

Graph Clustering: Block-models and model free results

Authored by:

Marina Meila
Yali Wan

Abstract

Clustering graphs under the Stochastic Block Model (SBM) and extensions are well studied. Guarantees of correctness exist under the assumption that the data is sampled from a model. In this paper, we propose a framework, in which we obtain "correctness" guarantees without assuming the data comes from a model. The guarantees we obtain depend instead on the statistics of the data that can be checked. We also show that this framework ties in with the existing model-based framework, and that we can exploit results in model-based recovery, as well as strengthen the results existing in that area of research.

1 Paper Body

: a framework for clustering with guarantees without model assumptions

In the last few years, model-based clustering in networks has witnessed spectacular progress. At the central of intact are the so-called block-models, the Stochastic Block Model (SBM), DegreeCorrected SBM (DC-SBM) and Preference Frame Model (PFM). The understanding of these models has been advanced, especially in understanding the conditions when recovery of the true clustering is possible with small or no error. The algorithms for recovery with guarantees have also been improved. However, the impact of the above results is limited by the assumption that the observed data comes from the model. This paper proposes a framework to provide theoretical guarantees for the results of model based clustering algorithms, without making any assumption about the data generating process. To describe the idea, we need some notation. Assume that a graph G on n nodes is observed. A modelbased algorithm clusters G , and outputs clustering C and parameters $M(G, C)$.

The framework is as follows: if $M(G, C)$ fits the data G well, then we shall prove that any other clustering C' of G that also fits G well will be a small perturbation of C . If this holds, then C with model parameters $M(G, C)$ can be said to capture the data structure in a meaningful way. We exemplify our approach

by obtaining model-free guarantees for the SBM and PFM models. Moreover, we show that model-free and model-based results are intimately connected.

2

Background: graphs, clusterings and block models

Graphs, degrees, Laplacian, and clustering Let G be a graph on n nodes, described by its adjacency matrix $A = [A_{ij}]$ the degree of node i , and $D = \text{diag}\{d_i\}$ the diagonal degree matrix. A is symmetric. In matrix of the node degrees. The (normalized) Laplacian of G is defined as $L = D - A$ [10]. Rigorously speaking, the normalized graph Laplacian is $I - L$.

30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

extension, we define the degree matrix D and the Laplacian L associated to any matrix $A \in \mathbb{R}^{n \times n}$, with $A_{ij} = A_{ji} \geq 0$, in a similar way.

Let C be a partitioning (clustering) of the nodes of G into K clusters. We use the shorthand notation $i \in k$ for "node i belongs to cluster k ". We will represent C by its $n \times K$ indicator matrix Z , defined by $Z_{ik} = 1$ if $i \in k$, 0 otherwise, for $i = 1, \dots, n$, $k = 1, \dots, K$. (1) $Z^T Z = \text{diag}\{n_k\}$ with n_k counting the number of nodes in cluster k , and $Z A Z^T = [n_{kl}]_{k,l=1}^K$ with n_{kl} counting the edges in G between clusters k and l . Moreover, for two indicator matrices Z, Z' for clusterings C, C' , $(Z^T Z')_{kk}$ counts the number of points in the intersection of C and C' , $(Z^T Z')_{kl}$ counts the number of points in the intersection of C and C' , and $(Z^T D Z')_{kk}$ counts the volume of the same intersection.

Block models for random graphs (SBM, DC-SBM, PFM) This family of models contains Stochastic Block Models (SBM) [1, 18], Degree-Corrected SBM (DC-SBM) [17] and Preference Frame Models (PFM) [20]. Under each of these model families, a graph G with adjacency matrix A over n nodes is generated by sampling its edges independently following the law $A_{ij} \sim \text{Bernoulli}(A_{ij})$, for all $i \neq j$. The symmetric matrix $A = [A_{ij}]$ describing the graph is the edge probability matrix. The three model families differ in the constraints they put on an acceptable A . Let C be a clustering. The entries of A are defined w.r.t C as follows (and we say that A is compatible with C). SBM : $A_{ij} = B_{kl}$ whenever $i \in k, j \in l$, with $B = [B_{kl}] \in \mathbb{R}^{K \times K}$ symmetric and nonnegative. DC-SBM : $A_{ij} = w_i w_j B_{kl}$ whenever $i \in k, j \in l$, with B as above and w_1, \dots, w_n non-negative weights associated with the graph nodes. PFM : A satisfies $D = \text{diag}(A \mathbf{1})$, $D^{-1} A Z = Z R$ where $\mathbf{1}$ denotes the vector of all ones, Z is the indicator matrix of C , and R is a stochastic matrix ($R \mathbf{1} = \mathbf{1}$, $R_{kl} \geq 0$), the details are in [20] While perhaps not immediately obvious, the SBM is a subclass of the DC-SBM, and the latter a subclass of the PFM. Another common feature of block-models, that will be significant throughout this work is that for all three, Spectral Clustering algorithms [15] have been proved to work well estimating C .

3

Main theorem: blueprint and results for PFM, SBM

Let M be a model class, such as SBM, DC-SBM, PFM, and denote $M(G, C)$ to be a model that is compatible with C and is fitted in some way to graph G (we do not assume in general that this fit is optimal). Theorem 1 (Generic

(6)

3. Change basis in $R(YZ)$ to align with Y . $Y = YZ U V T$.

Complete Y to an orthonormal basis $[Y B]$ of R_n .

(7)

4. Construct Laplacian L and edge probability matrix A . $T \leftarrow T + (BB^T - L)(BB^T - L = Y^T Y)$,

$\frac{1}{2} LD \leftarrow \frac{1}{2} A = D$

(8)

D, L, Y, Z be defined as above, and $(A, L) = PFM(G, C)$.

Then, Proposition 2 Let G, A, L , or A define a PFM with degrees $d_{1:n}$.
 1. $D \leftarrow \frac{1}{K}$. 2. The columns of Y are eigenvectors of L with eigenvalues λ_i .
 $\lambda_1 = 1$. $\frac{1}{2} LD \leftarrow \frac{1}{2} A$ is an eigenvector of both L and L^T with eigenvalue $\frac{1}{2}$. 3. D The proof is relegated to the Supplement, as are all the omitted proofs. $PFM(G, C)$ is an estimator for the PFM parameters given the clustering. It is evidently not the Maximum Likelihood estimator, but we can show that it is consistent in the following sense. 3

Proposition 3 (Informal) Assume that G is sampled from a PFM with parameters D, L and compatible with C , and let $L = PFM(G, C)$. Then, under standard recovery conditions for PFM (e.g [20]) $L \rightarrow L$ $\frac{1}{n} \rightarrow 0$ w.r.t. n . $L \rightarrow L$. Assumption 2 (Goodness of fit for PFM) $\rightarrow PFM(G, C)$ instantiates $M(G, C)$, and Assumption 2 instantiates the goodness of fit measure. It remains to prove an instance of Generic Theorem 1 for these choices. $\frac{1}{n}$ as defined, and $L \rightarrow L$ satisfy Theorem 4 (Main Result (PFM)) Let G be a graph with $d_{1:n}$, D , satisfy Assumption 1. Let C, C^T be two clusterings with K clusters, and L, L^T be their corresponding 2 and Laplacians, defined as in (8), and satisfy Assumption 2 respectively. Set $\epsilon = (\frac{1}{K} - \frac{1}{K+1})^2$

$K+1$

$\epsilon_0 = \min_k C_{kk} / \max_k C_{kk}$ with C defined as in (5), where k indexes the clusters of C . Then, whenever $\epsilon \leq \epsilon_0$, $\max_k C_{kk} \leq \epsilon \cdot \text{dist}_w(C, C^T) \leq k C_{kk}$

(9)

with dist_w being the weighted ME distance (3).

In the remainder of this section we outline the proof steps, while the partial results of Proposition 5, 6, 7 are proved in the Supplement. First, we apply the perturbation bound called the Sinus Theorem of Davis and Kahan, in the form presented in Chapter V of [19]. $\frac{1}{n}$, Y be defined as usual. If Assumptions 1 and 2 hold, then Proposition 5 Let $Y^T, Y \rightarrow \text{diag}(\sin \theta_{1:K}(Y^T, Y))$ $\epsilon = \frac{1}{K+1} - \frac{1}{K} = \frac{1}{K(K+1)}$

(10)

where $\theta_{1:K}$ are the canonical (or principal) angles between $R(Y^T)$ and $R(Y)$ (see e.g [8]). The next step concerns the closeness of Y, Y^T in Frobenius norm. Since Proposition 5 bounds the sines of the canonical angles, we exploit the fact that the cosines of the same angles are the singular values of $F = Y^T Y^T$ of (6). $\epsilon = Y^T Y^T T$ and F, ϵ as above. Assumptions 1 and 2 imply that Proposition 6 Let $M = Y Y^T, M \leftarrow T \leftarrow K \leftarrow (K-1)^2$. 1. $\rightarrow F \rightarrow 2F = \text{trace } M M^T \rightarrow 2 \leftarrow 2(K-1)^2$. 2. $\rightarrow M \leftarrow M F$ Now we show that

all clusterings which satisfy Proposition 6 must be close to each other in the weighted ME distance. For this, we first need an intermediate result. Assume we have two clusterings C, C' , with K clusters, for which we construct YZ, Y, L, M , respectively YZ', Y', L', M' as above. Then, the subspaces spanned by Y and Y' will be close. C, C' satisfy Assumption 1 and let C, C' represent two clusterings for which L, L' Proposition 7 Let L satisfy Assumption 2. Then, $\|YZ - YZ'\|_F \leq \sqrt{2} \sqrt{K} \sqrt{4(K-1)} = 2\sqrt{K}$. The main result now follows from Proposition 7 and Theorem 9 of [13], as shown in the Supplement. This proof approach is different from the existing perturbation bounds for clustering, which all use counting arguments. The result of [13] is a local equivalence, which bounds the error we need in terms of ϵ defined above (ϵ local meaning the result only holds for small ϵ). 4

3.2

Main Theorem for SBM

In this section, we offer an instantiation of Generic Theorem 1 for the case of the SBM. As before, we start with a model estimator, which in this case is the Maximum Likelihood estimator. SBM Estimation Algorithm clustering C with indicator matrix Z . Input Graph with A , Output $A = \text{SBM}(G, C)$

1. Construct an orthogonal matrix derived from Z : $YZ = ZC^{1/2}$ with $C = Z^T Z$. 2. Estimate the edge probabilities: $B = C^{1/2} Z^T A Z C^{1/2}$.
3. Construct A from B by $A = ZBZ^T$. $B = C^{1/2} B C^{1/2}$ and denote the eigenvalues of B , λ_i ordered by decreasing magnitude. Proposition 8 Let B be $B = U \Lambda U^T$, with U an orthogonal matrix nitude, by $1:K$. Let the spectral decomposition of B and $\Lambda = \text{diag}(\lambda_1:K)$. Then 1. A is a SBM. 2. $\lambda_1:K$ are the K principal eigenvalues of A . The remaining eigenvalues of A are zero. 3. $A = Y Y^T$ where $Y = YZ U$. Assumption 3 (Eigengap) B is non-singular (or, equivalently, $\lambda_K \gg 0$). Assumption 4 (Goodness of fit for SBM) $\|A - \hat{A}\|_F \leq \epsilon$. With the model (SBM), estimator, and goodness of fit defined, we are ready for the main result. Theorem 9 (Main Result (SBM)) Let G be a graph with incidence matrix A , K singular value of A . Let C, C' be two clusterings with K clusters, satisfying Assumptions 3 and 4. 2 Set $\epsilon = \frac{4K}{\lambda_K - \lambda_{K+1}}$ and $\epsilon_0 = \min_k \lambda_k / \max_k \lambda_k$, where k indexes the clusters of C . Then, whenever $\epsilon \leq K$

$\epsilon \leq \epsilon_0$, $\text{dist}(C, C') \leq \max_k \lambda_k / n$, where dist represents the ME distance (2).

Since the SBM is less Note that the eigengap of A , K exible than the PFM, we expect that for the same data G , Theorem 9 will be more restrictive than Theorem 4.

4.1

The results in perspective Cluster validation

Theorems like 4, 9 can provide model free guarantees for clustering. We exemplify this procedure in the experimental Section 6, using standard spectral clustering as described in e.g [18, 17, 15]. What is essential is that all the quantities such as ϵ and ϵ_0 are computable from the data. Moreover, if Y is available, then the bound in Theorem 4 can be improved. $\epsilon, M \leq F + (K-1)(\epsilon^2 + \epsilon)$ Proposition 10 Theorem 4 holds when ϵ is replaced by $\epsilon Y = K \epsilon M$.

simply replace the model assumption with the assumption that there is a C for which L (or A) satisfies Assumptions 1 and 2 (or 3 and 4).

5

Related work

To our knowledge, there is no work of the type of Theorem 1 in the literature on SBM, DC-SBM, PFM. The closest work is by [6] which guarantees approximate recovery assuming G is close to a DC-SBM. Spectral clustering is also used for loss-based clustering in (weighted) graphs and some stability results exist in this context. Even though they measure clustering quality by different criteria, so that the ϵ values are not comparable, we review them here. The recent paper of [16], Theorem 1.2 states $\epsilon = O(K^{-1/3})$ then the clustering error is $O(K^{-1/3})$ that if the K -way Cheeger constant of G is $\lambda_2(G) \geq \lambda_1(G) + c$ then the clustering error is $O(K^{-1/3})$. In the current proof, the constant $C = 2 \times 10^5$; moreover, $\lambda_2(G)$ cannot be computed tractably. In [14], the bound ϵ_{MSE} depends on $\lambda_2(G)$, the Normalized Cut scaled by the eigengap. Since both bounds refer to the result of spectral clustering, we can compare the relationship between ϵ_{MSE} and ϵ_{MSE} ; for [14], this is $\epsilon_{\text{MSE}} = 2\epsilon_{\text{MSE}} [1 + \epsilon_{\text{MSE}} / (K - 1)]$, 2 The results is stronger, bounding the perturbation of each cluster individually by ϵ_{MSE} , but it also includes a factor larger than 1, bounding the error of K -means algorithm.

6

which is about $K - 1$ times larger than ϵ when $\epsilon = \epsilon_{\text{MSE}}$. In [5], $\text{dist}(C, C^*)$ is defined in terms of $\|Y - YZ\|_F^2$, and the loss is (closely related) to $\|A - \hat{A}\|_F^2$. The bound does not take into account the eigengap, that is, the stability of the subspace Y itself. Bootstrap for validating a clustering C was studied in [11] (see also references therein for earlier work). In [3] the idea is to introduce a statistics, and large deviation bounds for it, conditioned on sampling from a SBM (with covariates) and on a given C .

6

Experimental evaluation

Experiment Setup Given G , we obtain a clustering C_0 by spectral clustering [15]. Then we calculate clustering C by perturbing C_0 with gradually increasing noise. For each C , we construct PFM (C, G) and SBM (C, G) model, and further compute ϵ , ϵ_{MSE} and ϵ_{MSE} . If $\epsilon \leq \epsilon_{\text{MSE}}$, C is guaranteed to be stable by the theorems. In the remainder of this section, we describe the data generating process for the simulated datasets and the results we obtained. **PFM Datasets** We generate from PFM model with $K = 5$, $n = 10000$, $\theta_{1:K} = (1, 0.875, 0.75, 0.625, 0.5)$. eigengap = 0.48, $n_{1:K} = (2000, 2000, 2000, 2000, 2000)$. The stochastic matrix R and its stationary distribution π are shown below. We sample an adjacency matrix A from A (shown below).

$A = \begin{bmatrix} 0.25 & 0.12 & 0.17 & 0.18 & 0.28 & 0.79 & 0.02 & 0.06 & 0.03 & 0.10 & 0.71 & 0.23 & 0.00 & 0.02 & 0.03 & 0.16 & 0.69 & 0.00 \\ 0.06 & 0.09 & 0.04 & 0.00 & 0.00 & 0.80 & 0.16 & 0.10 & 0.01 & 0.03 & 0.11 & 0.76 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \end{bmatrix}$

$\pi =$

$\pi =$

$A =$

$\pi \pi \pi \pi \pi$

Perturbed PFM Datasets A is obtained from the previous model by perturbing its principal subspace (details in Supplement). Then we sample A^\dagger from A . Lancichinetti-Fortunato-Radicchi (LFR) simulated matrix [12] The LFR benchmark graphs are widely used for community detection algorithms, due to heterogeneity in the distribution of node degree and community size. A LFR matrix is simulated with $n = 10000$, $K = 4$, $n_k = (2467, 2416, 2427, 2690)$ and $\gamma = 0.2$, where γ is the mixing parameter indicating the fraction of edges shared between a node and the other nodes from outside its community. γ of hyperlinks between weblogs on US politics, Political Blogs Dataset A directed network A compiled from online directories by Adamic and Glance [2], where each blog is assigned a political leaning, liberal or conservative, based on its blog content. The network A contains 1490 blogs. $\gamma(A) \approx 0.3$, which is a smoothed After erasing the disconnected nodes, $n = 983$. We study $A^\dagger = (A + A^T)/2$ undirected graph. For A we find no guarantees. The first two data sets are expected to fit the PFM well, but not the SBM, while the LFR data is expected to be a good fit for a SBM. Since all bounds can be computed on weighted graphs as well, we have run the experiments also on the edge probability matrices A used to generate the PFM and perturbed PFM graphs. The results of these experiments are summarized in Figure 1. For all of the experiments, the clustering C is ensured to be stable by Theorem 4 as the unweighted error grows to a breaking point, then the assumptions of the theorem fail. In particular, the C_0 is always stable in the PFM framework. 7

Comparing γ from Theorem 9 to that from Theorem 4, we find that Theorem 9 (guarantees for SBM) is much harder to satisfy. All γ values from Theorem 9 are above 1, and not shown. In particular, for the SBM model class, the C cannot be proved stable even for the LFR data.

Note that part of the reason why with the PFM model very little difference from the clustering C_0 can be tolerated for a clustering to be stable is that the large eigengap makes $PFM(G, C)$ differ from $PFM(G, C_0)$ even for very small perturbations. By comparing the bounds for A^\dagger with the bounds for the weighted graphs A , we can evaluate that the sampling noise on A^\dagger is approximately equal to that of the clustering perturbation. Of course, the sampling noise varies with n , decreasing for larger graphs. Moreover, from Political Blogs data, we see that “smoothing” a graph, by e.g. taking powers of its adjacency matrix, has a stability inducing effect.

Figure 1:

γ denotes a simple graph, while A denotes a Quantities γ , γ , γ_0 from Theorem 4 plotted vs $\text{dist}(C, C_0)$ for various datasets: A

weighted graph (i.e. a non-negative matrix). For the Political Blogs: Truth means C_0 is true clustering of [2], spectral means C_0 is obtained from spectral clustering. For SBM, γ is always greater than γ_0 .

7

Discussion

This paper makes several contributions. At a high level, it poses the problem of model free validation in the area of community detection in networks. The stability paradigm is not entirely new, but using it explicitly with model-based

clustering (instead of cost-based) is. So is “turning around” the model-based recovery theorems to be used in a model-free framework. All quantities in our theorems are computable from the data and the clustering C , i.e. do not contain undetermined constants, and do not depend on parameters that are not available. As with distribution-free results in general, making fewer assumptions allows for less confidence in the conclusions, and the results are not always informative. Sometimes this should be so, e.g. when the data does not “fit” the model well. But it is also possible that the “fit” is good, but not good enough to satisfy the conditions of the theorems as they are currently formulated. This happens with the SBM bounds, and we believe tighter bounds are possible for this model. It would be particularly interesting to study the non-spectral, sharp thresholds of [1] from the point of view of model-free recovery. A complementary problem is to obtain negative guarantees (i.e. that C is not unique up to perturbations). At the technical level, we obtain several different and model-specific stability results, that bound the perturbation of a clustering by the perturbation of a model. They can be used both in model-free and in existing or future model-based recovery guarantees, as we have shown in Section 3 and in the experiments. The proof techniques that lead to these results are actually simpler, more direct, and more elementary than the ones found in previous papers. 3

We also computed Ψ_{RCY} but the bounds were not informative.

8

2 References

- [1] Emmanuel Abbe and Colin Sandon. Community detection in general stochastic block models: fundamental limits and efficient recovery algorithms. arXiv preprint arXiv:1503.00609, 2015.
- [2] Lada A Adamic and Natalie Glance. The political blogosphere and the 2004 us election: divided they blog. In Proceedings of the 3rd international workshop on Link discovery, pages 36–43. ACM, 2005.
- [3] Edoardo M. Airolidi, David S. Choi, and Patrick J. Wolfe. Confidence sets for network structure. Technical Report arXiv:1105.6245, 2011.
- [4] Pranjali Awasthi. Clustering under stability assumptions. In Encyclopedia of Algorithms, pages 331–335. 2016.
- [5] Francis Bach and Michael I. Jordan. Learning spectral clustering with applications to speech separation. Journal of Machine Learning Research, 7:1963–2001, 2006.
- [6] Maria-Florina Balcan, Christian Borgs, Mark Braverman, Jennifer Chayes, and Shang-Hua Teng. Finding endogenously formed communities. In Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, pages 767–783. SIAM, 2013.
- [7] Shai Ben-David. Computational feasibility of clustering under clusterability assumptions. CoRR, abs/1501.00437, 2015.
- [8] Rajendra Bhatia. Matrix analysis, volume 169. Springer Science & Business Media, 2013.
- [9] Yonatan Bilu and Nathan Linial. Are stable instances easy? CoRR, abs/0906.3162, 2009.
- [10] Fan RK Chung. Spectral graph theory, volume 92. American Mathematical Soc., 1997.
- [11] Brian Karrer, Elizaveta Levina, and M. E. J. Newman. Robustness of community structure in networks. Phys. Rev. E, 77:046119, Apr 2008.
- [12] An-

drea Lancichinetti, Santo Fortunato, and Filippo Radicchi. Benchmark graphs for testing community detection algorithms. *Physical review E*, 78(4):046110, 2008. [13] Marina Meil’a. Local equivalence of distances between clusterings – a geometric perspective. *Machine Learning*, 86(3):369–389, 2012. [14] Marina Meil’a, Susan Shortreed, and Liang Xu. Regularized spectral learning. In Robert Cowell and Zoubin Ghahramani, editors, *Proceedings of the Artificial Intelligence and Statistics Workshop(AISTATS 05)*, 2005. [15] A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, Cambridge, MA, 2002. MIT Press. [16] Richard Peng, He Sun, and Luca Zanetti. Partitioning well-clustered graphs with k-means and heat kernel. In *Proceedings of the Annual Conference on Learning Theory (COLT)*, pages 1423–1455, 2015. [17] Tai Qin and Karl Rohe. Regularized spectral clustering under the degree-corrected stochastic blockmodel. In *Advances in Neural Information Processing Systems*, pages 3120–3128, 2013. [18] Karl Rohe, Sourav Chatterjee, and Bin Yu. Spectral clustering and the high-dimensional stochastic blockmodel. *The Annals of Statistics*, pages 1878–1915, 2011. [19] Gilbert W Stewart and Ji-guang Sun. *Matrix perturbation theory*, volume 175. Academic press New York, 1990. [20] Yali Wan and Marina Meila. A class of network models recoverable by spectral clustering. In Daniel Lee and Masashi Sugiyama, editors, *Advances in Neural Information Processing Systems (NIPS)*, page (to appear), 2015.