimal after  $T \ge \sum_{j \ne 1,2} \frac{1}{\Delta_i^2}$  samples. From the smoothness of rewards on the similarity graph N we know that,

$$-\sqrt{\epsilon} \le \mu_1 - \mu_j \le \sqrt{\epsilon} \quad \forall j \in [2, 3, \dots, n]$$
(52)

This information does not help us identify or even reduce the number of samples required to identify optimal arm between arm 1 and arm 2. Thus no  $\delta$ -PAC algorithm,  $\delta \leq 0.1$ , can determine the optimal arm from arm 1 and arm 2 without an additional  $\frac{4 \log 5\sigma^2}{\Delta_2^2}$  samples for determining the best arm.

Using above two steps, we construct the proof for lower bound as follows:

Now consider the graph side information as defined in the theorem, and let  $C_G$  denote the set of connected components of graph G and  $C^* \in C_G$  be the component containing the optimal arm. Finding the best arm in this setup requires elimination of the suboptimal arms with in the connected component containing optimal arm  $j \in C^*$  and elimination of the other connected components with suboptimal arms  $j \in C_G/C^*$ . Hence, the sample complexity lower bounds [26, 27] for any  $\delta$ -PAC algorithm with  $\delta \leq 0.1$  to eliminate these arms using the tools developed in step 1 and step 2, is given by

$$T \geq \sum_{j \in C^*/1} \frac{4\sigma^2 \log 5}{\Delta_j^2} + \sum_{C \in \mathcal{C}_G/C^*} \min_{j \in C} \left\{ \frac{4\sigma^2 \log 5}{(\Delta_j - \sqrt{\epsilon})^2} \right\}$$
(53)

# G (-GRUB Sample complexity proof

**Definition G.1.** Fix  $\mu \in \mathbb{R}^n$ , graph *D*, confidence parameter  $\delta$ , noise variance  $\sigma$ , and relaxation parameter  $\zeta$ . We define  $\mathcal{H}$  to be the set of competitive arms and  $\mathcal{N}$  to be the set of non-competitive arms for  $\zeta$ -GRUB as follows:

$$\mathcal{H}(D,\boldsymbol{\mu},\delta,\zeta) = \left\{ j \in [n] \middle| \Delta_i^{\zeta} \le 2\sqrt{\frac{2}{\rho \Im(i)}} \left( 2\sigma \sqrt{14 \log\left(\frac{2a_0 n \rho^2 \Im(i)^2}{\delta}\right)} + \rho \epsilon \right) \right\},\$$
$$\mathcal{N}(D,\boldsymbol{\mu},\delta,\zeta) \triangleq [n] \setminus \mathcal{H}(D,\boldsymbol{\mu},\delta,\zeta)$$

where  $\Delta_i^{\zeta} \triangleq \max{\{\Delta_i, \zeta\}}.$ 

**Lemma G.2.** Consider *n*-armed bandit problem with mean vector  $\mu \in \mathbb{R}^n$ . Let G be a given connected similarity graph on the vertex set [n], and further suppose that  $\mu$  is  $\epsilon$ -smooth. Define

$$T_{sufficient} \triangleq \sum_{i \in \mathcal{H}} \frac{1}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] + \max_{i \in \mathcal{N}} \left\{ \frac{2}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] \right\}$$
(54)

where  $\Delta_i^{\zeta} \triangleq \max{\{\Delta_i, \zeta\}}$ . Then, with probability at least  $1 - \delta$ , GRUB: (a) terminates in no more than  $T_{sufficient}$  rounds, and (b) returns a  $\zeta$ -best arm

*Proof.* With out loss of generality, assume that  $a^* = 1$ . Let  $\{t_i\}_{i=1}^n$  denote the number of plays of each arm up to time T. By Lemma 3.2, we can state that,

$$\mathbb{P}\left(|\hat{\mu}_T^i - \mu_i| \ge \gamma_i(\boldsymbol{\pi}_T)\right) \le \frac{2\delta}{a_0 n t_{\text{eff},i}^2}$$
(55)

where, 
$$\gamma_i(\boldsymbol{\pi}_T) = \beta_i(\boldsymbol{\pi}_T) \sqrt{t_{\text{eff},i}^{-1}}$$
 and  $\beta_i(\boldsymbol{\pi}_T) = \left(2\sigma \sqrt{14 \log\left(\frac{2a_0 n t_{\text{eff},i}^2}{\delta}\right)} + \rho \|\boldsymbol{\mu}\|_G\right)$ 

As is reflected in the elimination policy (4), at any time t, arm 1 can be mistakenly eliminated in GRUB only if  $\hat{\mu}_t^i > \hat{\mu}_t^1 + \gamma_i(\boldsymbol{\pi}_t) + \gamma_1(\boldsymbol{\pi}_t)$ . Let  $T_s$  be the stopping time of GRUB, then the total failure probability for GRUB can be upper-bounded as,

$$\mathbb{P}(\text{Failure}) \leq \sum_{t=2}^{T_s} \sum_{i=2}^n \mathbb{P}\left(\hat{\mu}_t^i \geq \hat{\mu}_t^1 + \gamma_i(\boldsymbol{\pi}_t) + \gamma_1(\boldsymbol{\pi}_t)\right)$$

Note that  $\mathbb{P}\left(\hat{\mu}_t^i \geq \hat{\mu}_t^1 + \gamma_i(\boldsymbol{\pi}_t) + \gamma_1(\boldsymbol{\pi}_t)\right) \leq \left[\mathbb{P}\left(\hat{\mu}_t^i \geq \mu^i + \gamma_i(\boldsymbol{\pi}_t)\right) + \mathbb{P}\left(\hat{\mu}_t^1 \leq \mu^1 - \gamma_1(\boldsymbol{\pi}_t)\right)\right]$ , provided that  $\gamma_i(\boldsymbol{\pi}_t), \gamma_1(\boldsymbol{\pi}_t) \leq \frac{\Delta_i^{\varsigma}}{2}$ . Hence the failure probability can be upperbounded as,

$$\mathbb{P}(\text{Failure}) \leq \sum_{i=2}^{n} \sum_{t=2}^{T_s} \left[ \mathbb{P}\left( \hat{\mu}_t^i \geq \mu^i + \gamma_i(\boldsymbol{\pi}_t) \right) + \mathbb{P}\left( \hat{\mu}_t^1 \leq \mu^1 - \gamma_1(\boldsymbol{\pi}_t) \right) \right]$$
(56)

conditioned on  $\gamma_i(\boldsymbol{\pi}_T), \gamma_1(\boldsymbol{\pi}_T) \leq \frac{\Delta_i^{\zeta}}{2}$ .

Let  $a_0 \ge 4 \sum_{t=1}^{\infty} t_{\text{eff},i}^{-2}$ , then from Lemma 3.2,

$$\mathbb{P}(\text{Failure}) \leq \sum_{i=2}^{n} \sum_{t=2}^{T_s} \frac{2\delta}{a_0 n t_{\text{eff},i}^2} \leq \delta$$
(57)

The finiteness of the infinite sum of  $t_{\text{eff},i}^{-2}$  can be found in Lemma I.13.

Thus, in order to keep  $\mathbb{P}(\text{Failure}) \leq \delta$ , it is sufficient if, at the time of elimination of arm *i*, we have enough samples to ensure,

$$\gamma_{i}(\boldsymbol{\pi}_{T}) \leq \frac{\Delta_{i}^{\zeta}}{2}$$

$$\sqrt{\frac{1}{t_{\text{eff},i}}} \left( 2\sigma \sqrt{14 \log\left(\frac{2a_{0}nt_{\text{eff},i}^{2}}{\delta}\right)} + \rho\epsilon \right) \leq \frac{\Delta_{i}^{\zeta}}{2}$$
(58)

Rewriting the above equation,

$$\frac{\log\left(a_{i}\right)}{a_{i}} \leq \sqrt{\frac{\delta}{d_{1}}} \frac{(\Delta_{i}^{\zeta})^{2}}{d_{0}}$$
(59)

where  $d_0 = 64 \times 14\sigma^2$ ,  $d_1 = 2na_0 e^{\frac{\rho^2 \epsilon^2}{4 \times 14\sigma^2}}$  and  $a_i = \sqrt{\frac{d_1}{\delta}} t_{\text{eff},i}$ . The following bound on  $a_i$  is sufficient to satisfy eq. (59),

$$a_i \geq 2\sqrt{\frac{d_1}{\delta}} \frac{d_0}{(\Delta_i^{\zeta})^2} \log\left(\sqrt{\frac{d_1}{\delta}} \frac{d_0}{(\Delta_i^{\zeta})^2}\right)$$

Resubstituting  $t_{eff,i}$ , we obtain the sufficient number of plays required to eliminate arm i as,

$$t_{\text{eff},i} \ge \frac{c_1}{(\Delta_i^{\zeta})^2} \left[ \log\left(\frac{c_2}{\delta^{\frac{1}{2}}(\Delta_i^{\zeta})^2}\right) + c_3 \right]$$
(60)

where  $c_1 = 2 \times 64 \times 14\sigma^2$ ,  $c_2 = 64 \times 14\sigma^2 \sqrt{2na_0}$  and  $c_3 = \frac{\rho^2 \epsilon^2}{8 \times 14\sigma^2}$ .

The further part of the proof depends crucially on the following bound on  $t_{\text{eff},i}$  for all  $i \in [n]$  from Theorem D.1 as follows:

$$t_{\text{eff},i} \ge t_i + \frac{1}{2} \min\left\{\rho \Im(i), T - t_i\right\}$$
(61)

Hence a sufficiency condition for the GRUB to produce the  $\zeta$ -best arm with probability  $1 - \delta$  is given when both the following conditions are satisfied,

$$t_i + \frac{\rho \Im(i)}{2} \ge \frac{1}{(\Delta_i^{\zeta})^2} \left[ c_1 \log\left(\frac{c_2}{\delta \Delta_i^{\zeta}}\right) + \frac{\rho \epsilon}{2} \right]$$
(62)

and,

$$T + t_i \ge T \ge \frac{2}{(\Delta_i^{\zeta})^2} \left[ c_1 \log\left(\frac{c_2}{\delta \Delta_i^{\zeta}}\right) + \frac{\rho \epsilon}{2} \right]$$
(63)

From the Definition G.1 we have the set of competitive arms  $\mathcal{H}$  and non-competitive arms  $\mathcal{N}$  as follows:

$$\mathcal{H} = \left\{ j \in [n] \middle| \Delta_i^{\zeta} \le 2\sqrt{\frac{2}{\rho \Im(i)}} \left( 2\sigma \sqrt{14 \log\left(\frac{2a_0 n \rho^2 \Im(i)^2}{\delta}\right)} + \rho \epsilon \right) \right\}$$
(64)

After the first  $\max_{i \in \mathcal{N}} \left\{ \frac{2}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] \right\}$  samples, all arms in  $\mathcal{N}$  are eliminated. Further, let  $k_1$  be the index of the first arm to be eliminated (in  $\mathcal{H}$ ) and  $t_{k_1}^*$  be the number of samples of arm  $k_1$  before getting eliminated then the total number of additional time steps played until the arm  $k_1$  is eliminated is at most  $|\mathcal{H}|t_{k_1}^*$ . Let  $k_2$  be the index of the next arm in  $\mathcal{H}$  to be eliminated. The number of additional plays until the next arm is eliminated is given by  $(|\mathcal{H}| - 1)[t_{k_2}^* - t_{k_1}^*]$  and so on.

Summing up all the samples required to converge to the optimal arm is given by, (let  $t_{k_0}^* = 0$ )

$$\sum_{h=1}^{|\mathcal{H}|} (|\mathcal{H}| - h))[t_{k_h}^* - t_{k_{h-1}}^*] = \sum_{h=1}^{|\mathcal{H}|-1} t_{k_h}^* = \sum_{i \in \mathcal{H}/1} t_i^*$$
(65)

Hence the final sample complexity can be computed as follows:

• Number of plays required for arms in  $\mathcal{H}$ :

$$\sum_{i \in \mathcal{H}/1} t_i^* \ge \sum_{i \in \mathcal{H}/1} \frac{1}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right]$$
(66)

• Number of plays required for all the arms in  $\mathcal{N} := [n]/\mathcal{H}$  to be eliminated:

$$T \ge \max_{i \in \mathcal{N}} \left\{ \frac{2}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] \right\}$$
(67)

Hence the final sample complexity can be given by,

$$T_{\text{sufficient}} \triangleq \max_{i \in \mathcal{N}} \left\{ \frac{2}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] \right\} + \sum_{i \in \mathcal{H}/1} \frac{1}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right]$$
(68)

We extend Lemma G.2 to the case when graph G has disconnected clusters.

Note: The following theorem stated in the main paper has a typographical error in the equation for  $T_{\text{sufficient}}$  in place of arg min it is supposed to be min.

**Theorem G.3.** Consider *n*-armed bandit problem with mean vector  $\boldsymbol{\mu} \in \mathbb{R}^n$ . Let  $\mathcal{G}$  be the set of subgraphs given similarity graph G on the vertex set [n], and further suppose that  $\boldsymbol{\mu}$  is  $\epsilon$ -smooth. Define

$$T_{sufficient} \triangleq \min_{D \in \mathcal{G}} \sum_{C \in \mathcal{C}_D} \left[ \sum_{i \in C \cap \mathcal{H}_D} \frac{1}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] + \max_{i \in C \cap \mathcal{N}_D} \left\{ \frac{2}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] \right\} \right]$$
(69)

#### Algorithm 2 $\zeta$ -GRUB

**Input:** Regularization parameter  $\rho$ , Smoothness parameter  $\epsilon$ , Error bound  $\delta$ , Total arms n, Laplacian  $L_G$ , Sub-gaussianity parameter  $\sigma$  $t \leftarrow 0$  $A = \{1, 2, \dots, n\}$ t = 0 $V_0 \leftarrow \rho L_G$  $\mathcal{C}(G) \leftarrow \text{Cluster-Identification}(L_G)$ for  $C \in \mathcal{C}(G)$  do  $t \leftarrow t + 1$ Pick random arm  $k \in C$  to observe reward  $r_{t,k}$  $V_t \leftarrow V_{t-1} + \mathbf{e}_k \mathbf{e}_k^T$ , and  $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + r_{t,k} \mathbf{e}_k^T$ end for while |A| > 1 do  $t \leftarrow t + 1$  $\beta(t) \leftarrow 2\sigma \sqrt{14 \log\left(\frac{2n(t+1)^2}{\delta}\right)} + \rho \epsilon$  $k \leftarrow \text{Sampling-Policy}(t, V_t, A, \mathcal{C}(G))$ Sample arm k to observe reward  $r_{t,k}$  $V_t \leftarrow V_{t-1} + \mathbf{e}_k \mathbf{e}_k^T$  $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + r_{t,k} \mathbf{e}_k$  $\hat{\boldsymbol{\mu}}_t \leftarrow V_t^{-1} \mathbf{x}_t$  $\mu_t \leftarrow v_t \quad x_t$   $a_{\max} \leftarrow \underset{i \in A}{\operatorname{arg\,max}} \left[ \hat{\mu}_t^i - \beta(t_i) \sqrt{[V_t^{-1}]_{ii}} \right]$   $A \leftarrow \left\{ \mathbf{a} \in A \mid \hat{\mu}_t^{a_{\max}} - \hat{\mu}_t^a \leq \beta(t_a) \sqrt{[V_t^{-1}]_{aa}} \right\}$  $+\beta(t_{a_{\max}})\sqrt{[V_t^{-1}]_{a_{\max}a_{\max}}}$  $A \leftarrow A / \left\{ a \in A \mid \beta(t_a) \sqrt{[V_t^{-1}]_{aa}} \leq \frac{\zeta}{2} \right\}$ end while return  $\arg \max \left\{ \mu_i | i \in \{a \in [n] | \beta(t_a) \sqrt{[V_t^{-1}]_{aa}} \le \frac{\zeta}{2} \} \cup A \right\}$ 

where  $\Delta_i^{\zeta} = \max{\{\Delta_i, \zeta\}}$  for all suboptimal arms,  $\mathcal{H}_D$  and  $\mathcal{N}_D$  are as in Definition G.1,  $\mathcal{C}_D$  is the set of connected components of subgraph  $D \in \mathcal{G}$  and  $c_1, c_2$  are constants independent of system parameters. Then, with probability at least  $1 - \delta$ , GRUB: (a) terminates in no more than  $T_{sufficient}$  rounds, and (b) returns a  $\zeta$ -best arm

*Proof.* From Lemma G.2, the sample complexity for each connected component  $C \in C$  can be given as,

$$T_{\text{sufficient}} = \left[\sum_{i \in C \cap \mathcal{H}} \frac{1}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] + \max_{i \in C \cap \mathcal{N}} \left\{ \frac{2}{(\Delta_i^{\zeta})^2} \left[ c_1 \log \frac{c_2}{\delta \Delta_i^{\zeta}} + \frac{\rho \epsilon}{2} \right] \right\} \right]$$
(70)

where, summing it over all the components  $C \in C$ , gives us the sample complexity for GRUB while considering graph G.

Any subgraph D of graph G satisfies,

$$\langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle \le \epsilon \Rightarrow \langle \boldsymbol{\mu}, L_D \boldsymbol{\mu} \rangle \le \epsilon$$
 (71)

As seen in Definition G.1, the influence factor is instrumental in deciding the competitive and non-competitive sets, which further dictates the sample complexity bounds. Further, notice from Lemma I.8 that the influence factor  $\Im(i, D)$  is not monotonic when considering subgraph D of graph G. Hence considering a subgraph of G could potentially increase the number of non-competitive arms and provide us with a tighter bound on the performance for GRUB.

Hence  $T_{\text{sufficient}}$  can be made tighter by considering the minimum value over the entire set of subgraphs  $\mathcal{G}$ .

Note that, as in the case of GRUB, the  $\zeta$ -GRUB algorithm's performance *automatically* adapts to the best possible subgraph in  $\mathcal{G}$ .

### H The Incomparability of the Graph Bandits problem with Linear Bandits

In this section, we demonstrate an example graph bandit problem that is cast as a linear bandit to reveal the incomparability of these frameworks. A typical linear bandit problem is defined as follows: Consider an *n*-armed linear bandit problem, each arm  $i \in [n]$  is associated with a feature vector  $\mathbf{x}_i \in \mathbb{R}^d$ , where *d* can be lower than *n*. In each round *t*, the learner chooses an action  $\mathbf{a}_t = \mathbf{x}_i$  for some  $i \in [n]$  and observes the reward  $y_t = \langle \mathbf{a}_t, \boldsymbol{\theta} \rangle + \eta_t$ , where  $\boldsymbol{\theta} \in \mathbb{R}^d$  is an unknown parameter and the  $\eta_t$  is a subgaussian random noise with  $\sigma^2$  variance. Denote the arm with the best mean reward with  $i^*$ , i.e.  $i^* = \arg \max_{i \in [n]} \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle$ . The goal of the learner is to to output the index of the arm  $i^*$  with probability  $1 - \delta$ ,  $\delta > 0$  in as few samples as possible.

Firstly, a *n*-armed bandit problem without any graph can be easily seen as linear bandits by associating the canonical basis for  $\mathbb{R}^n \{\mathbf{e}_i\}_{i=1}^n$  as the feature vectors and the mean vector  $\boldsymbol{\mu} \in \mathbb{R}^n$  as the unknown reward vector. This provides up with the mean reward function for arm  $i \in [n]$  as  $\langle \mathbf{e}_i, \boldsymbol{\mu} \rangle = \mu_i$ .

In order to cast the graph bandit problem in a linear bandit framework, we need to associate every arm index *i* with a feature vector  $\mathbf{x}_i$  and identify the unknown feature vector  $\boldsymbol{\theta}$  for the problem. We achieve this by modifying the feature vectors  $\{\mathbf{e}_i\}_{i=1}^n$  and the reward vector  $\boldsymbol{\mu}$  based on the graph Laplacian  $L_G$ .

Following is the information available at hand in the current graph bandit problem: we are provided with an *n*-armed bandit with an unknown mean vector  $\boldsymbol{\mu}$  smooth on a graph *G*, i.e.  $\langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle \leq \epsilon$ . For this toy problem, we consider the graph *G* to be connected.

Let  $\{\boldsymbol{\nu}_i\}_{i=1}^n$  and  $0 = \lambda_1 < \cdots < \lambda_n$  denote the eigenvectors and eigenvalues of the Laplacian  $L_G$  respectively. It can be easily seen that  $\boldsymbol{\mu} = \sum_{i=1}^n a_i \boldsymbol{\nu}_i$  for some  $a_i \ge 0 \quad \forall i \in [n]$ . The reward function of arm j is

$$\langle \mathbf{e}_j, \boldsymbol{\mu} \rangle = \sum_{i=1}^n a_i \langle \mathbf{e}_j, \boldsymbol{\nu}_i \rangle = a_1 + \sum_{i=2}^n a_i \langle \mathbf{e}_j, \boldsymbol{\nu}_i \rangle$$

the second equality follows from the properties of graph Laplacian we know that  $\nu_1 = \mathbb{1}_n$ , is the only eigenvector associated to 0 eigenvalue in a connected graph.

Without loss of generality we can assume  $a_1 = 0$  as  $a_1$  does not depend on the arm index j. Notice that letting  $a_1 = 0$  is equivalent to having  $\sum_{i=1}^{n} \mu_i = 0$ . Also, the graph constraint can be rewritten as follows:

$$\langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle \leq \epsilon \Rightarrow \sum_{i=1}^n \lambda_i a_i^2 = \langle \boldsymbol{\theta}, \boldsymbol{\theta} \rangle = \|\boldsymbol{\theta}\|_2^2 \leq \epsilon$$

where  $\boldsymbol{\theta} = (\sqrt{\lambda_1}a_1, \dots, \sqrt{\lambda_n}a_n).$ 

Using the above we can cast the graph bandit problem as the linear bandit problem with the mean reward function of arm j expressed as

$$\langle \mathbf{e}_j, \boldsymbol{\mu} \rangle = \sum_{i=2}^n \frac{\theta_i}{\sqrt{\lambda_i}} \langle \mathbf{e}_j, \nu_i \rangle = \langle \mathbf{x}_j, \boldsymbol{\theta} \rangle$$

Hence, the new linear bandit problem is such that the set of arms is  $\{\mathbf{x}_j\}_{j=1}^n$ , the unknown parameter is a vector  $\boldsymbol{\theta}$ , the expected reward of an arm is  $\langle \mathbf{x}_j, \boldsymbol{\theta} \rangle$  and the unknown parameter satisfies the constraint  $\|\boldsymbol{\theta}\|_2^2 \leq \epsilon$ .

We discuss below the drawbacks of casting a graph bandit problem into a linear bandit framework:

• The original best-arm identification is an n-armed problem and the recasted linear bandit problem still has feature vectors with dimensionality n and hence no low-dimensional

benefit of linear bandits is completely lost. Having a performance bound for any algorithm for linear bandits which scales in n, the number of arms gives us no additional advantage.

- The above conversion to linear bandit setup only works when the graph G is connected. Recasting problem setup with disconnected components require assumption of  $\sum_{i \in C} \mu_i = 0$  on individual connected components, which is unrealistic. The results of GRUB holds with or without this assumption.
- Consider the corner case of ε = 0, the linear bandit problem setup derived becomes that of arg max<sub>i</sub>(x<sub>i</sub>, θ) such that ||θ|| ≤ 0 which is only possible if ||θ|| = 0 and in this case we can observe two interesting facts:
  - If the graph G is completely connected then the problem is trivial, since

$$\epsilon = 0 \Rightarrow \langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle = 0 \Rightarrow (\mu_i - \mu_j)^2 = 0 \ \forall i, j \in [n], i \neq j$$

This implies all arms are equal and optimal and the solution is trivial. Here the mean reward function of all arms i is  $\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle = 0$  since  $\theta = 0$  and hence gives the correct output (any arm i).

Suppose graph G has two connected components C<sub>1</sub>, C<sub>2</sub>, where C<sub>k</sub> indicates the arm indices in the connected component k. Further assume that μ<sub>i</sub> = 1 ∀i ∈ C<sub>1</sub>, μ<sub>i</sub> = -1 ∀i ∈ C<sub>2</sub>. Considering the case of ε = 0 here gives us the following :

$$\epsilon = 0 \Rightarrow \langle \boldsymbol{\mu}, L_G \boldsymbol{\mu} \rangle = 0 \Rightarrow (\mu_i - \mu_j)^2 = 0 \quad \forall i \neq j, i, j \in C_k, k = 1, 2$$

Here the mean reward function of all arms *i* is  $\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle = 0$  since  $\theta = 0$  but this is incorrect since not all arms are optimal.

Our graph bandit setup and the performance of GRUB is independent of all of these drawbacks and provides us with a better sample complexity than vanilla best arm identification algorithms.

### I Supporting Results

This appendix is devoted to providing supporting results for many of the theorems and lemmas in the paper.

### I.1 Notation and Definition

Let  $\{t_i(T)\}_{i=1}^n$  (denoted as  $\{t_i\}_{i=1}^n$  for ease of reading) indicate the number of plays of each arm until time T. Let  $X \in \mathbb{R}^{n \times n}$  be a matrix, then  $\{\lambda_i(X)\}_{i=1}^n$  indicate the eigenvalues of matrix X in an increasing order.

Let  $N(\boldsymbol{\pi}_T) = \sum_{t=1}^T \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^T$  be the diagonal counting matrix. Note that  $N(\boldsymbol{\pi}_T)$  can be written as  $N(\{t_i\}_{i=1}^n)$  since the diagonal counting matrix only depends on the number of plays of each arm, rather than the each sampling sequence  $\boldsymbol{\pi}_T$ .

We next establish some properties of the influence function  $\Im$ .

**Lemma I.1.** Let D be an arbitrary graph with n nodes and let  $\{t_i\}_{i=1}^n$  be the number of times all arms are sampled till time T. For each node  $j \in [n]$ , the following are equivalent:

$$\frac{1}{\Im(j,D)} = \max_{\sum_{i \in D_j, i \neq j} t_i = T} \{ [K(i,D)]_{jj} \} \quad (A)$$
$$= \max_{k \in D_j, \sum_{i \in D_j, i \neq j} t_i = T} \{ [V_j(\{t_i\}_{i \in D_j},D)^{-1}]_{jj} - [V_j(\{t_i\}_{i \in D_j},D)^{-1}]_{kk} \} \quad (B)$$

$$= \max_{\sum_{i \in D_j, i \neq j} t_i = T} \left\{ [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{jj} - \min_{k \in D_j} [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{kk} \right\}$$
(C)

$$= \max_{\sum_{i \in D_j, i \neq j} t_i = T} \left\{ [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{jj} - \frac{1}{T} \right\}$$
(D) (72)

where K(i, D) be defined as in Definition 4.2

*Proof.* Let  $f(\cdot, \cdot)$  denote the following:

$$f(i,D) = \max_{\sum_{i \in D_j, i \neq j} t_i = T} \left\{ [K(i,D)]_{jj} \right\}$$

We prove the rest by showing equivalence between (A), (B), (C) and (D).

•  $(A) \Leftrightarrow (D)$ : A simple extension of Lemma I.3 to the case of disconnected clustered graph  $D, \forall \pi_T \in \mathcal{U}(T, D_j)$  we obtain,

$$V_j(\boldsymbol{\pi}_T, D)^{-1} = \frac{1}{T} \mathbb{1}\mathbb{1}^T + K(\pi_1, D)$$
(73)

where  $K(\pi_1, D)$  is as defined in Definition 4.2. Thus, we have the equivalence by explicitly writing the diagonal element of eq (73),

$$[V_j(\boldsymbol{\pi}_T, D)^{-1}]_{jj} - \frac{1}{T} = [K(\pi_1, D)]_{jj}$$
(74)

Hence we have the equivalence as,

$$f(i,D) = \max_{\sum_{i \in D_j, i \neq j} t_i = T} \left\{ [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{jj} - \frac{1}{T} \right\}$$
(75)

•  $(C) \Leftrightarrow (D)$ : Let  $\{t_i^*\}_{i \in D_j}$  denote the following:

$$\{t_i^*(j)\}_{i \in D_j} \in \operatorname*{arg\,max}_{\sum_{i \in D_j, i \neq j} t_i = T} \left\{ [V_j(\{t_i\}_{i \in D_j}, D)]_{jj}^{-1} - \frac{1}{T} \right\}$$
(76)

From Lemma I.2, the optimal  $\{t_i^*(j)\}_{i \in D_j}$  occurs in  $\mathcal{U}(j,T)$ , i.e.  $\exists \{t_i^*(j)\}_{i \in D_j}$  such that  $t_l^*(j) = T$  and  $t_k^*(j) = 0 \quad \forall k \neq l$  for some  $l \in D_j$ . Further by Lemma I.4,

$$\min_{k \in D_j} [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{kk} = \frac{1}{T}$$
(77)

Hence  $\{t_i^*(j)\}_{i \in D_j}$  is also a solution for the following problem:

$$\{t_{i}^{*}(j)\}_{i \in D_{j}} \in \arg \max_{\sum_{i \in D_{j}, i \neq j} t_{i} = T} \left\{ [V_{j}(\{t_{i}\}_{i \in D_{j}}, D)]_{jj}^{-1} - \min_{k \in D_{j}} [V_{j}(\{t_{i}\}_{i \in D_{j}}, D)^{-1}]_{kk} \right\}$$

$$(78)$$

Hence we can conlcude that,

$$f(i, D) = \max_{\sum_{i \in D_j, i \neq j} t_i = T} \left\{ [V_j(\{t_i\}_{i \in D_j}, D)]_{jj}^{-1} - \min_{k \in D_j} [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{kk} \right\}$$
(79)

 (B) ⇔ (C) : Note that max <sub>k∈D<sub>j</sub>,∑<sub>i∈D<sub>j</sub>,i≠j</sub> t<sub>i</sub>=T [V<sub>j</sub>({t<sub>i</sub>}<sub>i∈D<sub>j</sub></sub>, D)<sup>-1</sup>]<sub>jj</sub>] does not depend on arm node index k ∈ D<sub>j</sub>. Hence, the equivalence follows.
</sub>

The resistance distance r(i, j) Definition 4.1 is independent of  $\delta$  for all  $i, j \in [n]$  (The addition of diagonal elements and subtraction of off diagonal elements removes the dependence on  $\delta$  [5]).

Note that  $V_T = N_T + \rho L_G$ , hence  $V_T^{-1}$  gives the psuedo-inverse of the Laplacian matrix for graph G. We show in Lemma I.2 that the matrix R (denoting as  $R(\delta)$  to explicitly show dependence on  $\delta$ ) linked with  $V_T^{-1}$  is independent of number of samples T. Since both matrix R and  $V_T$  are psuedo-inverse of the Laplacian  $L_G$ . Thus we can conclude the following :

$$\lim_{\delta \to 0} [R(\delta)]_{ij} - \frac{1}{\delta} = \lim_{T \to 0} [V(\{t_i\}_{i=1}^n, G)^{-1}]_{ij} - \frac{1}{T}$$
(80)

where  $T \to 0$  implies  $t_i \to 0 \quad \forall i \in [n]$ . Further,

$$\lim_{\delta \to 0} R(\delta)_{ii} + R(\delta)_{jj} - R(\delta)_{ij} - R(\delta)_{ji}$$

$$= \lim_{T \to 0} [V(\{t_i\}_{i=1}^n, G)^{-1}]_{ii} + [V(\{t_i\}_{i=1}^n, G)^{-1}]_{jj}$$

$$- [V(\{t_i\}_{i=1}^n, G)^{-1}]_{ij} - [V(\{t_i\}_{i=1}^n, G)^{-1}]_{ji}$$

$$\to 0 \quad \forall i \in [n].$$
(81)

where  $T \to 0$  implies  $t_i$ 

Since the equation (81) holds for  $t_i \to 0$  for all  $i \in [n]$ , computing the value of limit for one trajectory should suffice for finding the value of the limit. Thereby, we provide an alternate equation for obtaining the resistance distance r(i, j) by

$$r(i,j) = [K(\pi_1 = i, D)]_{ij}$$
(82)

Note that  $[K(\pi_1 = i, D)]_{ii} = [K(\pi_1 = i, D)]_{ij} = [K(\pi_1 = i, D)]_{ij} = [K(\pi_1 = i, D)]_{ji} = 0$  from Lemma I.3). Thus we can say from Definition 4.2, (82)

$$f(i,D) = \frac{1}{\Im(j,D)}$$

Hence proved.

**Lemma I.2.** Let D be a given graph with n nodes. For every node  $j \in D$ , let  $\{t_i^*(j)\}_{i \in D_i}$  denote *the following:* 

$$\{t_i^*(j)\}_{i \in D_j} \in \underset{\sum_{i \in D_j, i \neq j} t_i = T}{\arg\max} \left\{ [V_j(\{t_i\}_{i \in D_j}, D)]_{jj}^{-1} - \frac{1}{T} \right\}$$
(83)

Then  $\exists \{t_i^*(j)\}_{i \in D_i}$ ,  $l \in D_j$  such that  $t_l^*(j) = T$  and  $t_k^*(j) = 0 \quad \forall k \neq l$ .

*Proof.* To simplify our proof, let graph D be connected. The proof for the case of disconnected components is an extension of the connected graph case, by analysing each individual connected component together.

If graph D is connected then  $D_i = D$ . For the rest of the proof we sometimes denote  $V(\boldsymbol{\pi}_T, D)$  as  $V(\{t_i\}_{i=1}^n, D)$  to make it more context relevant.

Let  $q: \mathbb{R}^n \to \mathbb{R}^{n \times n}$  be a partial function of  $V(\boldsymbol{\pi}_T, D)$  as follows:

$$g(\{t_i\}_{i=1}^n) = V(\{t_i\}_{i=1}^n, D)$$
(84)

For all  $i \in [n]$ , let  $t_i = \alpha_i T$  such that  $\sum_{i=1}^n \alpha_i = 1$ . Then we can say that,  $g(\{t_i\}_{i=1}^n) = g(\{\alpha_i T\}_{i=1}^n)$ 

$$= \sum_{i=1}^{n} \alpha_i g(\{0, 0, \dots, t_i = T, \dots, 0\})$$
(85)

Using convexity of matrix invertibility [42]  $V(\boldsymbol{\pi}_T, G)^{-1}$  satisfies,

$$g(\{t_i\}_{i=1}^n)^{-1} \preceq \sum_{i=1}^n \alpha_i g(\{0, 0, \dots, t_i = T, \dots, 0\})^{-1}$$
(86)

Hence  $g(\cdot)^{-1}$  is a convex function. Since we have the restriction as  $\sum_{i=1, i\neq j}^{n} t_i = T$ . We can say that,

$$\arg \max_{\sum_{i \in D_{j}, i \neq j} t_{i} = T} \left\{ [V(\{t_{i}\}_{i=1}^{n}, D)]_{jj}^{-1} - \frac{1}{\sum_{i=1}^{n} t_{i}} \right\}$$

$$= \arg \max_{\sum_{i \in D_{j}, i \neq j} t_{i} = T} [V(\{t_{i}\}_{i=1}^{n}, D)]_{jj}^{-1}$$

$$= \arg \max_{\sum_{i \in D_{j}, i \neq j} t_{i} = T} \langle \mathbf{e}_{j}, [V(\{t_{i}\}_{i=1}^{n}, D)]^{-1} \mathbf{e}_{j} \rangle$$

$$= \arg \max_{\sum_{i \in D_{j}, i \neq j} t_{i} = T} \langle \mathbf{e}_{j}, g(\{t_{i}\}_{i=1}^{n})^{-1} \mathbf{e}_{j} \rangle$$
(87)

Since  $g(\cdot)^{-1}$  is convex, for a convex function the maximization over a simplex happens at one of the vertices. Hence the max happens when  $t_i = T$  and  $t_k = 0 \quad \forall k \neq i$ .

Hence proved.

**Lemma I.3.** Let G be a given connected graph of n nodes and  $t_i$  be the number of samples of each arm i. Then  $\forall \pi_T \in \mathcal{U}(T)$ ,

$$V(\boldsymbol{\pi}_T, G)^{-1} = \frac{1}{T} \mathbb{1} \mathbb{1}^T + K(\pi_1, G)$$
(88)

where,  $\mathbb{1} \in \mathbb{R}^n$  is a vector or all ones and  $K(\pi_1, G) \in \mathbb{R}^{n \times n}$  is the matrix defined in Definition 4.2.

*Proof.* Let I be an identity matrix of dimension  $n \times n$ . We prove the result by showing that,  $\forall \pi_T \in \mathcal{U}(T), V(\pi_T, G)^{-1}V(\pi_T, G) = I$ ,

$$V(\boldsymbol{\pi}_{T}, G)^{-1}V(\boldsymbol{\pi}_{T}, G)$$

$$= \left(\frac{1}{T}\mathbb{1}\mathbb{1}^{T} + K(\pi_{1}, G)\right) \left(\sum_{t=1}^{T} \mathbf{e}_{\pi_{t}} \mathbf{e}_{\pi_{t}}^{T} + \rho L_{G}\right)$$

$$= \left(\frac{1}{T}\mathbb{1}\mathbb{1}^{T} + K(\pi_{1}, G)\right) \left(T\mathbf{e}_{\pi_{1}}\mathbf{e}_{\pi_{1}}^{T} + \rho L_{G}\right)$$

$$= \mathbb{1}\mathbf{e}_{\pi_{1}}^{T} + TK(\pi_{1}, G)\mathbf{e}_{\pi_{1}}\mathbf{e}_{\pi_{1}}^{T} + \rho K(\pi_{1}, G)L_{G}$$
(89)

From Definition 4.2,  $K(\pi_1, G)\mathbf{e}_{\pi_1}\mathbf{e}_{\pi_1}^T = 0$  and  $\mathbb{1}\mathbf{e}_{\pi_1}^T + \rho K(\pi_1, G)L_G = I$  implying that  $V(\boldsymbol{\pi}_T, G)^{-1}V(\boldsymbol{\pi}_T, G) = I$ .

Hence proved.

**Lemma I.4.** Let G be any connected graph and  $\pi_T \in \mathcal{U}(T, G)$ . Then,

$$\min_{j \in [n]} \{ [V(\boldsymbol{\pi}_T, G)^{-1}]_{jj} \} = \frac{1}{T}$$
(90)

*Proof.* From Definition 4.2,  $K(\pi_1, G)$  satisfies

$$K(\pi_1, G)L_G = \frac{1}{\rho} \left( I - \mathbb{1}\mathbf{e}_{\pi_1}^T \right)$$

Observe that  $\mathbb{1}\mathbf{e}_i^T$  is a rank 1 matrix with eigenvalue 1 and eigenvector  $\mathbf{e}_i$  and Identity matrix I is of rank n with all eigenvalues 1 and eigenvectors  $\{\mathbf{e}_i\}_{i=1}^n$ . Hence  $(I - \mathbb{1}\mathbf{e}_{\pi_1}^T)$  is a rank n - 1 matrix with rest nonzero eigenvalues as 1. Since the graph G is connected,  $\lambda_1(L_G) = 0$  and  $\lambda_2(L_G) > 0$ . The eigenvector corresponding to  $\lambda_1(L_G)$  is  $\mathbb{1}$ , the all 1 vector.

Given  $\rho > 0$ , we can conclude,

$$K(\pi_1, G)L_G \succeq 0$$
 s.t.  $\operatorname{rank}(K(\pi_1, G)L_G) = n - 1$  (91)

Hence, in order to satisfy eq. (91),  $K(\pi_1, G) \succeq 0$  and  $\operatorname{rank}(K(\pi_1, G)) \ge n - 1$ . By lower bounds on Rayleigh quotient we can conclude,

$$\langle \mathbf{e}_j, K(\pi_1, G) \mathbf{e}_j \rangle = [K(\pi_1, G)]_{jj} \ge 0 \quad \forall j \in [n]$$
(92)

From Lemma I.3,  $[K(\pi_1, G)]_{jj} = [V(\boldsymbol{\pi}_T, G)^{-1}]_{jj} - \frac{1}{T}$  implying that  $[V(\boldsymbol{\pi}_T, G)^{-1}]_{jj} \ge \frac{1}{T}$ . From Definition 4.2 it can be seen that  $[K(\pi_1, G)]_{\pi_1\pi_1} = 0$  and hence  $[V(\boldsymbol{\pi}_T, G)^{-1}]_{\pi_1\pi_1} = \frac{1}{T}$  which concludes the proof.

**Lemma I.5.** Given a connected graph G, the following bound holds for all the diagonal entries of  $[V(\pi_T, G)^{-1}]_{ii}$  for  $i \in [n]$ :

$$[V(\boldsymbol{\pi}_T, G)^{-1}]_{ii} \le \mathbb{1} \ (t_i = 0) \left(\frac{1}{\rho \Im(i, \mathcal{G})} + \frac{1}{T}\right) + \mathbb{1} \ (t_i > 0) \max\left\{\frac{1}{t_i + \frac{\rho \Im(i, G)}{2}}, \frac{1}{t_i + \frac{T}{2}}\right\}$$
(93)

*Proof.* From Definition 4.2 of  $\mathfrak{I}(\cdot, \mathcal{G})$  and Lemma I.1, Breaking the lemma statement into cases:

• Unsampled Arms : From Lemma I.1

$$\frac{1}{\Im(j,G)} = \max_{\sum_{i \in G_j, i \neq j} t_i = T} \left\{ [V_j(\{t_i\}_{i \in G_j}, G)^{-1}]_{jj} - \frac{1}{T} \right\} \quad \forall j \in [n]$$
(94)

Thus for any unsampled arm j,

$$[V(\boldsymbol{\pi}_T, G)]_{jj}^{-1} \le \left(\frac{1}{\mathfrak{I}(j, G)} + \frac{1}{T}\right)$$
(95)

• Sampled Arms : Since the matrix  $V(\pi_T, G)$  depends only on the final sampling distribution  $\{t_i\}_{i=1}^n$  rather than the sampling path  $\pi_T$ . Consider a sampling path such that  $\pi_t \neq j$  for  $t \leq T - t_j$  and  $\pi_t = j$  for  $T - t_j \leq t \leq T$ .

Assuming such a sampling path  $\pi_T$ , after  $\pi_{T-t_i}$  samples,

$$[V(\boldsymbol{\pi}_{T-t_j}, G)^{-1}]_{jj} \le \frac{1}{T} + \frac{1}{\Im(j, G)}$$
(96)

Then by the Sherman-Morrison rank 1 update identity<sup>9</sup>,

$$\frac{1}{[V(\boldsymbol{\pi}_T, G)^{-1}]_{jj}} = \frac{1}{[V(\boldsymbol{\pi}_{T-t_j}, G)^{-1}]_{jj}} + t_j$$
$$[V(\boldsymbol{\pi}_T, G)^{-1}]_{jj} = \frac{1}{t_j + \frac{1}{[V(\boldsymbol{\pi}_{T-t_j}, G)^{-1}]_{jj}}}$$
$$\leq \frac{1}{t_j + \frac{1}{\left(\frac{1}{\mathcal{I}(j, G)} + \frac{1}{T-t_j}\right)}}$$

Hence we have the bound on  $[V(\boldsymbol{\pi}_T, G)^{-1}]_{jj}$  as follows:

$$[V(\boldsymbol{\pi}_T, G)^{-1}]_{jj} \le \max\left\{\frac{1}{t_j + \frac{\Im(j, G)}{2}}, \frac{1}{t_j + \frac{T - t_j}{2}}\right\}$$
(97)

Hence proved.

**Lemma I.6.** Let D be a graph with n nodes and k disconnected components. If each of the connected components  $\{C_i(D)\}_{i=1}^k$  is a complete graph then  $\forall j \in [n]$ ,

$$\Im(j,D) = \frac{|\mathcal{C}(j,D)|}{2} \tag{98}$$

*Proof.* Let D be a complete graph (k = 1),  $\pi_T \in \mathcal{U}(T)$  and  $\rho = 1$ . Then,

$$V(\boldsymbol{\pi}_T, G)^{-1} = \frac{1}{T} \mathbb{1} \mathbb{1}^T + K$$
(99)

where  $\mathbb{1} \in \mathbb{R}^n$  is a vector or all ones and  $K \in \mathbb{R}^{n \times n}$  is a matrix given by,

$$K_{\pi_1\pi_1} = 0, \quad K_{jj} = \frac{2}{n} \quad \forall j \in [n] / \{\pi_1\}$$
  
$$K_{k\pi_1} = 0, \quad K_{\pi_1 j} = 0, \quad K_{jk} = \frac{1}{n} \quad \forall j, k \in [n] / \{\pi_1\}, \ j \neq k$$

The form of  $V(\boldsymbol{\pi}_T, G)^{-1}$  in eq.(99) can be verified by  $V(\boldsymbol{\pi}_T, G)^{-1}V(\boldsymbol{\pi}_T, G) = I$ .

The final statement of the lemma can be obtained by considering this analysis to just the nodes within a connected component of a diconnected graph G and Lemma I.1.

<sup>&</sup>lt;sup>9</sup>Hager, W. (1989). Updating the Inverse of a Matrix. SIAM Rev., 31, 221-239.

**Lemma I.7.** Let D be a graph with n nodes and k disconnected components. If each of the connected components  $\{C_i(D)\}_{i=1}^k$  is a line graph then  $\forall j \in [n]$ ,

$$\Im(j,D) > \frac{1}{|\mathcal{C}(j,D)|}$$
(100)

*Proof.* Let D be a complete graph (k = 1),  $\pi_T \in \mathcal{U}(T)$  and  $\rho = 1$ . Then,

$$V(\boldsymbol{\pi}_T, G)^{-1} = \frac{1}{T} \mathbb{1} \mathbb{1}^T + K$$
(101)

where  $\mathbbm{1}\in \mathbb{R}^n$  is a vector or all ones and  $K\in \mathbb{R}^{n\times n}$  is a matrix given by,

$$K_{\pi_1\pi_1} = 0, \quad K_{jj} = d(\pi_1, j) \quad \forall j \in [n] / \{\pi_1\}, K_{k\pi_1} = 0, \quad K_{\pi_1 j} = 0, K_{jk} = \min\{d(\pi_1, j), d(\pi_1, k)\} \quad \forall j, k \in [n] / \{\pi_1\}, \ j \neq k$$

The form of  $V(\boldsymbol{\pi}_T, G)^{-1}$  in eq.(101) can be verified by  $V(\boldsymbol{\pi}_T, G)^{-1}V(\boldsymbol{\pi}_T, G) = I$ .

The final statement of the lemma can be obtained by considering this analysis to just the nodes within a connected component of a diconnected graph G and Lemma I.1.

**Lemma I.8.** Let A = ([n], E) be any graph and let  $e \in E$  be an edge of graph A. Let  $B = ([n], E - \{e\})$  be a subgraph of A with one edge removed. Then the following holds for all non-isolated nodes i in B:

• If 
$$|\mathcal{C}(A)| = |\mathcal{C}(B)|$$
,  
 $\Im(i, A) \ge \Im(i, B)$   
• If  $|\mathcal{C}(A)| < |\mathcal{C}(B)|$ ,  
 $\Im(i, A) \le \Im(i, B)$ 

*Proof.* From Lemma I.1, for any graph  $D, \mathfrak{I}(\cdot, \cdot)$  satisfies,

$$\frac{1}{\Im(j,D)} = \max_{k \in D_j, \sum_{i \in D_j} t_i = T} \left\{ [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{jj} - [V_j(\{t_i\}_{i \in D_j}, D)^{-1}]_{kk} \right\} \quad \forall j \in [n]$$
(102)

Case I :  $|\mathcal{C}(A)| = |\mathcal{C}(B)|$ 

The edge set of B is smaller than edge set of A. Hence, from Lemma  $\Im(i, A) \ge \Im(i, B)$ 

**Case II**: C(A) < C(B) In this case,  $|B_i| \le |A_i|$ . Hence the max is over a smaller set of options, we can conclude that  $\Im(i, A) \le \Im(i, B)$ . Hence proved.

Given a graph D, we define a class of sampling policies  $\mathcal{U}(T, D)$  as follows,

**Definition I.9.** Let  $\mathcal{U}(T, D)$  denote the set of sampling policies,

 $\mathcal{U}(T,D) = \{ \boldsymbol{\pi}_T | \exists l \in D \text{ s.t. } \boldsymbol{\pi}_t = l \ \forall t \leq T \}$ 

**Lemma I.10.** Let G be the given graph and sampling policy  $\pi_T$  has been played for T time steps, then  $V_T$  satisfies the following structure,

$$V(\pi_T, D) = diag([V_1, V_2, \dots, V_{k(G)}])$$
(103)

where  $V_i$  depends on the connected component  $C_i \in C_D$  of the graph and the number of samples of the arms within the connected component  $\{t_j\}_{j \in C_i}$ .

*Proof.* Rewriting the definition of  $V(\boldsymbol{\pi}_T, D)$ ,

$$V(\boldsymbol{\pi}_T, D) \triangleq \sum_{t=1}^T \mathbf{e}_{\pi_t} \mathbf{e}_{\pi_t}^\top + \rho L_D$$
  
=  $N(\{t_i\}_{i=1}^n) + L_D$  (104)

Both component matrices  $N(\{t_i\}_{i=1}^n)$  (diagonal matrix) and  $L_D$  (Laplacian matrix of a graph) adhere to a block diagonal structure and hence  $V(\boldsymbol{\pi}_T, D)$  matrix also adheres to a block diagonal structure analogous to  $L_D$ . The block diagonal structure in  $L_D$  is dictated by connected components of graph D.

The following lemma establishes the invertibility of  $V(\pi_T, G)$  for a connected graph and T > 1: Lemma I.11. For a connected graph G,  $V(\pi_1, G)$  is invertible, but  $V(\pi_0, G)$  is not invertible.

*Proof.* Since the graph G is connected,  $\lambda_1(L_G) = 0$  and  $\lambda_2(L_G) > 0$ . The eigenvector corresponding to  $\lambda_1(L_G)$  is 1, the all 1 vector. At time T = 0,  $V(\boldsymbol{\pi}_T, G) = L_G$  and hence  $V(\boldsymbol{\pi}_T, G)$  is positive semi-definite matrix with one zero eigenvalues.

Let arm *i* be pulled at T = 1, i.e.  $\pi_1 = i$ , then the corresponding counting matrix is a positive semi definite matrix of rank one with the eigen value  $\lambda_n(N) = 1$  for the eigenvector  $e_i$ .

Observe that  $\mathbf{e}_i^T \mathbb{1} > 0$ . Also,  $N_T$  and  $L_G$  are positive semi-definite matrices with ranks 1 and n-1 respectively. The subspace without information (corresponding to the direction of zero eigenvalue) for matrix  $L_G$  is now provided by  $N(\boldsymbol{\pi}_1)$  and hence  $\lambda_{\min}(V(\boldsymbol{\pi}_1, G)) > 0$  making it invertible.  $\Box$ 

**Lemma I.12.** Let  $G = ([n], E_G, A), H = ([n], E_H, A)$  are two graphs with n nodes such that  $E_G \supseteq E_H$ . Then, assuming invertibility of  $[V(G, T)^{-1}]$  and  $[V(H, T)^{-1}]$ ,

$$[t_{eff,i}]_G \ge [t_{eff,i}]_H \quad \forall i \in [n], T > k(G)$$

$$(105)$$

where  $\forall i \in [n], [t_{eff,i}]_G, [t_{eff,i}]_H$  indicates the effective samples with graph G and H respectively.

*Proof.* Given graphs  $G = ([n], E_G), H = ([n], E_H)$  satisfy  $E_G \supseteq E_H$ .

The quadratic form of Laplacian for the graph G, H is given by,

$$\mathbf{x}L_G \mathbf{x} = \sum_{(i,j)\in E_G} (x_i - x_j)^2$$
$$\mathbf{x}L_H \mathbf{x} = \sum_{(i,j)\in E_H} (x_i - x_j)^2$$

Since  $E_G \supseteq E_H$ ,

$$\mathbf{x} L_G \mathbf{x} \ge \mathbf{x} L_H \mathbf{x} \quad \forall \ \mathbf{x} \in \mathbb{R}^n$$
$$\Rightarrow L_G \succeq L_H$$

Further, provided a sampling policy  $\pi_T$ , we can say that,

$$V(\boldsymbol{\pi}_T, G) \succeq V(\boldsymbol{\pi}_T, H)$$

For the number of samples T sufficient to ensure invertibility of  $V(\boldsymbol{\pi}_T, H)$ , we have

$$V(\boldsymbol{\pi}_T, G)^{-1} \leq V(\boldsymbol{\pi}_T, H)^{-1}$$
  

$$\mathbf{x}^T V(\boldsymbol{\pi}_T, G)^{-1} \mathbf{x} \leq \mathbf{x}^T V(\boldsymbol{\pi}_T H)^{-1} \mathbf{x} \quad \forall \mathbf{x} \in \mathbb{R}^n$$
  

$$[V(\boldsymbol{\pi}_T, G)^{-1}]_{ii} \leq [V(\boldsymbol{\pi}_T, H)^{-1}]_{ii} \quad (\text{taking } \mathbf{x} = \mathbf{e}_i)$$
  

$$\frac{1}{[V(\boldsymbol{\pi}_T, G)^{-1}]_{ii}} \geq \frac{1}{[V(\boldsymbol{\pi}_T, H)^{-1}]_{ii}}$$

Hence from the definition of effective samples 3.1, it is clear that for any  $i \in [n]$ ,

$$[t_{\text{eff},i}]_G \ge [t_{\text{eff},i}]_H \tag{106}$$

Hence proved.

**Lemma I.13.** Let effective samples  $t_{eff,i}$  be as is defined in Definition 3.1 and let  $\pi_T$  denote a cyclic sampling policy for T > k(G), then the infinite sum  $\sum_{T=k(G)+1}^{\infty} t_{eff,i}^{-2}$  is bounded. In fact,

$$\sum_{T=k(G)+1}^{\infty} t_{e\!f\!f,i}^{-2} < n \left(\frac{2(n-1)}{\rho}\right)^2 + \frac{n\pi^2}{6}$$
(107)

*Proof.* We first prove the lemma statement for connected graph G and then go towards a more general graph G. From Lemma D.1,

$$t_{\text{eff},i} \ge t_i + \min\{\rho \mathfrak{I}(i,G), T - t_i\}$$

if  $T - t_i \leq \rho \mathfrak{I}(i, G)$ , then  $t_{\text{eff},i} \geq \frac{T + t_i}{2} \geq \frac{T}{2}$ . For the reverse case of  $T - t_i \geq \rho \mathfrak{I}(i, G)$ ,  $t_{\text{eff},i} \geq t_i + \frac{\rho \mathfrak{I}(i, G)}{2} \geq t_i + \frac{\rho}{2(n-1)}$  (since  $\mathfrak{I}(i, G) \geq \frac{1}{n-1}$  by Remark ??).

Since  $\pi_T$  is a cyclic sampling policy, hence  $t_i$  increases by 1 at-least once every *n* samples. Thus, we can upperbound the infinite sum as,

$$\sum_{T=1}^{\infty} \frac{1}{t_{\text{eff},i}^2} \leq \sum_{T=1}^{\infty} \frac{1}{\left(t_i + \frac{\rho}{2(n-1)}\right)^2} \leq n \left(\frac{2(n-1)}{\rho}\right)^2 + n \sum_{t_i=1}^{\infty} \frac{1}{t_i^2} < n \left(\frac{2(n-1)}{\rho}\right)^2 + \frac{n\pi^2}{6}$$
(108)

Hence proved.

# J Better sampling strategies

Theorem 4.4 established a baseline w.r.t. sampling protocol by solving  $T_{\text{sufficient}}$  for naive cyclic sampling policy (a sampling policy which doesn't exploit the graph properties). Note that, even if the sampling policy doesn't utilize any graph properties, the similarity graph is still being utilized in computing the mean estimate and the confidence widths. For the safe elimination of suboptimal arms, the ultimate goal of GRUB is to shrink the confidence bounds  $\beta_i \sqrt{(t_{\text{eff},i})^{-1}}$  as quickly as possible. Accordingly, a few intelligent sampling policies that exploit the graph structure of the problem is given as follows:

- Marginal variance minimization (MVM): Since picking any arm impacts the confidence widths of all arms in it's connected component, we pick the arm with the maximum variance. Specifically,  $l = \underset{i \in A}{\operatorname{arg min}} t_{\text{eff},i} = \underset{i \in A}{\operatorname{arg max}} [V_T^{-1}]_{ii}$ , where A is the set of indices of the arms under consideration.
- Joint variance minimization nuclear (JVM-N): This variant is inspired from the concept of V-optimality [22]. This policy aims to select the arm that minimizes ℓ<sub>2</sub> regression loss of the estimated vector µ̂, i.e. the confidence interval across all remaining arms in A. Specifically, l = arg min<sub>i∈A</sub> ||(V<sub>T</sub> + e<sub>i</sub>e<sup>T</sup><sub>i</sub>)<sup>-1</sup>||<sub>\*</sub>
- Joint variance minimization operator (JVM-O). Taking inspiration from  $\Sigma$ optimality [40, 38] the next policy can be stated as,  $l = \underset{i \in A}{\operatorname{arg\,min}} ||(V_T + \mathbf{e}_i \mathbf{e}_i^T)^{-1}||_{\operatorname{op}} =$

$$\underset{i \in A}{\arg \max} \frac{\|\operatorname{Row}_{i}(V_{T}^{-1})\|_{2}^{2}}{1 + [(V_{T}^{-1})_{ii}]}$$

The main objective of sampling policies is to *decrease* the value of  $[V_T^{-1}]_{ii}$  for every arm *i* as fast as possible. The notion of *decrease* leads to different sampling policies for GRUB. The algorithm chooses the arm which maximizes this notion of decreases.

The objective of the sampling policy **Joint variance minimization – operator (JVM-O)** is equivalent to:

$$\max_{k \in A} \sum_{j \in [n]} |(V_T^{-1})_{k,j}| - |(V_T + \mathbf{e}_k \mathbf{e}_k^T)_{k,j}^{-1}|$$
(109)

Using Sherman-morrison rank 1 update we split the summation into different cases:



Figure 2: (Best seen in color) Performance of GRUB with using various sampling protocols for SBM ((p,q) = (0.9, 5e - 3)) [Left] and BA (m = 2) [Right]. The UCB method without graph information is significantly slower compared to the graph-based variants. Note that for these toy datasets, the sampling algorithm used does not alter the results too much.

• For j = k,

$$|\langle \mathbf{e}_{k} V_{T}^{-1} \mathbf{e}_{k} \rangle| - |\langle \mathbf{e}_{k} \left( V_{T} + \mathbf{e}_{k} \mathbf{e}_{k}^{T} \right)^{-1} \mathbf{e}_{k} \rangle| = \frac{\|\mathbf{e}_{k}\|_{V_{T}^{-1}}^{4}}{1 + \|\mathbf{e}_{k}\|_{V_{T}^{-1}}^{2}}$$
(110)

• For all connected-nodes of  $j \in \mathcal{N}_k$ ,

$$|\langle \mathbf{e}_{k} V_{T}^{-1} \mathbf{e}_{j} \rangle| - |\langle \mathbf{e}_{k} \left( V_{T} + \mathbf{e}_{k} \mathbf{e}_{k}^{T} \right)^{-1} \mathbf{e}_{j} \rangle| = \frac{\langle \mathbf{e}_{j}, \mathbf{e}_{k} \rangle_{V_{T}^{-1}}^{2}}{1 + \|\mathbf{e}_{k}\|_{V_{T}^{-1}}^{2}}$$
(111)

• For all other non-connected  $j \notin \mathcal{N}_k, i \neq k$ ,

$$\left|\left\langle \mathbf{e}_{k}V_{T}^{-1}\mathbf{e}_{i}\right\rangle\right|-\left|\left\langle \mathbf{e}_{k}\left(V_{T}+\mathbf{e}_{k}\mathbf{e}_{k}^{T}\right)^{-1}\mathbf{e}_{i}\right\rangle\right|=0$$
(112)

Hence the sampling policy decides on the arm to sample based on the following optimization problem,

$$\sum_{j \in [n]} |(V_T^{-1})_{k,j}| - |(V_T + \mathbf{e}_k \mathbf{e}_k^T)_{k,j}^{-1}| = \frac{\|\mathbf{e}_k\|_{V_T^{-1}}^4 + \sum_{j \in \mathcal{N}_k} \langle \mathbf{e}_i, \mathbf{e}_k \rangle_{V_T^{-1}}^2}{1 + \|\mathbf{e}_k\|_{V_T^{-1}}^2}$$
$$= \frac{\|\operatorname{Row}_k(V_T^{-1})\|_2^2}{1 + [(V_T^{-1})_{kk}]}$$

Hence we try to find the arm k within the remaining arms in consideration which maximizes  $\frac{\|\text{Row}_k(V_T^{-1})\|_2^2}{1+[(V_T^{-1})_{kk}]}.$ 

# **K** Additional Experiments

Synthetic Data : Setup 2 We consider an *n*-armed bandit setup with the aim of finding the best arm. The number of arms scale from n = 50 to 200 in steps of 50. We consider 2 cases: G is a Stochastic Block model(SBM) with parameters  $(p,q) = (0.9, 1e^{-4})$  and G is a Barabási–Albert(BA) graph with parameter m = 2, both containing 10 clusters. We run every setup for 20 runs and record the stopping time for all runs.

As can be seen in 1, all graph algorithm (GRUB with different sampling policies) the standard UCB based best-arm identification algorithm. Within the different GRUB, different sampling policies exploit the graph infromation in distinct ways, leading to a different in their performance. GRUB (cyclic sampling based) is outperformed by all other sampling based GRUB methods.

Synthetic Data: Setup 3 We consider the setup where G with n = 100 arms consists of 10 connected components with 10 arms per cluster. We consider 2 cases: G is a Stochastic Block model(SBM)



Figure 3: (Best seen in color) Performance of GRUB with using various sampling protocols for SBM ((p,q) = (0.9, 5e - 3)) [Left] and BA (m = 2) [Right]. The UCB method without graph information is significantly slower compared to the graph-based variants. Note that for these toy datasets, the sampling algorithm used does not alter the results too much.



Figure 4: (Best seen in color) Performance of GRUB using different sampling protocols for Github social graph (left) and LastFM graph (right). With no graph information, UCB requires orders of magnitude more samples compared to policies that use explicitly graph information. The cycic sampling policy is not as competitive on real world datasets

with parameters (p,q) = (0.9, 0.0005) and G is a Barabási–Albert(BA) graph with parameter m = 2. Results are provided in Figure 3

**Real Data:** We use graphs from SNAP [34] for the experiments involving real world graphs. We sub-sample the graphs using Breadth-First Search (to retain connected components) to generate the graphs for our experiments. We use the LastFM [46], subsampled to 229 nodes and Github Social [45] subsampled to 242 nodes.

In all the experiments, it is evident that GRUB with any of the sampling policies outperform UCB algorithm [32], which does not leverage the graph. Further within the various sampling policies, MVM sampling policy seems to outperform other sampling policies (Figure 4). For both Github and LastFM datasets, the MVM policy obtains the best arm in  $\sim 300$  rounds compared to traditional UCB that takes  $\sim 4500$  rounds. A rigorous theoretical characterization of the above sampling policies is an exciting avenue for future research. We refer the reader to Appendix A for a discussion on the results of the paper, potential extensions, and broader impacts.

## L Code Availability

The full code used for conducting experiments can be found at the following Github repository.