

Benchmarks for spectralGraphTopology

Vinícius

January 28, 2019

1 Clustering: warm-up

To warm-up, we consider the classical problem of clustering artificial nodes randomly distributed in \mathbb{R}^2 according to some geometry. More precisely, we use the package `clusterSim` to generate 100 nodes per cluster distributed according to structures colloquially known as *two-moons*, *two-circles*, *three-circles*, and *worms*.

Figure 1 depicts the results of learning the clusters structures using the algorithm *Spectral Graph Learning*, hereafter denoted as (SGL). As we can note, SGL is able to perfectly distinguish the cluster membership of all the nodes for all datasets. Additionally, Figure 2 depicts the convergence trend of the terms in the objective function for the *worms* dataset.

On what concerns hyperparameter selection, we use $K = 3$ for the *three-circles* dataset, otherwise we use $K = 2$. We fix $\beta = 0.25$ and $\alpha = 0$ for all the datasets.

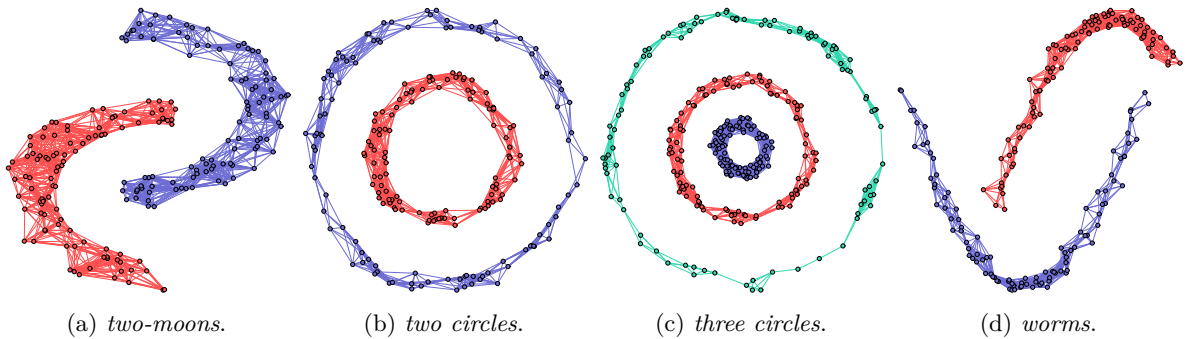


Figure 1: Clustering results learned by the SGL algorithm for synthetic datasets.

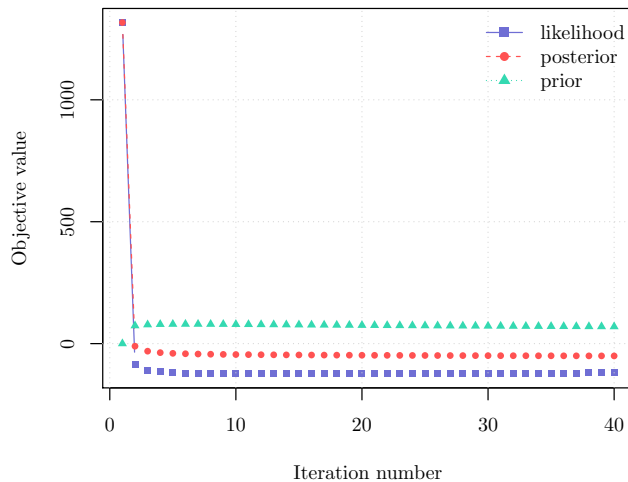


Figure 2: Convergence trend of the SGL algorithm for the *worms* dataset.

2 Grid graphs

A comprehensive performance evaluation of the SGL algorithm was conducted considering grid graph models with 64 nodes denoted as $\mathcal{G}_{\text{grid}}^{(64)}$. We compare the performance of the SGL algorithm against state-of-the-art algorithms, namely CGL, $\text{CGL}(\mathbf{A})$. We recollect that $\text{CGL}(\mathbf{A})$ stands for the CGL algorithm equipped with knowledge of the connectivity matrix \mathbf{A} , which gives the exact information of which nodes are connected to which ones.

The experimental setup is as follows. The edges of the graph model are sampled from $\text{Uniform}(0.1, 3)$. The Laplacian matrix estimation is carried out on the basis of T samples distributed according to $\mathcal{N}(\mathbf{0}, \mathbf{L}_{\text{grid}}^\dagger)$. We repeat that experiment 20 times for every value of T and we average out the relative errors and F-scores¹.

Some hyperparameter tuning is required. For the SGL algorithm, we fix $\beta = 10$ for the values of T such that $T/N > 5$. Otherwise, we start with $\beta = 10^{-2}$, and we exponentially increase it up to $\beta = 4$. Additionally, we fix $\alpha = 0$.

For the CGL and $\text{CGL}(\mathbf{A})$ algorithms, we follow the recipe by Eglimez, i.e., we choose α from $\{0\} \cup \{0.75^r (s_{\max} \sqrt{\log(N)/T}) \mid r = 1, 2, \dots, 14\}$, such that the relative error between the estimated Laplacian and the ground truth is minimized, where s_{\max} is the maximum value among the off-diagonal elements of the sample covariance matrix.

Figure ?? compares the performance of the algorithms for different sample size regimes for the grid graph model. As it can be noted, the SGL algorithm outperforms the CGL in both relative error and F-score senses. More precisely, for the case when the sample size is equal to the number of nodes, the difference in relative error and in F score are around 12% and 23%, respectively.

As expected, with the additional prior knowledge of the connectivity matrix \mathbf{A} , the $\text{CGL}(\mathbf{A})$ algorithm basically attains a perfect F-score for $T/N \geq 10$. However, the connectivity matrix is not always available in practical problems, especially in clustering tasks where the goal is precisely to understand the connectivity membership among the nodes.

Nonetheless, the proposed SGL algorithm presents a comparable performance against $\text{CGL}(\mathbf{A})$. For instance, at $T/N = 5$, the difference in relative error is only around 2.5%, and it keeps decreasing even further until virtually equal performance after $T/N = 100$, where the difference in relative error is around 1.2%. Additionally, we noted that the SGL algorithm requires far less tuning than CGL. At last, we include the relative error and the F-score for the Moore-Penrose inverse of the sample covariance matrix (ISCM) for completeness.

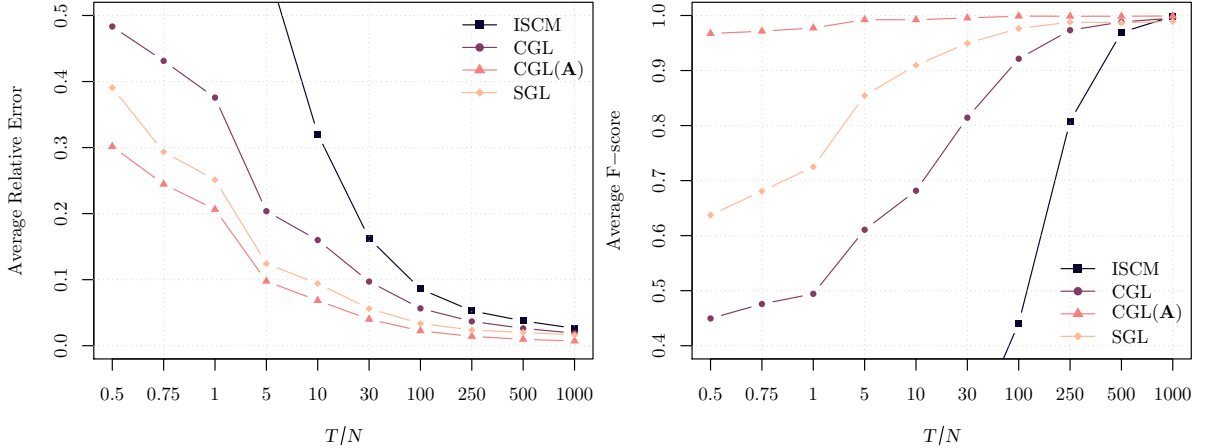


Figure 3: Average performance results for learning Laplacian matrix of a $\mathcal{G}_{\text{grid}}^{(64)}$.

Figure ?? pictorially compares the graph structures learned by SGL and CGL when $T/N = 100$.

¹For the computation of the F-score, we ignore edge weight values which are less than 10^{-1} .

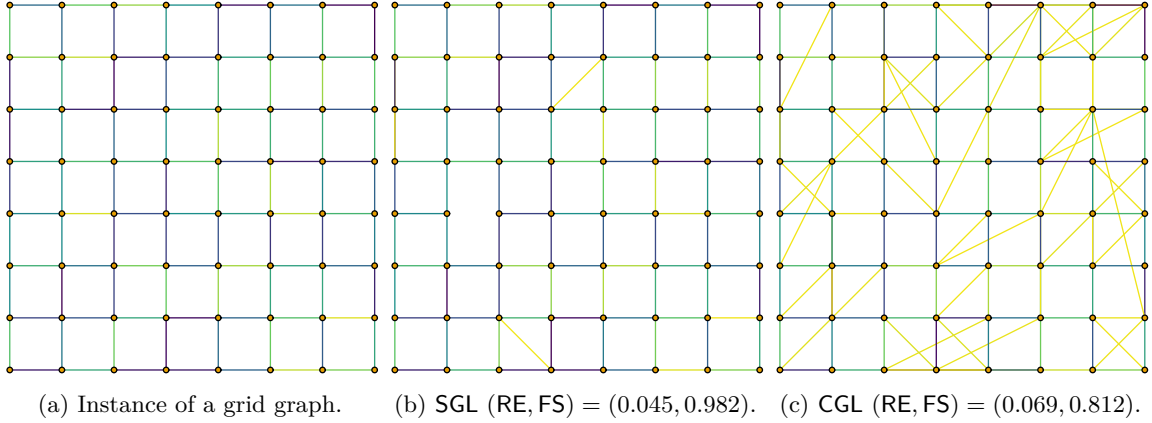


Figure 4: Sample results of learning $\mathcal{G}_{\text{grid}}^{(64)}$ for $T/N = 100$. Edges smaller than 0.05 were removed. For SGL, we fix $\beta = 10$, whereas for CGL we perform a grid search in α to find α that maximizes the F-score. We found $\alpha = 3.6 \cdot 10^{-3}$.

3 Erdos-Renyi graphs

The experiments involving Erdos-Renyi graph model is carried out in a similar fashion as for the grid graph. We consider an Erdos-Renyi graph with 64 nodes, denoted as $\mathcal{G}_{\text{ER}}^{(64,p)}$, with $p = 0.1$, where p is the probability that a particular node is connected to any other node.

On what concerns hyperparameter tuning, we fix $\alpha = 1.3 \cdot 10^{-2}$ and start with $\beta = 10^{-1}$ and we exponentially increase it up to $\beta = 1$ for the values of T such that $T/N \leq 1$. For $1 < T/N \leq 30$, we set $\alpha = 0$ and $\beta = 1.5$. Finally, for $T/N > 30$, we set $\alpha = 0$ and $\beta = 10$.

The hyperparameter tuning for the CGL and CGL(A) algorithms is done as described for the grid graph model.

Figure ?? reveals that the algorithms SGL and CGL obtain a similar performance across the sample size regimes. Unsurprisingly, the algorithm CGL(A) attains nearly perfect F-score for $T/N > 10$, because, again, it assumes that the connectivity information is fully available, which may be a non trivial assumption in practice.

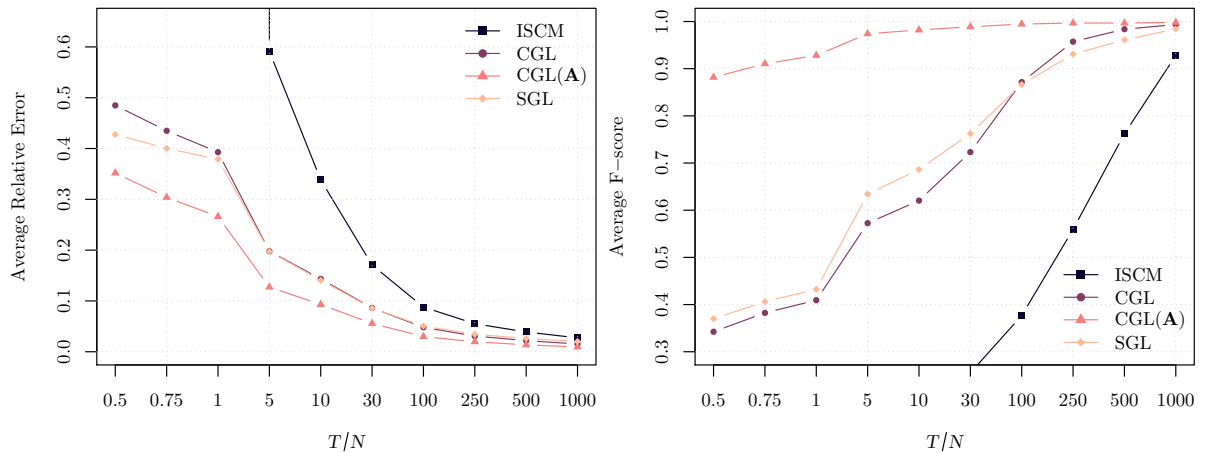


Figure 5: Average performance results for learning Laplacian matrix of a $\mathcal{G}_{\text{ER}}^{(64,0.1)}$.

4 Modular graphs

We consider a modular graph with 64 nodes and four modules, denoted as $\mathcal{G}_M^{(64, p_1, p_2)}$, in which $p_1 = 0.01$, and $p_2 = 0.5$, where p_1 is the probability that a particular node of a given module is connected to any other node of a different module, and p_2 is the probability that a particular node is connected to any other node of the same module.

On what concerns hyperparameter tuning, we fix $\beta = 4$ for all values of T/N .

Figure 3 shows the performance of the proposed algorithms and the benchmarks. As we can note, the SGL algorithm outperforms the CGL algorithm in terms of F-score, while having a comparable performance in terms of relative error.

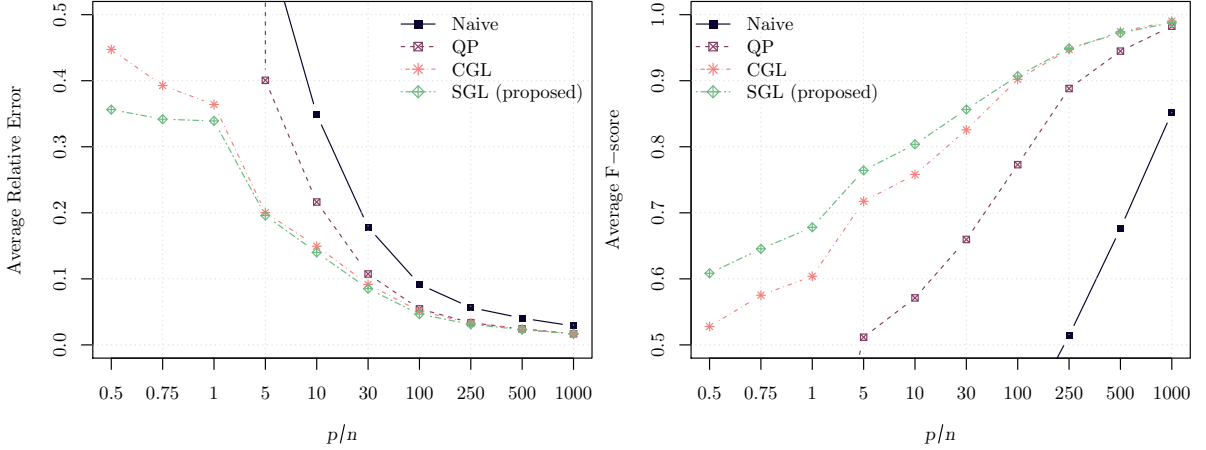


Figure 6: Average performance results for learning Laplacian matrix of a $\mathcal{G}_M^{(64,0.01,0.5)}$ with four modules.

5 Animals dataset

The **animals** dataset consists of binary values which are the answers to questions such as "is warm-blooded?", "has lungs?", etc. There are in total 102 such questions, which make up the features, and 33 animal categories, which are the nodes. Figure 4 shows the results of estimating the graph Laplacian of the **animals** dataset using the GGL algorithm with $\alpha = 0.1$ and the proposed SGL algorithm with $\beta = 1/2$ and $\alpha = 10^{-1}$. The input for all the algorithms is the sample covariance matrix plus an identity matrix scaled by $1/3$ (see Elgimez paper for an explanation of the input).

The evaluation of the estimated graphs is done by natural intuition, i.e., we expect that similar animals such as (*ant*, *cockroach*), (*bee*, *butterfly*), and (*trout*, *salmon*), would be clustered together, while presently virtually zero connection to other (group of) animals. Under this sense, it can be seen that the SGL algorithm yields a more intuitive graph than the ones learned by GGL and Glasso.

6 Does a better initialization affect the final estimation?

The standard initialization for the SGL algorithm is to take \mathbf{w}_0 as the off-diagonal elements of the pseudoinverse of the sample covariance matrix, denoted as $\hat{\Theta}_{\text{scm}}$. This is tantamount to solving the following optimization problem

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \|\hat{\Theta}_{\text{scm}} - \mathcal{L}\mathbf{w}\|_{\text{off},F}^2 \\ & \text{subject to} && \mathbf{w} \geq \mathbf{0}. \end{aligned}$$

In order to improve the initial estimate, we consider the following optimization problem

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \|\hat{\Theta}_{\text{scm}} - \mathcal{L}\mathbf{w}\|_F^2 \\ & \text{subject to} && \mathbf{w} \geq \mathbf{0}. \end{aligned} \tag{1}$$

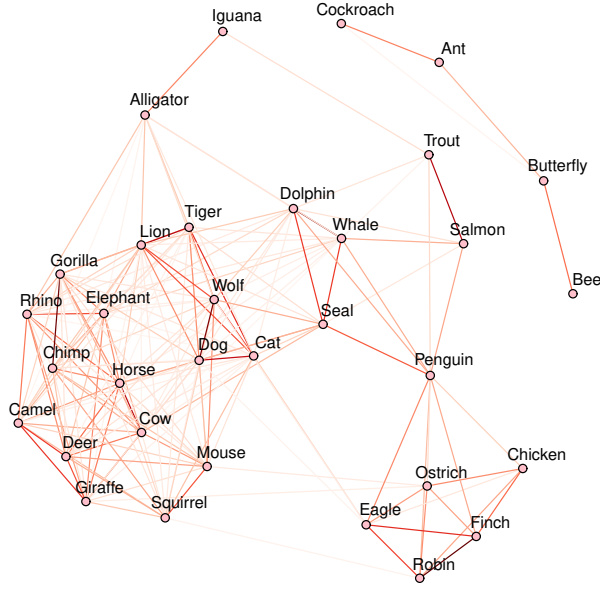


Figure 8: Learning the connectivity of the animals dataset with GLasso. Regularization parameter $\alpha = 0.05618$.

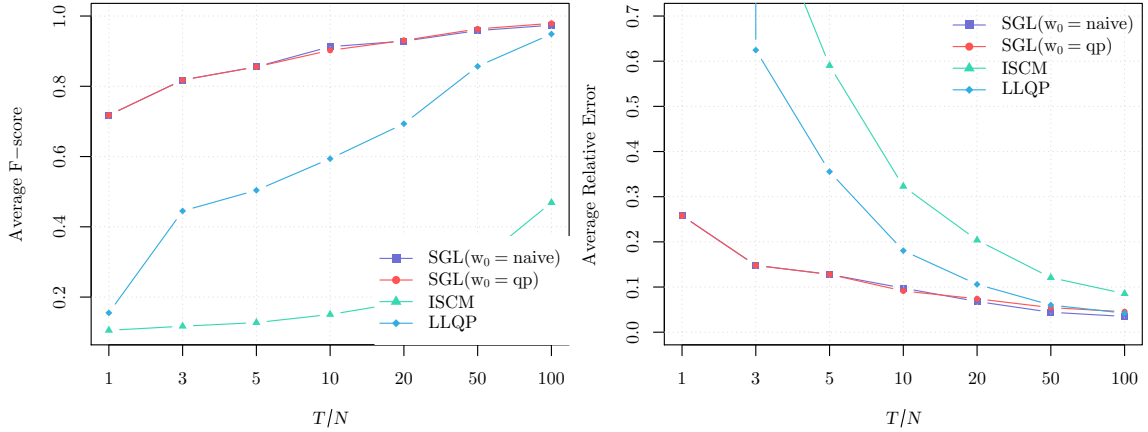


Figure 9: Relative error and F-score attained by SGL with different initialization strategies. $\text{SGL}(\mathbf{w}_0 = \text{naive})$ and $\text{SGL}(\mathbf{w}_0 = \text{qp})$ stand for the initialization with ISCM and LLQP, respectively.

7 4-component graph

An experiment with synthetic data generated by a 4-component graph was conducted. We consider the following model:

$$\mathbf{L}_{\text{noisy}} = \mathbf{L}_{\text{true}} + \mathbf{L}_{\text{ER}}, \quad (3)$$

where \mathbf{L}_{true} represents the Laplacian matrix of a K -component graph (for this example, $K = 4$) denoted as $\mathcal{G}_K^{(p_1, p_2)}$, in which p_1 and p_2 represent the probabilities of node connections across components and within components, respectively; \mathbf{L}_{ER} represents the Laplacian of an Erdos-Renyi graph $\mathcal{G}_{\text{ER}}^{(p)}$, in which p is the probability of a node connecting to any other node. The weighted edges of $\mathcal{G}_K^{(p_1, p_2)}$ were drawn from $\text{Uniform}(0, 1)$, whereas the ones of $\mathcal{G}_{\text{ER}}^{(p)}$ were drawn from $\text{Uniform}(0, \kappa)$, in which $\kappa \in (0, 1)$ controls how much noise is added to the true graph. Finally, we set $p_1 = 0$, $p_2 = 1$, $p = 0.35$, and $\kappa = 0.45$.

Then, data were sampled in the form of $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{L}_{\text{noisy}}^\dagger)$, where \mathbf{A}^\dagger denotes the generalized inverse of the matrix \mathbf{A} . The total number of nodes N and the number of drawn samples T were set to $N = 20$ (5 nodes per component) and $T/N = 30$.

Figure 6 illustrates the ground truth model, its noisy version, and the model learned by our spectral topology algorithm with $\beta = 20 \cdot N$ and $\alpha = 0.1$. We compute the performance of the learning process by means of the relative error (RE) and the F-score (FS). For this example, our algorithm achieves $(\text{RE}, \text{FS}) = (0.210, 1)$ which means a perfect clustering accuracy even in a noisy model that heavily suppress the ground truth weights.

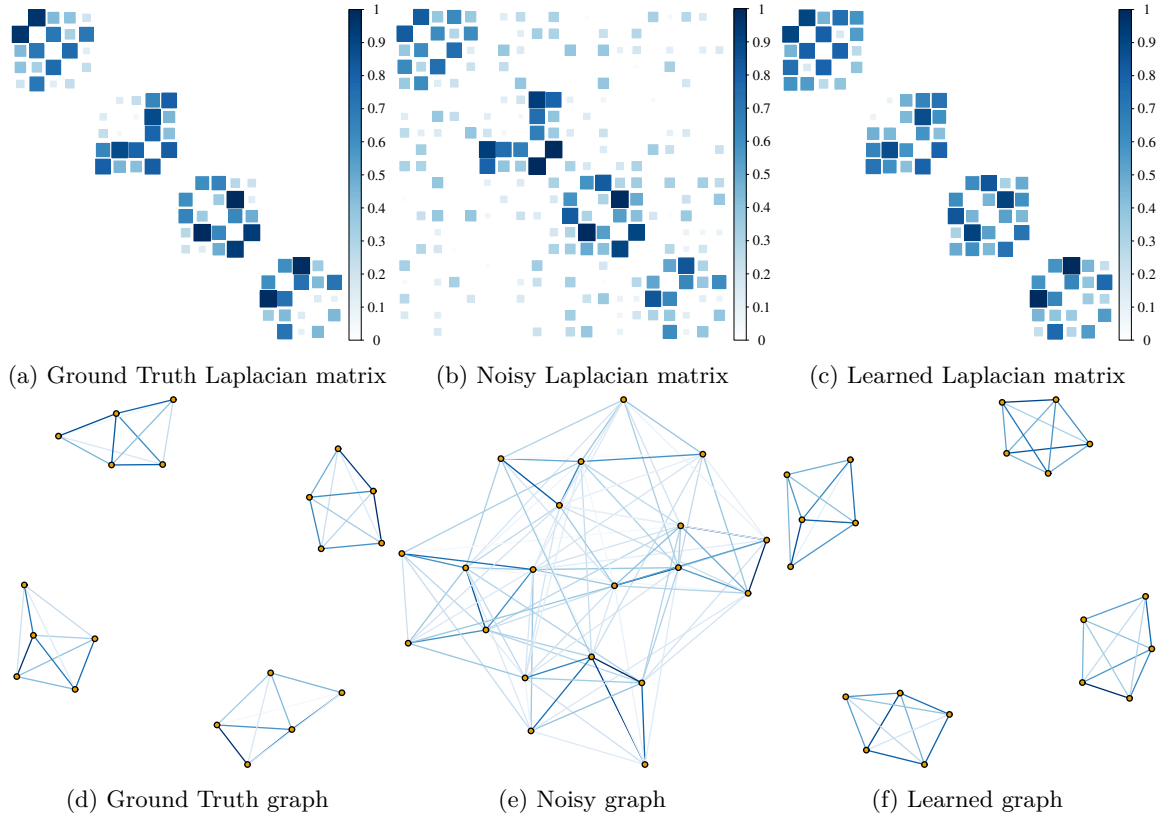


Figure 10: An example of estimating a 4-component graph. (a) the ground truth graph Laplacian matrix (\mathbf{L}_{true}), (b) \mathbf{L}_{true} after being corrupted by noise, (c) the learned graph Laplacian with a performance of $(\text{RE}, \text{FS}) = (0.210, 1)$. The panels (d), (e), and (f) correspond to the graphs represented by the Laplacian matrices in (a), (b), and (c), respectively.

Figure 8 shows the performance of our algorithm, using the same settings as the experiment above, but for different noise regimes.

8 Performance

To infer the performance as a function of the ratio T/N , we consider a 4-component graph in which every component has four nodes connected with probability 1. The edge weights among nodes are drawn from $\text{Uniform}(0.1, 3)$. Additionally, we perform a grid search on the sparsity parameter α , just like the one performed by Elgimez.

9 Model mismatch: k -component graph

In many applications, specially those in fully automated and unsupervised settings, the number of clusters is either unknown or it needs to be inferred from the data. In this section we consider an experiment involving model mismatch: the underlying Laplacian matrix that generates the data has K number of components but we actually use J , $J < K$, number of components to estimate it. Additionally, we consider that the generative

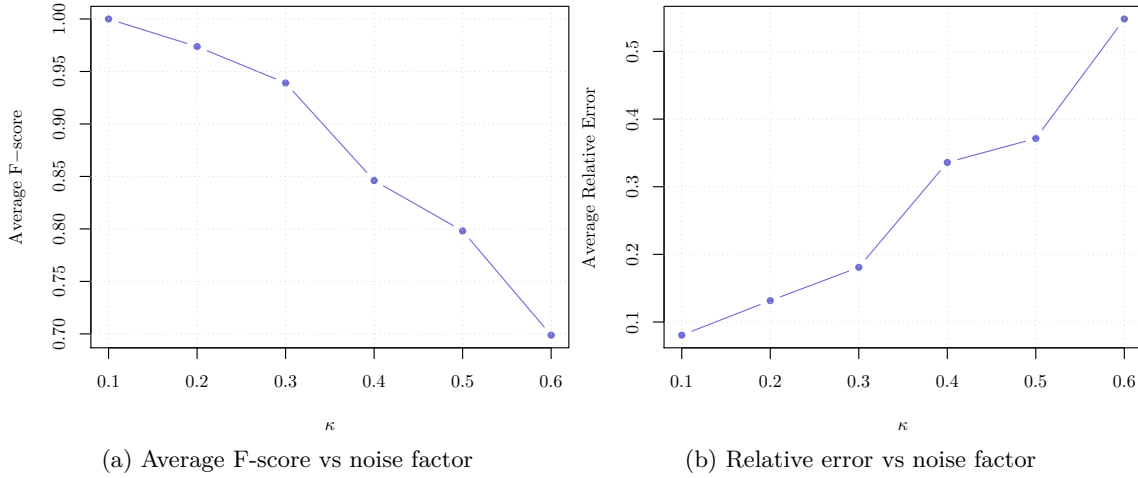


Figure 11: Average performance results, as a function of the noise factor, for learning Laplacian matrix of a modular graph embedded in noise.

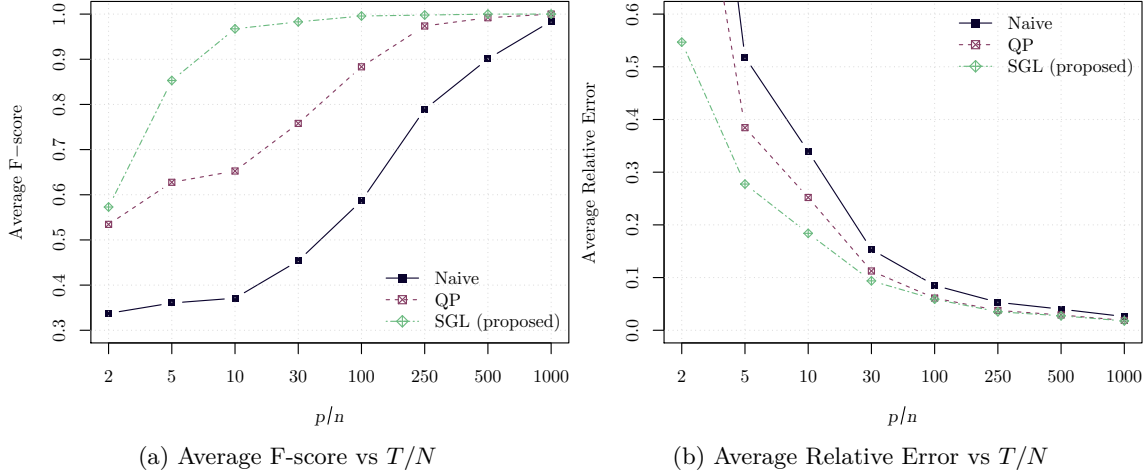


Figure 12: Average performance results as, a function of the number of samples, for learning Laplacian matrix of a 4-component graph.

model is noisy, i.e., $\mathbf{L}_{\text{noisy}} = \mathbf{L}_{\text{bm}} + \mathbf{L}_{\text{ER}}$, where the parameters of the Erdos-Renyi model that acts as noise are $p = 0.25$ and $\kappa = 0.45$. The block-model \mathbf{L}_{bm} is constructed with seven blocks, each of which with seven nodes that are fully connected among themselves. The weights on the edges are drawn from $\text{Uniform}(0, 1)$.

Figure ?? shows an example where the underlying graph has seven components, and we apply the SGL algorithm with $K = 2$. As we can see, even though the number of components is mismatched and the data is noisy, the SGL algorithm is still able to indentify the true structure with a resonable performance in terms of F-score and average relative error.

Additionally, we investigate the performance as a function of the number of components, as depicted in Figure ?. As we can note, the performance is monotonically increasing and eventually reaches a perfect F-score when $K = 7$.

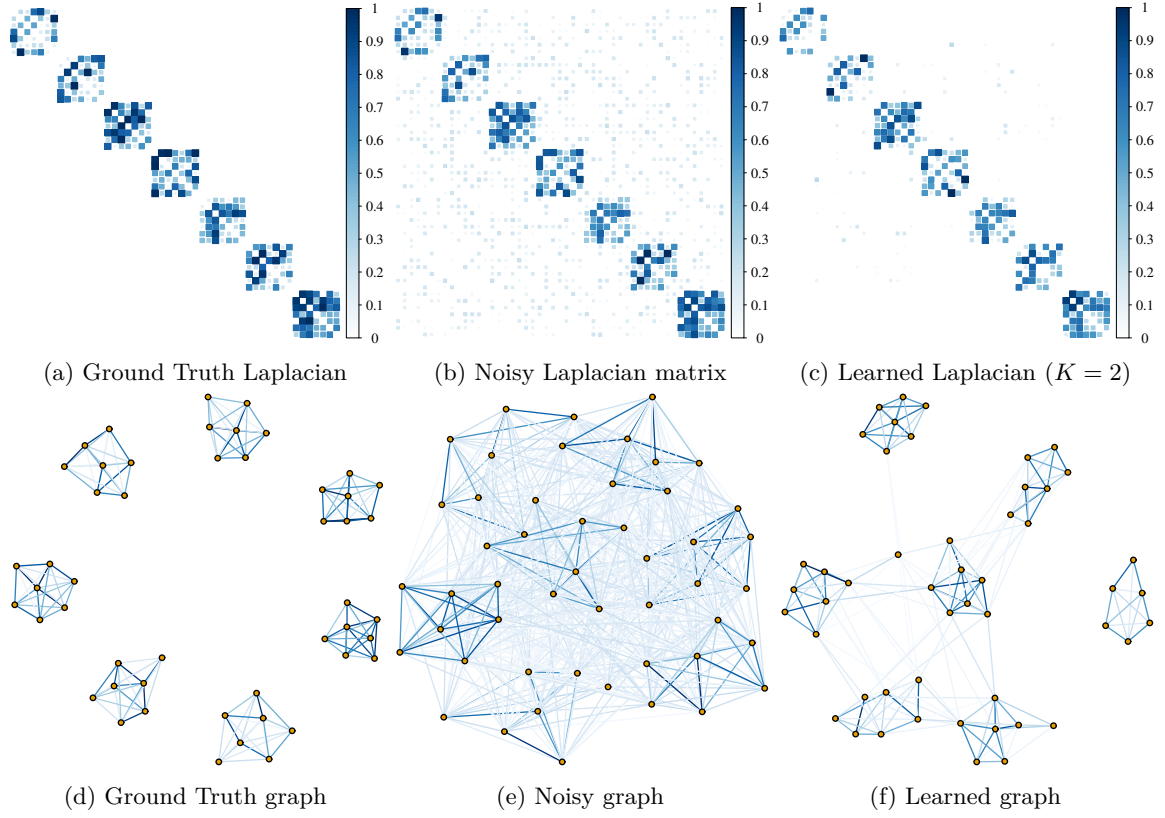


Figure 13: Example of model mismatching. (a) the ground truth graph Laplacian of a seven-component graph (\mathbf{L}_{true}), (b) \mathbf{L}_{true} after being corrupted by noise, (c) the learned graph Laplacian with a performance of $(\text{RE}, \text{FS}) = (0.18, 0.81)$. The panels (d), (e), and (f) correspond to the graphs represented by the Laplacian matrices in (a), (b), and (c), respectively.

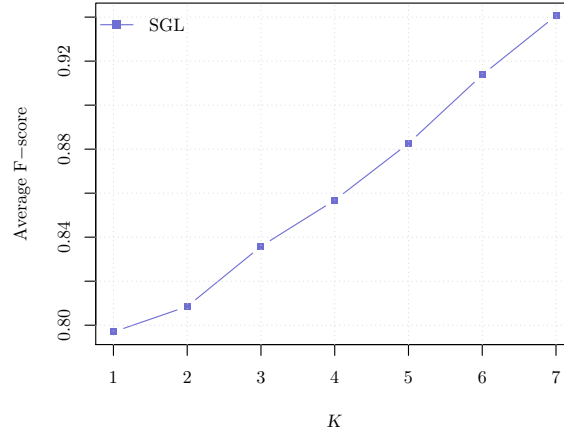


Figure 14: F-score as a function of the number of components K .

10 Cancer RNA gene dataset

As an example of clustering on another real dataset, we consider the RNA-Seq PANCAN dataset available at the UC-Irvine Machine Learning Database².

²<https://archive.ics.uci.edu/ml/datasets/gene+expression+cancer+RNA-Seq>

This dataset is a subset of a larger dataset designed by the Cancer Genome Atlas Research Network³. It consists of genetic features which map twelve types of cancer, namely: glioblastoma multiforme (GBM), lymphoblastic acute myeloid leukemia (LAML), head and neck squamous carcinoma (HNSC), lung adenocarcinoma (LUAD), lung squamous carcinoma (LUSC), breast carcinoma (BRCA), kidney renal clear-cell carcinoma (KIRC), ovarian carcinoma (OV), bladder carcinoma (BLCA), colon adenocarcinoma (COAD), uterine cervical and endometrial carcinoma (UCEC) and rectal adenocarcinoma (READ).

In the subset available at UCI, however, only five types of cancer are available. The dataset consists of 801 labeled samples, in which every sample has 20531 genetic features. The goal in this dataset is to classify and cluster the samples according to their tumor type on the basis of those genetic features.

In our experiments, we first apply the SGL algorithm to the first 70 nodes using $K = 5$ and $\beta = 5$. Figure 11 shows a perfect clustering, which indicates that SGL may be used to automatically predict types of cancer tumors based on gene features.

We also compare the classification performance against the GGL algorithm. As depicted in Figure 10, the GGL algorithm correctly classifies most of the nodes, but fails to segregate the clusters as it is evidenced by the vast number of mistaken edges placed across clusters.

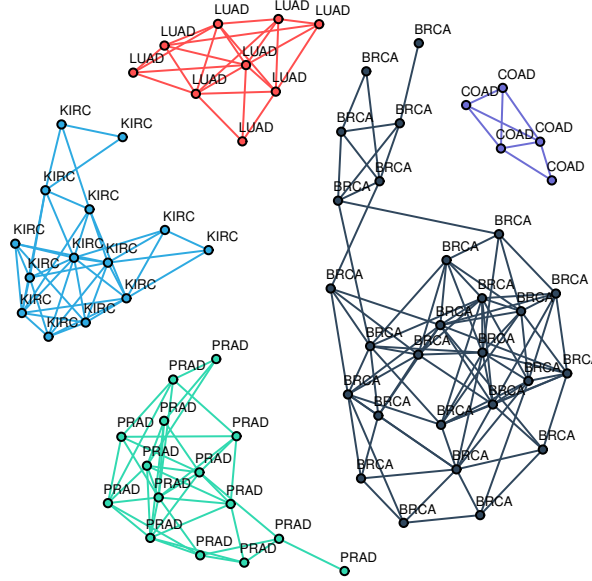


Figure 15: Clustering the first 70 nodes of the PANCAN dataset with $SGL(K = 5, \beta = 5)$.

Finally, we perform clustering for the full dataset. The result is depicted in Figure ??.

³<https://www.nature.com/articles/ng.2764>

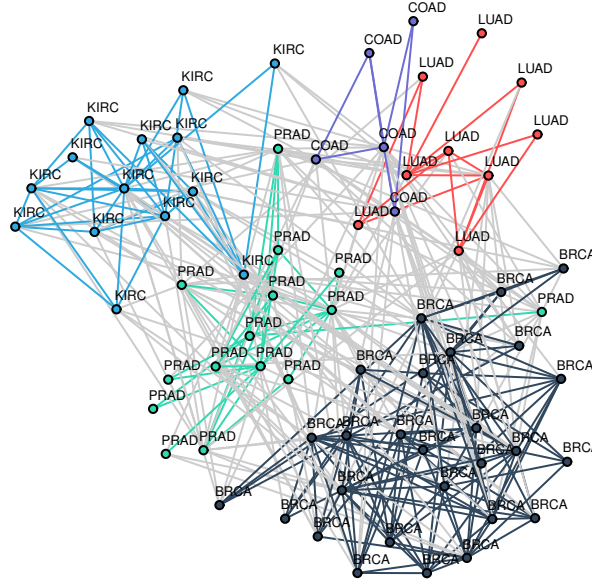


Figure 16: Clustering the first 70 nodes of the PANCAN dataset with $GGL(\alpha = 0.1)$. Grey colored edges indicate wrong connections.

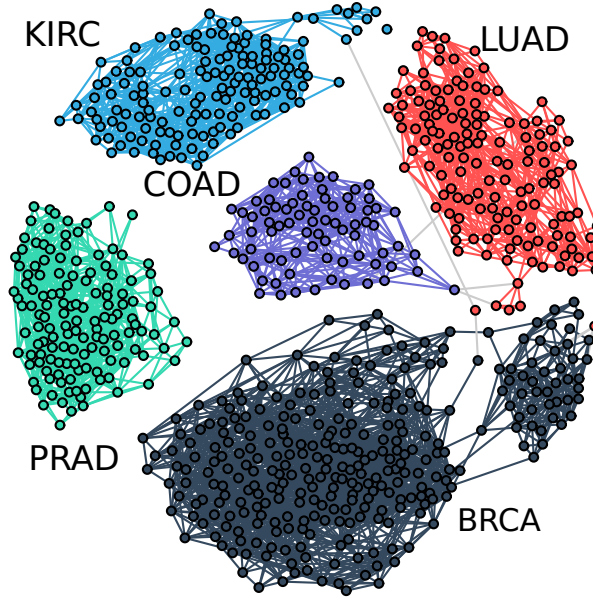


Figure 17: Clustering the full 801 nodes of the PANCAN dataset with $SGL(K = 5, \beta = 10)$.