

MIXTURE MODELS FOR GRAPH CLUSTERING

Sofiane Ezzehi¹, Bastien Le Chenadec¹, and Theïlo Terrisse¹

¹École des Ponts ParisTech

CONTRIBUTION STATEMENT

The Stochastic Block Model-based variational EM [1] was implemented in Numpy and Pytorch by Bastien Le Chenadec. The Newman variant [2] was added by Sofiane Ezzehi, as well as the spectral clustering algorithm [3]. The experiments were designed and performed by Theïlo Terrisse. The report was written by all three authors.

1 INTRODUCTION

Graph clustering, also referred to as “community detection” [4], is the problem of revealing clusters within a graph, where a “cluster” is often defined as a set of nodes that are more densely connected between them than with the rest of the network. More precisely, a *homophilic* cluster C is defined as a group of nodes with an internal density $\delta_{int}(C)$ larger than its inter-community density $\delta_{ext}(C)$, where

$$\delta_{int} = \frac{\# \text{ internal edges of } C}{n_C(n_C - 1)/2} \quad (1)$$

and

$$\delta_{ext} = \frac{\# \text{ inter-cluster edges of } C}{n_C(n - n_C)} \quad (2)$$

with n the number of nodes of the graph and n_C the number of nodes in the cluster C . Some applications may also be interested in *heterophilic* clusters, defined such that $\delta_{int}(C) < \delta_{ext}(C)$. Graph clustering has tremendous applications in the study of real networks such as social [5] or biological [6] ones. This problem is of importance to perform mesoscopic studies of graphs, where clusters can be seen as meta-nodes to reveal interactions within a complex system.

Graph clustering is known to be an NP-hard problem, as testing all possible partitions of a graph has complexity $O(B_n)$ with B_n the Bell number of order n . This is particularly problematic as real networks easily grow to thousands or even millions of nodes. Several classes of algorithms are documented in the literature [4]. To name a few, modularity-based approaches solve an optimization problem that maximizes modularity (defined in Section 3.3); hierarchical algorithms rely on similarity measures to aggregate or divide groups of nodes; spectral algorithms also rely on a similarity matrix, but operate in the space spanned by the eigenvectors of its Laplacian matrix. Yet, defining a similarity measure between nodes is not always trivial.

Model-based methods are an alternative to algorithmic methods, that aim to fit a mathematical model capable of explaining the observed connectivity, and from which descriptive properties of a graph may be extracted or computed. The Stochastic Block Model (SBM) introduced in [7] is a renowned example of model for graphs. In [1], J.-J. Daudin *et al.* study this model and propose a variational Expectation-

Maximization (EM) algorithm to find parameters that best fit a given graph. We will refer to this method as “Daudin-EM”.

In this report, we detail the Daudin-EM method as well as a variant from the literature in Section 2. In Section 3, we present our implementation of the method and the experiments performed to evaluate it, with a discussion of the results in Section 4. Lastly, we conclude in Section 5.

2 PROPOSED METHODS AND VARIANTS

2.1 A mixture model for random graphs

In “A mixture model for random graphs” [1], Daudin *et al.* propose a bayesian model for graphs, and an EM approach to fit this model to a given graph.

We will follow the same notations as in [1]. We consider an undirected graph with n nodes and no self-loops. We denote X the adjacency matrix of this graph. As such, $X_{ij} \in \{0, 1\}$ denotes the existence of an edge between nodes i and j , and $\forall i \in \{1, \dots, n\}, X_{ii} = 0$.

As explained by the authors, the SBM draws inspiration both from mixture models for distribution of degrees, that have the disadvantage of not dealing with the probability for two given nodes of being connected, and from the foundational Erdős-Rényi model, which is known to fit real graphs poorly. To define the SBM, we consider a mixture model that spreads the vertices into Q classes, with an a priori repartition $\{\alpha_1, \dots, \alpha_Q\}$ between classes. We introduce the random variables $Z_{iq} \in \{0, 1\}$ for $i \in \{1, \dots, n\}$ and $q \in \{1, \dots, Q\}$, that represent the membership of node i to class q . We have the following prior distribution on Z :

$$\forall i \in \{1, \dots, n\}, \quad \sum_{q=1}^Q Z_{iq} = 1 \quad (3)$$

and

$$\forall q \in \{1, \dots, Q\}, \quad \mathbb{P}(Z_{iq} = 1) = \alpha_q. \quad (4)$$

We introduce priors on the existence of edges between nodes of different classes. We denote π_{ql} the probability of an edge between a node of class q and a node of class l . Because the graph is undirected, we have $\pi_{ql} = \pi_{lq}$. This prior writes : $\forall q, l \in \{1, \dots, Q\}, \quad \forall i \neq j \in \{1, \dots, n\},$

$$\mathbb{P}(X_{ij} = 1 | Z_{iq} = 1, Z_{jl} = 1) = \pi_{ql}. \quad (5)$$

Figure 1 shows the graphical model associated with this mixture model.

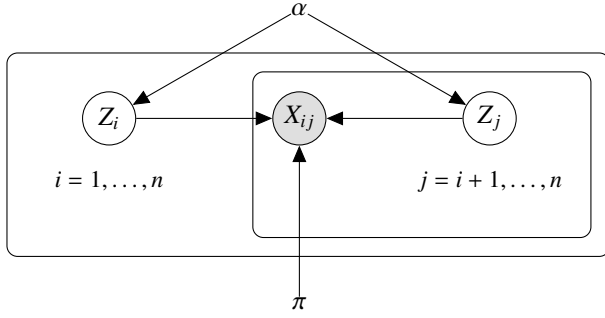


Figure 1: Graphical model of the SBM

2.1.1 Variational Expectation-Maximization algorithm

The log-likelihood of the model is given by :

$$\begin{aligned} \log \mathcal{L}(X, Z) = & \sum_i \sum_q Z_{iq} \log \alpha_q \\ & + \frac{1}{2} \sum_{i \neq j} \sum_{q,l} Z_{iq} Z_{jl} \pi_{ql}^{X_{ij}} (1 - \pi_{ql})^{1-X_{ij}} \end{aligned} \quad (6)$$

Because the likelihood $\mathcal{L}(X)$ is not tractable, the authors propose to use an Expectation-Maximization (EM) algorithm to fit the model. However the E-step of the EM algorithm is not tractable either because of the posterior distribution of Z given X . Instead the authors propose to optimize (7) which is a lower bound of $\log \mathcal{L}(X)$ obtained using the Kullback-Leibler divergence between the posterior distribution of Z given X and an approximated distribution R_X .

$$\mathcal{J}(R_X) = \log \mathcal{L}(X) - \text{KL}[R_X(\cdot), P(\cdot|X)] \quad (7)$$

By choosing the approximated distribution R_X to be a product of independant multinomial distributions (8), the authors obtain a fixed point relation (9) between the parameters of the model and the parameters of the approximated distribution maximizing the lower bound $\mathcal{J}(R_X)$. This fixed point relation is used in the E-step of the algorithm.

$$\begin{aligned} R_X(Z; \tau) = & \prod_{i=1}^n h(Z_i, \tau_i) \\ \forall i \in \{1, \dots, n\}, \quad h(Z_i, \tau_i) = & \prod_{q=1}^Q \tau_{iq}^{Z_{iq}}. \end{aligned} \quad (8)$$

The fixed point is, $\forall i \in \{1, \dots, n\}, \forall q \in \{1, \dots, Q\}$,

$$\hat{\tau}_{iq} \propto \alpha_q \prod_{j \neq i} \prod_l [\pi_{ql}^{X_{ij}} (1 - \pi_{ql})^{1-X_{ij}}]^{\hat{\tau}_{jl}} \quad (9)$$

This fixed point relation does not assure the theoretical convergence of the algorithm. We will see later that the convergence of the algorithm is not guaranteed in practice either. Finally we have the following updates for the M-step maximizing $\mathcal{J}(R_X)$:

$$\forall q, l \in \{1, \dots, Q\}, \quad \hat{\alpha}_q = \frac{1}{n} \sum_i \tau_{iq} \quad (10)$$

and

$$\hat{\pi}_{ql} = \frac{\sum_{i \neq j} \tau_{iq} \tau_{jl} X_{ij}}{\sum_{i \neq j} \tau_{iq} \tau_{jl}}. \quad (11)$$

These relations let us define an EM algorithm that is guaranteed to increase the lower bound $\mathcal{J}(R_X)$ at each iteration (conditionnal on solving the fixed point relation (9)).

2.1.2 Selection of Q

In practice the number of clusters Q has to be estimated. The authors propose to use a criterion based on the Integrated Classification Likelihood (ICL) which writes :

$$\begin{aligned} \text{ICL}(X, Q) = & \max_{\theta} \log \mathcal{L}(X, Z | \theta, Q) - \frac{Q-1}{2} \log n \\ & - \frac{1}{2} \times \frac{Q(Q+1)}{2} \log \frac{n(n-1)}{2} \end{aligned} \quad (12)$$

This criterion requires the knowledge of the parameters θ of the model, which we will estimate using the EM algorithm.

2.2 An alternative mixture-model by Newman et. al.

A different mixture model has been proposed the same year as Daudin-EM, by Newman et al. [2]. The main difference lies in the way a cluster is defined. In the previous model, the probability of connection between a node in cluster A and a node in cluster B is the same, regardless of the considered nodes.

In the Newman model, we build directed edges between nodes. In the undirected case, we simply consider that two nodes are connected if there are directed edges between them in both directions. By construction, two nodes belong to the same cluster if the probability of a directed edge between either of them and any other node in the graph is the same. Therefore, for any given pair of nodes (i, j) belonging to the same cluster C and any other node k , we have,

$$\mathbb{P}(i \rightarrow k) = \mathbb{P}(j \rightarrow k) = \Pi_{C,k}. \quad (13)$$

Therefore, we see that matrix Π , which is different from the π of the previous model, describes the probability of a directed edge existing between a cluster and a node. Furthermore, contrary to the previous model, the probability of a directed edge existing between a node in cluster A and a node in cluster B is not necessarily invariant with respect to the considered nodes.

2.2.1 Model

As in the previous model, we denote X the adjacency matrix of the input graph which is assumed to be undirected and containing n nodes. We also denote $(\alpha_1, \dots, \alpha_Q)$ the probabilities for a node of belonging to a cluster $q \in \{1, \dots, Q\}$. We introduce the random variables $g_{i \in \{1, \dots, n\}} \in \{1, \dots, Q\}$ which represent the cluster of node i ; these are analogous to the random variables Z_{iq} previously introduced. Therefore,

$$\forall i \in \{1, \dots, n\}, \quad \mathbb{P}(g_i = q) = \alpha_q. \quad (14)$$

Finally, we introduce the matrix $\Pi \in [0, 1]^{Q \times n}$ defined earlier. Since the input graph is undirected, the probability of an undirected edge existing between two nodes i and j is,

$$\mathbb{P}(X_{ij} = 1) = \mathbb{P}(i \rightarrow j) \mathbb{P}(j \rightarrow i) = \Pi_{g_i, j} \Pi_{g_j, i}. \quad (15)$$

Two normalization conditions are imposed,

$$\sum_{q=1}^Q \alpha_q = 1 \quad \text{and} \quad \forall q \in \{1, \dots, Q\}, \sum_{i=1}^n \Pi_{q,i} = 1 \quad (16)$$

2.2.2 Expectation-Maximization algorithm

The expected log-likelihood (with respect to the random variables g) is,

$$\mathcal{L}(X) = \sum_{i=1}^n \sum_{q=1}^Q \tau_{iq} \left[\log(\alpha_q) + \sum_{j=1}^n X_{ij} \log(\Pi_{q,j}) \right] \quad (17)$$

where $\tau_{iq} = \mathbb{P}(g_i = q | X, \alpha, \Pi)$.

Remark : The authors of the article made the surprising choice of not taking into account, in the likelihood of X given the parameters, the probability of the non-existing edges. Taking them into account would have added in the last sum of the log-likelihood a term of the form $(1 - X_{ij}) \log(1 - \Pi_{q,j})$.

The main advantage of this model is that the expectation step is much simpler than in the previous model as it doesn't require to solve a fixed-point problem. Indeed, we have an explicit expression for τ_{iq} ,

$$\hat{\tau}_{iq} = \frac{\alpha_q \prod_{j=1}^n \Pi_{q,j}^{X_{ij}}}{\sum_{s=1}^Q \alpha_s \prod_{j=1}^n \Pi_{s,j}^{X_{ij}}} \quad (18)$$

The M-step updates are,

$$\hat{\alpha}_q = \frac{1}{n} \sum_{i=1}^n \tau_{iq} \quad \text{and} \quad \hat{\Pi}_{qj} = \frac{\sum_{i=1}^n X_{ij} \tau_{iq}}{\sum_{i=1}^n k_i \tau_{iq}} \quad (19)$$

where $k_i = \sum_{j=1}^n X_{ij}$ is the degree of node i . We verify that the normalization conditions are satisfied.

2.3 Spectral clustering

A third method that we will use as a comparison baseline is spectral clustering, first introduced in [3]. The idea is to use the eigenvectors of the Laplacian matrix of the graph to cluster the nodes. The Laplacian matrix is defined as $L = D - X$ where X is the adjacency matrix of the graph and D is the diagonal matrix of degrees of the nodes.

The algorithm is simple: stacking the k first eigenvectors of L in a matrix $U \in \mathbb{R}^{n \times k}$, we perform a k -means clustering on the rows of U .

3 EXPERIMENTS

The experiments we have carried out can be separated into three parts, all based on a re-implementation of the methods covered in Section 2. A first experiment illustrates the lack of robustness of the fixed-point algorithm used in the E-step of the Daudin-EM algorithm. Second, Daudin-EM is tested isolately on synthetic graphs generated from the SBM introduced in Section 2.1, both in terms of model-fitting capacity and graph clustering power. Lastly, the performance of the mixture model is evaluated on classical real datasets, namely Zachary's karate club and the Cora dataset. This performance is compared to the baselines presented in Sections 2.2 and 2.3.

3.1 Implementation

All of the experiments were performed using the code attached with this report. The Daudin-EM method (simply referred to as "SBM" in the code) was implemented in Python using Numpy, and later in Pytorch to distribute the computations on a GPU. A second Pytorch implementation was added to handle memory issues occurring for large graphs. The Newman variant was directly implemented in Pytorch. As an indication, a speed comparison between the versions of the method is available in appendix A. The spectral clustering algorithm was implemented in Numpy in a dedicated notebook. Lastly, the algorithm proposed by Lancichinetti, Fortunato, and Kertesz in [8] to deal with overlapping communities was implemented in Numpy, but was not exploited in experiments for conciseness of this report.

The code comes with a battery of tests to validate the implementations. Scripts and notebooks to reproduce and visualize the experiments are also available; please refer to the README.md file in the attached code to know in which notebook each experiment is run. All experiments used Pytorch implementations when available, especially the PytorchLogImplementation for Daudin-EM, except for the smallest graphs of the second set of experiments that used the Numpy version as it was faster in this case. For visualizations, we used the Networkx library.

3.2 Robustness of the fixed-point algorithm

Early stages of the implementation of the method revealed a lack of robustness of the fixed-point algorithm used in the E-step of the method. In particular, the fixed-point algorithm was found to be sensitive to the initialization of the parameters α and π , and does not converge in all cases.

To illustrate this, we perform 2 experiments:

1. In the first experiment, we initialize 50 random graphs and fix a random $\hat{\tau}$ to get 50 initial states.
2. In the second experiment, we fix a random graph and initialize 50 random $\hat{\tau}$ to get 50 initial states.

In each experiment, the graphs are generated using the SBM described in Section 2.1, with $Q = 3$ and 100 nodes. More precisely, α and π are initialized randomly, then the latent variable Z is sampled from (4), and finally the adjacency matrix X is sampled from (5). Random $\hat{\tau}$ s are sampled as normalized uniform variables. In each experiment, 100 fixed-point algorithm iterations are run in parallel for the 50 initial states, resulting in figures 2a and 2b which plot the norm of the difference of $\hat{\tau}$ between two successive iterations as a function of the iteration, for all graphs, and for the two experiments.

3.3 Clustering metrics

The literature offers various metrics to assess the quality of a clustering. For our experiments, we re-implemented some supervised and unsupervised metrics.

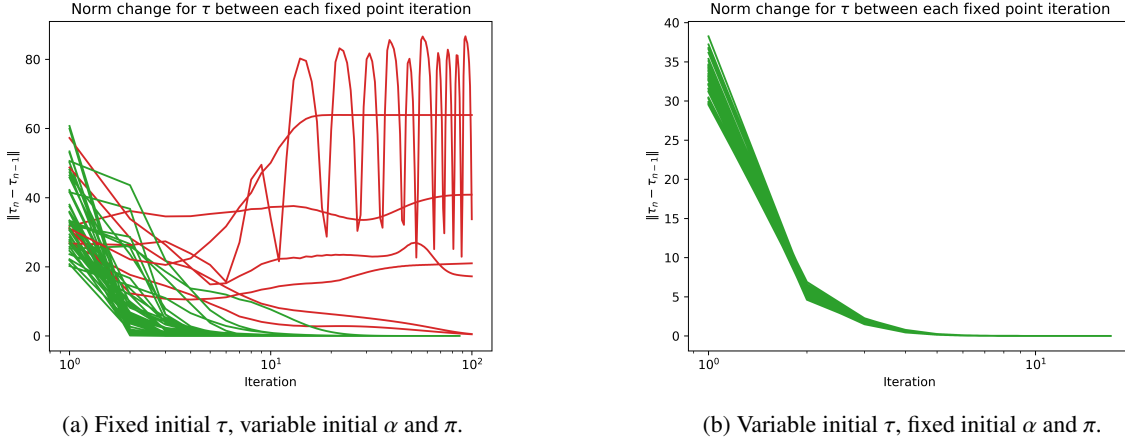


Figure 2: Convergence of the fixed point algorithm. Green curves converged within 100 iterations; red curves did not.

Supervised metrics By “supervised metrics”, we designate metrics that compare the labelling of nodes C^P as predicted by a clustering method, with a ground truth labelling C^T that is supposed to be given. We have used the Rand index (RI) and the normalized mutual information (NMI).

RI measures the similarity of two clusterings by measuring how consistently the two clusterings identify edges as connecting nodes of a same cluster, or of different clusters:

$$RI = \frac{a + b}{m} \quad (20)$$

with a the number of pair of nodes that are in a same cluster in both clusterings, b the number of pairs of nodes that are in different clusters in both clusterings and m the number of pairs of nodes in the graph.

MI is a measure that evaluates the mutual dependence between two variables. For two clusters $C = (C_i)_{1 \leq i \leq |C|}$ and $C' = (C'_j)_{1 \leq j \leq |C'|}$, it is defined as

$$MI(C, C') = \sum_{i=1}^{|C|} \sum_{j=1}^{|C'|} P(i, j) \log \left(\frac{P(i, j)}{P(i)P'(j)} \right) \quad (21)$$

with $P(i) = \frac{|C_i|}{n}$, $P'(j) = \frac{|C'_j|}{n}$ and $P(i, j) = \frac{|C_i \cap C'_j|}{n}$.

Both measures are insensitive to permutation of the clusters.

Unsupervised metrics Unsupervised metrics are defined without resorting to ground truth labelling.

The modularity M of a clustering is a common metric that compares the observed connectivity within clusters to its expected value for a random graph with the same degree distribution as the considered graph. It is defined as

$$M(C) = \sum_{i=1}^{|C|} \left(\frac{e_i}{m} - \left(\frac{d_i}{2m} \right)^2 \right) \quad (22)$$

with m the total number of edges, and e_i and d_i respectively the number of edges and total degree of nodes in cluster C_i .

Additionally, we use the clustering coefficient defined as the ratio of all triangles within a set of nodes C , over the

total number of triplets of nodes within this set; using the adjacency matrix, this writes:

$$CC(C) = \frac{\sum_{i,j,k \in C} X_{ij} X_{jk} X_{ki}}{\sum_{i \in C} k_i(k_i - 1)} \quad (23)$$

with k_i the degree of node i . We will use this metric both at the scale of clusters (to determine their individual qualities) and of whole graphs (to measure their clustering tendency). To close this section on metrics, note that the SBM model defined in Section 2.1 allows a computation of the expected clustering coefficient associated to graphs generated by a given set of parameters (α, π) :

$$CC_{SBM}(\alpha, \pi) = \frac{\sum_{q,l,m \in \{1, \dots, Q\}} \alpha_q \alpha_l \alpha_m \pi_{ql} \pi_{qm} \pi_{lm}}{\sum_{q,l,m \in \{1, \dots, Q\}} \alpha_q \alpha_l \alpha_m \pi_{ql} \pi_{qm}}. \quad (24)$$

3.4 Testing on SBM datasets

The aim of this second set of experiments performed on Daudin-EM is 3-fold:

1. assessing the capacity of the algorithm to recover the parameters of the model that generated a given graph;
2. assessing the performance of the algorithm in terms of graph clustering;
3. experimenting with ICL to select the number of clusters.

Building the dataset Using the SBM, we generate 6 batches of 10 graphs to test the method on. Each batch is generated using a specific set of parameters with the intent of modelling a specific type of graphs, as indicated by the name of the associated experiment. The parameters of the SBM for each batch are given in Table 1, with Dir the Dirichlet law, $\varepsilon_1 = 0.1$, $\varepsilon_2 = 0.01$, $a = 0.7$ and $b = 0.8$. Note that π was made symmetric, normalized and scaled in experiments 1 and 2. Sample graphs from each batch are visualized on Figure 4 in appendix.

#	Experiment Name	Hyper-parameters		Parameters	
		n	Q	α	π
1	Random-small	30	3	$\alpha \sim \text{Dir}(1.5)$	$\pi_{ij} \sim \mathcal{U}([0, 1])$
2	Random-large	500	3	$\alpha \sim \text{Dir}(1.5)$	$\pi_{ij} \sim \mathcal{U}([0, 1])$
3	Homophilic	150	3	$\alpha_i = (\frac{1}{Q})_i$	$\pi_{ii} = 1 - \varepsilon_1, \pi_{ij} = \varepsilon_2$
4	Homophilic-hard	150	5	$\alpha_i = (\frac{1}{Q})_i$	$\pi_{ii} \sim \mathcal{U}([0.5, 1]), \pi_{ij} = \varepsilon_2$
5	High-degree-minority	150	2	$\alpha^T = (\frac{9}{10}, \frac{1}{10})$	$\pi = \begin{pmatrix} \varepsilon_2 & a \\ a & b \end{pmatrix}$
6	Heterophilic	150	3	$\alpha_i = (\frac{1}{Q})_i$	$\pi_{ii} = \varepsilon_2, \pi_{ij} = 1 - \varepsilon_1$

Table 1: Parameters used for the generation of the SBM dataset.

Assessing parameter recovering performance Because the graphs are precisely generated using the same model as the one we are trying to fit, these experiments are primarily designed to estimate the performance of the proposed EM algorithm. In particular, a ground truth (α_T, π_T) is trivially given for the expected best-fitting model. Consequently, we could simply measure the distance between the laws induced by the predicted (α_P, π_P) and (α_T, π_T) . However, the intractability of the likelihood for SBM makes it hard to apply criteria like the Kullback-Leibler divergence or the χ^2 statistic. We have considered using the relative distance between $\mathcal{J}(R_X)$ computed for both sets of parameters, but ditched this metric due to its lack of interpretability across batches. Ultimately, we decided to simply sum the Euclidean distances between α_T and α_P and between π_T and π_P , normalized by the number of classes. More precisely, we retain the minimum such distance over all possible permutations of the predicted clusters, to avoid penalizing an indexing of clusters that would be different from the ground truth.

Assessing clustering performance To assess the clustering capacity of the algorithm, we use the metrics introduced in Section 3.3. In particular, the true and predicted clusterings are defined as $\forall i \in \{1, \dots, Q\}$,

$$C_i^T = \arg \max_q z_{iq} \quad \text{and} \quad C_i^P = \arg \max_q \hat{\tau}_{iq}. \quad (25)$$

with $\hat{\tau}$ the predicted value of τ output by the algorithm and z the latent variable sampled to generate the graph. Also, for more interpretability of modularity, instead of averaging over graphs, we choose the graph that produces largest modularity with the predicted parameters and compute ground truth and predicted modularity on this graph.

Testing the ICL criterion We test the ICL criterion separately, on a single graph from experiments 1, 4 and 5, testing values of Q respectively in $\{1, \dots, 4\}$, $\{1, 3, 5, 7\}$ and $\{1, \dots, 4\}$.

Running the experiments Each experiment is run for 100 iterations, and each time using 5 random initializations of α and π to limit the effects of the lack of robustness underlined in Section 3.2. For each experiment, excluding tests on ICL criterion, the true number of classes Q is considered given.

The results of the experiments are presented in Table 2 in appendix, along with sample predictions on Figure 5. Note that, due to the lack of robustness underlined in Section 3.2, the fixed-point algorithm did not converge within the budget

of 1000 iterations for all graphs. This is why the number of passed graphs is indicated in the results. Lastly, the ICL criterion wrongly predicts a best number of classes $Q = 1$ for experiment 1 and $Q = 3$ for experiment 3, but correctly predicts $Q = 2$ for experiment 5.

3.5 Comparison to alternatives

In the last set of experiments, the Daudin-EM method [1] described in Section 2.1 is tested on real graphs and compared with 2 other methods:

1. The alternative mixture model from [2] described in section 2.2;
2. The spectral clustering method from [3] described in section 2.3.

The results are then visualized and compared in terms of running time and clustering power, using the metrics defined in Section 3.3.

Datasets These methods are all tested on 2 famous datasets, commonly used as benchmarks for community detection algorithms:

1. Zachary’s Karate-club dataset: Published in an Anthropology journal in 1977 [9], it describes the friendship between 34 members of a karate club at a US university after an internal dispute that split the social relationships in two factions. The graph is undirected and contains 34 nodes and 78 edges. Figure 6a shows the ground truth of the dataset.
2. The Cora dataset: A citation network of scientific publications [10]. The graph is undirected and contains 2708 nodes(=papers) and 5429 edges(=citations). The nodes are labeled with 7 different classes. Figure 7a shows the dot-plot representation of the adjacency matrix of the ground truth, with associated labels.

A third dataset similar to that used in [1], the Escherichia-Coli metabolic network available on www.ecocyc.org [11], was also considered, but was eventually abandoned as it was not as interpretable as Cora. While the Karate-club dataset is a small graph with a clear community structure, the Cora dataset is a much larger graph with a more complex structure. Both datasets are constructed from real-world data and observations which makes them interpretable. They also have a

ground truth that we can use to evaluate the performance of the algorithms.

Running the experiments For Zachary’s karate club, Daudin-EM and Newman-EM are run for 100 iterations. For Daudin-EM, 5 parallel initializations of (α, π, τ) are used to increase the probability of converging. For Cora, they are run for 50 iterations, using 10 initializations of (α, π, τ) . In both cases, the spectral clustering algorithm is run once. For this set of experiments, the true number of classes Q was considered given.

Figures 6 and 7 show for both datasets the true clusterings and those proposed by the 3 methods, and tables 3 and 4 show the resulting metrics. For Cora, a reordering of the predicted clusters has been applied so as to coincide as much as possible with ground truth classes. To close this part, we also computed the estimation of the global clustering coefficient using (24) and obtained 0.211 for Zachary’s karate club and 0.006 for Cora.

4 RESULTS DISCUSSION

First set of experiments From the first set of experiments, we understand that the main weakness of Daudin-EM is the difficult convergence of the fixed-point algorithm of the E-step. In particular, as observed on Figure 2, this algorithm is particularly sensitive to the initialization on (α, π) , but not as much on the initialization of τ . Some ideas have been tested to solve this, such as applying a momentum-based formula of the form $\hat{\tau}_{n+1} = \hat{\tau}_n + \gamma \hat{\tau}_n^{\text{FP}}$ for the n^{th} E-step, with $\gamma \in (0, 1]$ an update rate and $\hat{\tau}_n^{\text{FP}}$ the result of the fixed-point algorithm; but this showed little improvement.

Second set of experiments In the SBM experiments, we observe that Daudin-EM performs generally rather poorly, both for fitting the SBM and clustering graphs. In particular, it has a tendency to predict only one cluster (resulting in null NMI), especially for clear-cut clusters (both homophilic and heterophilic). On the other hand, it brings better model-fitting performance for random graphs, especially large ones, and also yields higher supervised metrics in this case. Lastly, it achieves near-perfect results for the “High-degree-minority” case of experiment 5 according to both criteria, with a modularity that is rightly negative in this instance of graph type.

To sum up, Daudin-EM seems bad at detecting clusters, especially rather homogeneous ones. Our main observation is that this method seems very sensitive to the distribution of degrees of nodes, which would explain bad performance on experiments 3, 4 and 6 as rather homogeneous clusters are hardly distinguished by this quantity. This would also explain the best results observed for experiment 5. One way to alleviate this problem might be to offer semi-supervised guidance to enforce a favourable initialization.

Third set of experiments For both datasets, Daudin-EM brings mostly dissatisfying performance, especially when compared with Newman which always gives the best results. As for the spectral clustering method, it behaves well on Zachary’s karate club, but very badly on Cora, probably due to its large size. Daudin-EM is also particularly slow, notably

due to the number of iterations required for convergence of the E-step. Nevertheless, the clustering coefficients estimated using (24) fall close to their empirical values, although they are a bit underestimated.

For Zachary’s karate club, Daudin-EM yields low NMI and RI, as well as a negative modularity. On the other hand, the resulting per-cluster clustering coefficients are higher than that of the ground truth clusters in average, thus revealing clusters of better quality from that perspective.

On the other hand, Daudin-EM brought interesting results for Cora. Indeed, the supervised metrics are comparable to those for Newman-EM. Furthermore, although the final modularity is quite low with Daudin-EM, and despite some clusters having a clustering coefficient of 0, the method reveals some good-quality clusters, especially clusters 2 and 4. Cluster 2 is particularly well preserved, as “genetic algorithms” might be somewhat isolated from the other considered scientific subjects. Additionally, cluster 5 is interesting as it gathers papers that cite (or are cited by) many papers from many different scientific topics: this cluster probably gathers survey articles. Thereby, the method succeeds to bring a new perspective on the dataset that does not restrict to homophily-driven clustering.

5 CONCLUSION

The SBM-based variational EM method introduced by Daudin *et al.* in [1] has been presented, along with another EM-method from Newman & Leicht [2] and with the spectral-clustering baseline [3]. Daudin-EM relies on an approximate maximization of the log-likelihood of the law induced by the SBM model, so as to circumvent the intractability of this law - and of the conditional law given the membership of nodes to classes. We have illustrated that the main flaw of this algorithm is the unproven and unreliable convergence of the fixed-point-based E-step, a flaw not shared by Newman-EM which benefits from a closed-form solution. This is at the root of a lack of robustness of the method as well as to its lengthy execution due to numerous iterations being required. Testing on several datasets revealed poor performance of the method against Newman-EM, but better performance against spectral clustering on Cora. In particular, the experiments hinted towards a large sensitivity of the method to node degree distributions and a tendency to gather nodes into clusters as large as possible. On the other hand, it behaves better when applied to large graphs, and may reveal clustering properties that are not revealed by other clustering methods.

Further work may include trying other strategies to maximize the lower bound of (7) with respect to τ , for instance using a projected gradient descent; testing Daudin-EM on semi-supervised graph problems; and reflecting on how to extend this method to overlapping clusters, a class of problems over-viewed in [12].

REFERENCES

- [1] J.-J. Daudin, F. Picard, and S. Robin. A mixture model for random graphs. *Statistics and Computing*, 18(2):173–183, June 2008. ISSN 1573-1375. doi: 10.1007/s11222-007-9046-7. URL <https://doi.org/10.1007/s11222-007-9046-7>.
- [2] M. E. J. Newman and E. A. Leicht. Mixture models and exploratory analysis in networks. *Proceedings of the National Academy of Sciences*, 104(23):9564–9569, June 2007. ISSN 1091-6490. doi: 10.1073/pnas.0610537104. URL <http://dx.doi.org/10.1073/pnas.0610537104>.
- [3] Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001. URL https://proceedings.neurips.cc/paper_files/paper/2001/file/801272ee79cfde7fa5960571fee36b9b-Paper.pdf.
- [4] Santo Fortunato. Community detection in graphs. *Physics Reports*, 486(3-5):75–174, February 2010. ISSN 03701573. doi: 10.1016/j.physrep.2009.11.002. URL <https://linkinghub.elsevier.com/retrieve/pii/S0370157309002841>.
- [5] Petr Chunaev. Community detection in node-attributed social networks: A survey. *Computer Science Review*, 37:100286, August 2020. ISSN 15740137. doi: 10.1016/j.cosrev.2020.100286. URL <https://linkinghub.elsevier.com/retrieve/pii/S1574013720303865>.
- [6] Pratha Sah, Lisa O. Singh, Aaron Clauset, and Shweta Bansal. Exploring community structure in biological networks with random graphs. *BMC Bioinformatics*, 15(1):220, June 2014. ISSN 1471-2105. doi: 10.1186/1471-2105-15-220. URL <https://doi.org/10.1186/1471-2105-15-220>.
- [7] Tom A.B. Snijders and Krzysztof Nowicki. Estimation and Prediction for Stochastic Blockmodels for Graphs with Latent Block Structure. *Journal of Classification*, 14(1):75–100, January 1997. ISSN 1432-1343. doi: 10.1007/s003579900004. URL <https://doi.org/10.1007/s003579900004>.
- [8] Andrea Lancichinetti, Santo Fortunato, and János Kertész. Detecting the overlapping and hierarchical community structure in complex networks. *New Journal of Physics*, 11(3):033015, March 2009. ISSN 1367-2630. doi: 10.1088/1367-2630/11/3/033015. URL <https://iopscience.iop.org/article/10.1088/1367-2630/11/3/033015>.
- [9] Wayne Zachary. An information flow model for conflict and fission in small groups¹. *Journal of anthropological research*, 33, 11 1976. doi: 10.1086/jar.33.4.3629752.
- [10] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI Magazine*, 29(3):93, Sep. 2008. doi: 10.1609/aimag.v29i3.2157. URL <https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/2157>.
- [11] Ingrid M. Keseler, Socorro Gama-Castro, Amanda Mackie, Richard Billington, César Bonavides-Martínez, Ron Caspi, Anamika Kothari, Markus Krummenacker, Peter E. Midford, Luis Muñoz-Rascado, Wai Kit Ong, Suzanne Paley, Alberto Santos-Zavaleta, Pallavi Subhraveti, Víctor H. Tierrafria, Alan J. Wolfe, Julio Collado-Vides, Ian T. Paulsen, and Peter D. Karp. The EcoCyc Database in 2021. *Frontiers in Microbiology*, 12:711077, 2021. ISSN 1664-302X. doi: 10.3389/fmicb.2021.711077.
- [12] Jierui Xie, Stephen Kelley, and Boleslaw K. Szymanski. Overlapping community detection in networks: The state-of-the-art and comparative study. *ACM Computing Surveys*, 45(4):43:1–43:35, August 2013. ISSN 0360-0300. doi: 10.1145/2501654.2501657. URL <https://doi.org/10.1145/2501654.2501657>.

Appendix

A IMPLEMENTATIONS SPEED COMPARISON

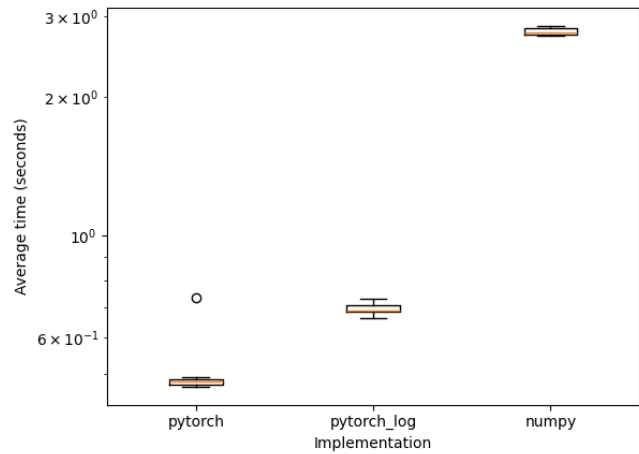


Figure 3: Speed comparison (in seconds) between implementations of Daudin-EM. Each implementation was run 10 times, for 10 runs on a graph randomly generated using SBM for 100 nodes and 3 classes.

B SECOND SET OF EXPERIMENTS

B.1 Samples

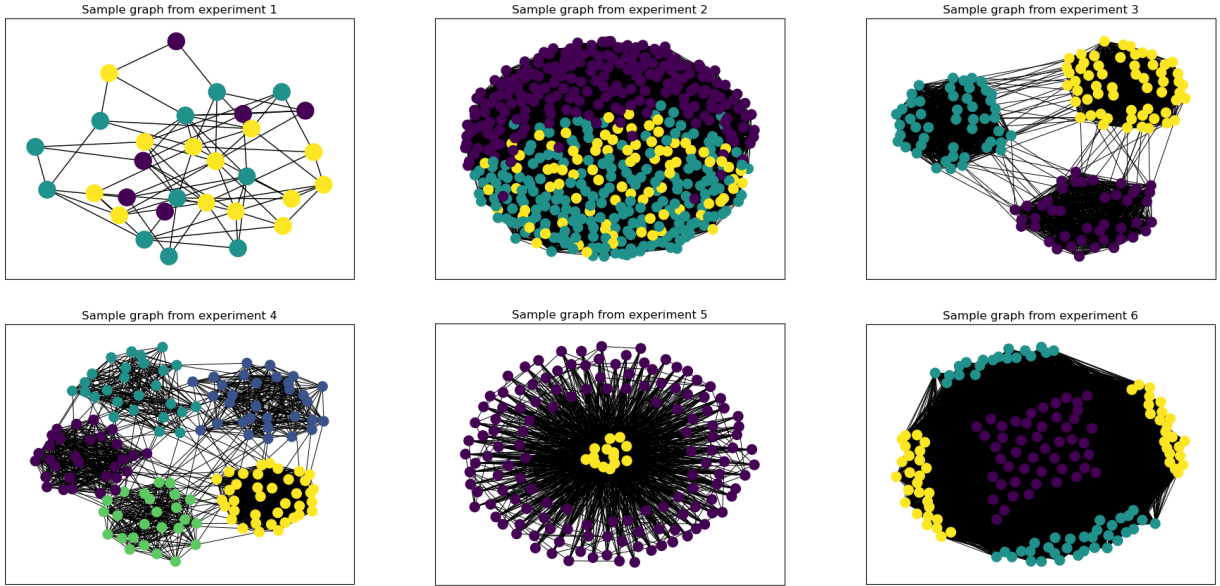


Figure 4: Sample graphs from the second set of experiments

B.2 Results

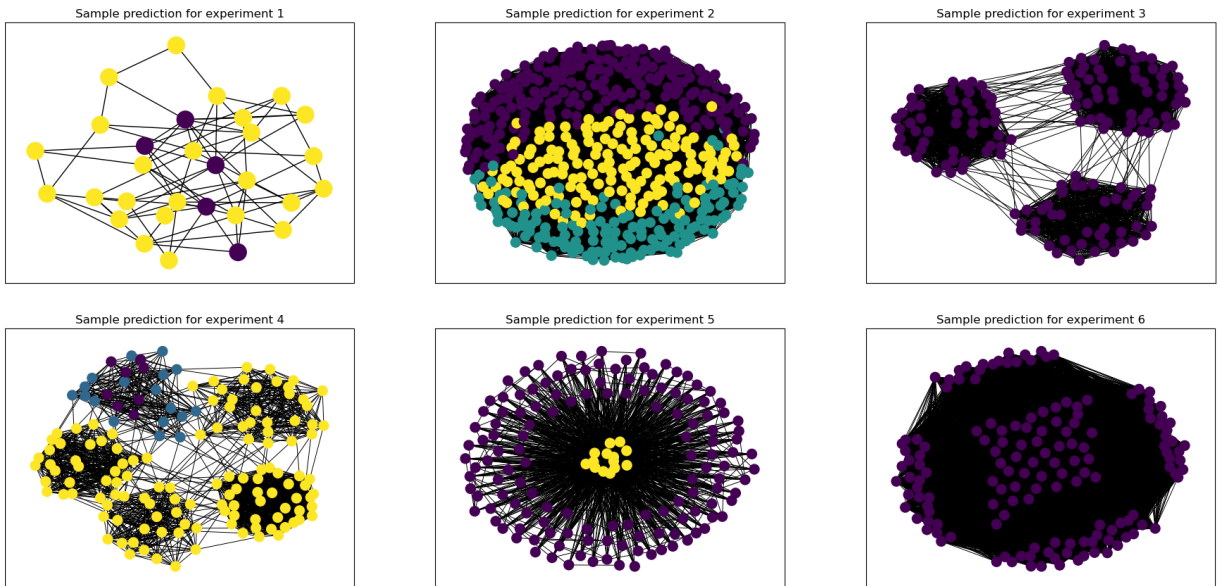


Figure 5: Some predicted graphs from the second set of experiments

Experiment #	# graphs passed	Param. dist.	NMI	RI	G.t. M	Pred. M
1	10	0.13 ± 0.06	0.07 ± 0.07	0.52 ± 0.10	0.08	0.02
2	8	0.08 ± 0.03	0.33 ± 0.13	0.66 ± 0.13	0.07	0.07
3	10	0.22 ± 0.12	0 ± 0	0.33 ± 0.00	0.63	0.00
4	4	0.22 ± 0.02	0.11 ± 0.18	0.27 ± 0.12	0.71	0.15
5	10	0.02 ± 0.01	1 ± 0	1 ± 0	-0.38	-0.38
6	9	0.18 ± 0.08	0 ± 0	0.33 ± 0.00	-0.33	0.00

Table 2: Results of the second set of experiments.

C THIRD SET OF EXPERIMENTS

C.1 Zachary’s Karate club

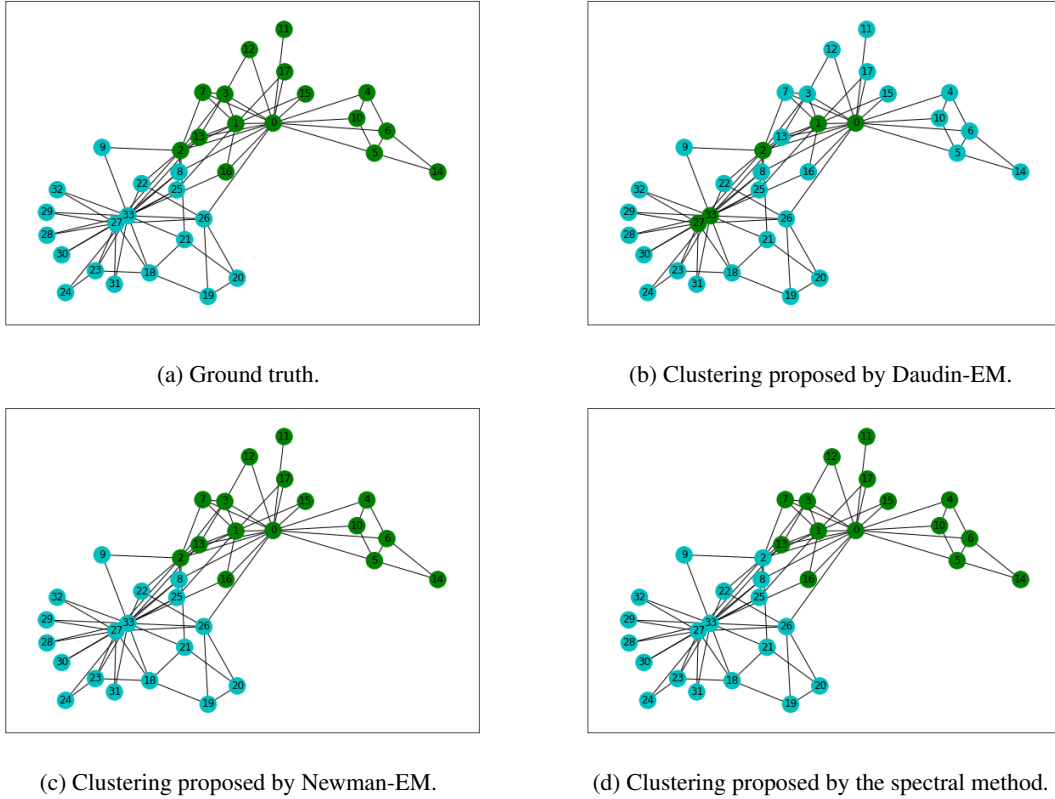


Figure 6: Clustering results for Zachary’s Karate Club.

Method	Time (s)	NMI	RI	M	CC	CC (per cluster)	
						1	2
Ground truth	-	-	-	0.37	0.26	0.42	0.26
Daudin-EM	137	0.01	0.49	-0.21	-	0.50	0.23
Newman-EM	2	1.00	1.00	0.37	-	0.41	0.26
Spectral	0	0.84	0.94	0.36	-	0.35	0.24

Table 3: Results on Zachary’s Karate Club

C.2 Cora

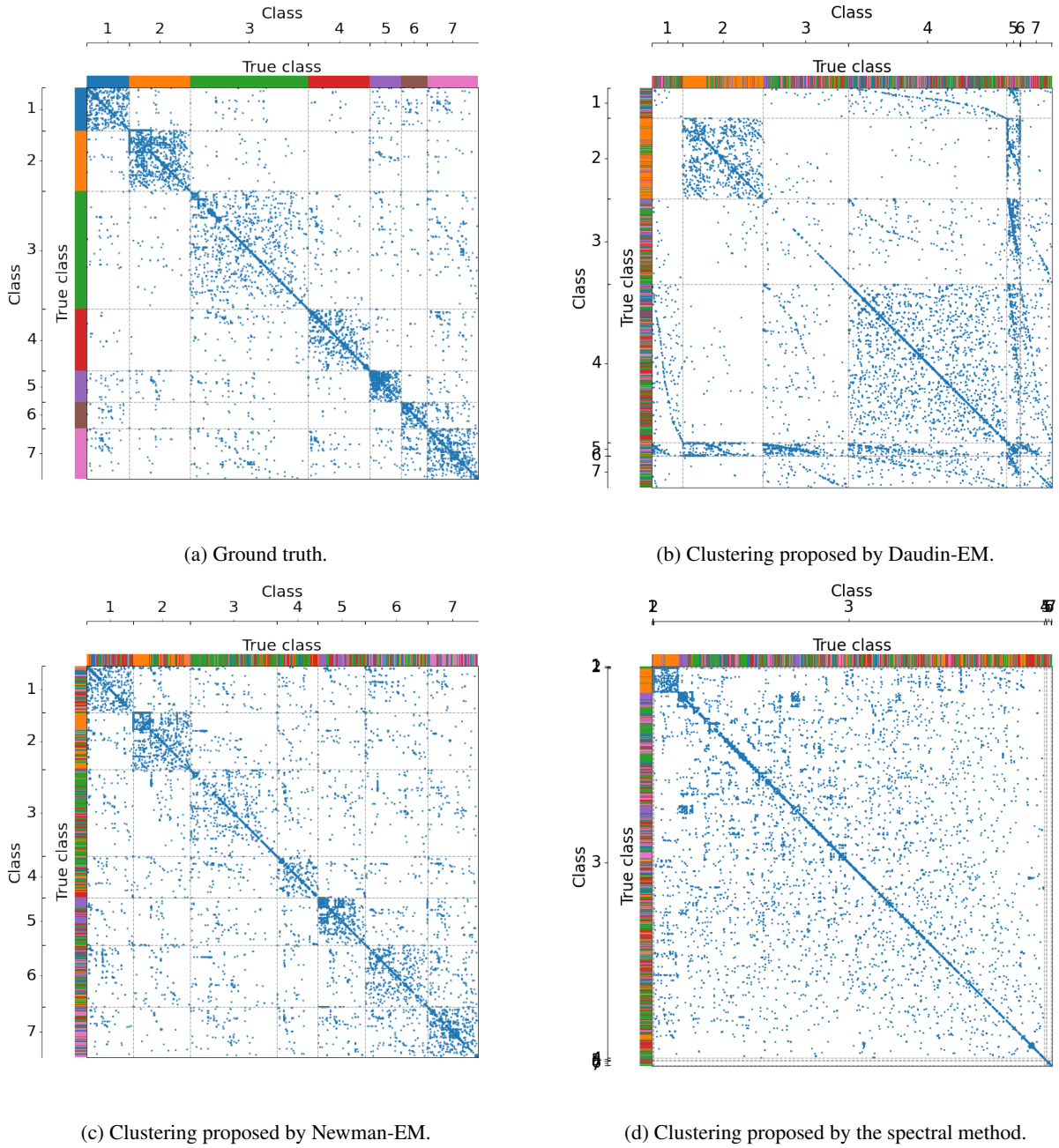


Figure 7: Dot-plot of the results for Cora.

G.T.: 1: Case Based; 2: Genetic Algorithms; 3: Neural Networks; 4: Probabilistic Methods; 5: Reinforcement Learning; 6: Rule Learning; 7: Theory.

Method	Time (s)	NMI	RI	M	CC	CC (per cluster)						
						1	2	3	4	5	6	7
Ground truth	-	-	-	0.64	0.09	0.19	0.06	0.12	0.23	0.10	0.22	0.16
Duadin-EM	3802	0.15	0.70	0.22	-	0	0.16	0	0.29	0.25	0	0.86
Newman-EM	27	0.18	0.76	0.53	-	0.25	0.06	0.16	0.28	0.12	0.29	0.18
Spectral	1	0.01	0.21	0.02	-	0	0	0.09	0	0	0.66	0

Table 4: Results on Cora.