



Towards Cohesion-Fairness Harmony: Contrastive Regularization in Individual Fair Graph Clustering

Siamak Ghodsi¹(✉) , Seyed Amjad Seyedi² , and Eirini Ntoutsi³

¹ L3S Research Centre, Leibniz University Hannover, Hannover, Germany
ghodsi@l3s.de

² University of Kurdistan, Sanandaj, Iran
amjadseyedi@uok.ac.ir

³ RI CODE, University of the Bundeswehr Munich, Munich, Germany
eirini.ntoutsi@unibw.de

Abstract. Conventional fair graph clustering methods face two primary challenges: i) They prioritize balanced clusters at the expense of cluster cohesion by imposing rigid constraints, ii) Existing methods of both individual and group-level fairness in graph partitioning mostly rely on eigen decompositions and thus, generally lack interpretability. To address these issues, we propose *iFairNMTF*, an individual Fairness Nonnegative Matrix Tri-Factorization model with contrastive fairness regularization that achieves balanced and cohesive clusters. By introducing fairness regularization, our model allows for customizable accuracy-fairness trade-offs, thereby enhancing user autonomy without compromising the interpretability provided by nonnegative matrix tri-factorization. Experimental evaluations on real and synthetic datasets demonstrate the superior flexibility of *iFairNMTF* in achieving fairness and clustering performance.

Keywords: Fair Graph Clustering · Fair-Nonnegative Matrix Factorization · Fair Unsupervised Learning · Individual Fairness

1 Introduction

Graph-structured data is ubiquitous in various real-world applications including recommender systems, e-commerce, social networks, and neural networks. Graph clustering is essential for identifying meaningful patterns within graphs. Despite the advancements in algorithmic fairness for supervised learning scenarios which are mostly tailored for independent and identically distributed (i.i.d.) data [20], the topic of fairness is less explored in the unsupervised learning domain and especially for graphs. A motivating example comes from the educational domain [22]: how to divide students in a classroom into smaller groups for collaborative assignments. It is demanded to diversify group members from different genders or races while respecting existing friendship networks and maintaining connections. Graphs comprise non-i.i.d. data; thus, the broad literature

on fairness for i.i.d. data is generally not applicable to graphs [6]. However, some approaches mitigate bias by converting graph data into tabular form and leveraging existing methods. Additionally, there exist bias mitigation approaches that transform tabular data into hypergraphs based on dataset similarities, e.g. [8].

In the realm of fairness in i.i.d. clustering, the pioneering work of [3] introduced balance score, a fairness measure rooted in statistical parity [7], aiming at clusters of balanced demographic subgroups given a sensitive feature. Inspired by this work, [13] proposed a spectral graph clustering (SC) framework promoting group fairness that was later extended in [26] to scaled networks. However, there is not much literature on clustering with individual fairness, which prioritizes treating similar individuals (nodes in our context) similarly. A spectral model based on PageRank was proposed in [12] introducing a notion of individual fairness but for supervised node-classification tasks, whereas [27] introduces an individual-fair model for multi-view graph clustering. Only in [10], an (unsupervised) graph partitioning method employed an individual fairness approach which constrains a spectral clustering with a representation graph constructed solely based on sensitive information of individuals. SC methods are based on minimizing either the Ratio-cut or the Normalized-cut heuristic that generally tend to minimize the number of links pointing outside each cluster [17]. These cut-based heuristics do not guarantee to discover the optimal graph partitioning. Thus, incorporating (hard) fairness constraints into these rigid frameworks, which is usually also not a trivial and straightforward process, makes achieving the optimal solution challenging such that usually a relaxed form of the problem is being solved as in [10, 13]. In addition, since the solution to these hard-constrained spectral approaches is based on the eigen-decomposition of the graph, it lacks interpretability.

To address the identified issues, we introduce a versatile fairness-aware model for graph clustering, the so-called individually-Fair Symmetric Nonnegative Matrix Tri-Factorization (*iFairNMTF*) model with contrastive regularization. Building on the symmetric NMF [14, 16], a model tailored for graph clustering, the NMTF [21] extends its capabilities inheriting its intrinsic interpretability through non-negativity and direct clustering, while other models require steps like graph and/or node embedding [4, 5], representation learning [25], or eigen-decomposition [12, 13, 26] before performing the final clustering. Additionally, NMTF provides better clustering and also introduces an explicit interpretability factor for inter-cluster interactions. We integrate these capabilities with a novel soft individual fairness regularization in *iFairNMTF* with an adjustable parameter λ for balancing both fairness and clustering objectives. Our key contributions include: i) A flexible joint learning framework with adjustable fairness regularization, accommodating customization of fairness enforcement in relation to clustering quality. The framework supports the linear integration of fairness and other problem-specific constraints via a customizable cost function. ii) Introduction of a contrastive fairness regularization, promoting the distribution of similar individuals across clusters based on sensitive attribute membership while ensuring distinct representation of dissimilar individuals within each cluster. iii) Reten-

tion of SNMF advantages, providing an interpretable data representation due to non-negativity and direct clustering. iv) Integration of an explicit interpretability factor, exposing inter-cluster relationships. v) Extensive experiments demonstrating the efficacy of our model with soft-fairness constraints and emphasizing the significance of the adjustable trade-off optimization.

To the best of our knowledge, our proposed joint learning contrastive framework is the first attempt to integrate an NMF model into a fairness-aware learning framework. The rest of this paper is organized as follows: In Sect. 2, we review related work. Our method is introduced in Sect. 3. The experimental evaluation is presented in Sect. 4. Conclusions and outlook are discussed in Sect. 5.

2 Background and Related Works

Problem Formulation. Let us assume an undirected graph $\mathcal{G} = (V, E)$ where $V = \{v_1, v_2, \dots, v_n\}$ is the set of n nodes and $E \subseteq V \times V$ is the set of edges. The adjacency matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ encodes the edge information; the existence or non-existence of an edge between two nodes v_i, v_j is modeled as $a_{ij} = 1$ and $a_{ij} = 0$, respectively. Also, we assume no self-loops (edge connecting a node to itself), so $a_{ii} = 0$ for all $i \in [n]$. Let us further assume that the set of vertices constitutes m disjoint groups identified based on a sensitive attribute e.g., gender or race, such that $V = \dot{\cup}_{s \in [m]} V_s$. The goal is to find a non-overlapping clustering of V into $k \geq 2$ clusters $V = \{C_1 \dot{\cup} \dots \dot{\cup} C_k\}$ which is subject to individual fairness.

Individual Fairness. Individual fairness primarily formalized in [7] identifies a model f to be fair if, for any pair of inputs v_i, v_j which are sufficiently close (as per an appropriate metric), the model outputs $f(v_i), f(v_j)$ should also be close (as per another appropriate metric). In other terms, pairwise node distances in the input space and output space should satisfy the Lipschitz continuity Condition. Specifically, it requires the distance of any node pairs in the output space to be smaller or equal to their corresponding distance in the input space (usually re-scaled by a scalar). Given a pair of nodes v_i and v_j , the Lipschitz condition is:

$$D(f(v_i), f(v_j)) \leq L \cdot d(v_i, v_j) \quad (1)$$

where $f(\cdot)$ is the predictive model producing the node-level outputs (e.g., embeddings). $D(\cdot, \cdot)$ and $d(\cdot, \cdot)$ are the distance metrics of output and input space and L is the Lipschitz constant that re-scales the input distance between nodes v_i, v_j . In order to measure individual fairness based on L , [29] proposed consistency on non-graph data with the intuition to measure the average distance of the output between each individual and its k -nearest neighbors such that:

$$1 - \frac{1}{n \cdot k} \sum_{i=1}^n \left| f(x_i) - \sum_{j \in kNN(x_i)} f(x_j) \right| \quad (2)$$

where $f(x_i)$ is the probabilistic classification output for node features x_i of node v_i and $kNN()$ is the neighborhood of node v_i . In general, a larger average distance indicates a lower level of individual fairness.

Individual Fairness for Graph Clustering. The notion of individual fairness in graph mining [6] can be divided into three categories by application: i) node pair distance-based fairness, ii) node ranking-based fairness, and iii) individual fairness in graph clustering. The core idea in the first category is the investigation of achieving individual fairness in node representation and node embedding problems based on pairwise node distances. For example, in [15] a notion of consistency is proposed based on a similarity matrix S that characterizes node similarity in input space and can be derived from node attributes, graph topology, or domain experts. Moreover, in [12] a measure is proposed that calculates the similarity-weighted output discrepancy between nodes to measure unfairness. This metric calculates the weighted sum of pairwise node distance in the output space, where the weighting score is the pairwise node similarity. Hence for any graph mining algorithm, a smaller value of the similarity-weighted discrepancy typically implies a higher level of individual fairness.

The second category aims to achieve individual fairness by establishing node rankings. This involves creating two ranking lists in the input and output space, R_1 and R_2 based on a pairwise similarity matrix S in the input space. The satisfaction of individual fairness is determined by the alignment of these ranking lists, ensuring that R_1 and R_2 are identical for each individual [5].

The third category which remains relatively less explored and is the focus of our work, surveys individual-level fairness for graph clustering. In essence, if all neighbors of each node in a graph, are proportionally distributed to each cluster, individual fairness is then fulfilled [9]. One of the pioneering recent works in this direction is the work of Gupta, et.al., [10] according to which a clustering algorithm satisfies individual fairness for node v_i if:

$$\frac{|\{v_j : \mathbf{A}_{i,j} = 1 \wedge v_j \in C_k\}|}{|C_k|} = \frac{|\{v_j : \mathbf{A}_{i,j} = 1\}|}{|V|} \quad (3)$$

for all clusters C_k . The key intuition is that for each node, the ratio occupied by its one-hop neighbors in its cluster should be the same as the ratio occupied by its one-hop neighbors in the entire population (i.e. the main graph).

3 The iFairNMTF Model

Inspired by the individual fairness of [10], we propose a novel individual fairness regularization for graph clustering. It constitutes a contrastive graph regularization that incorporates positive and negative elements, signifying the attraction and repulsion of individuals towards their similar and dissimilar neighbors, based on a sensitive (node) attribute. By integrating this regularization into a flexible clustering framework, we introduce a unique **individually Fair Non-negative Matrix Tri-Factorization** joint learning model (iFairNMTF).

3.1 The iFairNMTF Model Formulation

Symmetric NMF (SNMF) [14] is an extension of the traditional NMF that transforms it into a versatile graph clustering model. This model factorizes an adjacency matrix $\mathbf{A} \in \mathbb{R}_+^{n \times n}$ and is based on the assumption that similar samples

($A_{ij} > 0$) should have similar representations ($\mathbf{h}_i \mathbf{h}_j^\top > 0$) and dissimilar samples ($A_{ij} = 0$) should have opposite representations ($\mathbf{h}_i \mathbf{h}_j^\top = 0$), where H can be interpreted as the *node-to-cluster membership matrix*. More formally:

$$\min_{H \geq 0} \|\mathbf{A} - \mathbf{H}\mathbf{H}^\top\|_F^2, \quad (4)$$

An extended form of the SNMF is the SNM-Tri-Factorization [21] (we omit the “S” and refer NMTF hereafter) which has been tailored to address graph clustering tasks [1, 11]. It takes into account the *cluster-cluster interactions matrix* using an additional factor \mathbf{W} such that, $A_{ij} \approx \mathbf{h}^{(i)} \mathbf{W} \mathbf{h}^{(j)\top}$. More formally:

$$\min_{H, W \geq 0} \|\mathbf{A} - \mathbf{H}\mathbf{W}\mathbf{H}^\top\|_F^2, \quad (5)$$

where \mathbf{W} can be interpreted as the cluster interactions. We build upon this model and extend it into an individual fairness joint learning framework using a contrastive regularization. More formally:

$$\min_{H, W \geq 0} \|\mathbf{A} - \mathbf{H}\mathbf{W}\mathbf{H}^\top\|_F^2 + \lambda \mathcal{R}_C(\mathbf{H}), \quad (6)$$

Schematically, the model block diagram is illustrated in Fig. 1. The left term comes from Eq. (5) and $\mathcal{R}_C(\mathbf{H})$ is a contrastive regularization constraining the cluster indicator \mathbf{H} relatively adjusted by the magnitude of a flexible λ parameter ensuring its alignment with group demographics. The contrastive term $\mathbf{C} = \mathbf{P} - \mathbf{N}$ consists of a positive and a negative component:

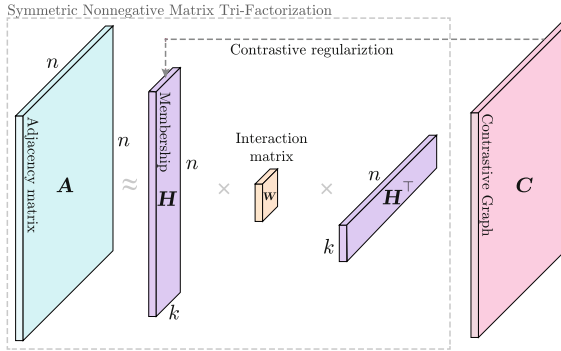


Fig. 1. Schematic representation of the iFairNMTF model with contrastive regularization.

$$\begin{aligned} \mathcal{N}_{i,j} &= \begin{cases} 1, & \text{if } g_i = g_j \\ 0, & \text{otherwise.} \end{cases} & \mathcal{P}_{i,j} &= \begin{cases} 1, & \text{if } g_i \neq g_j \\ 0, & \text{otherwise.} \end{cases} \\ N_{ij} &= \mathcal{N}_{ij} / \sum_{r=1}^n \mathcal{N}_{ir}, & P_{ij} &= \mathcal{P}_{ij} / \sum_{r=1}^n \mathcal{P}_{ir}, \end{aligned} \quad (7)$$

which can be enforced to apply the attraction of different demographic groups into the same cluster, and repulsion of same-group members to ensure diversity of their distribution into different clusters according to Eq. (8):

$$\min_{\mathbf{H}} \mathcal{R}_C = \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{h}^{(i)} - \mathbf{h}^{(j)}\|^2 C_{ij} = \text{Tr}(\mathbf{H}^\top \mathbf{L} \mathbf{H}). \quad (8)$$

where $\mathbf{L} = \mathbf{D} - \mathbf{C}$ is the graph Laplacian and $D_{ii} = \sum_{j=1}^n C_{ij}$. By adding the contrastive regularization \mathcal{R}_C to the NMTF (5), we derive the final objective function (loss function) of iFairNMTF, $\mathcal{L} = \mathcal{L}_{\mathcal{F}} + \lambda \mathcal{R}_C$ as follow:

$$\min_{\mathbf{H}, \mathbf{W} \geq 0} \|\mathbf{A} - \mathbf{H} \mathbf{W} \mathbf{H}^\top\|_F^2 + \lambda \text{Tr}(\mathbf{H}^\top \mathbf{L} \mathbf{H}), \quad (9)$$

The objective function in Eq. (9) is a combination, trading-off between the clustering loss and the contrastive regularization term to ensure individual fairness. The hyper-parameter $\lambda \in [0, +\infty)$ controls the compromise between clustering performance and fairness. Smaller λ implies a higher importance of the clustering performance and prompts the model to prioritize generating strong and cohesive clusters. Conversely, a higher λ prioritizes fairness, prompting the model to create diversified clusters that fairly represent groups of V_s .

3.2 The iFairNMTF Model Optimization

In this section, we focus on solving the iFairNMTF model. The objective function in Eq. (9) is a fourth-order non-convex function with respect to the entries of \mathbf{H} and has multiple local minima. For these types of problems, it is difficult to find a global minimum; thus a good convergence property we can expect is that every limit point is a stationary point. Therefore, we adopt multiplicative updating rules to update the membership matrix \mathbf{H} and introduce two Lagrangian multiplier matrices of $\boldsymbol{\Theta}$, and $\boldsymbol{\Phi}$ to enforce the nonnegative constraints on \mathbf{H} , and \mathbf{W} respectively, resulting in the following equivalent objective function:

$$\min_{\mathbf{H}, \mathbf{W}} \mathcal{L} = \|\mathbf{A} - \mathbf{H} \mathbf{W} \mathbf{H}^\top\|_F^2 + \lambda \text{Tr}(\mathbf{H}^\top \mathbf{L} \mathbf{H}) - \text{Tr}(\boldsymbol{\Theta}^\top \mathbf{H}) - \text{Tr}(\boldsymbol{\Phi}^\top \mathbf{W}),$$

which can be further rewritten as follows:

$$\begin{aligned} \min_{\mathbf{H}, \mathbf{W}} \mathcal{L} = & \text{Tr}(\mathbf{A}^\top \mathbf{A} - 2\mathbf{A}^\top \mathbf{H} \mathbf{W} \mathbf{H}^\top + \mathbf{H} \mathbf{W}^\top \mathbf{H}^\top \mathbf{H} \mathbf{W} \mathbf{H}^\top) \\ & + \lambda \text{Tr}(\mathbf{H}^\top \mathbf{L} \mathbf{H}) - \text{Tr}(\boldsymbol{\Theta}^\top \mathbf{H}) - \text{Tr}(\boldsymbol{\Phi}^\top \mathbf{W}). \end{aligned} \quad (10)$$

The partial derivative of \mathcal{L} with respect to \mathbf{H} is

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{H}} = & -2\mathbf{A}^\top \mathbf{H} \mathbf{W} - 2\mathbf{A} \mathbf{H} \mathbf{W}^\top + 2\mathbf{H} \mathbf{W}^\top \mathbf{H}^\top \mathbf{H} \mathbf{W} \\ & + 2\mathbf{H} \mathbf{W} \mathbf{H}^\top \mathbf{H} \mathbf{W}^\top + 2\lambda \mathbf{L} \mathbf{H} - \boldsymbol{\Theta}. \end{aligned} \quad (11)$$

Algorithm 1. Individual Fair Nonnegative Matrix Tri-Factorization (iFairNMTF)**Input:** adjacency matrix \mathbf{A} , group set g , latent factor k , trade-off parameter λ ;**Output:** cluster assignment M ;

- 1: Construct the contrastive graph \mathbf{C} according to (7);
- 2: **while** convergence not reached **do**
- 3: Update cluster-membership matrix \mathbf{H} according to (13);
- 4: Update cluster-interaction matrix \mathbf{W} according to (16);
- 5: **end while**
- 6: Calculate cluster assignment $M_i \leftarrow \arg \max(\mathbf{h}^{(i)}), \forall i \in \{1, \dots, n\}$
- 7: **return** cluster-membership matrix \mathbf{H} and cluster-interaction matrix \mathbf{W} ;

By setting the partial derivative $\frac{\partial \mathcal{L}}{\partial \mathbf{H}}$ to 0, we have:

$$\Theta = -2\mathbf{A}^\top \mathbf{H} \mathbf{W} - 2\mathbf{A} \mathbf{H} \mathbf{W}^\top + 2\mathbf{H} \mathbf{W}^\top \mathbf{H}^\top \mathbf{H} \mathbf{W} + 2\mathbf{H} \mathbf{W} \mathbf{H}^\top \mathbf{H} \mathbf{W}^\top + 2\lambda \mathbf{L} \mathbf{H}. \quad (12)$$

From the Karush-Kuhn-Tucker complementary slackness conditions (KKT), we obtain $\mathbf{H} \odot \Theta = \mathbf{0}$ where \odot denotes the element-wise product. This is the fixed point equation that the solution must satisfy at convergence. By solving this equation, we derive the following updating rule for \mathbf{H} :

$$\mathbf{H} \leftarrow \mathbf{H} \odot \left(\frac{\mathbf{A}^\top \mathbf{H} \mathbf{W} + \mathbf{A} \mathbf{H} \mathbf{W}^\top + \lambda \mathbf{L}^\top \mathbf{H}}{\mathbf{H} \mathbf{W}^\top \mathbf{H}^\top \mathbf{H} \mathbf{W} + \mathbf{H} \mathbf{W} \mathbf{H}^\top \mathbf{H} \mathbf{W}^\top + \lambda \mathbf{L} + \mathbf{H}} \right)^{\frac{1}{4}}. \quad (13)$$

To guarantee the nonnegativity, we separate the positive and negative elements as $\mathbf{L} = \mathbf{L}^+ - \mathbf{L}^-$. Similarly, we differentiate \mathcal{L} with respect to \mathbf{W} such that:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = -2\mathbf{H}^\top \mathbf{A} \mathbf{H} + 2\mathbf{H}^\top \mathbf{H} \mathbf{W} \mathbf{H}^\top \mathbf{H} - \Phi \quad (14)$$

By setting the partial derivative $\frac{\partial \mathcal{L}}{\partial \mathbf{W}}$ to 0, we obtain Φ as:

$$\Phi = -2\mathbf{H}^\top \mathbf{A} \mathbf{H} + 2\mathbf{H}^\top \mathbf{H} \mathbf{W} \mathbf{H}^\top \mathbf{H}. \quad (15)$$

From the complementary slackness KKT conditions we obtain $\mathbf{W} \odot \Phi = \mathbf{0}$. This is another fixed point equation that the solution must satisfy at convergence. Finally, by solving this equation, we derive the following updating rule for \mathbf{W} :

$$\mathbf{W} \leftarrow \mathbf{W} \odot \frac{\mathbf{H}^\top \mathbf{A} \mathbf{H}}{\mathbf{H}^\top \mathbf{H} \mathbf{W} \mathbf{H}^\top \mathbf{H}}. \quad (16)$$

4 Experimental Evaluation

4.1 Experimental Setup

Datasets. In the paper, six real-world and three synthetic networks are used for benchmarking the performance of the proposed method against competitors. Our synthetic networks are generated according to a generalized Stochastic Block Model (SBM) [13] with equal-sized clusters $|C_l| = n/k$ and groups

$|V_s| = n/g$ randomly distributed among the clusters. We generate three SBM networks of 2K, 5K, and 10K nodes with $k = 5$ clusters and $g = 5$ groups. Real datasets include three high school friendship networks [18]: *Facebook*, *Friendship*, and *Contact-Diaries* which represent connections among a group of French high school students. *DrugNet* [28] is a network encoding acquaintanceship between drug users in Hartford, CT. *LastFMNet* [24] contains mutual follower relations among users of Last.fm, a recommendation-based online radio and music community in Asia. Lastly, *NBA* is a network containing relationships between around 400 NBA basketball players [4]. A detailed description of both real and synthetic datasets, as well as instructions on generating the SBM networks are provided in the supplementary material¹. Likewise for dataset statistics including size and number of sensitive groups and also details on cleaning the real datasets.

Competitors. We compare *iFairNMTF* with four state-of-the-art graph clustering methods, namely, with two group-fair models: i) *Fair-SC* [13], and its scalable version (ii) *sFair-SC* [26], iii) an individual-fairness model (*iFair-SC*) [10] and iv) a deep graph neural network (*DMoN*) [25]. The three former models are fairness-aware and have been already discussed. The latter model is one of the very few DNNs developed for pure graph-partitioning problems, but does not consider fairness. This model extends the general graph neural network (GNN) architecture into a deep modularity optimization GNN. It operates on attributed graphs, thus we pass the sensitive attribute as node-attribute to it. The number of layers and learning rate are set according to the official source code provided by the authors (layers = 64 or 512 for small and large networks, $\alpha = 0.001$). The number of epochs for *DMoN* and our method is 500. To produce reliable results, all experiments are averaged over 10 independent runs.

Evaluation Measures. We use accuracy for measuring clustering assignment quality on synthetic networks. For real-world networks, since the ground truth cluster structures are unknown, we use Newman’s modularity (Q) measure [2, 19] which analyzes the homogeneity of clusters by calculating the proportion of internal links in each cluster for a given partitioning compared to the expected proportion of edges in a null graph with the same degree distribution. Modularity is preferable over cut-based measures due to its robustness against imbalanced cluster sizes. We measure the fairness of clustering in terms of the popular average balance (B) measure [13, 26]: $B = \frac{1}{k} \sum_{l=1}^k \text{Balance}(C_l)$, where $\text{Balance}(C_l)$ calculates the minimum group proportion of C_l according to Eq. (17):

$$\text{Balance}(C_l) = \min_{s \neq s' \in [m]} \frac{|V_s \cap C_l|}{|V_{s'} \cap C_l|}, \quad (17)$$

where $l \in [1, k]$ iterates over all the k clusters and V_s identifies each sensitive group of the sensitive attribute. The minimum balance of each cluster can range between $[0, 1]$, thus their average also ranges between $[0, 1]$.

Parameters. Our model has an adjustable hyper-parameter λ to trade-off between the degree of fairness and clustering efficiency (Eq. (9)). The range of λ includes 50 values from $[0, 100]$ with a median of 3 for small and from

¹ Link to supplemental file and source codes: [Github.com/SiamakGhodsiiFairNMTF](https://github.com/SiamakGhodsiiFairNMTF).

$[0, 3500]$ for large datasets (must be set separately for each dataset.). The effect of λ is discussed in Sect. 4.3. The trade-off parameter λ can be set based on user preferences between fairness and clustering quality. A practical way is to select the best value according to the intersection point of B and Q , see Fig. 3.

4.2 Clustering Quality vs Fairness

In real datasets, the ground truth partitioning of the networks is unknown, therefore we report the performance for various number of clusters. Figure 2 illustrates the comparison of our method’s results in terms of Q (clustering quality/ modularity) and B (fairness/balance) with those of other models on two datasets, for various numbers of clusters. Dataset balance, highlighted by the yellow dashed line, identifies the proportion of the smallest to the largest group of the sensitive attribute, calculated according to Eq. (17). For iFairNMTF, the best λ values for each k are used. They are selected based on the intersection of Q and B charts as in Fig. 3: $\lambda = 2$ for DrugNet, and $\lambda = 100$ for LastFM.

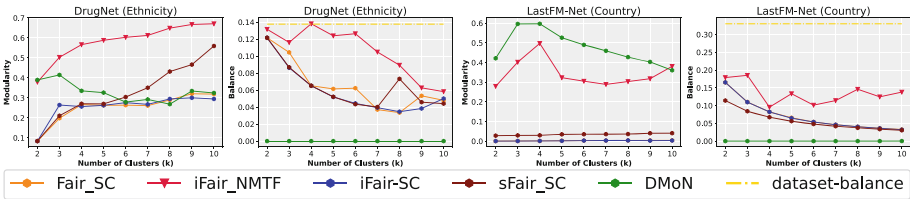


Fig. 2. Performance comparison w.r.t. clustering quality/modularity Q and cluster fairness B (higher values are better for both measures) on DrugNet, and LastFM for different number of clusters $k \in [2, 10]$. $k = 10$ is the convergence point of all models.

As we can see from this figure, our model outperforms the SC-based models in terms of both measures on both datasets. It reports a lower clustering quality Q on LastFM than DMoN which is a neural model primarily focusing on identifying the most modular partitioning of the graph through modularity optimization. DMoN’s Q outcomes reveal varied patterns on LastFM and DrugNet, attributed to differences in network size and density. Neural models typically excel in data-intensive learning cycles, yielding better performance on larger datasets. For instance, LastFM, a substantially larger graph with 5k nodes and 20k edges, showcases this advantage compared to the 200-node DrugNet. However, in terms of fairness, DMoN fails to generate diverse clusters w.r.t. the sensitive attribute, as evidenced by low balance (B). In contrast, our model consistently achieves well-distributed clusters, boasting the highest balance scores among all competitors.

Next, we compare all the models, on all the datasets with a fixed number of clusters $k = 5$, the median of our selected number of clusters. The results are presented in Table 1. We distinguish between real and SBM networks based on

the measurable accuracy of partitioning quality in SBM networks, where ground-truth clusters are known. Additionally, we present the average modularity (Q) and balance (B) across all clusters for real datasets.

Table 1. Results illustrating modularity (Q) and average balance (B) of real networks, and accuracy (Acc) and average balance (B) results on SBM networks for $k = 5$ clusters. (**Bold-underline**) and underline indicate best and second best B results. Best Q , Acc are highlighted with **boldfaced gray**.

Network	FairSC		sFairSC		iFairSC		DMoN		iFairNMTF	
	B	Q	B	Q	B	Q	B	Q	B	Q
Diaries	0.708	0.612	0.809	0.684	0.699	0.647	0.263	0.145	0.648	0.640
Facebook	0.327	0.449	0.602	0.500	0.330	0.448	0.268	0.048	<u>0.514</u>	0.509
Friendship	0.391	0.483	<u>0.485</u>	0.627	0.374	0.392	0.183	0.140	0.631	0.669
DrugNet	0.052	0.263	0.052	0.270	<u>0.061</u>	0.263	0.000	0.326	0.124	0.588
NBA	0.083	0.000	0.323	0.113	0.072	0.000	0.036	0.057	<u>0.286</u>	0.150
LastFM	0.065	0.003	0.056	0.035	<u>0.066</u>	0.002	0.000	0.526	0.069	0.600
	B	Acc	B	Acc	B	Acc	B	Acc	B	Acc
SBM-2K	<u>0.575</u>	0.588	–	–	0	0.799	–	–	0.953	0.958
SBM-5K	0.995	0.998	–	–	0	0.799	–	–	<u>0.941</u>	0.962
SBM-10K	<u>0.999</u>	0.999	–	–	0	0.600	–	–	<u>1</u>	1

The results on real networks indicate the superiority of our proposed iFairNMTF model while reporting the best Q values on 5/6 (meaning 5 out of 6) datasets and 3/6 w.r.t. B . Similarly, on SBM networks iFairNMTF stands the best with 2/3 best accuracy and balance scores. It is worth noting that, in the SBM experiment, DMoN and sFairSC failed to deliver the required number of clusters resulting in empty clusters implying inconsistency in accuracy calculation since the cluster assignments need to be masked to be comparable to true labels.

4.3 Parameter Analysis

This section studies the effect of the λ hyper-parameter on the iFairNMTF model’s performance in terms of Q and B for $k = 5$ clusters in comparison with the performance of other models. In this experiment, we also provide the results of the vanilla SC and vanilla NMTF (the same as iFairNMTF with $\lambda = 0$) models. The results are illustrated in Fig. 3. Based on the results, a comparably good value for the λ parameter can be selected in the intersection of the two measures. These twin charts provide a nice opportunity to visualize the distribution of results and make it easy to select. For instance, values in the range $[0.1, 4]$ for Drugnet and $[55, 200]$ for LastFM are suggested. It gives the end-user a desirable autonomy and depends on the user’s demands on how to select values for this parameter. Consider that, since LastFM is a much larger network than Drugnet, we increase the range of λ with 100 values from 0 to $\lambda = 5000$. Complementary results can be found in the supplementary material (see footnote 4).

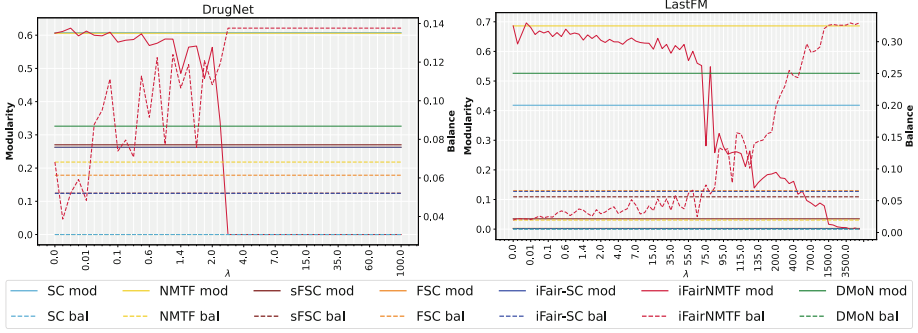


Fig. 3. Parameter λ analysis of the iFairNMTF on Drugnet and LastFM-Net datasets with $k = 5$ in terms of Q and B for $\lambda \in [0, 100]$. Solid lines depict modularity and dashed lines represent balance. Only the behavior of FairSNMF depends on λ .

4.4 Interpretability Analysis

In this section, in addition to the inherent model interpretability through the direct clustering given by the \mathbf{H} factor, we investigate the explicitly interpretable intermediary factor $\mathbf{W} \in \mathbb{R}_+^{k \times k}$ of the iFairNMTF model introduced in Eq. (5). This factor is a symmetric square matrix consisting of non-negative scores representing the strength of cluster-cluster interactions. Diagonal elements reflect intra-cluster connectivity such that the score for dense clusters is expected to be higher. An illustrative example of a graph with 40 nodes distributed between 4 clusters and an imbalanced group distribution of 35% (square shape) to 65% (triangle shape) is shown in Fig. 4. We apply Algorithm 1 to this graph with $\lambda = 1$, and the model identifies the true clusters. Entries corresponding to clusters like I–II, which have no interactions (no links) together, are assigned a value of 0. Furthermore, the score for clusters IV–I is notably lower compared to IV–II, reflecting the difference in the number of connecting links between these clusters.

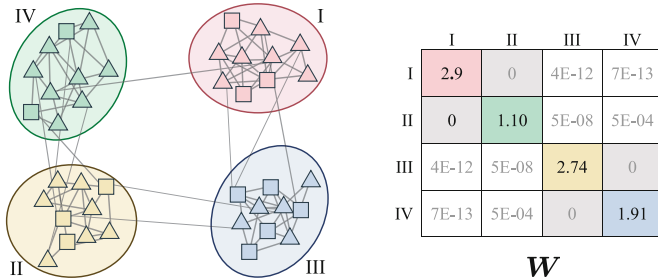


Fig. 4. Interpretability of \mathbf{W} factor for a 40-node graph divided to 4 clusters. Shapes indicate groups.

5 Conclusion and Outlook

In this paper, we introduce the iFairNMTF model, an individually fair flexible approach for graph clustering that takes sensitive (node) attributes into account. iFairNMTF modifies the NMTF model's objective function by incorporating a contrastive penalty term, ensuring that clustering outcomes align with sensitive demographic information and thereby promoting individually fair cluster representations through the attraction and repulsion advantage of the proposed contrastive regularization term. The trade-off regularization parameter λ empowers users to customize the balance between clustering performance and fairness based on their specific needs. Our experiments on both real and synthetic datasets demonstrate that adjusting the trade-off parameter allows for achieving a desired equilibrium between maximizing clustering cohesion and promoting fairness. Promising directions for future research include exploring multi-objective techniques to effectively balance fairness and cohesion objectives, particularly in complex, multi-dimensional discrimination scenarios [23]. Additionally, developing NMF tailored for group fairness, with an emphasis on integrating both individual and group notions into algorithmic design. Finally, evaluating fair clustering methods, esp. for individual fairness remains an ongoing challenge.

Acknowledgements. This work has received funding from the European Union's Horizon 2020 research and innovation programme under Marie Skłodowska-Curie Actions (grant agreement number 860630) for the project "NoBIAS - Artificial Intelligence without Bias". This work reflects only the authors' views and the European Research Executive Agency (REA) is not responsible for any use that may be made of the information it contains. The research was also supported by the EU Horizon Europe project MAMMOth (GrantAgreement 101070285).

References

1. Abdollahi, R., Amjad Seyedi, S., Reza Noorimehr, M.: Asymmetric semi-nonnegative matrix factorization for directed graph clustering. In: ICCKE, pp. 323–328 (2020)
2. Chakraborty, T., Dalmia, A., Mukherjee, A., Ganguly, N.: Metrics for community analysis: A survey. *ACM Comput. Surv.* **50**(4), 1–37 (2017)
3. Chierichetti, F., Kumar, R., Lattanzi, S., Vassilvitskii, S.: Fair clustering through fairlets. In: *Advances in NeurIPS*, pp. 5029–5037 (2017)
4. Dai, E., Wang, S.: Say no to the discrimination: learning fair graph neural networks with limited sensitive attribute information. In: *WSDM*, pp. 680–688 (2021)
5. Dong, Y., Kang, J., Tong, H., Li, J.: Individual fairness for graph neural networks: a ranking based approach. In: *KDD*, pp. 300–310. ACM (2021)
6. Dong, Y., Ma, J., Wang, S., Chen, C., Li, J.: Fairness in graph mining: a survey. *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–22 (2023)
7. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., Zemel, R.S.: Fairness through awareness. In: *Proceedings of the 3rd ITCS Conference*, pp. 214–226 (2012)
8. Ghodsi, S., Ntoutsis, E.: Affinity clustering framework for data debiasing using pairwise distribution discrepancy. In: *EWAF. CEUR Proceedings*, vol. 3442 (2023)

9. Gupta, S., Dukkipati, A.: Protecting individual interests across clusters: Spectral clustering with guarantees. CoRR abs/2105.03714 (2021)
10. Gupta, S., Dukkipati, A.: Consistency of constrained spectral clustering under graph induced fair planted partitions. In: *Advances in NeurIPS*, pp. 13527–13540 (2022)
11. Hajiveisheh, A., Seyedi, S.A., Tab, F.A.: Deep asymmetric nonnegative matrix factorization for graph clustering. *Pattern Recognit.* **148**, 110179 (2024)
12. Kang, J., He, J., Maciejewski, R., Tong, H.: Inform: individual fairness on graph mining. In: *KDD*, pp. 379–389. ACM (2020)
13. Kleindessner, M., Samadi, S., Awasthi, P., Morgenstern, J.: Guarantees for spectral clustering with fairness constraints. In: *ICML*, vol. 97, pp. 3458–3467 (2019)
14. Kuang, D., Park, H., Ding, C.H.Q.: Symmetric nonnegative matrix factorization for graph clustering. In: *SDM*, pp. 106–117 (2012)
15. Lahoti, P., Gummadi, K.P., Weikum, G.: Operationalizing individual fairness with pairwise fair representations. *Proc. VLDB Endow.* **13**(4), 506–518 (2019)
16. Li, T., Ding, C.: Nonnegative matrix factorizations for clustering: a survey. In: *Data Clustering*, pp. 149–176. Chapman and Hall/CRC (2018)
17. von Luxburg, U.: A tutorial on spectral clustering. *Stat. Comput.* **17**(4), 395–416 (2007)
18. Mastrandrea, R., Fournet, J., Barrat, A.: Contact patterns in a high school: a comparison between data collected using wearable sensors, contact diaries and friendship surveys. *PLoS ONE* **10**(9), e0136497 (2015)
19. Newman, M.E.J.: Modularity and community structure in networks. *Proc. Natl. Acad. Sci.* **103**(23), 8577–8582 (2006)
20. Ntoutsis, E., et al.: Bias in data-driven artificial intelligence systems-an introductory survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* **10**(3), e1356 (2020)
21. Pei, Y., Chakraborty, N., Sycara, K.P.: Nonnegative matrix tri-factorization with graph regularization for community detection in social networks. In: *IJCAI*, pp. 2083–2089. AAAI Press (2015)
22. Quy, T.L., Friege, G., Ntoutsis, E.: Multi-fair capacitated students-topics grouping problem. In: *PAKDD* (1). LNCS, vol. 13935, pp. 507–519. Springer (2023)
23. Roy, A., Horstmann, J., Ntoutsis, E.: Multi-dimensional discrimination in law and machine learning - A comparative overview. In: *FAccT*. pp. 89–100. ACM (2023)
24. Rozemberczki, B., Sarkar, R.: Characteristic functions on graphs: Birds of a feather, from statistical descriptors to parametric models. In: *CIKM*, pp. 1325–1334 (2020)
25. Tsitsulin, A., Palowitch, J., Perozzi, B., Müller, E.: Graph clustering with graph neural networks. *J. Mach. Learn. Res.* **24**, 127:1–127:21 (2023)
26. Wang, J., Lu, D., Davidson, I., Bai, Z.: Scalable spectral clustering with group fairness constraints. In: *AISTATS*, pp. 6613–6629 (2023)
27. Wang, Y., Kang, J., Xia, Y., Luo, J., Tong, H.: ifg: Individually fair multi-view graph clustering. In: *IEEE Big Data*, pp. 329–338. IEEE (2022)
28. Weeks, M.R., Clair, S., Borgatti, S.P., Radda, K., Schensul, J.J.: Social networks of drug users in high-risk sites: finding the connections. *AIDS Behav.* **6**, 193–206 (2002)
29. Zemel, R.S., Wu, Y., Swersky, K., Pitassi, T., Dwork, C.: Learning fair representations. In: *ICML* (3). JMLR Workshop and Conference Proceedings, vol. 28, pp. 325–333. JMLR.org (2013)

Supplementary Material for Towards Cohesion-Fairness Harmony: Contrastive Regularization in Individual Fair Graph Clustering

Anonymous Authors

No Institute Given

S1 The iFairNMTF complementary details

An illustrative example is shown in Figure S1, comparing iFairNMTF’s clustering results with three different degrees of enforced fairness regularization on a toy graph with 45 nodes divided into three clusters and two groups. Cluster and group divisions are described in the caption. In Figure S1(a), the model with no fairness penalty prioritizes the clustering objective and obtains a lower Balance (potentially biased clusters). The model in Figure S1(b), enforces an intermediate fairness regularization and well equilibrates the objectives Modularity (indicating clustering cohesion) and Balance (indicating fairness). However, in Figure S1(c), a large penalty is enforced and the model prioritizes fairness leading to well-distributed but less modular clusters.

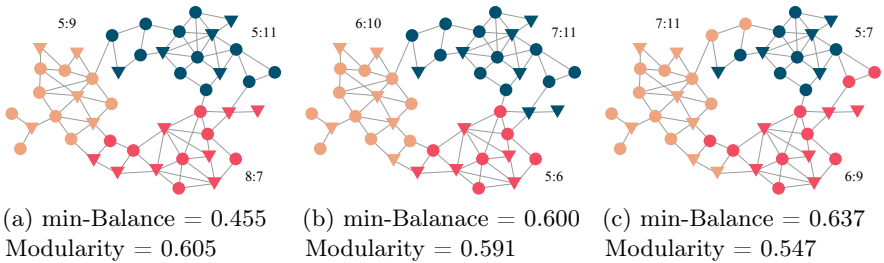


Fig.S1: Comparison of Modularity (Eq. (S1)) and minimum Balance ($\min\text{-Balance}(C) = \min\{\text{Balance}(C_1), \text{Balance}(C_2), \text{Balance}(C_3)\}$) where Balance is computed according to Eq. (S3)) for three different partitionings of a graph, constituting 45 nodes with 27 \circ and 18 ∇ representing sensitive groups. Colors identify clusters. The effect of fairness penalty: a) No penalty: the model prioritizes clustering performance, potentially leading to biased clusters. b) An intermediate penalty: equilibrates clustering and fairness objectives, producing more equitable representations. c) A large penalty: prioritizes fairness leading to well-distributed clusters but potentially with lower performance.

S2 Experimental Evaluation

S2.1 Datasets: Full Details

In this section, the complementary details of the datasets of Table S1 are provided. Table S1 summarizes the statistics of the datasets. It illustrates the raw and cleaned numbers of nodes, and edges with $|V|$, and $|E|$ respectively. The cardinality of $|g|$ indicates the number of sensitive groups and $|c|$ indicates the ground-true number of classes (of the SBM networks).

Table S1: Dataset statistics. $|V|$, $|E|$, $|g|$ and $|c|$ indicate the number of nodes, edges, groups, and ground-truth classes respectively.

Network	$ V $		$ E $		Sensitive Attribute	$ g $	$ c $
	raw	clean	raw	clean			
SBM	2,000	-	267,430	-	attribute	5	5
	5,000	-	978,959	-	attribute	5	5
	10,000	-	2,603,190	-	attribute	5	5
Friendship	134	127	406	396	gender	2	-
Facebook	156	155	1,437	1,437	gender	2	-
DrugNet	293	193	284	273	ethnicity	3	-
NBA	403	403	8,285	8,285	nationality	2	-
LastFM	7,624	5,576	27,806	19,587	country	6	-

Synthetic (SBM): The Synthetic networks used in the paper experiments are generated according to an extension of the popular Stochastic Block Model (SBM) model [3]. The SBM is a random graph model that has been widely used to study the performance of clustering algorithms. According to [4], this extension of the SBM can generate two or more meaningful ground-truth partitionings such that only one of these ground-truths are fair (distributed w.r.t. the sensitive attribute). In the traditional SBM, there is a ground-truth clustering of the vertex set $V = [n]$ into k clusters, and in a random graph generated from the model, two vertices i and j are connected with a probability that only depends on which clusters i and j belong to.

In the extended SBM-generator that we use, we assume that $V = [n]$ comprises m groups $V = V_1 \dot{\cup} \dots \dot{\cup} V_m$ and is partitioned into k ground-truth non-overlapping clusters $V = C_1 \dot{\cup} \dots \dot{\cup} C_k$ such that $|V_s \cap C_l|/|C_l| = \eta_s$, $s \in [m]$, and $l \in [k]$ for some $\eta_1, \dots, \eta_m \in (0, 1)$ with $\sum_{s=1}^m \eta_s = 1$. Hence, in every cluster each group is represented with the same fraction as in the whole data set V and this ground-truth clustering is fair. The edges that connect the set of vertices V on the graph are defined between every arbitrary pair of vertices i and j with a certain probability $\Pr(i, j)$ that depends only on whether i and j are in the same

cluster (or not) and on whether i and j are in the same group (or not). More specifically:

$$\Pr(i, j) = \begin{cases} a, & i \text{ and } j \text{ in same cluster and in same group,} \\ b, & i \text{ and } j \text{ not in same cluster, but in same group,} \\ c, & i \text{ and } j \text{ in same cluster, but not in same group,} \\ d, & i \text{ and } j \text{ not in same cluster, not in same group,} \end{cases}$$

and assume that $a > b > c > d$. In our experiments, similar to [4, 8] we choose the probabilities proportional to the number of nodes of the network such that $a = 10p$, $b = 7p$, $c = 4p$, and $d = 1p$ where $p = \left(\frac{\log(n)}{n}\right)^{2/3}$.

Real: The three high school friendship networks; *Facebook*, *Friendship*, and *Contact-Diaries* datasets [5] are collected from a French high school based on three different strategies but from the same statistical population. Vertices correspond to students and are split into two groups of males and females. The Contact diaries network is constructed based on students’ face-to-face contacts measured through their contact diaries. The Friendship network constitutes self-reported surveys of students’ friendship connections and Facebook is collected from their online social network profiles.

The *DrugNet* [9], which is a network encoding acquaintanceship between drug users in Hartford, CT, can be either used with ethnicity as a sensitive attribute constituting three ethnic groups of African Americans, Latinos, and others or gender (i.e. male and female). In our experiments, we have used the ethnicity feature because it has three imbalanced groups and better challenges algorithms.

The *LastFMNet* [7] contains mutual follower relations among users of Last.fm, a recommendation-based online radio and music community in Asia. LastFMNet was collected from public API in 2020 and used to study the distribution of vertex features on graphs.

NBA is an extension of a Kaggle¹ dataset containing relationships between around 400 NBA basketball players [2]. It has demographics including nationality, age, salary, and a number of statistical performance indicators of players in the 2016-2017 season. The nationality of players is used as a sensitive attribute categorizing them to US and non-US players.

S2.2 Metrics

In the paper, accuracy is used for measuring clustering assignment quality on synthetic networks. For real-world networks, since the ground truth cluster structures are unknown, Newman’s modularity measure [1, 6] is adopted which basically analyzes the homogeneity of clusters by calculating the difference in the proportion of internal links in each cluster for a given partitioning compared to

¹ <https://www.kaggle.com/datasets/noahgift/social-power-nba>

the expected proportion of edges in a null graph with the same degree distribution as of the original graph:

$$Q = \frac{1}{|E|} \sum_{i,j} \left(A_{ij} - \frac{\deg(i)\deg(j)}{|E|} \right) \delta(c_i, c_j), \quad (\text{S1})$$

where $|E|$ is the total number of edges in the graph, $\deg(i)$ and $\deg(j)$ are degrees of each arbitrary pair of i and j nodes, c_i and c_j are cluster assignments of nodes i and j , respectively and $\delta(c_i, c_j)$ is the Kronecker delta function, indicating 1 if its arguments are equal, and 0 otherwise. In other terms, $\delta(c_i, c_j)$ equals 1 only if nodes i and j belong to the same cluster (community). Modularity can range between $[-1, 1]$. We measure the fairness of clustering in terms of the average balance (B) over k clusters such that:

$$B = \frac{1}{k} \sum_{l=1}^k \text{Balance}(C_l), \quad (\text{S2})$$

Where $\text{Balance}(C_l)$ calculates the minimum group proportion of C_l according to Equation (S3):

$$\text{Balance}(C_l) = \min_{s \neq s' \in [m]} \frac{|V_s \cap C_l|}{|V_{s'} \cap C_l|}, \quad (\text{S3})$$

The average balance is a common measure that was also previously used by [4, 8]. The minimum balance of each cluster can range between $[0, 1]$, thus their average also ranges between $[0, 1]$. For both modularity and average balance, higher values indicate better results.

S2.3 Convergence Analysis

Figure S2 illustrates how the proposed model's loss converges w.r.t. the number of iterations. In particular, the objective value (i.e. the loss) of the model is calculated according to Equation (S4) in each iteration of the optimization. As demonstrated in the plots, for all datasets, the loss values converge immediately in the initial 50 iterations. Note that, loss values in Figure S2 are shown in scientific notation with a power of $1e2$.

$$\min_{\mathbf{H}, \mathbf{W} \geq 0} \|\mathbf{A} - \mathbf{H}\mathbf{W}\mathbf{H}^\top\|_F^2 + \lambda \text{Tr}(\mathbf{H}^\top \mathbf{L} \mathbf{H}), \quad (\text{S4})$$

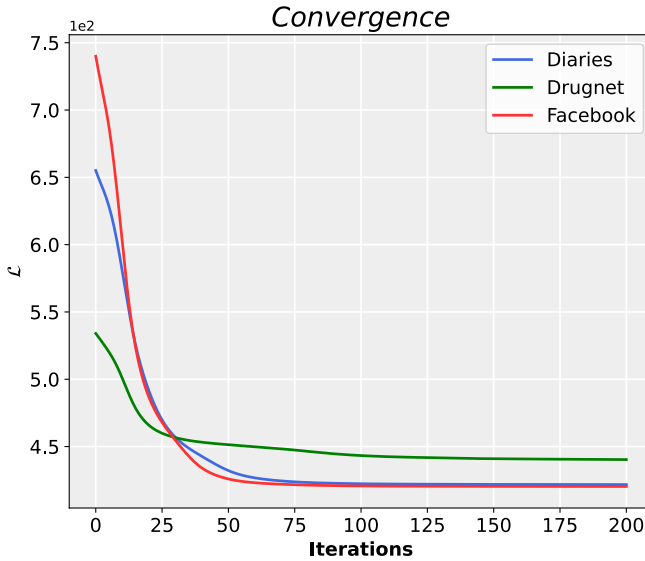


Fig. S2: iFairNMTF loss convergence as a function of the number of iterations on three datasets Contact Diaries, Drugnet, Facebook for a parameter value of $\lambda = 1$ and $k = 5$ clusters.

References

1. Chakraborty, T., Dalmia, A., Mukherjee, A., Ganguly, N.: Metrics for community analysis: A survey. *ACM Computing Surveys* **50**(4), 1–37 (2017)
2. Dai, E., Wang, S.: Say no to the discrimination: Learning fair graph neural networks with limited sensitive attribute information. In: *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. pp. 680–688 (2021)
3. Holland, P.W., Laskey, K.B., Leinhardt, S.: Stochastic blockmodels: First steps. *Social networks* **5**(2), 109–137 (1983)
4. Kleindessner, M., Samadi, S., Awasthi, P., Morgenstern, J.: Guarantees for spectral clustering with fairness constraints. In: *Proceedings of the 36th International Conference on Machine Learning*. vol. 97, pp. 3458–3467 (2019)
5. Mastrandrea, R., Fournet, J., Barrat, A.: Contact patterns in a high school: a comparison between data collected using wearable sensors, contact diaries and friendship surveys. *PloS one* **10**(9), e0136497 (2015)
6. Newman, M.E.J.: Modularity and community structure in networks. *Proceedings of the National Academy of Sciences* **103**(23), 8577–8582 (2006)
7. Rozemberczki, B., Sarkar, R.: 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*. pp. 1325–1334 (2020)
8. Wang, J., Lu, D., Davidson, I., Bai, Z.: Scalable spectral clustering with group fairness constraints. In: *International Conference on Artificial Intelligence and Statistics*. pp. 6613–6629 (2023)
9. Weeks, M.R., Clair, S., Borgatti, S.P., Radda, K., Schensul, J.J.: Social networks of drug users in high-risk sites: Finding the connections. *AIDS and Behavior* **6**, 193–206 (2002)