

# Spectral Graph Forge: Graph Generation Targeting Modularity

Luca Baldesi  
DISI Dept  
University of Trento  
luca.baldesi@unitn.it

Carter T. Butts  
Sociology Dept  
University of California, Irvine  
butts@uci.edu

Athina Markopoulou  
EECS Dept, CPCC  
University of California, Irvine  
athina@uci.edu

**Abstract**—Community structure is an important property that captures inhomogeneities common in large networks, and modularity is one of the most widely used metrics for such community structure. In this paper, we introduce a principled methodology, the Spectral Graph Forge, for generating random graphs that preserves community structure from a real network of interest, in terms of modularity. Our approach leverages the fact that the spectral structure of matrix representations of a graph encodes global information about community structure. The Spectral Graph Forge uses a low-rank approximation of the modularity matrix to generate synthetic graphs that match a target modularity within user-selectable degree of accuracy, while allowing other aspects of structure to vary. We show that the Spectral Graph Forge outperforms state-of-the-art techniques in terms of accuracy in targeting the modularity and randomness of the realizations, while also preserving other local structural properties and node attributes. We discuss extensions of the Spectral Graph Forge to target other properties beyond modularity, and its applications to anonymization.

## I. INTRODUCTION

Generating random graphs with certain prescribed properties is an increasingly important field [1], [2], [3], [4], [5]. Random graphs are widely used to model structure in social, computer, and biological systems, with applications ranging from importance sampling to the generation of synthetic data sets. One use of synthetic graphs is as a proxy of real-world data sets. For example, in the measurement community, researchers collect Internet [6], [7] or online social network [8], [9] topologies, making them available to other researchers for use in studying their properties or evaluating the performance of graph algorithms. In this context, it is often desirable to not release the original graph itself, but to generate synthetic graphs that (exactly or approximately) specify properties of interest of the real (or “target”) graph, while varying with respect to other properties. This can be useful for anonymization purposes (*i.e.*, to obfuscate individual or other sensitive information) and/or for simulation purposes (*i.e.*, to ensure acceptable performance over a range of network structures, while preserving features known to be present in real-world use cases).

This work was conducted when Luca Baldesi visited UC Irvine in 2016. All authors were partially supported by NSF Award III-1526736. A. Markopoulou is a member of the Center for Pervasive Communications and Computing (CPCC).

Prior work has primarily focused on generating synthetic graphs with certain local structural properties, such as degree distribution, degree correlations, subgraph counts etc. [1], [3], [4], [5], [10], while global properties have received less attention. In this paper, we are interested in an important global property exhibited by many social and communications networks, namely *community structure*, which captures the complex pattern of inhomogeneities that characterize them. Real-world networks are typically irregular, containing numerous and often nested or overlapping subsets of vertices that are internally cohesive while being less well-connected to other subsets. This *community* [11] or *cohesive subgroup* structure [12] has a profound impact on phenomena such as diffusion or robustness to failure; it can also affect the performance of sparse-graph algorithms for tasks such as path or cycle counting that are sensitive to the presence of small, densely connected clusters of vertices. Graph generation methods that focus on local features (such as subgraph counts, the degree distribution, clustering etc.) do not necessarily preserve community structure, and the resulting graphs may be poor proxies for real-world networks in applications that are sensitive to it, such as diffusion. Therefore, distinct methods are needed to target it.

In this paper, we present a novel random graph generation framework that is well-suited to the above applications. In particular, we target community structure, within desired accuracy while allowing other features to vary. We use *modularity*, a well-known property of matrices to capture the community structure [11]. Starting from a given real graph as input, we use the original graph’s high-level structure captured by the eigenstructure of its modularity matrix to derive a class of random graphs sharing the same modularity value, within a level of accuracy.

Our approach can be summarized as follows and is depicted in Fig. 1. Given the adjacency matrix of a graph  $A$ , we capture its community structure via its modularity matrix ( $B$ ). Because  $B$  is symmetric and real, it admits a spectral decomposition [13] which expresses the matrix in terms of a set of eigenvectors and eigenvalues. We apply a low-rank spectral approximation [14] on this eigenstructure to yield an approximated matrix reflecting the features most strongly associated with global structure; this matrix is then transformed and normalized to produce an edge probability matrix

( $A^\dagger$ ), which is then passed to a sampling algorithm to yield a synthetic graph (with adjacency matrix  $A'$ ). The synthetic graphs produced will preserve the higher-order properties (modularity) of the input, while allowing other details to vary.

Spectral Graph Forge has several strengths. First, our approach is principled, based on spectral approximation of matrices, thus the name Spectral Graph Forge (SGF), and provides a tunable parameter to tradeoff between accuracy of the targeted modularity and the randomness of the produced realizations. Second, Spectral Graph Forge outperforms state-of-the-art baseline approaches in several aspects and for a range of graphs: (i) it achieves the target modularity more accurately than baselines (Section V-B) (ii) it preserves some metrics of local structure, beyond modularity (Section V-C) (iii) it can maintain the node and edge attributes, hence it can preserve properties beyond just pure topology (Section V-D). Third, the analysis of the entropy of graphs produced shows that they vary substantially on dimensions other than those deliberately targeted (Section V-E); we discuss how to exploit this property and use SGF as anonymization technique for social networks. Last but not least, SGF is not limited to targeting modularity, which is the focus of this paper. The general SGF framework can be used to target other matrix representations of a graph, such as the adjacency matrix or one of its derived matrices (including but not limited to modularity, laplacian, clique co-membership).

The structure of the rest of the paper is as follows. In Section II, we review related work. In Section III, we discuss community structure and its relation to modularity. In Section IV, we present our *Spectral Graph Forge* Framework, and how to apply it specifically to generate graphs with a prescribed modularity. In Section V, we evaluate our algorithms for various datasets (both synthetic and real-world social networks), against state-of-the-art algorithms targeting modularity. Section VI concludes the paper.

## II. RELATED WORK

There is a significant amount of work on generation of graphs that preserve *local* structural properties (e.g., degree distribution, small subgraph counts). Graphlets [5], motifs [4] and dK-series [1], [3] graph generators are all intended to produce graphs that preserve such properties, starting from an initial realization, or from a vector of input statistics. Exponential family random graph models (ERGMs) [15] can be employed to simulate graphs having the same *expected* values for a series of targeted features as the original, and most applications to date have likewise focused on local properties (including properties that depend on nodal or edgewise covariates); among other differences from the above, ERGMs preserve only expected statistics (not exact values), but provide maximum entropy guarantees.

In contrast, for global properties such as the graph modularity [11], there exist only a few approaches so far that target them by design. Trajanovski et al. [16] propose an algorithm for generating graphs with a prescribed modularity. The Karrer and Newman stochastic block model [17] generates

graphs with a given community structure and modularity value. In principle, other latent structure models such as the latent space models [18] and latent stochastic block models [19] could also be used to simulate graphs with related properties, although applications to date have been focused on inference. In Section V-A2, we describe the modularity-targeting algorithms by Trajanovski et al. and Karrer and Newman, and we used them as baselines for comparison. The relation between modularity and eigenstructure has been investigated mainly for applications such as community detection [11], [20], visualization [21], and analysis [18]. To the best of our knowledge, our work is the first to design and implement a general framework for generation of synthetic graphs with specified high-level structure (modularity being the focus in this paper) by means of spectrally transformed inputs.

## III. MODULARITY AND COMMUNITY STRUCTURE

### A. Definitions and Intuition

*Modularity* is one of the most widely used metrics for characterizing community structure [22]. It expresses the extent to which a partition of vertices divides the graph such that within-group densities are higher than between-group densities; the maximum modularity obtainable over all partitions is hence a global measure of the extent to which a graph approximates a union of dense structures that are minimally connected to one another. Modularity is closely related to *cohesive subgroup* structure, a topic of historical sociological interest [23], [24], and has been the object of extensive study within the large literature on *community detection* [22], [11], [20].

Formally, modularity is defined as follows. Given a graph  $G = (V, E)$ , where  $|V| = n$ , with adjacency matrix  $A$  and degree vector  $K = (k_1, \dots, k_n)$ , the modularity matrix  $B$  is defined as [11]

$$B_{ij} = A_{i,j} - \frac{k_i k_j}{|K|} \quad (1)$$

where  $|K| = \sum_{i=1}^n k_i$  is the degree sum.  $B$  effectively captures the number of edges within nodes of the same community compared to a random graph with the same expected degree sequence. Given a node partition  $\{C_1, \dots, C_m\}$  with nodes assigned to communities  $c_1, \dots, c_n \in \{C_1, \dots, C_m\}$ , the modularity value measuring this concentration over all node pairs is:

$$Q = \frac{1}{|K|} \sum_{i,j} \left( A_{i,j} - \frac{k_i k_j}{|K|} \right) \delta(c_i, c_j). \quad (2)$$

where  $\delta(a, b) = 1$  iff  $a = b$  and 0 otherwise.

The notion of modularity is used in three ways. First, given a partition of nodes into communities  $c_1, \dots, c_n$ ,  $Q$  provides a metric for how well the partitioning matches the actual topology. Second,  $Q$  can be used to search for community structure, by seeking a partition of nodes into  $m^*$  communities,  $c_1^*, \dots, c_n^*$ , such that  $Q^*$  is maximized. Finally, the maximized value of  $Q^*$  over all partitions can be taken as an indicator of the extent to which the graph naturally divides into distinct communities. This quantity is sometimes referred to generically as the *modularity of the graph*.

## B. Modularity and Spectral Structure

The eigenstructure of an adjacency matrix  $A$  associated to a given graph is intimately related to the structure of the associated graph.  $A$  can be decomposed in a set of eigenvalues  $\lambda_i$  and eigenvectors  $v_i$ ; the eigenvectors associated with positive eigenvalues describe *core-periphery* structures [25] and the signs of their entries describe possible graph bi-partitions. The eigenvectors associated with negative eigenvalues relate instead to memberships in *bi-partitions*, so that, the sign of their entries indicate a likelihood of node connection.

Hence, there is an obvious affinity between the information described by the eigenvectors and modularity but for the leading eigenvector, associated to the main core of the network (the whole graph) and it would be natural to consider the eigenstructure of the residual matrix  $A - \lambda_1 v_1 v_1^T$  (where  $\lambda_1, v_1$  are respectively the leading eigenvalue and eigenvector) as a basis for community detection.

However, in practice,  $v_1$  is typically correlated with  $K$ , and hence the modularity matrix  $B = A - (KK^T)/|K|$  is very close to the residual matrix. In keeping with this observation, it is common practice (e.g., [11]) to work with the eigenstructure of  $B$  to identify partitions, recursively subdividing nodes into classes based on the signs of their eigenvector elements associated with high-modulus, positive eigenvalues of  $B$ . For example, in [22], Newman showed that the leading eigenvector of the modularity matrix  $B$  determines the best bipartition of nodes into two communities.

**Beyond modularity:** While in this paper, we focus on targeting modularity, our spectral decomposition framework is more general and can be tuned to capture desired structural elements of interest, while randomizing out other features. More generally, we could also apply SGF to matrices other than  $B$  (namely  $A$  and its derived matrices), and use it to target other global properties. For example, while the leading eigenvector defines the best bipartition [22], positive/negative eigenvalues capture core-periphery/bipartition structures, respectively.

## IV. SPECTRAL GRAPH FORGE

The Spectral Graph Forge takes an undirected graph  $G$  with  $n$  nodes as input and, for each run, produces a random output graph  $G'$  with the same number of nodes. We represent  $G$  and  $G'$  with their respective adjacency matrices  $A, A' \in \{0, 1\}^{n \times n}$ . The Spectral Graph Forge is designed to be as modular and extensible as possible; its graph construction process can be viewed, hence, as a pipeline composed by several nearly independent steps, each of which can be varied as needed by the user. A representation of this pipeline is depicted in Fig. 1.

### A. Transformation

To allow the Spectral Graph Forge to preserve structure associated with either  $A$  or a derived matrix, we work with a user-selected real-valued, symmetry preserving transformation of  $A$ , denoted  $M$ . The main use case of  $M$  in this paper is the modularity matrix  $M = B = A - (KK^T)/|K|$ . Other examples of  $M$  include the adjacency matrix itself  $M = A$ ,

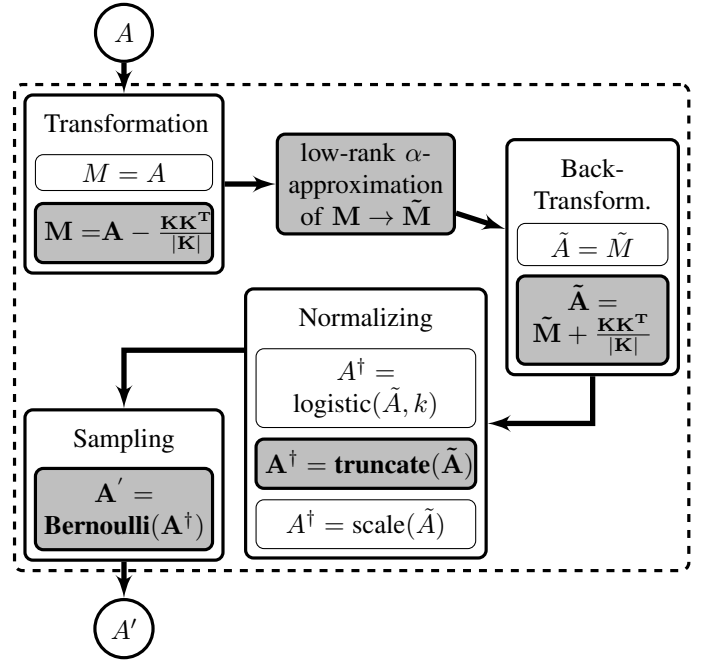


Fig. 1. The pipeline of the Spectral Graph Forge (SGF) Framework. Given an undirected graph adjacency matrix  $A$  as input, SGF outputs a “similar” one  $A'$  from which we can build the corresponding graph, called  $\text{SGF}(\alpha)$ . Sub-blocks indicate mutually exclusive options for each step. The focus of this paper is on using SGF to target modularity, by setting  $M = B = A - (KK^T)/|K|$ , and the corresponding blocks are highlighted in grey.

and the clique co-membership matrix (whose  $i, j$  cells are the counts of cliques containing both  $i$  and  $j$ ).

### B. Low-rank $\alpha$ approximation.

Since  $M$  is real and symmetric, it can be expressed as

$$M = V\Lambda V^T = \sum_{i=1}^n \lambda_i v_i v_i^T, \quad (3)$$

where  $\lambda_i, v_i$  are respectively the  $i$ th eigenvalue and the  $i$ th eigenvector of  $M$  (all of which are real). W.l.o.g., we scale  $v_i$  so that  $\|v_i\|_2 = 1$  and sort the eigenvectors/eigenvalues such that  $|\lambda_1| \geq \dots \geq |\lambda_n|$ . Hence the first components of Eq. (3) contribute to  $M$  more than the last ones. We can perform a lossy compression of  $M$  by removing the least important frequencies. We use a parameter  $\alpha \in [0, 1]$  to control the fraction of the eigenvectors we want to keep from the original  $M$  and our approximation is

$$\tilde{M} = \sum_{i=1}^{\lceil \alpha n \rceil} \lambda_i v_i v_i^T \quad (4)$$

This is essentially a low-pass filter on  $M$ , which retains global characteristics captured by the first  $\alpha n$  eigenvectors of the matrix, while removing the remaining ones, associated with idiosyncratic local structures.<sup>1</sup> The euclidean norm of the error,  $M - \tilde{M}$ , due to transformation can be computed as

$$\|M - \tilde{M}\|_2 = \left\| \sum_{i=\lceil \alpha n \rceil + 1}^n \lambda_i v_i v_i^T \right\|_2 = \lambda_{\lceil \alpha n \rceil + 1} \quad (5)$$

<sup>1</sup>We could also filter on phase (a *graph polarizing filter*) to extract core-periphery (positive eigenvalues) or bipartite (negative eigenvalues) structures.

Hence, SGF can target the reproduction of  $M$  within arbitrary precision, by tuning  $\alpha$  (Appendix A of extended version [26]).

### C. Back-Transformation.

Once the low-rank approximation  $\tilde{M}$  is obtained, we need to transform it to obtain an approximated adjacency matrix  $\tilde{A} \in \mathbb{R}^{n \times n}$ . Where the back-transformation is the inverse of the transformation  $M$  to  $\tilde{M}$ . For the modularity matrix  $M = B = A - (KK^T)/|K|$ , the back transformation is intuitive:  $\tilde{A} = \tilde{M} + (KK^T)/|K|$ . Other choices of back-transformation may be employed where  $M$  is not invertible (e.g., in the case of the co-clique matrix), resulting in a different approximation. Since  $\tilde{A}$  is not, in general, an adjacency matrix, our SGF fills the gap between  $\tilde{A}$  and the final  $A'$ , by normalizing the values to be within  $[0, 1]$  ( $A^\dagger$ ) and by using a stochastic process that forces the elements to be either 1 or 0 ( $A'$ ), as described next.

### D. Normalization

We first normalize the values of the latter in  $[0, 1]$  obtaining:  $A_{i,j}^\dagger = \text{norm}(\tilde{A}_{i,j})$ , where *norm* is a user-specified normalization function. With modularity in mind, we investigated three variants for this normalization function:

- $\text{logistic}(\tilde{A}_{i,j}, k) = \frac{1}{1 + e^{(0.5 - \tilde{A}_{i,j}) \cdot k}}$ , where  $k \in [2, 10]$
- $\text{truncate}(\tilde{A}_{i,j}) = \begin{cases} 1 & \text{if } \tilde{A}_{i,j} \geq 1 \\ 0 & \text{if } \tilde{A}_{i,j} \leq 0 \\ \tilde{A}_{i,j} & \text{otherwise} \end{cases}$
- $\text{scale}(\tilde{A}_{i,j}) = \frac{\tilde{A}_{i,j} - \min_{s,t} \tilde{A}_{s,t}}{\max_{s,t} \tilde{A}_{s,t} - \min_{s,t} \tilde{A}_{s,t}}$

where  $k \in [2, 10]$  is a parameter we use to tune the inclination of the *logistic()* function.

We have extensively evaluated the above normalization functions and we found that the truncation rule gives the cleanest entropy and distance relationships; in fact, *scale* introduces a distortion highly dependent on the difference between the minimum and maximum value while *logistic* (for which we found with  $k = 6$  provides the best performance for our case of study) shows poor sensitivity to  $\alpha$ , preserving relatively little structure throughout. Details on this investigation are deferred to Appendix B of the extended version [26].

### E. Sampling

Taking  $A^\dagger$  as a matrix of expectations for  $A'$ , we complete our process by drawing  $A'$  from an inhomogeneous Bernoulli graph distribution:  $A'_{i,j} = \text{bernoulli}(A_{i,j}^\dagger)$ ,  $\forall j > i$ . Note that to generate a simple undirected graph  $A'$ , we set:

- $A'_{i,i} = 0, \forall 1 \leq i \leq n$
- $A'_{j,i} = A'_{i,j}, \forall j > i$

As with other aspects of the pipeline, this can be generalized (e.g., to allow for degree or other constraints). In the following, we denote by  $\text{SGF}(\alpha)$  as the graph resulting from  $A'$ , drawn using our method on fraction  $\alpha$  of the eigenvectors.

*Entropy of the Synthetic Graphs.* One advantage of the use of an inhomogeneous Bernoulli graph in our pipeline is that we can compute exactly the (Shannon) entropy of our synthetic distribution. Entropy indicates the degree of variation in the output graphs, and is an indicator of the extent to

which a graph generation method is able to provide a range of structures compatible with its arguments; at constant graph order  $n$ , it necessarily varies between  $n(n-1)/2$  bits for a uniform random graph and 0 bits for a generator that produces only a single realization with certainty. Given a fixed level of preservation of targeted features, generators with higher entropy are typically preferable (as they vary the distribution of structures as widely as possible, given the desired constraints). The ability to easily compute output entropy is a useful feature of our approach; by contrast, most other approaches provide no known way of computing output entropy, making it difficult to assess the extent to which they are able to produce a diverse range of structures.

Since  $A'$  is drawn from a Bernoulli graph distribution with parameter matrix  $A^\dagger$ , the entropy is immediate:

$$H(A^\dagger) = \frac{1}{S} \sum_{i=1}^n \sum_{j=i+1}^n -A_{i,j}^\dagger \log_2[A_{i,j}^\dagger] - (1 - A_{i,j}^\dagger) \log_2[1 - A_{i,j}^\dagger],$$

which is readily computed from  $A^\dagger$ , where  $S = \frac{-n(n-1)}{2} (\log_2 \delta + \log_2(1 - \delta))$  is the normalization factor and  $\delta = \frac{2}{n^2 - n} \sum_{i,j} A_{i,j}^\dagger$  is the expected graph density under  $A^\dagger$ . This index can be interpreted as the fraction of the maximum possible entropy that  $A'$  attains, given its expected density.

### F. Computational Complexity

To analyze the complexity of SGF, we consider the complexity of each block in Fig. 1. The pipeline starts with a graph  $G = (V, E)$ ,  $|V| = n$ ,  $|E| = m$ . The *Transformation* and *Back-Transformation* blocks can clearly be performed in  $O(n^2)$ . The *Normalizing* and *Sampling* blocks operate element-wise on  $n \times n$  matrices and apply a float operation so their cost is  $O(n^2)$ . The *Low-rank  $\alpha$ -approximation* block performs a spectral decomposition on a  $n \times n$  matrix which can be done in  $O(n^3)$  or  $O(nm)$  in case of an input sparse matrix<sup>2</sup>. Overall, the bottleneck of SGF is the spectral decomposition, which, for sparse input matrices  $M$ , has hence a complexity of  $O(n^2)$  [11]. We mainly target social graphs where the community structure analysis is of main interest; these graphs are generally quite sparse with a low number of nodes (each of the network in our real-world dataset has less than  $10^5$  nodes). Furthermore, SGF can gracefully tradeoff complexity for accuracy by tuning the parameter  $\alpha$ : keeping a few eigenvectors in the approximation (Eq. (4)), increases the error (Eq. (5)) but decreases the complexity of the spectral decomposition (e.g. using the moment method).

## V. EVALUATION

### A. Evaluation Setup

To validate our approach, we test the Spectral Graph Forge against state-of-the art baselines (Section V-A2) on several datasets (Section V-A1), and in terms of several metrics.

<sup>2</sup>Since  $M$  symmetric, we can find an orthogonal matrix  $Q$  such that the similarity transform  $T = Q^T A Q$  gives a tridiagonal matrix [11], which belongs to  $O(nm)$  with the Lanczos algorithm [27]. The eigenvectors  $v_i = Q w_i$  of  $M$  are hence directly computable from the eigenvectors  $w_i$  of  $Q$ . Since  $Q$  is tridiagonal, its eigenvectors can be computed with the *QL algorithm* with a cost of  $O(n)$  [11].

1) *Datasets*: We use both synthetic and real-world datasets; the former for a fair comparison on well-known state-of-the-art datasets in the field of community detection, and the latter to show that our results hold with real data as well.

**I. Synthetic Datasets.** A random graph generator that controls modularity is the following: fix the number of communities and then control the intra-community vs. inter-community probability of edge: fixing the inter-community edge probability, higher values of intra-community edge probability leads to higher modularity value. Although stylized, this model allows for explicit control of the inherent modularity of the synthetic graph. This idea has been widely used since the early days of community detection. For example, Girvan and Newman [28] generated graphs with 128 nodes, communities of the same size and essentially constant node degree. The idea has been extended by Lancichinetti et al. [29] to produce networks of 1000 nodes and geometrically distributed node degrees and community sizes. We use both generators to produce ten networks from each, and refer to them as **Girvan** or **Lancichinetti** datasets, respectively.

**II. Real World Datasets.** In the real world, community structure can be found in nearly any network representing interpersonal relationships. Some of the most striking group structures are seen in the context of friendship networks, particularly for populations such as adolescents (who are frequently segregated by age, among other factors). The National Longitudinal Study of Adolescent to Adult Health (referred to as **Add-Health** in the following) is a U.S. national study on adolescents in grades 7-12. It is a large, longitudinal study of students drawn from a sample of American secondary schools, with data being collected during an initial in-school survey followed by four in-home interviews spanning from 1994 to 2008 [30]. In this work we employ the public-use network data from the first wave, in which students were asked to nominate peers they considered to be friends; the data used here comprises 16 different networks, in which each network represents a school (or pair of linked schools in which graduates of the first proceed to the second), each vertex represents a student, and each edge represents a friendship nomination. For each student, we are also provided with data on gender, race/ethnicity, and grade. Overall, friendships are strongly segmented by grade, with additional divisions by other demographic characteristics in most schools.

Finally, we also consider smaller networks frequently encountered when analyzing large-scale network data. Given a large network and a given node, an *egocentric network* (or *ego net*) for this node is the induced subgraph formed by the node (“ego”) and its first-order neighborhood (“alters”). An ego net hence includes all ego-alter ties, together with the alter-alter relationships; cohesive subgroups in the ego net (with ego subtracted) thus mark the *local* communities of alters with whom ego interacts. Here, we employ the **Facebook** dataset of Mcauley and Leskovec [31], which contains a set of ego-networks derived from the Facebook social network comprising of 10 different networks of size ranging from 52 to 1034 nodes (ego removed for analysis).

2) *Baseline Modularity-Targeting Algorithms*: There are two classes of prior state-of-the-art graph generation algorithms that have been used to explicitly target modularity. We briefly present here these algorithms to facilitate better understanding of our experimental results.

**I. Trajanovski Algorithm** Trajanovski et al. [16] present three different rewiring operations capable of varying the modularity value of a graph with a given partition. The input to the algorithm is a target modularity value  $Q_t^*$ , a number of communities  $m$  and a total number of links  $L$ . In the initialization phase, it creates the communities of nodes in the form of a tree with no more than one link between two communities. This configuration has the maximum achievable modularity for the given input. Then, rewiring transformations are applied to lower the modularity value to the target one.

A limitation of this approach is that it assumes the community graph partitioning  $c_1^*, \dots, c_n^*$  never changes. Since the best partition is dictated by the topology, the rewiring process may alter them and, consequently, yield an apparent modularity that is lower than the optimum. Hence, even if the rewiring process successfully obtains the desired modularity value  $Q_t^*$  with respect to the initial partitioning, the output graph may have topology with a slightly different inner community structure, leading to a completely different modularity value  $Q_o^*$ .

**II. Degree Corrected - Stochastic Block Model (DC-SBM)** The stochastic block model [32] is a generative algorithm for groups or communities in networks. The related algorithms work by placing each edge  $e_{i,j}$  between the vertices  $v_i, v_j$  with a probability which is function of  $c_i$  and  $c_j$ . Since  $c_i, c_j \in \{C_1, \dots, C_m\}$ , these probabilities form a matrix of inter and intra community (block) connection probabilities.

Generally, simple stochastic block models have difficulty in building graphs with community structure matching real-world data sets [17]. A limitation of the most basic models in this class is that they do not include effects for other forms of structure found in real-world networks, e.g. degree heterogeneity. The approach proposed by Karrer and Newman [17] rectifies this last limitation by correcting for the degree distribution, yielding improved performance on realistic topologies.

### B. Results on Targeting Modularity

Here, we evaluate the main objective of this paper, *i.e.*, how well Spectral Graph Forge achieves the target modularity. When comparing the modularity value of an input graph  $Q_i$  (or  $Q_i^*$ ) with the modularity value of an output one  $Q_o$  (or  $Q_o^*$ ), we consider the (maximum) *modularity ratio*  $\frac{Q_o}{Q_i}$  (or  $\frac{Q_o^*}{Q_i^*}$ ): the closer this ratio to 1, the closer the output graph matches the input in terms of modularity. To compute the modularity value we use the well-known algorithm by Lefebvre et al. [33].<sup>3</sup>

1) *Insights into SGF itself*: We first consider a trivial dataset for which our graphs are made of 128 nodes and 2 or 8 communities, and we vary parameters of SGF ( $\alpha$ ) as well as the inherent community structure of the target network to get

<sup>3</sup>There is a large body of work on community detection algorithms to maximize  $Q$  [34]. In this work, we compute  $c_1^*, \dots, c_n^*$  and the associated modularity values using the well-known algorithm by Lefebvre et al. [33].

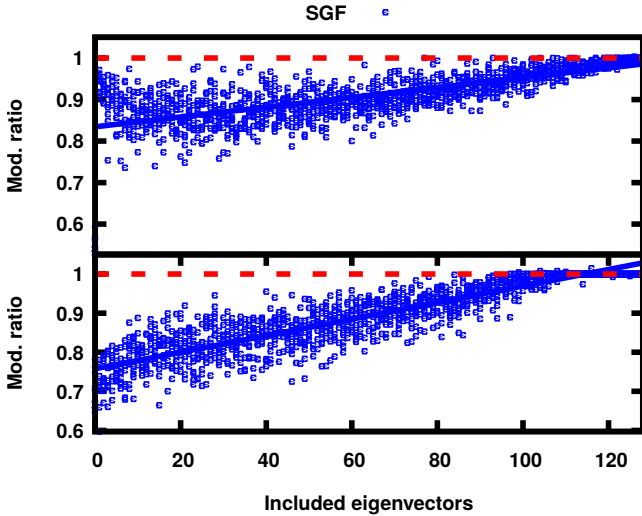


Fig. 2. Modularity ratio for Spectral Graph Forge varying  $\alpha$  on our simple graphs with 2 communities (upper plot) and 8 communities (lower plot).

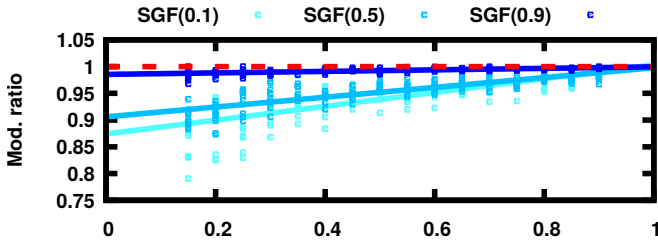


Fig. 3. Modularity ratio for Spectral Graph Forge on our simple graphs varying their connectivity with 2 communities.

insights into the behavior of SGF. Fig. 2 shows the maximum modularity ratio  $\frac{Q^*}{Q_i^*}$  varying the  $\alpha$  parameter with networks of two and eight communities respectively. As expected, the as  $\alpha \rightarrow 1$ , the more eigenvectors are used and the better the performance ( $\frac{Q^*}{Q_i^*} \rightarrow 1$ ).

The SGF performance, in terms of how well it targets modularity, is also influenced by the inherent community structure. In Fig. 2, we see that SGF requires more eigenvectors for 8 than for 2 communities, in order to obtain similar results. In Fig. 3, we see that the higher the modularity of the input graph (controlled by the intra-community edge probability) the better the performance of SGF.

2) *Evaluating SGF against baselines.: Experiments.* In the SGF pipeline, we target the modularity matrix as the  $M = B$ , the truncation as the normalization function, and a value  $\alpha = 0.9$ . We apply SGF to each input network ten different times, computing the modularity ratio for each simulation; the same process is carried out for the two baseline algorithms. In particular, for the Trajanovski et al. algorithm, we compute the modularity value  $Q_o^*$ , the number of nodes  $n$ , the number of edges and the number of communities from each input graph and we use them for the Trajanovski approach. As input

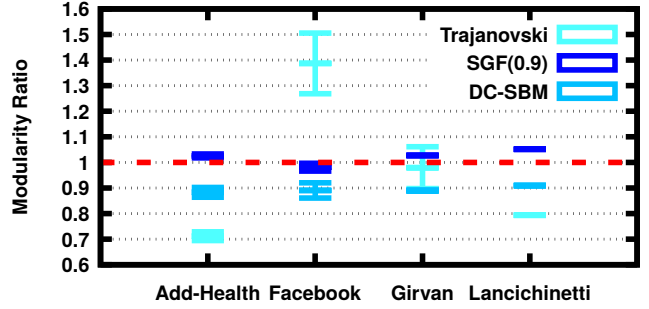


Fig. 4. Means and 99% confidence intervals for the modularity ratio, by method and data set.

TABLE I  
MODULARITY RATIO FOR ALL THE STRATEGIES ON ALL THE DATASETS.

Strategy	Dataset	Mean	Std	Dataset	Mean	Std
SGF(0.1)	Add-Health	0.76038	0.10479	Facebook	0.80551	0.07143
SGF(0.1)	Girvan	0.74602	0.01709	Lancichinetti	0.58268	0.00799
SGF(0.2)	Add-Health	0.71985	0.06740	Facebook	0.82773	0.08925
SGF(0.2)	Girvan	0.70728	0.01033	Lancichinetti	0.54349	0.00576
SGF(0.3)	Add-Health	0.75379	0.07496	Facebook	0.87341	0.08013
SGF(0.3)	Girvan	0.70690	0.01553	Lancichinetti	0.55989	0.00530
SGF(0.4)	Add-Health	0.80250	0.10278	Facebook	0.86587	0.06847
SGF(0.4)	Girvan	0.73735	0.00793	Lancichinetti	0.59954	0.00445
SGF(0.5)	Add-Health	0.82452	0.08711	Facebook	0.89184	0.06929
SGF(0.5)	Girvan	0.77007	0.01134	Lancichinetti	0.64635	0.00339
SGF(0.6)	Add-Health	0.86057	0.08472	Facebook	0.91149	0.06688
SGF(0.6)	Girvan	0.83302	0.01151	Lancichinetti	0.71476	0.00564
SGF(0.7)	Add-Health	0.88024	0.05368	Facebook	0.93420	0.07223
SGF(0.7)	Girvan	0.88050	0.01007	Lancichinetti	0.78831	0.00714
SGF(0.8)	Add-Health	0.93563	0.03820	Facebook	0.96249	0.02018
SGF(0.8)	Girvan	0.92902	0.00797	Lancichinetti	0.87126	0.00553
SGF(0.9)	Add-Health	1.02524	0.03708	Facebook	0.98131	0.05789
SGF(0.9)	Girvan	1.02722	0.00429	Lancichinetti	1.05135	0.00338
DC-SBM	Add-Health	0.88304	0.09229	Facebook	0.89041	0.11375
DC-SBM	Girvan	0.88989	0.00892	Lancichinetti	0.90944	0.00462
Trajanovski	Add-Health	0.71150	0.08293	Facebook	1.38715	0.45120
Trajanovski	Girvan	0.97946	0.31149	Lancichinetti	0.79338	0.00097

to the DC-SBM algorithm, instead, we use the node degree sequence, the node group assignment  $c_1^*, \dots, c_n^*$ , the group degree sequence on this partition and the number of edges between each pair of groups, computed on the original graph.

*Comparison.* Table I reports the mean and the standard deviation of obtained modularity ratios comparing the methods with different values of  $\alpha$  on every dataset and it shows the Spectral Graph Forge requires  $\alpha \sim 0.7, 0.8$  in order to obtain results similar to the one of DC-SBM and that Trajanovski obtains results with a very high standard deviation from the target value. Fig. 4 reports the mean and 99% confidence intervals for the resulting modularity ratios, for which SGF with  $\alpha = 0.9$  consistently produces graphs with a maximum modularity nearly equal to the input data set (modularity ratio appx. 1). In contrast, the Trajanovski et al. algorithm frequently produces graphs with too high or too low modularity, and does so with high variance. The DC-SBM performs more consistently, but tends to achieve modularity lower than the target, and generally shows higher variance. Interestingly, both the Trajanovski et al. and DC-SBM algorithms perform better on the synthetic data sets than in their real-world counterparts. Across all four data sets, the Spectral Graph Forge is consistently closer to the target.

Fig. 5 looks at a related but different metric: it presents means and 99% confidence intervals for the ratio of *number*

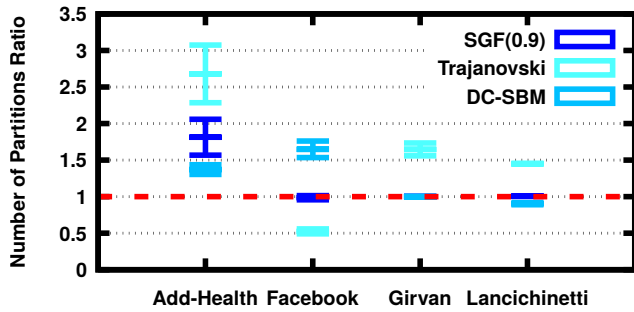


Fig. 5. Means and 99% confidence intervals for the partition number ratio, by data set. Note that the SGF(0.9) and DC-SBM algorithm results overlap with the Girvan data set graphs.

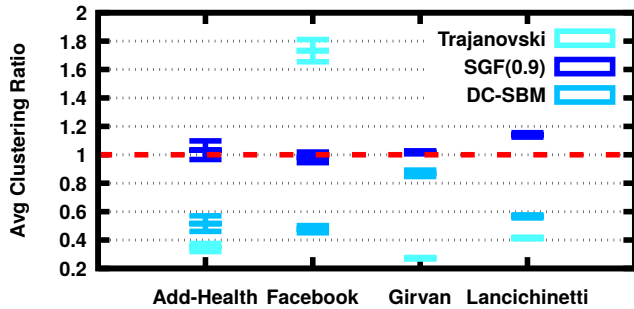


Fig. 6. Means and 99% confidence intervals for average clustering ratio, by data set.

of detected output partitions (under modularity maximization) vs. the number of input partitions; this expresses the extent to which each algorithm preserves the number of communities found in the input graph. The target value is again 1, as our objective is to generate graphs with the same number of communities. Spectral Graph Forge consistently reproduces the number of input partitions for three out of the four data sets, but tends to split communities within AddHealth. None of the other algorithms consistently performs as well as the Spectral Graph Forge, with some being tied on some data sets; the only case of better performance by a competing algorithm is DC-SBM on AddHealth, which overproduces communities by a slightly smaller margin.<sup>4</sup>

### C. Results on Other (Local) Structural Properties

In addition to accurately targeting modularity, our approach also generates output graphs  $A' = SGF(\alpha)$  maintaining other important (local structural) properties of the input  $A$ .

*Clustering* is a local network property, indicating the extent of triadic closure. Fig. 6 presents means and confidence intervals for the average clustering ratio between the output and input graphs. Although it is not apparent that the Spectral Graph Forge would do well at capturing this local property -

<sup>4</sup>None of the three algorithms perform optimally on this data set. That can be due to input partition solutions not necessarily stable to small perturbations, and a generated graph could have different numbers of identified subgroups than the source graph even if the two are similar in most respects.

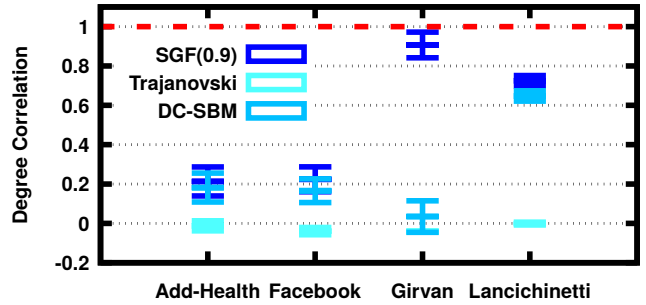


Fig. 7. Means and 99% confidence intervals for degree sequence correlation, by data set. Note that the DC-SBM result lower bound overlaps the Trajanovski results with the Girvan dataset.

since it targets global structure - it in fact does well for all four data sets. By contrast, the other modularity-targeting methods do not tend to preserve clustering. This is particularly evident for the Trajanovski et al. algorithm, whose rewiring strategy can easily alter the structure of triads.

*Degree sequence.* We also examined the correlation of the degree sequence of the input and output graphs. (Recall that, in the presence of covariates, nodes are non-exchangeable.) While the eigenvectors of the modularity matrix are only loosely connected with degree, we show better performance on this metric than the competing methods. In particular, Trajanovski doesn't attempt to preserve the degree distribution at all, leading (reasonably enough) to a nearly complete loss of correlation between degree sequences. The relatively poor performance of the DC-SBM method is more surprising, given that it attempts to preserve degree information. Overall, it performs only slightly worse than the Spectral Graph Forge for three out of four data sets, the lone exception being the Girvan synthetic networks (where the Spectral Graph Forge has a correlation of nearly 1, vs. almost 0 for DC-SBM).

### D. Results on Attribute Modularity Preservation

We find that our approach not only successfully preserves the maximum modularity value  $Q^*$  but also preserves the modularity of significant partitions other than the topologically best one. If we consider partitions  $c_1, \dots, c_n$  indicating e.g. gender, race/ethnicity, or school grade provided by the Add-Health dataset, we can compute the resulting modularity values  $Q_i$ , on the original graphs, and  $Q_o$  on the output graphs after applying the same attribute sequence, summarizing the similarity by the resulting modularity ratio. The closer this ratio is to 1, the better maintained the community structure for partitions associated with important covariates (whether or not they correspond to modularity-maximizing partitions).

Fig. 8 shows the comparison. The Add-Health dataset is a real-world dataset with attributes associated with the nodes. Attribute labeling for the output graphs of Trajanovski and DC-SBM is done trivially, by assigning attributes to the nodes in the same order as the input ones. As expected, the Spectral Graph Forge closely preserves modularity on all

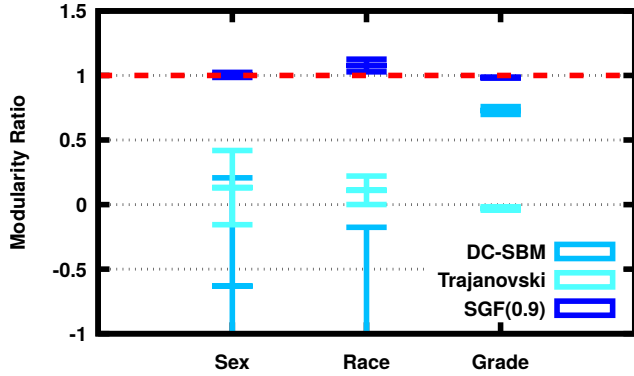


Fig. 8. 99% confidence interval on attribute modularity ratio for the Add-Health dataset.

attribute partitions, while the other methods typically perform poorly, entirely losing the graph structure. This is not surprising as the other methods were not designed to preserve other community structure but the topological one and their performance are shown only for completeness. On contrary, SGF successfully reproduce all the possible (even overlapping) community structures from the original ones, being, to the best of our knowledge, the first suitable generator for creating graph proxies with complex community structures.

#### E. Randomness of Realizations and (De)anonymization

The closer  $\alpha$  is to 1, the better our approach targets the intended modularity. However, two natural questions arise from this relationship: (1) as  $\alpha \rightarrow 1$ , how concentrated is the distribution of graphs  $A' = SGF(\alpha)$  we produce? (2) as  $\alpha \rightarrow 1$ , how distinct is our output  $A'$  from the input  $A$ ? Ideally, we would like a relatively low level of concentration and a reasonable distance from the input graph, to grant a sufficient degree of randomness in our realizations. We use the normalized Shannon entropy to measure the concentration of the SGF distribution in the space of graphs; this is important for the use case of simulation. To measure the distinctness of a realized graph from the input network, we consider a pragmatic metric that maps to another important potential use case: *resistance to de-anonymization attacks*. We envision that Spectral Graph Forge could be used not only as a graph generator for simulating graphs resembling real-world networks, but also as an anonymization technique (in conjunction with others) to anonymize  $A$  to  $A'$ , while preserving several global and local properties.

De-anonymization attacks attempt to identify the nodes in a partially labeled graph, exploiting similarities between the two graphs and potentially auxiliary information. To the extent that a generated network cannot be readily de-anonymized with respect to the original graph structure, the two networks are clearly distinct; this is important for applications such as anonymization of real network datasets, where concealment of node identities is a major priority. For example, a Facebook dataset may be collected and made publicly available after

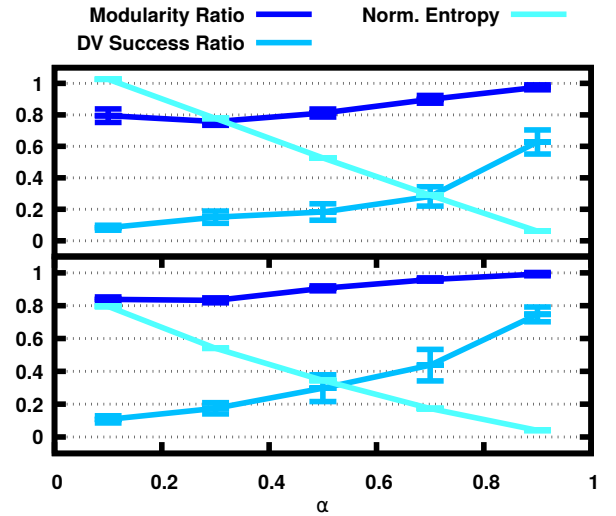


Fig. 9. Means and 99% confidence intervals for graph distribution entropy and DV de-anonymization success rates on one Add-Health graph (upper plot) and one Facebook graph (lower plot), by  $\alpha$ .

proper anonymization (removing node ids, and perturbing the network structure); de-anonymization attacks can infer the node ids based on graph structure and auxiliary information. A survey of state of the art anonymization and de-anonymization techniques can be found in [35]. We pick one of the state-of-the-art de-anonymization attacks: the Distance Vector attack [36], which is proven to be scalable, robust and exploitative of global graph characteristics [35] preserved by the SGF. At each run, we feed the Distance Vector attack with a seed of 5% of nodes as ground truth and we report the mean fraction of nodes the attack successfully identified. Examining the fraction of nodes identified as a function of  $\alpha$  allows us to examine the tradeoff between privacy preservation and synthetic data quality.

The upper plot of Fig. 9 presents the results for one connected graph from the Add-Health dataset. We note that the modularity ratio remains close to the target 1 even for low values of  $\alpha$ , implying that it is possible to preserve much of the community structure with a fairly small fraction of eigenvectors. Entropy is close to 0 for very large  $\alpha$ , but climbs steadily as  $\alpha$  is reduced, indicating the extent of reduction in concentration. More striking is the success rate of the DV attack. Even at  $\alpha = 0.9$ , the attack can identify only around 60% of nodes from the original graph given the synthetic graph and a 5% seed sample. This falls apace, with a success rate of only around 20% once  $\alpha = 0.5$  (corresponding to a modularity ratio  $\sim 0.82$  from Table I), and a rate not significantly different from baseline once  $\alpha = 0.25$  (modularity ratio  $\in (0.71, 0.75)$ ). Given the strength of this attack, this implies a high degree of discrepancy between the initial and target graph (while maximum modularity is preserved). Similar results are reported in the lower plot of Fig. 9 for a connected graph from the Facebook dataset. Again, even for high values of  $\alpha$  the de-anonymization attack cannot fully align the graphs.



For both datasets, when  $\alpha = 0.1$ , the modularity ratio is still close to the target 1 and the de-anonymization attack can identify only approximately 10% of nodes (while knowing 5% in the first place). This suggests that the SGF has the potential to be used as a tool (in combination with others [35]) for anonymizing sensitive graph datasets.

## VI. CONCLUSION

We propose Spectral Graph Forge (SGF) - a framework for generating graphs that resemble real-world networks in terms of global properties, via spectral decomposition and approximation. Our focus in this paper was the use of SGF specifically to generate graphs with a target modularity; we showed that it succeeds in doing so and it outperforms baselines. At the same time, SGF preserves other important local structure properties (such as clustering and the degree sequence) and node attribute structure. It also does well at preserving modularity associated with such node attributes, even when they are not specifically targeted; this allows SGF to generate reliable graph proxies for real-world networks with arbitrarily complex and overlapping community structures. The parameter  $\alpha$  can control not only the error in modularity ( $\alpha \rightarrow 1$  brings us arbitrarily close to the target modularity) but also the entropy of graph distribution. Lowering  $\alpha$  leads to higher entropy of the realizations, thus making the node identities in the input graph increasingly difficult to reveal via de-anonymization attacks. This makes SGF a potentially useful tool (in combination with existing anonymization techniques) for generating synthetic data sets from sensitive data. Our codebase is available on-line [37].

## REFERENCES

- [1] P. Mahadevan, D. Krioukov, K. Fall, and A. Vahdat, "Systematic topology analysis and generation using degree correlations," *Acm Sigcomm Computer Communication Review*, vol. 36, pp. 135–146, 2006.
- [2] C. Orsini, M. M. Dankulov, P. Colomer-de Simón, A. Jamakovic, P. Mahadevan, A. Vahdat, K. E. Bassler, Z. Toroczka, M. Boguñá, G. Caldarelli *et al.*, "Quantifying randomness in real networks," *Nature communications*, 2015.
- [3] M. Gjoka, B. Tillman, and A. Markopoulou, "Construction of simple graphs with a target joint degree matrix and beyond." 2015 IEEE Conference on Computer Communications (INFOCOM), 2015.
- [4] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon, "Network motifs: simple building blocks of complex networks," *Science*, 2002.
- [5] Ö. N. Yaveroglu, N. Malod-Dognin, D. Davis, Z. Levnajic, V. Janjic, R. Karapandza, A. Stojmirovic, and N. Przulj, "Revealing the hidden language of complex networks," *Scientific reports*, 2014.
- [6] B. Donnet and T. Friedman, "Internet topology discovery: a survey," *IEEE Communications Surveys Tutorials*, vol. 9, no. 4, pp. 56–69, 2007.
- [7] M. H. Gunes and K. Sarac, "Resolving Anonymous Routers in Internet Topology Measurement Studies," in *IEEE INFOCOM 2008 - The 27th Conference on Computer Communications*, April 2008.
- [8] A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, and B. Bhattacharjee, "Measurement and Analysis of Online Social Networks," in *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement*, ser. IMC '07. New York, NY, USA: ACM, 2007, pp. 29–42. [Online]. Available: <http://doi.acm.org/10.1145/1298306.1298311>
- [9] F. Benevenuto, T. Rodrigues, M. Cha, and V. Almeida, "Characterizing User Behavior in Online Social Networks," in *IMC*, ser. IMC '09. New York, NY, USA: ACM, 2009, pp. 49–62. [Online]. Available: <http://doi.acm.org/10.1145/1644893.1644900>

- [10] R. Aldecoa, C. Orsini, and D. Krioukov, "Hyperbolic graph generator," *Computer Physics Communications*, vol. 196, pp. 492 – 496, 2015. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0010465515002088>
- [11] M. E. J. Newman, *Networks: an introduction*. Oxford university press, 2010.
- [12] S. Wasserman and K. Faust, *Social Network Analysis: Methods and Applications*. Cambridge: Cambridge University Press, 1994.
- [13] S. J. Axler, *Linear algebra done right*. Springer, 1997.
- [14] L. Mirsky, "Symmetric gauge functions and unitarily invariant norms," *The quarterly journal of mathematics*, vol. 11, pp. 50–59, 1960.
- [15] G. Robins, P. Pattison, Y. Kalish, and D. Lusher, "An introduction to exponential random graph (p\*) models for social networks," *Social Networks*, vol. 29, no. 2, pp. 173 – 191, 2007, special Section: Advances in Exponential Random Graph (p\*) Models. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0378873306000372>
- [16] S. Trajanovski, F. A. Kuipers, J. Martín-Hernández, and P. Van Mieghem, "Generating graphs that approach a prescribed modularity," *Computer Communications*, vol. 36, no. 4, 2013.
- [17] B. Karrer and M. E. J. Newman, "Stochastic blockmodels and community structure in networks," *Physical Review E*, vol. 83, no. 1, 2011.
- [18] P. D. Hoff, A. E. Raftery, and M. S. Handcock, "Latent Space Approaches to Social Network Analysis," *Journal of the American Statistical Association*, vol. 97, no. 460, pp. 1090–1098, 2002.
- [19] K. Nowicki and T. A. B. Snijders, "Estimation and Prediction for Stochastic Blockstructures," *Journal of the American Statistical Association*, vol. 96, no. 455, pp. 1077–1087, 2001.
- [20] M. E. J. Newman, "Spectral methods for community detection and graph partitioning," *Phys. Rev. E*, vol. 88, p. 042822, Oct 2013.
- [21] W. D. Richards and A. J. Seary, "Multinet. Version 4.38 for Windows," 2003.
- [22] M. E. J. Newman, "Modularity and community structure in networks," *PNAS*, vol. 103, pp. 8577–8582, June 2006.
- [23] J. S. Coleman, *Introduction to mathematical sociology*. London Free Press Glencoe, 1964.
- [24] L. Freeman, *The development of social network analysis: a Study in the Sociology of Science*. BookSurgePublishing, 2004.
- [25] S. P. Borgatti and M. G. Everett, "Models of core/periphery structures," *Social Networks*, vol. 21, pp. 375–395, 1999.
- [26] L. Baldesi, A. Markopoulou, and C. T. Butts, "Spectral Graph Forge - Extended Version," <https://arxiv.org/abs/1801.01715>, 2018.
- [27] C. D. Meyer, *Matrix analysis and applied linear algebra*. Siam, 2000.
- [28] M. Girvan and M. E. J. Newman, "Community structure in social and biological networks," *Proceedings of the national academy of sciences*, vol. 99, no. 12, pp. 7821–7826, 2002.
- [29] A. Lancichinetti, S. Fortunato, and F. Radicchi, "Benchmark graphs for testing community detection algorithms," *Physical review E*, vol. 78, no. 4, p. 046110, 2008.
- [30] K. M. Harris, C. T. Halpern, E. Whitsel, J. Hussey, J. Tabor, P. Entzel, and J. R. Udry, "The national longitudinal study of adolescence to adult health: Research design." See <http://www.epc.unc.edu/projects/addhealth/design> (accessed 9 April 2015), 2009.
- [31] J. J. McAuley and J. a. Leskovec, "Learning to discover social circles in ego networks," in *Advances in neural information processing systems*, 2012, pp. 539–547.
- [32] T. A. B. Snijders and K. Nowicki, "Estimation and prediction for stochastic block models for graphs with latent block structure," *Journal of Classification*, vol. 14, pp. 75–100, 1997.
- [33] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, p. P10008, 2008.
- [34] S. Fortunato, "Community detection in graphs," *Physics Reports*, vol. 486, no. 3, pp. 75 – 174, 2010. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0370157309002841>
- [35] S. Ji, W. Li, P. Mittal, X. Hu, and R. A. Beyah, "SecGraph: A Uniform and Open-source Evaluation System for Graph Data Anonymization and De-anonymization." in *Usenix Security*, 2015, pp. 303–318.
- [36] M. Srivatsa and M. Hicks, "Deanonymizing mobility traces: Using social network as a side-channel," in *Proceedings of the 2012 ACM conference on Computer and communications security*. ACM, 2012, pp. 628–637.
- [37] L. Baldesi, "Spectral Graph Forge - Source Code," [https://github.com/lucabaldesi/spectral\\_graph\\_forge](https://github.com/lucabaldesi/spectral_graph_forge), December 2017.