

# Densest Diverse Subgraphs: How to Plan a Successful Cocktail Party with Diversity

Atsushi Miyauchi  
atsushi.miyauchi@centai.eu  
CENTAI Institute  
Turin, Italy

Tianyi Chen  
ctony@bu.edu  
Boston University  
Boston, MA, USA

Konstantinos Sotiropoulos  
ksotirop@bu.edu  
Boston University  
Boston, MA, USA

Charalampos E. Tsourakakis  
ctsourak@bu.edu  
Boston University  
Boston, MA, USA

## ABSTRACT

Dense subgraph discovery methods are routinely used in a variety of applications including the identification of a team of skilled individuals for collaboration from a social network. However, when the network’s node set is associated with a sensitive attribute such as race, gender, religion, or political opinion, the lack of diversity can lead to lawsuits.

In this work, we focus on the problem of finding a densest diverse subgraph in a graph whose nodes have different attribute values/types that we refer to as colors. We propose two novel formulations motivated by different realistic scenarios. Our first formulation, called the *densest diverse subgraph problem* (DDSP), guarantees that no color represents more than some fraction of the nodes in the output subgraph, which generalizes the state-of-the-art due to Anagnostopoulos et al. (CIKM 2020). By varying the fraction we can range the diversity constraint and interpolate from a diverse dense subgraph where all colors have to be equally represented to an unconstrained dense subgraph. We design a scalable  $\Omega(1/\sqrt{n})$ -approximation algorithm, where  $n$  is the number of nodes. Our second formulation is motivated by the setting where any specified color should not be overlooked. We propose the *densest at-least- $\vec{k}$ -subgraph problem* (DalkS), a novel generalization of the classic DalkS, where instead of a single value  $k$ , we have a vector  $\mathbf{k}$  of cardinality demands with one coordinate per color class. We design a  $1/3$ -approximation algorithm using linear programming together with an acceleration technique. Computational experiments using synthetic and real-world datasets demonstrate that our proposed algorithms are effective in extracting dense diverse clusters.

## CCS CONCEPTS

• **Theory of computation** → **Graph algorithms analysis**; **Approximation algorithms analysis**.

## KEYWORDS

social network analysis, densest subgraph problem, diversity, fairness, approximation algorithms

## 1 INTRODUCTION

Dense subgraph discovery (DSD) is a fundamental graph-mining primitive, routinely used to mine social, financial, and biological

networks among others [22]. Applications include team formation [19, 40], detecting correlated genes [47], community and spam link farm detection in the Web graph [18, 21], finding experts in crowdsourcing systems [26], spotting money laundering in financial networks [15, 33, 44], assessing the statistical significance of motifs [14], and modeling real-world networks [11]. See the tutorial [22] for an extensive list of related applications.

Among various DSD formulations, the densest subgraph problem (DSP) stands out for various reasons [32]. For a given undirected graph  $G = (V, E)$  with  $n = |V|$  nodes and  $m = |E|$  edges, DSP is solvable in polynomial time using maximum flow [23, 38] or linear programming (LP) [12] and can be approximated within a factor of 2 using a greedy algorithm [12, 30], and more recently solved near-optimally by an iterative greedy algorithm over the input [9, 13]. Furthermore, DSP is solvable on massive graphs using distributed and streaming implementations [7] and admits useful variants, e.g., [27, 31, 46, 48].

In numerous real-world settings, we have extra information about the nodes. For example, in a graph database we may know each individual’s gender and race. On the Twitter follow network, there exist methods to infer from tweets whether a node is positive, neutral, or negative towards a controversial topic [17]. On brain networks neurons can play different functional roles [1, 34]. Consider also the problem of organizing a cocktail party [43] with diversity constraints as follows:

A number of computer scientists aim to organize a cocktail party to celebrate Turing’s legacy. They believe that the success of the event will be higher if they invite computer scientists who have collaborated in the past but also who span different research areas. Whom should they invite?

We will refer to the set of different attribute values/types of nodes as *colors*. The aforementioned settings motivate the problem of finding a *densest diverse subgraph*, namely a subset of nodes that induce many edges, but also are diverse in terms of colors. We will be referring to the concept of diversity as *fairness* when the attribute concerns sensitive information such as gender, race, or religion. Applying a DSD method does not guarantee that the extracted densest subgraph will be diverse. Actually on a variety of real data, we observe that this is the typical case, i.e., the densest subgraph often exhibits strong homophily. Suppose the output of such a method is used to select a group of individuals in a social

network. In that case, it will not be representative of the different races/religions/opinions that may co-exist in the network. This can be especially harmful in the context of selecting teams using dense subgraphs [40], recommending material to social media users that is not balanced in terms of opinions and hence increasing polarization [35] or even leading to lawsuits [24]. In such cases, it becomes of paramount importance to have algorithms that can extract a cluster with diversity guarantees.

With the exception of a recent paper by Anagnostopoulos et al. [4], very little attention has been given to dense diverse subgraph extraction, despite the extensive research on the DSP and its applications [32]. Although Anagnostopoulos et al. [4] have made progress, there are still many unanswered questions. For example, the methods they propose exclusively focus on scenarios involving two colors and strive to achieve the complete fairness in the output. Furthermore, their spectral approach offers theoretical guarantees, but these guarantees are contingent upon restrictive conditions for the spectrum of the input graph that are computationally burdensome to verify. Additionally, they offer heuristics for scenarios involving more than two colors, but without any assurances regarding their quality.

## 1.1 Our Contributions

We introduce two novel formulations for finding a densest diverse subgraph. The first one is called the *densest diverse subgraph problem* (DDSP), and the second one is called the *densest at-least- $k$ -subgraph problem* (DalkS). Informally, the first problem aims to offer diversity guarantees that concern the relative sizes of the color classes, while the second guarantees in terms of absolute counts.

Let  $C$  be the set of colors. Our first formulation guarantees the diversity of the output in the sense that it is not dominated by any single color. To this end, the formulation introduces a parameter  $\alpha \in [1/|C|, 1]$  that determines the maximum portion of any color in the output solution. It should be noted that our formulation is a substantial generalization of the *fair densest subgraph problem* introduced by Anagnostopoulos et al. [4], enabling us to deal with non-binary attributes and to freely adjust the degree of diversity. Interestingly, from a theory perspective, our formulation contains two important variants of DSP, the densest at-least- $k$ -subgraph problem (DalkS) and the densest at-most- $k$ -subgraph problem (DamkS), as special cases, cf. Sections 2 and 3.1 for details.

On the other hand, our second formulation guarantees the diversity of the output in the sense that it does not overlook any specified color. In particular, consider a graph where  $|C|$  is a very small constant and the minority colors appear only in a handful of nodes. Instead of imposing *relative* constraints on the sizes through ratios, we impose *absolute* constraints on the cardinalities of the nodes from each color. Specifically, the formulation requires the output to contain at least a given number of representatives from each color. The formulation is a novel generalization of DalkS, where instead of just demanding  $k$  nodes in the output, we have a vector  $\mathbf{k}$  of demands from each color, i.e., they are lower bounds on how many nodes we have to include from each possible color.

As both formulations are NP-hard, we design polynomial-time approximation algorithms: For the first problem, we provide an approximation algorithm for the case where  $V$  is already diverse

(i.e.,  $V$  is a feasible solution for the problem). Our algorithm has an approximation ratio of  $\gamma \cdot \max\left\{\frac{1}{\lceil 1/\alpha \rceil}, \frac{1}{\alpha n}\right\}$ , where  $\gamma$  is the best approximation ratio known for DalkS (currently equal to  $1/2$  [28]). By simple calculation, we observe that the above approximation ratio leads to an approximation ratio of  $\Omega(1/\sqrt{n})$ , irrespective of any parameter other than the number of nodes. Moreover, we can also see that the approximation ratio is lower bounded by  $1/|C|$ , meaning that the algorithm attains a constant-factor approximation for the case of  $|C| = O(1)$  and is a generalization of the  $1/2$ -approximation algorithm for the fair densest subgraph problem by Anagnostopoulos et al. [4]. Our algorithm is based on an approximation algorithm for DalkS with a carefully selected value of  $k$ , together with a simple postprocessing. The primary factor determining the time complexity of our algorithm is the time complexity of the approximation algorithm used for DalkS.

For the second problem, we devise a  $1/3$ -approximation algorithm, which runs in polynomial time for constant  $|C|$ . In the design and analysis of our algorithm, we generalize the existing  $1/2$ -approximation algorithm for DalkS and its approximation ratio analysis. As shown later, we can get a  $1/4$ -approximate solution directly using the  $1/2$ -approximation algorithm for DalkS. Our effort improves the approximation ratio from  $1/4$  to  $1/3$ , by sacrificing some degree of scalability. We also present an acceleration technique for the proposed algorithm with the aid of the well-known greedy peeling algorithm [12]. The running time of our original algorithm is  $O((n/|C|)^{|C|}T_{LP})$ , where  $T_{LP}$  is the time required to solve an LP used in the algorithm, while the accelerated version runs in  $O(|C|(n/|C|)^{|C|-1}T_{LP})$ . In the case where  $|C|$  is a small constant, the reduction of the running time due to the acceleration is drastic.

We evaluate our algorithms on real-life attributed datasets, including social networks with gender and profession attributes. We compare against Anagnostopoulos et al. [4], but we also develop a novel baseline that uses node embeddings [10, 37, 39] combined with advances in scalable fair clustering of points [6]. The algorithms we propose have the capability to extract dense and diverse subgraphs. We demonstrate that real-world networks contain dense subgraphs that exhibit significant homophily, emphasizing the importance of employing our tools in scenarios where diversity or fairness is essential.

## 2 RELATED WORK

**DSP and its variations.** Given an undirected graph  $G = (V, E)$ , we define for any non-empty subset of nodes  $S \subseteq V$  the (degree) density  $d(S) = |E(S)|/|S|$ , where  $E(S) = \{\{u, v\} \in E \mid u, v \in S\}$ . The DSP aims to maximize the degree density over all possible subsets. Notice that the degree density is just half of the average degree of an induced subgraph. The DSP is polynomial-time solvable using  $O(\log n)$  max-flow computations [23, 38],  $O(1)$  number of flows using parametric maximum flow [20], or LP [12]. There also exists a linear-time  $1/2$ -approximation algorithm that greedily removes a node of the smallest degree, and reports the maximum degree density seen among these subsets [12, 29]. This kind of algorithm is often called the greedy peeling algorithm. Recently, Boob et al. [9] proposed GREEDY++, an iterative peeling algorithm that generalizes the above and converges to a near-optimal solution

extremely fast without the use of flows. Very recently, Chekuri et al. [13] analyzed the performance of GREEDY++, and showed that its algorithmic idea could be generalized to general fractional supermodular maximization.

DSP has a lot of problem variants [32]. Unlike the original DSP, its size-restricted variants are known to be NP-hard. Indeed, the densest  $k$ -subgraph problem (DkS) that asks for the densest subgraph with exactly  $k$  nodes, is not only NP-hard, but also hard to approximate, with the best-known approximation ratio being  $\Omega(1/n^{1/4+\epsilon})$  for any  $\epsilon > 0$  [8]. This approximability result is far off from the best-known hardness result that assumes the Exponential Time Hypothesis (ETH). If ETH holds, then DkS cannot be approximated within a ratio of  $n^{1/(\log \log n)^c}$  for some  $c > 0$ . Another size-restricted variant, DamkS, aims to maximize the degree density over all subsets of nodes  $S$  such that  $|S| \leq k$  [5]. It is known that  $\alpha$ -approximation to DamkS leads to  $\alpha/4$ -approximation to DkS [28].

Closest to this work lies DalkS, which imposes the cardinality constraint  $|S| \geq k$  [5]. The problem is also known to be NP-hard [28]. Andersen and Chellapilla [5] designed a  $1/3$ -approximation algorithm that runs in linear time, using greedy peeling. Khuller and Saha [28] designed two different approximation algorithms, that both achieve  $1/2$ -approximation using either a small number of flows, or by solving an LP [28].

**Fair densest subgraph problem.** Despite the large amount of research on DSD, the problem of finding a densest diverse subgraph has not received attention with the single exception of Anagnostopoulos et al. [4] who introduced the fair densest subgraph problem for two colors. Assuming the graph is fair to begin with, i.e., the two colors are equally represented in  $V$ , they demand equal representation of each category in the output. For this case, they proposed a greedy  $1/2$ -approximation algorithm and a spectral approach. The spectral approach comes with guarantees only in limited cases (e.g., all degrees being almost equal), which are not typical on real data that tend to have a skewed degree distribution. This algorithm can be extended to the case of  $|C| > 2$ , but without quality guarantees. Finally, the authors studied the hardness of the problem and showed that their formulation is at least as hard as DamkS: Any  $\alpha$ -approximation to their formulation leads to  $\alpha$ -approximation to DamkS.

**Fairness and algorithms.** While DSD with diversity is not yet well studied, fair clustering of clouds of points has received much attention from the data mining community. Chierichetti et al. [16] initiated the problem of finding balanced clusters in a cloud of points, namely clusters where two groups are equally represented. They proposed a method called fairlet decomposition, that decomposes a dataset into minimal sets that satisfy fair representation, called fairlets. Afterwards, typical machine learning algorithms, like  $k$ -median, can be used to obtain fair clusters. Backurs et al. [6] provided a scalable algorithm for the fairlet decomposition. Later work has extended the problem of finding fair clusters to the case of correlation clustering [3] and hierarchical clustering [2].

### 3 PROBLEM STATEMENTS

In this section, we formally introduce our optimization problems. Let  $G = (V, E)$  be an undirected graph with  $n = |V|$  nodes and

$m = |E|$  edges. Let  $C$  be a set of colors. Without loss of generality  $|C| \leq n$ . Let  $\ell: V \rightarrow C$  be the coloring function that assigns a color to each node. Given the above as input, we aim to find a densest diverse subgraph. We mathematically formalize the notion of diversity in two ways found in Sections 3.1 and 3.2, respectively.

#### 3.1 Densest Diverse Subgraph Problem (DDSP)

Our first notion aims to ensure that no single color dominates the rest. Specifically, for  $S \subseteq V$ , we denote by  $c_{\max}(S)$  the maximum number of nodes in  $S$  that receive the same color, i.e.,  $c_{\max}(S) = \max_{c \in C} |\{v \in S \mid \ell(v) = c\}|$ . We also denote by  $\alpha(S)$  the maximum fraction of monochromatic nodes in  $S$ , i.e.,  $\alpha(S) = c_{\max}(S)/|S|$ . Our problem can be formulated as follows:

**PROBLEM 1 (DDSP).** *Given an undirected graph  $G = (V, E)$ ,  $\ell: V \rightarrow C$ , and  $\alpha \in [1/|C|, 1]$ , find a subset of nodes  $S \subseteq V$  that maximizes the degree density  $d(S)$  subject to the constraint  $\alpha(S) \leq \alpha$ .*

This problem is a major generalization of the fair densest subgraph problem introduced by Anagnostopoulos et al. [4], which is obtained for the special values  $|C| = 2$  and  $\alpha = 1/2$ . As the fair densest subgraph problem is NP-hard [4], Problem 1 is also NP-hard. Clearly when  $\alpha = 1$  we are oblivious to diversity and obtain (polynomial-time solvable) DSP. More interestingly, Problem 1 contains two totally different optimization problems, DalkS and DamkS, as special cases:

**PROPOSITION 1.** *There exist polynomial-time reductions from DalkS and DamkS to Problem 1.*

**PROOF.** The reductions are obtained by appropriately setting the number of colors and the parameter  $\alpha$ . For DalkS, it suffices to construct the instance of Problem 1 by coloring the nodes with  $n$  distinct colors and setting  $\alpha$  to be  $1/k$ . For DamkS, it suffices to construct the instance of Problem 1 by coloring the nodes with the single color and adding  $k$  isolated nodes (i.e., dummy nodes) with another color, and setting  $\alpha = 1/2$ .  $\square$

#### 3.2 Densest At-Least- $\vec{k}$ -Subgraph (DalkS)

Our second formulation diversifies the output by ensuring that it does not overlook any specified color. To this end, the formulation requires the output to contain at least a given number of representatives from each color. For  $S \subseteq V$  and  $c \in C$ , let  $S_c = \{v \in S \mid \ell(v) = c\}$ . Our problem formulation is as follows:

**PROBLEM 2 (DalkS).** *Given an undirected graph  $G = (V, E)$ ,  $\ell: V \rightarrow C$ , and  $\mathbf{k} = (k_c)_{c \in C} \in \mathbb{Z}_{\geq 0}^{|C|}$ , find a subset of nodes  $S \subseteq V$  that maximizes the degree density  $d(S)$  subject to  $|S_c| \geq k_c$  for any  $c \in C$ .*

Obviously the above problem is a generalization of DalkS. Thus, the problem is NP-hard. If  $\mathbf{k} = \mathbf{0}$ , the problem is reduced to the original DSP; therefore, throughout the paper, we assume that  $k_c \geq 1$  for some  $c \in C$ . As mentioned in the introduction, we can easily get a  $1/4$ -approximation algorithm for the problem:

**PROPOSITION 2.** *For Problem 2, there exists a polynomial-time  $1/4$ -approximation algorithm.*

**PROOF.** Let  $G = (V, E)$ ,  $\ell: V \rightarrow C$ ,  $\mathbf{k} = (k_c)_{c \in C} \in \mathbb{Z}_{\geq 0}^{|C|}$  be an instance of Problem 2, and OPT the optimal value of the instance.

**Procedure 1:** Diversify( $S$ )

---

```

1 while  $\alpha(S) > \alpha$  do
2   Find  $v_{\min} \in V \setminus S$  that satisfies
    $\ell(v_{\min}) \in \operatorname{argmin}_{c \in C} |S_c|$ ;
   /* In practice consider also the objective
   value. */
3    $S \leftarrow S \cup \{v_{\min}\}$ ;
4 return  $S$ ;
```

---

To get a feasible solution for Problem 2, we have to take at least  $k_c$  nodes for every color  $c \in C$ ; therefore, DalkS with  $k = \|\mathbf{k}\|_1 = \sum_{c \in C} k_c$  on  $G$  is a relaxation of Problem 2 on  $G$ . Let  $S \subseteq V$  be an  $\alpha$ -approximate solution for DalkS with  $k = \|\mathbf{k}\|_1$ . As DalkS with  $k = \|\mathbf{k}\|_1$  is a relaxation of Problem 2, we have  $d(S) \geq \alpha \cdot \text{OPT}$ . Note that  $S$  is not necessarily feasible for Problem 2, but we can make it feasible by adding at most  $k_c$  nodes for every color  $c \in C$ , resulting in adding at most  $\|\mathbf{k}\|_1 \leq |S|$  nodes in total. Letting  $S' \subseteq V$  be the resulting subset, we have  $d(S') = \frac{|E(S')|}{|S'|} \geq \frac{|E(S)|}{|S'|} \geq \frac{1}{2} \cdot d(S) \geq \frac{\alpha}{2} \cdot \text{OPT}$ , meaning that  $S'$  is an  $\alpha/2$ -approximate solution for Problem 2. As mentioned above, there is a polynomial-time  $1/2$ -approximation algorithm for DalkS [28]. Therefore, we can set  $\alpha = 1/2$  and have a polynomial-time  $1/4$ -approximation algorithm for Problem 2.  $\square$

## 4 ALGORITHM FOR PROBLEM 1

In this section, we design a polynomial-time  $\Omega(1/\sqrt{n})$ -approximation algorithm for Problem 1. A subset of nodes  $S \subseteq V$  is said to be diverse if  $\alpha(S) \leq \alpha$  holds. In what follows, we assume that  $V$  is diverse, i.e.,  $\alpha(V) \leq \alpha$ .

### 4.1 The Proposed Algorithm

Our algorithm first computes a constant-factor approximate solution (say  $\gamma$ -approximate solution) to DalkS on  $G$  with  $k = \lceil 1/\alpha \rceil$ . For example, we can use a  $1/2$ -approximation algorithm using LP [28] or a  $1/3$ -approximation algorithm using greedy peeling [5]. Then, the algorithm makes the solution feasible by adding an arbitrary node with a color of the lowest participation iteratively (Procedure 1). For reference, the entire procedure of our algorithm is summarized in Algorithm 2.

The time complexity of Algorithm 2 is dominated by the algorithm we use for DalkS. Even if we consider the objective value in Procedure 1, it still depends on the choice of the approximation algorithm for DalkS. If we employ a  $1/2$ -approximation algorithm using LP, the time complexity of Algorithm 2 is dominated by that required for solving the LP. If we use a  $1/3$ -approximation algorithm using greedy peeling, Algorithm 2 can be implemented to run in  $O(m + n \log n)$  time, by handling the nodes outside  $S$  using a Fibonacci heap for each color with key values being degrees to  $S$ .

### 4.2 Analysis

We analyze the approximation ratio of Algorithm 2. We have the following key lemma:

**Algorithm 2:**  $\Omega(1/\sqrt{n})$ -approximation algorithm

---

```

Input :  $G = (V, E)$ ,  $\ell: V \rightarrow C$ ,  $\alpha \in [1/|C|, 1]$ 
Output:  $S \subseteq V$ 
1  $S \leftarrow \gamma$ -approximate solution to DalkS on  $G$  with  $k = \lceil 1/\alpha \rceil$ ;
  /* See [28] and [5] for the algorithms achieving
   $\gamma = 1/2$  and  $\gamma = 1/3$ , respectively. */
2 return Diversify( $S$ );
```

---

LEMMA 1. Assume that  $\alpha(V) \leq \alpha$  holds. Then, for any  $S \subseteq V$  with  $|S| \geq \lceil 1/\alpha \rceil$ , it holds that  $|\text{Diversify}(S)| \leq \min\{\lceil 1/\alpha \rceil, \alpha n\} \cdot |S|$ .

PROOF. We first prove that  $c_{\max}(S) = c_{\max}(\text{Diversify}(S))$  holds. As  $S \subseteq \text{Diversify}(S)$ , we have  $c_{\max}(S) \leq c_{\max}(\text{Diversify}(S))$ . Therefore, it suffices to show  $c_{\max}(S) \geq c_{\max}(\text{Diversify}(S))$ . For  $S \subseteq V$  and  $c \in C$ , we denote by  $f(S, c)$  the fraction of the nodes in  $S$  that receive the color  $c$ , i.e.,  $f(S, c) = |S_c|/|S|$ . Let us focus on an arbitrary iteration of Procedure 1 and let  $S' \subseteq V$  be the subset kept at the beginning of the iteration. Then there exists  $c \in C$  that satisfies the condition  $f(S', c) < \alpha$ . Suppose, for contradiction, that there exist no such colors. Then, for any color  $c \in C$ , we have  $f(S', c) \geq \alpha$ . Moreover, as  $S'$  is not yet feasible, we see that there exists  $c' \in C$  that satisfies  $f(S', c') > \alpha$ . Therefore, we have

$$1 = \sum_{c \in C} f(S', c) = f(S', c') + \sum_{c \in C \setminus \{c'\}} f(S', c) > \alpha + (|C| - 1)\alpha = |C|\alpha \geq 1,$$

a contradiction. From the above, recalling the greedy rule of Procedure 1, we see that Procedure 1 only adds the nodes with the colors  $c \in C$  that satisfy  $f(S', c) < \alpha$ . To increase  $c_{\max}$  in this iteration, Procedure 1 needs to add a node with a color  $c \in C$  that satisfies  $f(S', c) > \alpha$ , but it does not happen. As we fixed an iteration arbitrarily,  $c_{\max}$  does not increase throughout Procedure 1. Therefore, we have  $c_{\max}(S) \geq c_{\max}(\text{Diversify}(S))$ .

Assume that in some iteration, which produces  $S'' \subseteq V$ , of Procedure 1,  $|S''| = \lceil |S|/\alpha \rceil$  holds. Then we have

$$\alpha(S'') = \frac{c_{\max}(S'')}{|S''|} = \frac{c_{\max}(S)}{\lceil |S|/\alpha \rceil} \leq \alpha \cdot \frac{c_{\max}(S)}{|S|} \leq \alpha,$$

where the second equality follows from  $c_{\max}(S) = c_{\max}(\text{Diversify}(S))$ . This means that the algorithm terminates at or before this iteration. Therefore, we have  $|\text{Diversify}(S)| \leq \lceil |S|/\alpha \rceil$ .

Similarly, assume that in some iteration, which produces  $S'' \subseteq V$ , of Procedure 1,  $|S''| = c_{\max}(S)|S|$  holds. Then we have

$$\alpha(S'') = \frac{c_{\max}(S'')}{|S''|} = \frac{c_{\max}(S)}{c_{\max}(S)|S|} = \frac{1}{|S|} \leq \frac{1}{\lceil 1/\alpha \rceil} \leq \alpha,$$

where the first inequality follows from the assumption of the lemma. This means that the algorithm terminates at or before this iteration. Therefore, we have  $|\text{Diversify}(S)| \leq c_{\max}(S)|S|$ .

From the above, we see that

$$\begin{aligned} |\text{Diversify}(S)| &\leq \min\{\lceil |S|/\alpha \rceil, c_{\max}(S)|S|\} \\ &\leq \min\{\lceil 1/\alpha \rceil, c_{\max}(V)\} \cdot |S| \\ &\leq \min\{\lceil 1/\alpha \rceil, \alpha n\} \cdot |S|, \end{aligned}$$

which completes the proof.  $\square$

Using the above lemma, we can prove the following:

**THEOREM 3.** *Algorithm 2 is a  $\left(\gamma \cdot \max\left\{\frac{1}{\lceil 1/\alpha \rceil}, \frac{1}{\alpha n}\right\}\right)$ -approximation algorithm for Problem 1 when  $\alpha(V) \leq \alpha$  holds. Here  $\gamma$  is the approximation ratio of the algorithm for DalkS used in Algorithm 2.*

**PROOF.** Let OPT be the optimal value of Problem 1. Let  $S \subseteq V$  be a constant-factor approximate solution (say  $\gamma$ -approximate solution) for DalkS on  $G$  with  $k = \lceil 1/\alpha \rceil$ . Note that DalkS with  $k = \lceil 1/\alpha \rceil$  is a relaxation of Problem 1, because even if we pick the nodes, all of which have different colors, we need at least  $k = \lceil 1/\alpha \rceil$  nodes to satisfy the diversity constraint. Thus, we have  $d(S) = \frac{|E(S)|}{|S|} \geq \gamma \cdot \text{OPT}$ . The output of the algorithm is Diversify( $S$ ) whose objective value can be evaluated as follows:

$$\begin{aligned} d(\text{Diversify}(S)) &= \frac{|E(\text{Diversify}(S))|}{|\text{Diversify}(S)|} \\ &\geq \frac{|E(S)|}{\min\{\lceil 1/\alpha \rceil, \alpha n\} \cdot |S|} \\ &\geq \gamma \cdot \max\left\{\frac{1}{\lceil 1/\alpha \rceil}, \frac{1}{\alpha n}\right\} \cdot \text{OPT}, \end{aligned}$$

where the first inequality follows from Lemma 1.  $\square$

It should be emphasized that the approximation ratio of Algorithm 2 is lower bounded by  $1/|C|$ , meaning that the algorithm is a constant-factor approximation algorithm for the case of  $|C| = O(1)$ . In fact, we see that  $\frac{1}{\lceil 1/\alpha \rceil} \geq \frac{1}{\lceil |C| \rceil} = \frac{1}{|C|}$ . Therefore, Algorithm 2 can be seen as a generalization of the 1/2-approximation algorithm for the fair densest subgraph problem by Anagnostopoulos et al. [4].

Moreover, the above analysis leads to an approximation ratio that is independent of any parameter other than  $n$ :

**COROLLARY 1.** *Algorithm 2 is an  $\Omega(1/\sqrt{n})$ -approximation algorithm for Problem 1 when  $\alpha(V) \leq \alpha$  holds.*

**PROOF.** We can lower bound the approximation ratio given in Theorem 3 as follows:

$$\gamma \cdot \max\left\{\frac{1}{\lceil 1/\alpha \rceil}, \frac{1}{\alpha n}\right\} \geq \gamma \cdot \sqrt{\frac{1}{\lceil 1/\alpha \rceil} \cdot \frac{1}{\alpha n}} = \Omega\left(\frac{1}{\sqrt{n}}\right).$$

$\square$

## 5 ALGORITHM FOR PROBLEM 2

In this section, we design a polynomial-time 1/3-approximation algorithm for Problem 2 with  $|C| = O(1)$ . Using a sophisticated LP, we can improve the approximation ratio of 1/4 given in Proposition 2. Recall that  $\mathbf{k} = (k_c)_{c \in C} \in \mathbb{Z}_{\geq 0}^{|C|}$  is the input lower-bound vector.

---

### Procedure 3: Make\_it\_feasible( $S$ )

---

```

1 for each  $c \in C$  do
2   if  $|S_c| < k_c$  then
3     Take arbitrary  $v \in V_c \setminus S_c$ ;
     /* In practice consider also the objective
     value. */
4      $S \leftarrow S \cup \{v\}$ ;
5 return  $S$ ;
```

---



---

### Algorithm 4: 1/3-approximation algorithm for Problem 2

---

```

Input :  $G = (V, E)$ ,  $\ell: V \rightarrow C$ ,  $\mathbf{k} = (k_c)_{c \in C} \in \mathbb{Z}_{\geq 0}^{|C|}$ 
Output :  $S \subseteq V$ 
1 for each  $\mathbf{p}$  such that  $\mathbf{k} \leq \mathbf{p} \leq (|V_c|)_{c \in C}$  do
2   Solve LP( $\mathbf{p}$ ) to obtain an optimal solution  $(\mathbf{x}^*, \mathbf{y}^*)$ ;
3   Construct Candidates( $\mathbf{p}$ ) := {Make_it_feasible( $S(r)$ ) |
    $r \in \{y_v^* \mid v \in V\} \cup \{0\}$ };
4   Take Best( $\mathbf{p}$ )  $\in \text{argmax}\{d(S) \mid S \in \text{Candidates}(\mathbf{p})\}$ ;
5 return  $\text{argmax}\{d(S) \mid S \in \{\text{Best}(\mathbf{p}) \mid \mathbf{k} \leq \mathbf{p} \leq (|V_c|)_{c \in C}\}\}$ ;
```

---

## 5.1 The Proposed Algorithm

Let  $\mathbf{p} = (p_c)_{c \in C} \in \mathbb{Z}_{\geq 0}^{|C|}$  be a vector that satisfies  $\mathbf{p} \geq \mathbf{k}$ . We consider the following LP:

$$\begin{aligned} \text{LP}(\mathbf{p}): \quad & \text{maximize} \quad \sum_{e \in E} x_e \\ & \text{subject to} \quad x_e \leq y_u, \quad x_e \leq y_v \quad \forall e = \{u, v\} \in E, \\ & \quad \sum_{v \in V_c} y_v = \frac{p_c}{\|\mathbf{p}\|_1} \quad \forall c \in C, \\ & \quad y_v \leq \frac{1}{\|\mathbf{p}\|_1} \quad \forall v \in V, \\ & \quad x_e, y_v \geq 0 \quad \forall e \in E, \forall v \in V. \end{aligned}$$

Note that this is a major generalization of the LP used in the 1/2-approximation algorithm for DalkS [28]. The LP for DalkS is parameterized by a single value, while our LP is parameterized by multiple values (i.e., a vector  $\mathbf{p}$ ). This modification is essential to address the generalization caused by DalkS.

Let  $(\mathbf{x}^*, \mathbf{y}^*)$  be an optimal solution to LP( $\mathbf{p}$ ) for  $\mathbf{p}$ . For  $(\mathbf{x}^*, \mathbf{y}^*)$  and  $r \geq 0$ , we define  $S(r) = \{v \in V \mid y_v^* \geq r\}$ . Now we are ready to present our algorithm. For each  $\mathbf{p} \in \mathbb{Z}_{\geq 0}^{|C|}$  such that  $\mathbf{k} \leq \mathbf{p} \leq (|V_c|)_{c \in C}$ , the algorithm conducts the following procedure: It first solves LP( $\mathbf{p}$ ) to obtain an optimal solution  $(\mathbf{x}^*, \mathbf{y}^*)$ . Then using the solution, the algorithm enumerates all possible  $S(r)$ 's by setting  $r \in \{y_v^* \mid v \in V\} \cup \{0\}$  each, makes them feasible by adding nodes with appropriate colors using Procedure 3, and takes the best subset among them, as a candidate for  $\mathbf{p}$ . After the above iterations, the algorithm finally outputs the best subset among all candidates for  $\mathbf{p}$  with  $\mathbf{k} \leq \mathbf{p} \leq (|V_c|)_{c \in C}$ . The entire procedure is summarized in Algorithm 4.

## 5.2 Analysis

We first give a lower bound on the optimal value of LP( $\mathbf{p}$ ).

LEMMA 2. For any  $S \subseteq V$  such that  $|S_c| = p_c$  for every  $c \in C$ , there exists a feasible solution of  $\text{LP}(\mathbf{p})$  whose objective function value is greater than or equal to  $d(S)$ .

PROOF. Construct a solution  $(\mathbf{x}, \mathbf{y})$  of  $\text{LP}(\mathbf{p})$  as follows:

$$x_e = \begin{cases} 1/\|\mathbf{p}\|_1 & \text{if } e \in E(S), \\ 0 & \text{otherwise,} \end{cases} \quad y_v = \begin{cases} 1/\|\mathbf{p}\|_1 & \text{if } v \in S, \\ 0 & \text{otherwise.} \end{cases}$$

Then we can see that  $(\mathbf{x}, \mathbf{y})$  is feasible for  $\text{LP}(\mathbf{p})$ . In fact,  $\sum_{v \in V_c} y_v = \sum_{v \in S_c} y_v = \frac{p_c}{\|\mathbf{p}\|_1}$ . The objective function value of  $(\mathbf{x}, \mathbf{y})$  is  $\sum_{e \in E} x_e = \sum_{e \in E(S)} x_e = \frac{|E(S)|}{\|\mathbf{p}\|_1} = d(S)$ . Thus, we have the lemma.  $\square$

For  $S \subseteq V$ , let  $C_{\text{sat}}(S) = \{c \in C \mid |S_c| \geq k_c\}$ , i.e., the set of colors for which the constraint is satisfied by  $S$ . We can prove the following key lemma (see Appendix A for the proof):

LEMMA 3. Let  $S^* \subseteq V$  be an optimal solution to Problem 2. For each  $c \in C$ , let  $p_c^* = |S_c^*|$ . Let  $(\mathbf{x}^*, \mathbf{y}^*)$  be an optimal solution to  $\text{LP}(\mathbf{p}^*)$  and  $\lambda$  its objective value. Then there exists  $S \in \{S(r) \mid r \in \{y_v^* \mid v \in V\} \cup \{0\}\}$  that satisfies (exactly) one of the following:

- (1)  $d(S) \geq \lambda/3$  and  $C_{\text{sat}}(S) = C$ ;
- (2)  $\frac{|E(S)|}{\sum_{c \in C_{\text{sat}}(S)} |S_c| + \sum_{c \in C \setminus C_{\text{sat}}(S)} k_c} \geq \lambda/3$  and  $C_{\text{sat}}(S) \neq \emptyset \neq C \setminus C_{\text{sat}}(S)$ ;
- (3)  $|E(S)| \geq |E(S^*)|/3$  and  $C_{\text{sat}}(S) = \emptyset$ .

Based on the above lemma, we prove the following:

THEOREM 4. Algorithm 4 is a 1/3-approximation algorithm for Problem 2 and runs in  $O((n/|C|)^{|C|} T_{\text{LP}})$  time, where  $T_{\text{LP}}$  is the time complexity required to solve the LP.

PROOF. As the time complexity analysis is straightforward, we show the approximation ratio of 1/3. Let  $S^* \subseteq V$  be an optimal solution to Problem 2 and for each  $c \in C$ , let  $p_c^* = |S_c^*|$ . Looking at Lemma 3, we see that when the algorithm solves  $\text{LP}(\mathbf{p}^*)$ , one of the possible  $S(r)$ 's itself is just the subset of nodes  $S \subseteq V$  whose existence is guaranteed in Lemma 3. Then it is easy to see that  $\text{Make\_it\_feasible}(S)$  is a 1/3-approximate solution. Therefore,  $\text{Best}(\mathbf{p}^*)$  is also a 1/3-approximate solution. Noticing that the algorithm outputs the best of  $\text{Best}(\mathbf{p})$ 's for all possible  $\mathbf{p}$  with  $\mathbf{k} \leq \mathbf{p} \leq (|V_c|)_{c \in C}$  (containing  $\mathbf{p}^*$ ), we are done.  $\square$

Here we mention a simple acceleration technique for Algorithm 4. As shown in Lemma 2,  $\text{LP}(\mathbf{p}^*)$  with  $p_c^* = |S_c^*|$  has the optimal value greater than or equal to  $d(S^*)$ . Therefore, if we solve  $\text{LP}(\mathbf{p})$  for some  $\mathbf{p}$  and its optimal value is less than the density of the current best feasible solution, the LP must not be  $\text{LP}(\mathbf{p}^*)$  and thus we can skip the procedure to be applied to its optimal solution (i.e., Lines 3 and 4). In the next section, we present an additional acceleration technique that can reduce the number of LPs to solve.

### 5.3 Acceleration via Greedy Peeling

Our acceleration technique for Algorithm 4 is based on greedy peeling. Algorithm 5 is a straightforward application of greedy peeling to Problem 2, where for  $S \subseteq V$  and  $v \in S$ ,  $\text{deg}_S(v)$  denotes the degree of  $v$  in  $G[S] = (S, E(S))$ .

This algorithm achieves an approximation ratio of 1/2 in a very specific case as follows:

---

#### Algorithm 5: Greedy peeling algorithm for Problem 2

---

**Input** :  $G = (V, E)$ ,  $\ell: V \rightarrow C$ ,  $\mathbf{k} = (k_c)_{c \in C} \in \mathbb{Z}_{\geq 0}^{|C|}$   
**Output** :  $S \subseteq V$

- 1  $S^{(n)} \leftarrow V, i \leftarrow n$ ;
- 2 **while**  $|S_c^{(i)}| > k_c$  for every  $c \in C$  with  $k_c \geq 1$  **do**
- 3      $v_{\min} \leftarrow \text{argmin}_{v \in S^{(i)}} \text{deg}_{S^{(i)}}(v)$ ;
- 4      $S^{(i-1)} \leftarrow S^{(i)} \setminus \{v_{\min}\}, i \leftarrow i - 1$ ;
- 5 **return**  $\text{argmax}\{d(S) \mid S \in \{S^{(n)}, \dots, S^{(i)}\}\}$ ;

---



---

#### Algorithm 6: Accelerated 1/3-approximation algorithm for Problem 2

---

**Input** :  $G = (V, E)$ ,  $\ell: V \rightarrow C$ ,  $\mathbf{k} = (k_c)_{c \in C} \in \mathbb{Z}_{\geq 0}^{|C|}$   
**Output** :  $S \subseteq V$

- 1 Run Algorithm 5 and obtain its output  $S_{\text{peel}}$ ;
- 2 Run Algorithm 4 after replacing " $\mathbf{k} \leq \mathbf{p} \leq (|V_c|)_{c \in C}$ " by " $\mathbf{k} \leq \mathbf{p} \leq (|V_c|)_{c \in C}$  with  $p_c = k_c$  for some  $c \in C$  with  $k_c \geq 1$ " in Lines 1 and 5, and obtain its output  $S_{\text{LP}}$ ;
- 3 **return**  $\text{argmax}\{d(S) \mid S \in \{S_{\text{peel}}, S_{\text{LP}}\}\}$ ;

---

LEMMA 4. Assume that there exists an optimal solution  $S^* \subseteq V$  that satisfies  $|S_c^*| > k_c$  for every  $c \in C$  with  $k_c \geq 1$ . Then Algorithm 5 outputs a 1/2-approximate solution for Problem 2.

PROOF. It is easy to see that for any  $v \in S^*$ ,  $S^* \setminus \{v\}$  is a feasible solution for Problem 2. Therefore, we have that for any  $v \in S^*$ ,

$$d(S^*) \geq d(S^* \setminus \{v\}).$$

Transforming the above inequality, we have that for any  $v \in S^*$ ,

$$\text{deg}_{S^*}(v) \geq d(S^*). \quad (1)$$

Let  $v^*$  be the node that is contained in  $S^*$  and removed first by the algorithm. Note that the existence of such a node is guaranteed due to the assumption of the lemma and the design of the algorithm. Let  $S' \subseteq V$  be the subset of nodes kept just before the removal of  $v^*$ . Obviously  $S^* \subseteq S'$  and thus  $S'$  is a feasible solution of Problem 2. We can evaluate the density of  $S'$  as follows:

$$d(S') = \frac{\frac{1}{2} \sum_{v \in S'} \text{deg}_{S'}(v)}{|S'|} \geq \frac{1}{2} \text{deg}_{S^*}(v^*) \geq \frac{1}{2} d(S^*),$$

where the first inequality follows from the greedy choice of  $v^*$  and the relation  $S' \supseteq S^*$ , and the second inequality follows from inequality (1). Noticing that  $S'$  is one of the output candidates of the algorithm, we have the lemma.  $\square$

Finally, our accelerated algorithm is described in Algorithm 6.

THEOREM 5. Algorithm 6 is a 1/3-approximation algorithm for Problem 2 and runs in  $O(|C|(n/|C|)^{|C|-1} T_{\text{LP}})$  time, where  $T_{\text{LP}}$  is the time complexity required to solve the LP.

PROOF. The time complexity analysis is again trivial. Hence, in what follows, we guarantee the approximation ratio. If there exists an optimal solution  $S^* \subseteq V$  that satisfies  $|S_c^*| > k_c$  for every  $c \in C$  with  $k_c \geq 1$ , then the output  $S_{\text{peel}}$  of Algorithm 5 is a 1/2-approximate solution; thus, we are done. Otherwise there exists an

**Table 1: Dataset statistics summary.**

Dataset	$ C $	$n$	$m$	$\alpha(V)$
Amazon Product	2	3.3K $\pm$ 4.7K	24.7K $\pm$ 47.3K	0.68 $\pm$ 0.12
Facebook100	2 or 4	10.7K $\pm$ 8.4K	399K $\pm$ 315K	0.48 $\pm$ 0.09
GitHub Developers	2	37,700	289,003	0.74
LastFM Asia	18	7,624	27,806	0.20
Deezer Europe	2	28,281	92,752	0.56
DBLP	6	25,176	151,670	0.32

optimal solution  $S^* \subseteq V$  that satisfies  $|S_c^*| = k_c$  for some  $c \in C$  with  $k_c \geq 1$ . The modified version of Algorithm 4 used in Algorithm 6 skips some  $p$ 's with  $k \leq p \leq (|V_c|)_{c \in C}$  but still tests  $p^*$  with  $p_c^* = |S_c^*|$  for any  $c \in C$ . Therefore, the analysis similar to that for Algorithm 4 still works and we see that  $S_{LP}$  is a 1/3-approximate solution. Therefore, we have the theorem.  $\square$

## 6 EXPERIMENTAL EVALUATION

In this section, we evaluate our proposed algorithms using a variety of synthetic and real-world networks.

### 6.1 Experimental Setup

**Datasets.** We use two collections of datasets and four single-graph datasets, with attributes for the nodes in the graphs. Table 1 summarizes the statistics of the following datasets, where for the first two datasets, we put the average values with the standard deviations.

- *Amazon Product Metadata* - 299 graphs [36]. This consists of a collection of 299 graphs, each with 2 colors, as curated by Anagnostopoulos et al. [4]. Classes are product categories, while edges indicate that products were co-purchased.
- *Facebook100* - 100 graphs [45]. This contains 100 anonymized networks of the Facebook social network from universities across the United States. Nodes have demographic attributes, like profession (faculty/student), gender (male/female), the year they joined the university, etc. For our purpose, we also create  $2^2 = 4$  categories that combine profession and gender. The largest graph in terms of  $m$  has 1,382,325 edges.
- *GitHub Developers* [41]. Nodes are developers in GitHub who have starred at least 10 repositories, and edges are mutual follower relationships between them. The attribute is whether the user is a web or a machine learning developer.
- *LastFM Asia Social Network* [42]. Nodes are LastFM users from Asian countries and edges are mutual follower relationships between them. The attribute is the location of a user.
- *Deezer Europe Social Network* [42]. Nodes are Deezer users from European countries and edges are mutual follower relationships between them. The attribute is the gender of a user.
- *DBLP Co-Authorship*. We also create a dataset from the DBLP database. We create a graph from authors of papers published in major conferences in 6 areas (Theory, Data Management, Data Mining, Learning, Networking, and Image & Video Processing) between 2003 and 2022. Nodes are authors with attributes being the area (s)he has most publications. Edges

indicate that the two authors have collaborated at least once. See Appendix B.1 for details.

**Baselines.** For Problem 1, we employ the following baselines, including a novel baseline we devise using node embeddings [39]:

- *DSP*. The (unconstrained) densest subgraph. Specifically, we use GREEDY++ by Boob et al. [9] for 5 iterations, which in practice finds an optimal solution to DSP in very few iterations. For  $S \subseteq V$ , we define the normalized density as its density divided by the optimal value to DSP.
- *PS* and *FPS*. The spectral algorithms by Anagnostopoulos et al. [4] that split the entries of the largest eigenvector of the adjacency matrix (*PS*) and “fairified” adjacency matrix (*FPS*), based on their color, and sort them in a similar spirit with spectral clustering.
- *Embedding+Fair Clustering*. We design a novel baseline that first embeds the graph using the node embedding method NetMF [39] and then clusters the nodes using the fair  $k$ -means implementation of Backurs et al. [6] for various values of  $k$ . We compute the density of each fair cluster for each value of  $k$ , and output the maximum among all.

For Problem 2, we run the following two baselines:

- *IP*. We implement an exact algorithm using integer programming and test its scalability. See Appendix B.2 for details.
- *Prop2*. We also implement the 1/4-approximation algorithm introduced in Proposition 2. To solve DalkS, we use the 1/2-approximation algorithm based on LP with a greedy peeling acceleration [28].

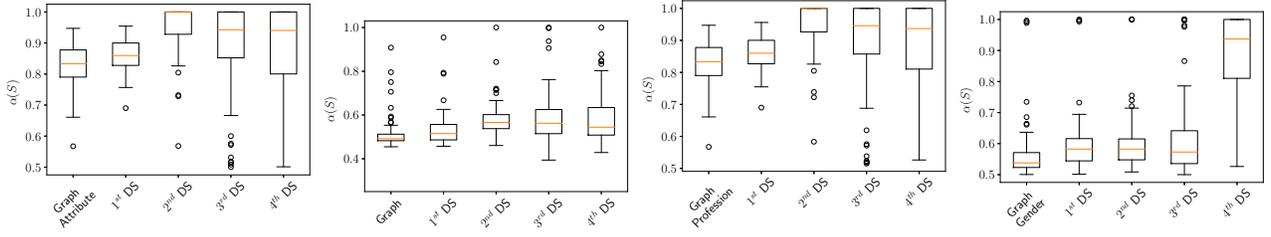
**Machine specs and code.** The experiments are performed on a single machine, with Intel i7-10850H CPU @ 2.70GHz and 32GB of main memory. We use C++ for all experiments. Linear programming and integer programming are implemented with SciPy and solved with HiGHS [25]. The code and datasets are available online.<sup>1</sup>

### 6.2 Evaluation of Algorithm 2 (for Problem 1)

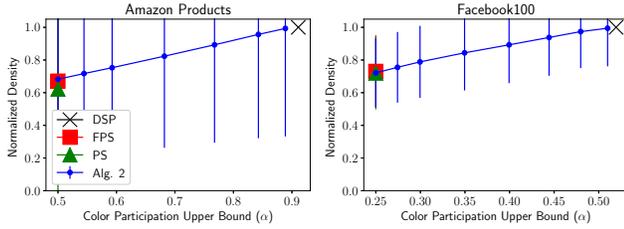
**Preliminary experiments.** Although we have created a reliable algorithmic toolbox for dense diverse clusters, we must assess their practical effectiveness. Specifically, when examining the densest subgraph or the top- $k$  densest subgraphs with sensitive attributes, do we observe diversity in the subgraphs? The answer is generally no, as observed in the vast majority of cases. This trend can result in unfair solutions, particularly in the context of sensitive attributes.

In Figure 1 we present the statistics of the entire graph and the top-4 densest subgraphs of the Amazon Product and Facebook100 datasets from Table 1, using standard box plots. As can be seen, the densest subgraphs tend to be not diverse. Especially for the Amazon Product dataset (leftmost figure), the densest subgraph most of the times is nearly monochromatic. This phenomenon is not restricted to the densest subgraph but also found for the top-4 densest subgraphs, implying that strong homophily is a key factor behind dense clusters. We observe a similar trend for the Facebook100 dataset (the latter three figures). Even if we restrict our attention to the specific attribute we consider, profession or gender, we see that again they are not represented equally in the densest subgraph. These findings motivate the study of Problem 1,

<sup>1</sup><https://github.com/tsourakakis-lab/densest-diverse-subgraphs>



**Figure 1: The diversity of the entire graph and the top-4 densest subgraphs of the first two datasets in Table 1. The first figure is for the Amazon Product dataset ( $|C| = 2$ ), while the latter three figures are for the Facebook100 dataset (left: Profession and Gender ( $|C| = 4$ ), middle: Profession ( $|C| = 2$ ), right: Gender ( $|C| = 2$ )). Note that  $\alpha(S) = 1/|C|$  means the complete diversity of  $S$ .**



**Figure 2: Performance of algorithms for the Amazon Products and Facebook100 datasets.**

where we can control the extent to which a specific color dominates the subgraph returned by an algorithm.

**Solution quality of algorithms.** We compare Algorithm 2 with the baseline algorithms with respect to the density and diversity of the returned subgraphs, using the six datasets of Table 1. We run Algorithm 2 by varying the value of  $\alpha$ . Note that when  $\alpha(V) > \alpha$ , Algorithm 2 may fail to find a feasible solution. In that case, we resort to iteratively peeling a node with the minimum degree from  $S$  with the most dominant color, until we obtain a feasible solution, or conclude that none exists. *PS* and *FPS* by Anagnostopoulos et al. [4] return by definition completely-balanced subgraphs, thus corresponding to the solutions of Problem 1 for a value of  $\alpha = 1/|C|$ .

Our results are depicted in Figure 2 for the two collections of datasets, Amazon Product and Facebook100. As these two datasets consist of a number of graphs, we present the average values with the standard deviations. We see that Algorithm 2 outperforms the baselines. Let us first focus on the case of  $\alpha = 1/|C|$ . Algorithm 2 returns a better solution than that of *PS* and *FPS* for the Amazon Product dataset, while *PS* and *FPS* are comparable to Algorithm 2 for the Facebook100 dataset. We observe similar trends for the single graph datasets, as seen in Figure 3. More precisely, when  $\alpha = 1/|C|$ , Algorithm 2 performs better than *PS* and *FPS* in 3 out of 4 cases. Moreover, we see that by varying the parameter  $\alpha$ , we can adjust the trade-off between getting a subgraph as dense as possible (when we let  $\alpha$  approach the diversity of the densest subgraph in the graph), or sacrifice density with the aim of getting more diverse subgraphs (completely balanced when we set  $\alpha = 1/|C|$ ).

**(Poor) performance of Embedding+Fair Clustering.** Here we use a random sample of 35 graphs from the Facebook100 dataset due to the cost of producing embeddings for each graph, as we did not

notice significant changes after including more datasets in terms of results. We also use only gender as our attribute (as fair clustering methods allow only two colors), while we range  $k$  in  $k$ -means from 2 to 32 and report the best result. In Figure 4, we see the results of this baseline in terms of the diversity of the clusters we get, as well their density. We see that simply performing  $k$ -means on the embedded graph results in clusters that are highly unbalanced. On the other hand, fair  $k$ -means allows us to get fair clusters by design (Figure 4 (left)), but the corresponding subgraph is not as dense as the output by Algorithm 2 (Figure 4 (right)). Notice that for the attribute we consider in this case, the ratio of the two colors in the graph is only slightly different from that of the densest subgraph (see the rightmost of Figure 1). Hence, Algorithm 2 finds subgraphs almost as dense as the densest.

**Running time of algorithms.** In Figure 5 (left) we report the runtime of algorithms on the Amazon Product dataset. We employed this dataset because the sizes of the graphs (in terms of  $n$  and  $m$ ) vary substantially so that we can easily see how the running time of algorithms grows. We see that all of them scale (almost) linearly with the number of edges; in particular, the result for Algorithm 2 is consistent with the theoretical analysis in Section 4.1. Note that for Algorithm 2, we report the runtime for the most difficult case, when  $\alpha = 0.5$ ; hence the subgraphs need to be completely balanced. As we relax the diversity constraint, the runtime drops significantly as evident on the synthetic graphs, as shown in Figure 5 (right). The synthetic graph we employed resembles the stochastic block model paradigm. We constructed a graph with 5 clusters of 40,000 nodes each, with higher intra-connectivity within clusters, than inter-connectivity across clusters. Most importantly, one of the clusters has a much higher density than the rest. Hence, the densest subgraph consists of 40,000 nodes. Initially, we assign nodes to one of five colors uniformly at random. Hence, the (original) densest subgraph is already diverse. Therefore, Algorithm 2 does not need to invoke Procedure 1. Thus, relaxing the diversity constraint has no impact in this case (blue dotted line). On the other hand, if we assign a unique color to each cluster, then the densest subgraph becomes monochromatic and in order to have a diverse representation we need to involve the whole graph to our solution, resulting in a longer runtime that gets decreased as we relax the diversity constraint (solid orange line). Finally we remark that on the Facebook100 dataset Algorithm 2 runs in 3s for  $\alpha = 0.5$  and 11s for  $\alpha = 0.25$  even on the aforementioned largest graph with more than one million edges.

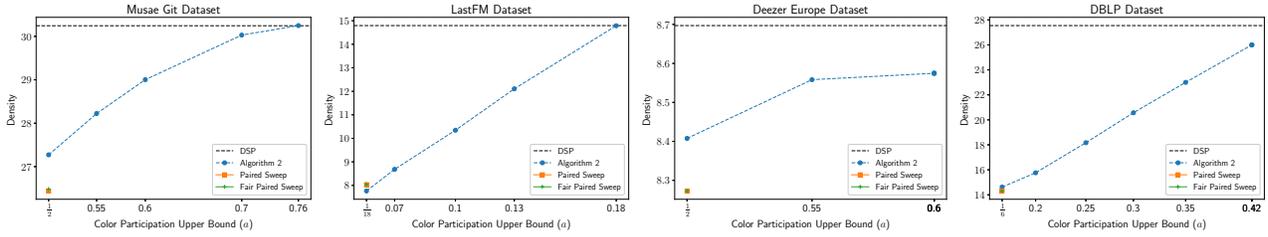


Figure 3: Performance of algorithms for the single-graph datasets.

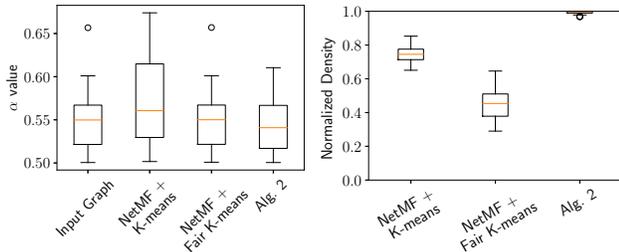


Figure 4: Diversity and normalized density of subgraphs obtained using *Embedding+Fair Clustering* (together with its unfair variant), and Algorithm 2.

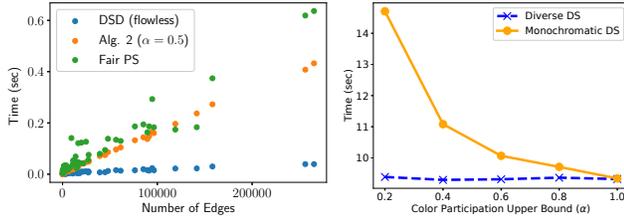


Figure 5: Left: Runtime as a function of the number of edges for all algorithms. Right: Runtime of Algorithm 2 as a function of the diversity parameter  $\alpha$ .

### 6.3 Evaluation of Algorithm 6 (for Problem 2)

We empirically study the accuracy and efficiency of Algorithm 6 to understand its practical performance. We apply our algorithm and the baseline *Prop2* on a part of the Amazon Product dataset, consisting of the graphs with at most 1,000 nodes and 5,000 edges. For each color  $c \in C$ , we set the lower bound as  $k_c = |V_c|/2$ . We also calculate the optimal solutions using *IP* and report the empirical approximation ratios of Algorithm 6 and *Prop2*.

The results are shown in Figure 6. In Figure 6 (left) we observe that in most graphs the solutions returned by Algorithm 6 are optimal or near-optimal, with empirical approximation ratios consistently higher than 0.95. This means that our proposed algorithm is of high accuracy in practice, despite the theoretical approximation guarantee being  $1/3$ . The baseline *Prop2* returns solutions with empirical approximation ratios greater than 0.8 but the performance is worse than that of ours. In Figure 6 (right) we show the running times of all algorithms applied. As expected, *Prop2* is the fastest as

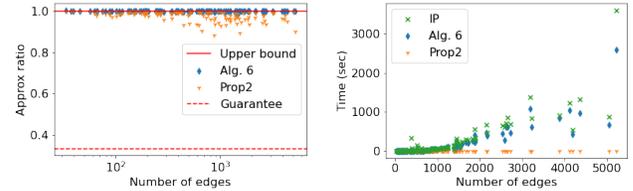


Figure 6: Empirical approximation ratios and running times of Algorithm 6 and baselines on the Amazon Product dataset with at most 1,000 nodes and 5,000 edges.

it solves only one LP. Algorithm 6 can scale to graphs in two colors with thousands of nodes, and it is faster than *IP* in general. Although from the figure the scalability of Algorithm 6 looks not much better than that of *IP*, we wish to remark that there are 7 graphs for which Algorithm 6 runs more than 10 times faster than *IP*. In the extreme example, where  $n = 30$  and  $m = 381$ , Algorithm 6 runs in 0.58s, while *IP* consumes 242.2s, meaning that Algorithm 6 runs more than 400 times faster than *IP*. We further discuss the scalability of Algorithm 6 particularly with respect to  $|C|$  in Appendix B.3.

## 7 CONCLUSIONS

In this work, we have focused on the problem of finding a densest diverse subgraph. We proposed novel formulations and approximation algorithms for two different notions of diversity. We performed various experiments on synthetic and real-world datasets, verifying that the densest subgraphs tend to be driven by homophily and that our tools provide the state-of-the-art methods.

Our work makes significant progress towards DSD with diversity and opens up several interesting problems. Can we improve the  $\Omega(1/\sqrt{n})$ -approximation for Problem 1 (in the case of  $V$  being diverse)? Can we design a better algorithm for Problem 2, in terms of both the approximation ratio and the runtime? Investigating the hardness of approximation is also an interesting direction.

## REFERENCES

- [1] Christoph Adami, Jifeng Qian, Matthew Rupp, and Arend Hintze. 2011. Information content of colored motifs in complex networks. *Artif. Life* 17, 4 (2011), 375–390.
- [2] Sara Ahmadian, Alessandro Epasto, Marina Knittel, Ravi Kumar, Mohammad Mahdian, Benjamin Moseley, Philip Pham, Sergei Vassilvitskii, and Yuyan Wang. 2020. Fair hierarchical clustering. In *NeurIPS '20*. 21050–21060.
- [3] Sara Ahmadian, Alessandro Epasto, Ravi Kumar, and Mohammad Mahdian. 2020. Fair correlation clustering. In *AISTATS '20*. 4195–4205.
- [4] Aris Anagnostopoulos, Luca Becchetti, Adriano Fazzone, Cristina Menghini, and Chris Schwiegelshohn. 2020. Spectral relaxations and fair densest subgraphs. In

- CIKM '20*, 35–44.
- [5] Reid Andersen and Kumar Chellapilla. 2009. Finding dense subgraphs with size bounds. In *WAW '09*, 25–37.
- [6] Arturs Backurs, Piotr Indyk, Krzysztof Onak, Baruch Schieber, Ali Vakilian, and Tal Wagner. 2019. Scalable fair clustering. In *ICML '19*, 405–413.
- [7] Bahman Bahmani, Ravi Kumar, and Sergei Vassilvitskii. 2012. Densest subgraph in streaming and MapReduce. In *Vldb '12*, 454–465.
- [8] Aditya Bhaskara, Moses Charikar, Eden Chlamtac, Uriel Feige, and Aravindan Vijayaraghavan. 2010. Detecting high log-densities: An  $O(n^{1/4})$  approximation for densest  $k$ -subgraph. In *STOC '10*, 201–210.
- [9] Digvijay Boob, Yu Gao, Richard Peng, Saurabh Sawlani, Charalampos E. Tsourakakis, Di Wang, and Junxing Wang. 2020. Flowless: Extracting densest subgraphs without flow computations. In *TheWebConf '20*, 573–583.
- [10] Sudhanshu Chanpuriya, Cameron Musco, Konstantinos Sotiropoulos, and Charalampos E. Tsourakakis. 2020. Node embeddings and exact low-rank representations of complex networks. In *NeurIPS '20*, 13185–13198.
- [11] Sudhanshu Chanpuriya, Cameron Musco, Konstantinos Sotiropoulos, and Charalampos E. Tsourakakis. 2021. On the power of edge independent graph models. In *NeurIPS '21*, 24418–24429.
- [12] Moses Charikar. 2000. Greedy approximation algorithms for finding dense components in a graph. In *APPROX '00*, 84–95.
- [13] Chandra Chekuri, Kent Quanrud, and Manuel R. Torres. 2022. Densest subgraph: Supermodularity, iterative peeling, and flow. In *SODA '22*, 1531–1555.
- [14] Tianyi Chen, Brian Matejek, Michael Mitzenmacher, and Charalampos E. Tsourakakis. 2022. Algorithmic tools for understanding the motif structure of networks. In *ECML PKDD '22*, 3–19.
- [15] Tianyi Chen and Charalampos Tsourakakis. 2022. AntiBenford subgraphs: Unsupervised anomaly detection in financial networks. In *KDD '22*, 2762–2770.
- [16] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. 2017. Fair clustering through fairlets. In *NIPS '17*, 5029–5037.
- [17] Abir De, Isabel Valera, Niloy Ganguly, Sourangshu Bhattacharya, and Manuel Gomez-Rodriguez. 2016. Learning and forecasting opinion dynamics in social networks. In *NIPS '16*, 397–405.
- [18] Yon Dourisboure, Filippo Geraci, and Marco Pellegrini. 2007. Extraction and classification of dense communities in the web. In *WWW '07*, 461–470.
- [19] Amita Gajewar and Atish Das Sarma. 2012. Multi-skill collaborative teams based on densest subgraphs. In *SDM '12*, 165–176.
- [20] Giorgio Gallo, Michael D. Grigoriadis, and Robert E. Tarjan. 1989. A fast parametric maximum flow algorithm and applications. *SIAM J. Comput.* 18, 1 (1989), 30–55.
- [21] David Gibson, Ravi Kumar, and Andrew Tomkins. 2005. Discovering large dense subgraphs in massive graphs. In *Vldb '05*, 721–732.
- [22] Aristides Gionis and Charalampos E. Tsourakakis. 2015. Dense subgraph discovery: KDD 2015 Tutorial. In *KDD '15*, 2313–2314.
- [23] Andrew V. Goldberg. 1984. *Finding a maximum density subgraph*. University of California Berkeley.
- [24] Julie Graber. 2020. Companies Hit With Lawsuits for Lack of Diversity at the Top. <https://www.linkedin.com/pulse/companies-hit-lawsuits-lack-diversity-top-julie-graber/>.
- [25] Qi Huangfu and Julian Hall. 2015. Parallelizing the dual revised simplex method. *Math. Program. Comput.* 10 (2015), 119–142.
- [26] Yasushi Kawase, Yuko Kuroki, and Atsushi Miyachi. 2019. Graph mining meets crowdsourcing: Extracting experts for answer aggregation. In *IJCAI '19*, 1272–1279.
- [27] Yasushi Kawase and Atsushi Miyachi. 2018. The densest subgraph problem with a convex/concave size function. *Algorithmica* 80, 12 (2018), 3461–3480.
- [28] Samir Khuller and Barna Saha. 2009. On finding dense subgraphs. In *ICALP '09*, 597–608.
- [29] Guy Kortsarz and Zeev Nutov. 2005. Approximating  $k$ -node connected subgraphs via critical graphs. *SIAM J. Comput.* 35, 1 (2005), 247–257.
- [30] Guy Kortsarz and David Peleg. 1994. Generating sparse 2-spanners. *J. Algorithms* 17, 2 (1994), 222–236.
- [31] Yuko Kuroki, Atsushi Miyachi, Junya Honda, and Masashi Sugiyama. 2020. Online dense subgraph discovery via blurred-graph feedback. In *ICML '20*, 5522–5532.
- [32] Tommaso Lanciano, Atsushi Miyachi, Adriano Fazzino, and Francesco Bonchi. 2023. A survey on the densest subgraph problem and its variants. *arXiv preprint arXiv:2303.14467* (2023).
- [33] Xiangfeng Li, Shenghua Liu, Zifeng Li, Xiaotian Han, Chuan Shi, Bryan Hooi, He Huang, and Xueqi Cheng. 2020. FlowScope: Spotting money laundering based on graphs. In *AAAI '20*, 4731–4738.
- [34] Brian Matejek, Donglai Wei, Tianyi Chen, Charalampos E. Tsourakakis, Michael Mitzenmacher, and Hanspeter Pfister. 2022. Edge-colored directed subgraph enumeration on the connectome. *Sci. Rep.* 12, 1 (2022), 11349.
- [35] Cameron Musco, Christopher Musco, and Charalampos E. Tsourakakis. 2018. Minimizing polarization and disagreement in social networks. In *TheWebConf '18*, 369–378.
- [36] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *EMNLP-IJCNLP '19*, 188–197.
- [37] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *KDD '14*, 701–710.
- [38] Jean-Claude Picard and Maurice Queyranne. 1982. A network flow solution to some nonlinear 0-1 programming problems, with applications to graph theory. *Networks* 12, 2 (1982), 141–159.
- [39] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. 2018. Network embedding as matrix factorization: Unifying DeepWalk, LINE, PTE, and node2vec. In *WSDM '18*, 459–467.
- [40] Syama Sundar Rangapuram, Thomas Bühler, and Matthias Hein. 2013. Towards realistic team formation in social networks based on densest subgraphs. In *WWW '13*, 1077–1088.
- [41] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. 2021. Multi-scale attributed node embedding. *J. Complex Netw.* 9, 2 (2021), 1–22.
- [42] Benedek Rozemberczki and Rik Sarkar. 2020. Characteristic functions on graphs: Birds of a feather, from statistical descriptors to parametric models. In *CIKM '20*, 1325–1334.
- [43] Mauro Sozio and Aristides Gionis. 2010. The community-search problem and how to plan a successful cocktail party. In *KDD '10*, 939–948.
- [44] Michele Starnini, Charalampos E. Tsourakakis, Maryam Zamanipour, André Panisson, Walter Allasia, Marco Fornasiero, Laura Li Puma, Valeria Ricci, Silvia Ronchiadin, Angela Ugrinoska, Marco Varetto, and Dario Moncalvo. 2021. Smurf-based anti-money laundering in time-evolving transaction networks. In *ECML PKDD '21*.
- [45] Amanda L. Traud, Peter J. Mucha, and Mason A. Porter. 2012. Social structure of Facebook networks. *Physica A* 391, 16 (2012), 4165–4180.
- [46] Charalampos E. Tsourakakis. 2015. The  $k$ -clique densest subgraph problem. In *WWW '15*, 1122–1132.
- [47] Charalampos E. Tsourakakis, Francesco Bonchi, Aristides Gionis, Francesco Gullo, and Maria Tsiarli. 2013. Denser than the densest subgraph: Extracting optimal quasi-cliques with quality guarantees. In *KDD '13*, 104–112.
- [48] Nate Veldt, Austin R. Benson, and Jon Kleinberg. 2021. The generalized mean densest subgraph problem. In *KDD '21*, 1604–1614.

## A PROOF OF LEMMA 3

First we consider the case where  $y_c \geq 1$  for every  $c \in C$ . Let  $a$  be the maximum number that satisfies  $C_{\text{sat}}(S(a)) = C$ . Let  $b$  be the infimum of the numbers  $b$  that satisfy  $C_{\text{sat}}(S(b)) = \emptyset$ . Let  $y_{\text{max}}^* = \max_{v \in V} y_v^*$ . Note that  $0 \leq a \leq b \leq y_{\text{max}}^*$  hold.

To prove the lemma, it suffices to show that (at least) one of the following cases occurs:

**Case (i)** There exists some  $r' \leq a$  such that  $d(S(r')) \geq \frac{\lambda}{3}$ ;

**Case (ii)** There exists some  $a < r' \leq b$  such that  $\frac{|E(S(r'))|}{\sum_{c \in C_{\text{sat}}(S(r'))} |S(r')_c| + \sum_{c \in C \setminus C_{\text{sat}}(S(r'))} k_c} \geq \lambda/3$ ;

**Case (iii)** There exists some  $r' > b$  such that  $|E(S(r'))| \geq \frac{|E(S^*)|}{3}$ .

To show that, suppose that none of the above cases occurs. We define indicator functions  $Z_v(r) : [0, y_{\text{max}}^*] \rightarrow \{0, 1\}$  for  $v \in V$  and  $Z_e(r) : [0, y_{\text{max}}^*] \rightarrow \{0, 1\}$  for  $e = \{u, v\} \in E$  as follows:

$$Z_v(r) = \begin{cases} 1 & \text{if } r \leq y_v^*, \\ 0 & \text{otherwise,} \end{cases} \quad Z_e(r) = \begin{cases} 1 & \text{if } r \leq \min\{y_u^*, y_v^*\}, \\ 0 & \text{otherwise.} \end{cases}$$

Since Case (i) does not hold, for any  $r' \leq a$ , we have  $d(S(r')) = \frac{|E(S(r'))|}{|S(r')|} < \frac{\lambda}{3}$ . Thus, we have

$$\begin{aligned} \int_0^a |E(S(r))| dr &< \frac{\lambda}{3} \int_0^a |S(r)| dr \\ &= \frac{\lambda}{3} \int_0^a \sum_{v \in V} Z_v(r) dr = \frac{\lambda}{3} \sum_{v \in V} \int_0^a Z_v(r) dr. \end{aligned}$$

Since Case (ii) does not occur, for any  $a < r' \leq b$ , we have  $|E(S(r'))| < \frac{\lambda}{3} \left( \sum_{c \in C_{\text{sat}}(S(r'))} |S(r')_c| + \sum_{c \in C \setminus C_{\text{sat}}(S(r'))} k_c \right)$ . Hence, we have

$$\int_a^b |E(S(r))| dr < \frac{\lambda}{3} \int_a^b \left( \sum_{c \in C_{\text{sat}}(S(r))} |S(r)_c| + \sum_{c \in C \setminus C_{\text{sat}}(S(r))} k_c \right) dr.$$

Now we see that

$$\begin{aligned} \int_a^b \sum_{c \in C_{\text{sat}}(S(r))} |S(r)_c| dr &\leq \int_a^b \sum_{c \in C} |S(r)_c| dr \\ &= \int_a^b \sum_{c \in C} \sum_{v \in V_c} Z_v(r) dr = \sum_{c \in C} \sum_{v \in V_c} \int_a^b Z_v(r) dr \\ &= \sum_{v \in V} \int_a^b Z_v(r) dr \end{aligned}$$

and

$$\int_a^b \sum_{c \in C \setminus C_{\text{sat}}(S(r))} k_c dr \leq \int_a^b \sum_{c \in C} p_c^* dr \leq \|\mathbf{p}^*\|_1 \cdot y_{\text{max}}^* \leq 1.$$

Therefore, we have

$$\int_a^b |E(S(r))| dr < \frac{\lambda}{3} \left( \sum_{v \in V} \int_a^b Z_v(r) dr + 1 \right).$$

Since Case (iii) does not hold, for any  $r' > b$ , we have  $|E(S(r'))| < \frac{|E(S^*)|}{3}$ . Thus, we have

$$\begin{aligned} \int_b^{y_{\text{max}}^*} |E(S(r))| dr &< \frac{|E(S^*)|}{3} \int_b^{y_{\text{max}}^*} dr \\ &\leq \frac{|E(S^*)|}{3} y_{\text{max}}^* \leq \frac{1}{3} \frac{|E(S^*)|}{\|\mathbf{p}^*\|_1} \leq \frac{\lambda}{3}, \end{aligned}$$

where the last inequality follows from Lemma 2. Thus, we have

$$\begin{aligned} \int_0^{y_{\text{max}}^*} |E(S(r))| dr &= \int_0^a |E(S(r))| dr + \int_a^b |E(S(r))| dr + \int_b^{y_{\text{max}}^*} |E(S(r))| dr \\ &< \frac{\lambda}{3} \left( \sum_{v \in V} \int_0^a Z_v(r) dr + \left( \sum_{v \in V} \int_a^b Z_v(r) dr + 1 \right) + 1 \right) \\ &\leq \frac{\lambda}{3} \left( \sum_{v \in V} y_v^* + 2 \right) = \frac{\lambda}{3} \left( \sum_{c \in C} \sum_{v \in V_c} y_v^* + 2 \right) = \lambda. \end{aligned}$$

On the other hand, by a simple calculation, we have

$$\begin{aligned} \int_0^{y_{\text{max}}^*} |E(S(r))| dr &= \int_0^{y_{\text{max}}^*} \sum_{e \in E} Z_e(r) dr \\ &= \sum_{e \in E} \int_0^{y_{\text{max}}^*} Z_e(r) dr = \sum_{e \in E} \min\{y_u^*, y_v^*\} \geq \sum_{e \in E} x_e^* = \lambda, \end{aligned}$$

which contradicts to the above inequality.

Next we consider the case where  $y_c = 0$  for some  $c \in C$ . Defining  $a, b$ , and  $y_{\text{max}}^*$  in the same way, we have  $0 \leq a \leq y_{\text{max}}^* < b = \infty$ . It suffices to show that (at least) one of the following cases occurs:

**Case (i)** There exists some  $r' \leq a$  such that  $d(S(r')) \geq \frac{\lambda}{3}$ ;

**Case (ii)** There exists some  $a < r' \leq y_{\text{max}}^*$  such that  $\frac{|E(S(r'))|}{\sum_{c \in C_{\text{sat}}(S(r'))} |S(r')_c| + \sum_{c \in C \setminus C_{\text{sat}}(S(r'))} k_c} \geq \lambda/3$ .

The proof is similar to the above and thus omitted.  $\square$

## B SUPPLEMENTARY MATERIAL FOR EXPERIMENTS

### B.1 DBLP Co-Authorship Dataset

We create this dataset from all authors that have published at least 3 papers in the following conferences between 2003 and 2022:

- **Theory:** COLT, FOCS, ICALP, SODA, STOC
- **Data Management:** CIDR, ICDE, PODS, SIGMOD, VLDB
- **Data Mining:** CIKM, ICDM, KDD, WSDM, WWW
- **Learning:** AAAI, ICLR, ICML, NeurIPS
- **Networking:** ICC, IMC, INFOCOM, MOBICOM, SIGCOMM
- **Image & Video Processing:** CVPR, ECCV, ICCV

### B.2 IP for Problem 2

Let us introduce a 0–1 variable  $x_e$  for each  $e \in E$  and 0–1 variable  $y_v$  for each  $v \in V$ . Then for each  $k_{\text{guess}} \in \{\|\mathbf{k}\|_1, \dots, n\}$ , we construct

the following (mixed) integer linear programming problem:

$$\begin{aligned}
& \text{maximize} && \sum_{e \in E} x_e / k_{\text{guess}} \\
& \text{subject to} && \sum_{v \in V} y_v = k_{\text{guess}}, \\
& && \sum_{v \in V_c} y_v \geq k_c \quad \forall c \in C, \\
& && x_e \leq y_u, x_e \leq y_v \quad \forall e = \{u, v\} \in E, \\
& && 0 \leq x_e \leq 1 \quad \forall e \in E, \\
& && y_v \in \{0, 1\} \quad \forall v \in V.
\end{aligned}$$

It is easy to see that this problem is equivalent to Problem 2 with size constraint  $|S| = k_{\text{guess}}$ . Note that the 0–1 constraints for  $x_e$ 's are relaxed because they are redundant.

The baseline *IP* computes an optimal solution to Problem 2 as follows: Solve the above integer programming problem for all  $k_{\text{guess}} \in \{\|\mathbf{k}\|_1, \dots, n\}$  and return a solution with the maximum density value.

### B.3 Scalability of Algorithm 6 on Synthetic Data

We run Algorithm 6 on a set of Erdős-Rényi graphs with size  $n$  ranged in  $\{18, 54, 90, 126\}$ , edge probability  $p = 5/n$ , and the number

of colors  $|C|$  ranged in  $\{2, 3, 6\}$ . Nodes are evenly split into colors, and the lower bound is set as  $k_c = \lfloor \frac{n}{2|C|} \rfloor$  for each  $c \in C$ . Random graphs are sampled ten times for each setting, and the averaged running times are reported in Figure 7. When we have only two colors, Algorithm 6 can scale up to a thousand nodes as shown in the main text. However, the algorithm becomes rather inefficient as the graph size goes beyond that or the number of colors increases.

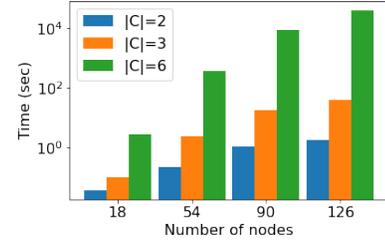


Figure 7: Running time of Algorithm 6.