
Modularity-based Sparse Soft Graph Clustering

Alexandre Hollocoou
INRIA, Paris

Thomas Bonald
Telecom Paristech

Marc Lelarge
INRIA & ENS, Paris

Abstract

Clustering is a central problem in machine learning for which graph-based approaches have proven their efficiency. In this paper, we study a relaxation of the modularity maximization problem, well-known in the graph partitioning literature. A solution of this relaxation gives to each element of the dataset a probability to belong to a given cluster, whereas a solution of the standard modularity problem is a partition. We introduce an efficient optimization algorithm to solve this relaxation, that is both memory efficient and local. Furthermore, we prove that our method includes, as a special case, the Louvain optimization scheme, a state-of-the-art technique to solve the traditional modularity problem. Experiments on both synthetic and real-world data illustrate that our approach provides meaningful information on various types of data.

1 Introduction

Modularity is a quality function defined on node partitions that has been widely used to tackle the problem of node clustering in graphs [9, 26]. Numerous algorithms have been proposed to solve the problem of maximizing this objective function [6, 13, 24]. The Louvain algorithm [2], which is built in a number of graph analysis software packages, is arguably the most popular approach to solve this problem. A solution to the modularity optimization problem is a partition of the graph nodes. In particular, a node cannot belong to more than one cluster. Yet, in many real-life applications, it makes more sense to allow some elements to belong to multiple clusters. For instance, if we consider

the problem of clustering the users of a social network such as Facebook or LinkedIn, we see that a person can belong to several social circles, such as her family, her colleagues and her friends, and we might be interested in recovering all these circles with a clustering algorithm, even if they do not form a partition.

In order to overcome this limit of the classical definition of modularity, we introduce a relaxation of the problem to perform a so-called *soft clustering* that allows partial membership like *fuzzy clustering* methods [1, 8]. A solution to this new problem gives for each node a degree of membership to each cluster, instead of returning a simple node partition. More precisely, we obtain a probability $p_{ik} \in [0, 1]$ for each node i to belong to a given cluster k , where the classical modularity-based techniques implicitly consider that this probability can be either 0 or 1.

The main contribution of this paper is to introduce an efficient algorithm to find an approximation of this *soft* version of the modularity maximization problem. This algorithm has four main advantages. First, the number of clusters does not need to be specified. Second, the algorithm is local, in the sense that each update of the membership information of a given node depends only on its direct neighbors in the graph. Third, the solution found by the algorithm is sparse, i.e. most of the membership probabilities p_{ik} returned by our algorithm are equal to 0, which guarantees an efficient storage of the solution. Finally, an update performed by our algorithm reduces to an update performed by the Louvain algorithm as soon as the unique parameter t of our algorithm is larger than a threshold w/δ that depends on the graph.

The remainder of the paper is organized as follows. In section 2, we present the related work on the topic of modularity optimization. We introduce the relaxation of the modularity maximization problem in section 3. In section 4, we present our optimization method and obtain theoretical guarantees about its convergence. In section 5, we take benefits of the specificities of our optimization technique to propose an algorithm that is both local and memory efficient. In section 6, we study the ties of our approach to the Louvain algorithm.

Finally, in section 7, we present experimental results on synthetic and real-world data, and section 8 concludes the paper. Complete proofs of the main results of the paper are presented in the supplementary material.

In the rest of the paper, we consider that we are given a weighted and undirected graph $G = (V, E, \mathbf{W})$. We use V to denote the set of nodes, E the set of edges, and $\mathbf{W} = (W_{ij})_{i,j \in V}$ the adjacency matrix of the graph. If $(i, j) \in E$, $W_{ij} > 0$ is the weight of edge (i, j) , and if $(i, j) \notin E$, $W_{ij} = 0$. We use w_i to denote the weighted degree of node i , $w_i = \sum_{j \in V} W_{ij}$, and w the total weight of the graph $w = \sum_{i \in V} w_i$. We use $\text{Nei}(i)$ to denote the set of the neighbors of node $i \in V$, and $j \sim i$ as a notation for $j \in \text{Nei}(i)$. Finally, we use $n = |V|$ to denote the number of nodes, and $m = |E|$ the number of edges of the graph.

2 Related work

2.1 Modularity optimization

The modularity function (1) has first been introduced by Newman and Girvan [26] to measure the quality of a partition of graph nodes. For a given node partition P , modularity is defined as

$$Q(P) = \frac{1}{w} \sum_{C \in P} \sum_{i,j \in C} \left(W_{ij} - \frac{w_i w_j}{w} \right). \quad (1)$$

Modularity Q can be interpreted as the difference between the probability to pick an edge in G between two nodes of the same cluster C , i.e. an *intra-cluster* edge, and the probability to pick such an intra-cluster edge in a random graph generated with the so-called *null model*, where the probability that an edge (i, j) exists is $w_i w_j / w$.

Clustering approaches based on modularity try to solve the modularity maximization problem

$$\max_{P \in \mathcal{P}} Q(P), \quad (2)$$

where \mathcal{P} is the set of the partitions of V .

We refer to the Newman and Girvan modularity as the *hard modularity*, to distinguish it from the *soft modularity* that we introduce in Section 3. The hard modularity optimization problem (2) cannot be solved exactly in real applications as it has been proven to be NP-hard [3]. However, many methods have been proposed to find good approximations of the modularity optimum in reasonable time. These techniques includes greedy algorithms [24], spectral approaches [25], and simulated annealing methods [13]. A thorough review of existing algorithms can be found in [9]. One of the most popular algorithm for finding a good approximation of the modularity maximum is the Louvain

algorithm [2], which is widely used in benchmarks for its ability to scale up to very large graphs (e.g. the Friendster social graph with 1.8 billions edges [29]). In this paper, we introduce an algorithm for optimizing the relaxation of the modularity maximization problem (4) that takes a parameter t , and we prove that each update performed by our algorithm reduces to the update performed by the Louvain algorithm if t is large enough. This constitutes a strong guarantee since the Louvain algorithm has proven its worth on multiple real-life applications.

2.2 Modularity relaxation

Relaxations of the hard modularity optimization problem have been tackled by few articles [4, 12, 14, 16, 20, 21, 27]. They all introduce membership matrices $\mathbf{p} \in \mathbb{R}^{n \times K}$, with $K \geq 1$, where $p_{ik} \geq 0$ represents the degree of membership of node $i \in V$ to the k^{th} cluster, $k \in \llbracket 1, K \rrbracket$. K is a given positive integer, that represents the maximum number of clusters. These relaxations all have similar forms and are very close to the one we introduce in Section 3. These works differ in the optimization techniques they use to solve the problem. In [4], Chang et al. use a softmax parameterization for \mathbf{p} , defining for each $i \in V$ and k $p_{ik} = \exp(\theta_{ik}) / (\sum_l \exp(\theta_{il}))$ where $\boldsymbol{\theta} \in \mathbb{R}^{n \times K}$ is the parameter to optimize. In [27], Nicosia et al. use a genetic algorithm to optimize a more general relaxation of the modularity problem. In [12], Griechisch et al. do not directly study the optimization of the relaxation problem, but they rely on an external quadratic solver. Finally, in [14], Havens et al. use a spectral approach that does not directly solve the relaxation of the modularity problem, but where modularity is used as a selection criterion. The main limitation of these methods lies in the maximum number of clusters K that must be specified. We could get around this issue by taking large values for K , but all the approaches cited above do not scale well to large K . Indeed, the solutions $\mathbf{p} \in \mathbb{R}^{n \times K}$ found by these methods are dense matrices. In other words, the number of parameters to store in memory and to optimize is in $O(nK)$, which quickly becomes prohibitive for large values of K . For instance, in the approach of [4], the matrix \mathbf{p} is the densest possible matrix, i.e. all its coefficients are positive. In the approach of [27], the genetic algorithm starts with dense random matrices of $\mathbb{R}^{n \times K}$, and its hybridation and mutation mechanisms do not lead to sparser matrices, so that all the coefficients of the solution \mathbf{p} found by this algorithm are positive with probability 1.

In this paper, we introduce an efficient algorithm to solve a relaxation of the modularity optimization problem with a membership matrix $\mathbf{p} \in \mathbb{R}^{n \times n}$. Thus, the maximum number of clusters does not need to be spec-

ified. Besides, unlike the approaches presented above, the updates performed by our algorithm preserve the sparsity of the solution \mathbf{p} and are *local*, in the sense that the membership vector \mathbf{p}_i of a node $i \in V$ is updated using only membership information from its neighbors in the graph. Thus, our algorithm can easily scale up to large datasets with large number of clusters, which is not the case of the algorithms listed above.

As proposed in [32], the problem of soft graph clustering can also be tackled by a combination of spectral embedding [10, 31] and soft-clustering in the vector space, with algorithms such as the fuzzy C-means algorithm [8] or the NEO k-means algorithm [32]. However, the number of clusters K has to be specified for these techniques. Besides, spectral embedding methods do not scale well to large size graphs [28]. For instance, they are unable to handle the subgraph of Wikipedia that we use in our experiments in Section 7. Finally, it is worth mentioning that several algorithms have been proposed for the related problem of overlapping community detection in networks [18, 34, 36]. Whereas these approaches return clusters $C \subset V$, soft-clustering methods give for each node the degrees of membership to all clusters, which represents a richer information. Note that some overlapping community detection algorithms also allow soft clustering [19]. In the survey of [33], our algorithm would fall into the class of local expansion algorithms, where all algorithms have worst case complexity $O(n^2)$ (except iCLD dealing with dynamic networks).

3 Soft modularity

Given an ordered partition $P = (C_1, \dots, C_K)$ of the nodes of V (i.e. a partition of V whose sets are given in a certain order), we define the membership matrix \mathbf{p} associated with P as follows: $\forall i \in V, \forall k \in \llbracket 1, K \rrbracket$, $p_{ik} = 1$ if $i \in C_k$, and 0 otherwise. It is easy to see that P is an optimal solution of the modularity maximization problem (2) if and only if its associated membership matrix \mathbf{p} is an optimal solution of

$$\begin{aligned} & \underset{\mathbf{p} \in \mathbb{Z}^{n \times n}}{\text{maximize}} && \frac{1}{w} \sum_{i,j \in V} \sum_{k=0}^n \left(W_{ij} - \frac{w_i w_j}{w} \right) p_{ik} p_{jk} \\ & \text{subject to} && \forall i \in V, \sum_{k=1}^n p_{ik} = 1, \forall k \in \llbracket 1, n \rrbracket, p_{ik} \geq 0. \end{aligned} \quad (3)$$

Given a membership matrix \mathbf{p} , we use \mathbf{p}_i to denote the vector that corresponds to the i^{th} line of \mathbf{p} and $\mathbf{p}_{\cdot k}$ to denote the vector that corresponds to its k^{th} column. From the new formulation (3), we introduce the natural relaxation of the problem where coefficients

of \mathbf{p} can take real values

$$\begin{aligned} & \underset{\mathbf{p} \in \mathbb{R}^{n \times n}}{\text{maximize}} && Q(\mathbf{p}) \\ & \text{subject to} && \forall i \in V, \mathbf{1}^T \mathbf{p}_i = 1, \mathbf{p}_i \geq 0. \end{aligned} \quad (4)$$

where $Q(\mathbf{p}) = \frac{1}{w} \sum_{i,j \in V} (W_{ij} - \frac{w_i w_j}{w}) \mathbf{p}_i^T \mathbf{p}_j$.

Note that if \mathbf{p} is a feasible solution of this problem, then p_{ik} can be interpreted as a probability. It corresponds to the probability for a node i to belong to cluster k . So, by solving this relaxation of the modularity optimization problem we can obtain nodes that belong to multiple clusters unlike in the classical modularity maximization problem. We refer to a solution to this problem as a *soft clustering* of the nodes of G . We refer to this new objective function $Q(\mathbf{p})$ as the *soft modularity*. Remark that, when $\mathbf{p}_i \geq 0$, the constraint $\forall i \in V, \mathbf{1}^T \mathbf{p}_i = 1$ is equivalent to the l_1 constraint $\|\mathbf{p}_i\|_1 = 1$, where $\|\cdot\|_1$ denotes the l_1 norm.

4 Alternating projected gradient descent

We can rewrite the relaxation of the modularity optimization problem as a classic minimization problem $\min_{\mathbf{p} \in \mathcal{X}} J(\mathbf{p})$, where $J : \mathbf{p} \mapsto -Q(\mathbf{p})$ is the cost function and $\mathcal{X} = \{\mathbf{p} \in \mathbb{R}^{n \times n} : \forall i, \mathbf{p}_i \geq 0, \mathbf{1}^T \mathbf{p}_i = 1\}$ is the set of constraints. We define the normalized weighted Laplacian of the graph, as $\mathcal{L} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$, where \mathbf{D} is the diagonal matrix whose entries are the weighted degrees of the nodes: $\mathbf{D} = \text{diag}(w_1, \dots, w_n)$. The Laplacian matrix is symmetric and semi-definite positive [5]. We use $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ to denote its eigenvalues.

Theorem 4.1. *The function J is convex if and only if the second lowest eigenvalue λ_2 of the normalized Laplacian \mathcal{L} verifies $\lambda_2 \geq 1$.*

Proof. We prove in the supplementary material that the hessian matrix \mathbf{H} of J is congruent to a matrix whose eigenvalues are $0, \lambda_2 - 1, \dots, \lambda_n - 1$. \square

The condition $\lambda_2 \geq 1$, which corresponds to a large spectral gap, is in general not satisfied. For instance, Lemma 1.7. in [5] proves that $\lambda_2 \leq 1$ if G is unweighted and not complete. Therefore, the loss function associated with our problem is in general non-convex. However, it is convex in \mathbf{p}_i for a graph with no self-loop as shown in the following proposition. In the rest of the paper, we assume that the graph G does not contain any self-loop.

Proposition 4.2. *If the graph does not contain any self loops, i.e. if $W_{ii} = 0$ for all node i , the function $\mathbf{p} \mapsto J(\mathbf{p})$ is convex with respect to variable \mathbf{p}_i for all $i \in V$.*

Proof. If $W_{ii} = 0$, the hessian matrix \mathbf{H}_i of J with respect to \mathbf{p}_i is $\mathbf{H}_i = 2(w_i/w)^2 \mathbf{I}$. \square

With this result, we can apply the methods of convex optimization with convex constraints to the problem $\min_{\mathbf{p}_i \in \mathcal{Y}} J(\mathbf{p})$, with fixed \mathbf{p}_j for $j \neq i$, where the set \mathcal{Y} corresponds to the probability simplex of \mathbb{R}^n : $\mathcal{Y} = \{\mathbf{q} \in \mathbb{R}^n : \mathbf{q} \geq 0, \mathbf{1}^T \mathbf{q} = 1\}$. We use $\pi_{\mathcal{Y}}$ to denote the euclidean projection onto the probability simplex \mathcal{Y} , $\pi_{\mathcal{Y}}(\mathbf{x}) = \arg \min_{\mathbf{y} \in \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|^2$. Then, using the projected gradient descent method [11, 22], we can define an update rule for variable \mathbf{p}_i such that modularity is non-decreasing at each step. The gradient of J with respect to \mathbf{p}_i is: $\nabla_i J = -\frac{2}{w} \sum_{j \in V} (W_{ij} - w_i w_j / w) \mathbf{p}_j$.

Theorem 4.3. *The soft modularity objective function $Q(\mathbf{p})$ is non-decreasing under the update rule*

$$\mathbf{p}_i \leftarrow \pi_{\mathcal{Y}} \left(\mathbf{p}_i + \frac{2t}{w} \sum_{j \in V} \left(W_{ij} - \frac{w_i w_j}{w} \right) \mathbf{p}_j \right) \quad (5)$$

for all node $i \in V$ if the step size t verifies $t < (w/w_i)^2$. $Q(\mathbf{p})$ is invariant under all these update rules iff \mathbf{p}_i maximizes $Q(\mathbf{p})$ with fixed \mathbf{p}_j , $j \neq i$, for all node $i \in V$.

Proof. In the supplementary material, we prove that $Q(\mathbf{p}^+) \geq Q(\mathbf{p}) + t \left(1 - t \left(\frac{w_i}{w}\right)^2\right) \|\mathbf{G}_i(\mathbf{p})\|^2$ where \mathbf{p}^+ is the matrix \mathbf{p} after the update, and $\mathbf{G}_i(\mathbf{p}) = (\mathbf{p}_i - \mathbf{p}_i^+)/t$. \square

We now consider the natural algorithm that cycles through the nodes of the graph to apply the update (5).

Theorem 4.4. *The algorithm converges to a local maximum of the soft modularity function $\mathbf{p} \mapsto Q(\mathbf{p})$ which is a fixed point of the updates of (5).*

Proof. This is a direct consequence of the analysis conducted in the proof of Theorem 4.3 provided in the supplementary material. \square

5 Soft clustering algorithm

In the previous section, we have presented updates rules that can be applied in an alternating fashion to find a local maximum of the soft modularity $Q(\mathbf{p})$. However, there are three major issues with a naive implementation of this method. First, the gradient descent update for each node i is computationally expensive because the update rule (5) involves a sum over all nodes of V . Then, the size of a solution \mathbf{p} is n^2 , which is prohibitive for large datasets. Finally, the computational cost of the projection $\pi_{\mathcal{Y}}$ onto the probability simplex is classically in $O(n \log n)$ [7], which can be a limiting factor in practice.

In the present section, we present an efficient implementation of the algorithm that solves these three problems. In particular, our implementation guarantees that the update step for node i only requires local computation, is memory efficient, and uses a fast mechanism for projection.

5.1 Local gradient descent step

The update of Theorem 4.3 relative to node $i \in V$ can be decomposed into two steps:

1. $\hat{\mathbf{p}}_i \leftarrow \mathbf{p}_i + \frac{2t}{w} \sum_{j \in V} (W_{ij} - \frac{w_i w_j}{w}) \mathbf{p}_j$.
2. $\mathbf{p}_i \leftarrow \pi_{\mathcal{Y}}(\hat{\mathbf{p}}_i)$.

At first sight, step 1 seems to depend on information from all nodes of V , and to require the computation of a sum over n terms. However, the gradient descent step can be written as

$$\hat{\mathbf{p}}_i \leftarrow \mathbf{p}_i + \frac{2t}{w} \sum_{j \sim i} W_{ij} (\mathbf{p}_j - \bar{\mathbf{p}}),$$

where $\bar{\mathbf{p}}$ is weighted average of vector \mathbf{p}_j : $\bar{\mathbf{p}} = \sum_{j \in V} \frac{w_j}{w} \mathbf{p}_j$. The vector $\bar{\mathbf{p}} \in \mathbb{R}^n$ can be stored and updated throughout the algorithm. Then, the computation of $\hat{\mathbf{p}}_i$ is purely local, in the sense that it only requires information from neighbors of node i . Moreover, the sum to compute contains only d_i terms, where $d_i = |\text{Nei}(i)|$ is the unweighted degree of node i . In most real-life graphs, the degree distribution follows a heavy tail distribution and most of nodes have a low degree d_i .

5.2 Cluster membership representation

The cluster membership variable \mathbf{p} is a $n \times n$ matrix. In a naive implementation, the memory cost of storing \mathbf{p} is therefore in $O(n^2)$. However, we will see below that the projection step can be written $\pi_{\mathcal{Y}}(\hat{\mathbf{p}}_i) = (\hat{\mathbf{p}}_i - \theta \mathbf{1})_+$ with $\theta \geq 0$, where $[(\mathbf{x})_+]_k = \max(x_k, 0)$. Thus, the projection acts as a thresholding step. We expect it to have the same effect as a Lasso regularization¹ and to lead to a sparse vector \mathbf{p}_i . To take benefit of this sparsity, we only store the non-zero coefficients of \mathbf{p} . If, for all node $i \in V$, the number of non-zero values in vector \mathbf{p}_i is lower than a given L , then the memory cost of storing \mathbf{p} is in $O(nL)$. We will see, in our experiments in Section 7, that the average number of positive components that we observe for vectors \mathbf{p}_i throughout the execution of our algorithm is lower than 2, and that the maximum number of positive components that we record is $L = 65$ for a graph with 731, 293 nodes.

¹ The l_1 proximal operator is $x \mapsto \text{sign}(x)(|x| - \lambda)_+$.

5.3 Projection onto the probability simplex

The classic algorithm to perform the projection $\pi_{\mathcal{Y}}(\hat{\mathbf{p}}_i)$ onto \mathcal{Y} is described in [7]. It starts by sorting the vector $\hat{\mathbf{p}}_i$ into $\boldsymbol{\mu}$: $\mu_1 \geq \mu_2 \geq \dots \geq \mu_n$. Then, it finds $\rho = \max \left\{ j \in [n], \mu_j - \frac{1}{j} \left(\sum_{r=1}^j \mu_r - 1 \right) > 0 \right\}$, and defines $\theta = \frac{1}{\rho} \left(\sum_{i=1}^{\rho} \mu_i - 1 \right)$ so that $\pi_{\mathcal{Y}}(\hat{\mathbf{p}}_i) = (\hat{\mathbf{p}}_i - \theta \mathbf{1})_+$. This algorithm runs in $O(n \log n)$ because of the sorting step.

In our case, the complexity can be reduced to $O(L_i \log L_i)$, with $L_i = |\text{supp}_i(\mathbf{p})|$, where $\text{supp}_i(\mathbf{p}) = \{k \in \llbracket 1, n \rrbracket, \exists j \in \text{Nei}(i) \cup \{i\}, p_{ik} > 0\}$ is the union of the supports of vectors \mathbf{p}_j . for $j \in \text{Nei}(i) \cup \{i\}$. L_i is upper-bounded by $L(d_i + 1)$ which is typically much smaller than the total number of nodes n .

Proposition 5.1. *In order to compute $\pi_{\mathcal{Y}}(\hat{\mathbf{p}}_i)$, we only need to sort the components $k \in \text{supp}_i(\mathbf{p})$ of $\hat{\mathbf{p}}_i$ to determine ρ and θ . All components $k \notin \text{supp}_i(\mathbf{p})$ of $\pi_{\mathcal{Y}}(\hat{\mathbf{p}}_i)$ are set to zero.*

Proof. We prove this result in the supplementary material. \square

5.4 MODSOFT

Finally, the algorithm can be described as follows.

- **Initialization:** $\mathbf{p} \leftarrow \mathbf{I}$ and $\bar{\mathbf{p}} \leftarrow \mathbf{w}/w$.
- **One epoch:** For each node $i \in V$,
 - $\forall k \in \text{supp}_i(\mathbf{p}), \hat{p}_{ik} \leftarrow p_{ik} + t' \sum_{j \sim i} W_{ij} (p_{jk} - \bar{p}_k)$
 - $\mathbf{p}_i^+ \leftarrow \text{project}(\hat{\mathbf{p}}_i)$
 - $\bar{\mathbf{p}} \leftarrow \bar{\mathbf{p}} + (w_i/w)(\mathbf{p}_i^+ - \mathbf{p}_i)$ and $\mathbf{p}_i \leftarrow \mathbf{p}_i^+$.

where \mathbf{w} is the vector $(w_i)_{1 \leq i \leq n}$, project is the adaptation of the algorithm of [7] presented above, and t' is the effective learning rate $t' = 2t/w$. One epoch of the algorithm is typically repeated until the increase of modularity falls below a given threshold $\epsilon > 0$, as in the Louvain algorithm. Note that we chose, in this implementation, to put each node in its own cluster at initialization, but our algorithm could also be initialized with the results of another algorithm (e.g. Louvain). We refer to this algorithm as MODSOFT (MODularity SOFT clustering).

6 Link with the Louvain algorithm

In this section, we show that if the learning rate t is larger than a graph-specific threshold, one update performed by our algorithm reduces to a local update performed by the Louvain algorithm. First, we observe that, if the membership matrix \mathbf{p} verifies $\mathbf{p} \in \{0, 1\}^{n^2}$, then the update rule (5) for a node $i \in V$ becomes

$$\mathbf{p}_i \leftarrow \pi_{\mathcal{Y}} \left(\left[\mathbf{1}_{\{i \in C_k\}} + \frac{2t}{w} \left(w_i(C_k) - w_i \frac{\text{Vol}(C_k)}{w} \right) \right]_k \right)$$

where $C_k = \{i \in V : p_{ik} = 1\}$ refers to the k^{th} cluster defined by \mathbf{p} , $w_i(C)$ denotes the degree of node i in cluster $C \subset V$, $w_i(C) = \sum_{j \in C} W_{ij}$, and $\text{Vol}(C)$ denotes the volume of cluster C , $\text{Vol}(C) = \sum_{j \in C} w_j$.

In the following, we assume that: $\forall C, C' \subset V, \forall i \in V$,

$$C \neq C' \Rightarrow w_i(C) - \frac{w_i \text{Vol}(C)}{w} \neq w_i(C') - \frac{w_i \text{Vol}(C')}{w}. \quad (\text{H})$$

The hypothesis (H) is non-binding in real cases, because, if the weights $(W_{ij})_{(i,j) \in E}$ are continuously distributed, then we will have (H) with probability 1; and if they follow a discrete distribution, we can always add a small noise to the weights, so that (H) is verified with probability 1. Moreover, note that this hypothesis does not hold only in specific symmetrical graphs. In real-life data, these symmetrical cases are rare.

Proposition 6.1. *Let $\mathbf{p} \in \mathcal{X}$ be a membership matrix. If $\mathbf{p} \in \{0, 1\}^{n^2}$, if the hypothesis (H) is verified, and if $t > w/\delta$ where*

$$\delta = \min_{\substack{C, C' \subset V \\ i \in V \\ C \neq C'}} \left| \left(w_i(C) - \frac{w_i \text{Vol}(C)}{w} \right) - \left(w_i(C') - \frac{w_i \text{Vol}(C')}{w} \right) \right|,$$

then the update rule (5) for node $i \in V$ reduces to $p_{ik} \leftarrow 1$ if $k = \arg \max_{l: j \in C_l, j \sim i} \left[w_i(C_l) - w_i \frac{\text{Vol}(C_l)}{w} \right]$, and

$p_{ik} \leftarrow 0$ otherwise, where C_k denotes the k^{th} cluster defined by \mathbf{p} .

Proof. We show in the supplementary material that assuming that the updated vector \mathbf{p}_i has more than one positive component leads to a contradiction with $t > w/\delta$. \square

In particular, this result shows that if $\mathbf{p} \in \{0, 1\}^{n^2}$ and t is large enough, then for each update of our algorithm, the updated membership matrix \mathbf{p}^+ verifies $\mathbf{p}^+ \in \{0, 1\}^{n^2}$, which corresponds to the maximum sparsity that can be achieved by a feasible solution $\mathbf{p} \in \mathcal{X}$.

The Louvain algorithm [2] is a popular heuristic to find a local maximum for the *hard* modularity maximization problem (2). One epoch of the main routine of the Louvain algorithm considers successively each node $i \in V$, and (1) removes i from its current cluster, (2) applies an update rule, that we refer to as LouvainUpdate(i), which consists in transferring i to the cluster that leads to the largest increase in hard modularity. Proposition 6.2 gives an explicit formula to pick the cluster C_{k^*} to which node i is transferred by LouvainUpdate(i).

Proposition 6.2. *The update performed by LouvainUpdate(i) is equivalent to transferring node i to the cluster C_{k^*} such that*

$$k^* = \arg \max_{k: j \in C_k, j \sim i} \left[w_i(C_k) - w_i \frac{\text{Vol}(C_k)}{w} \right].$$

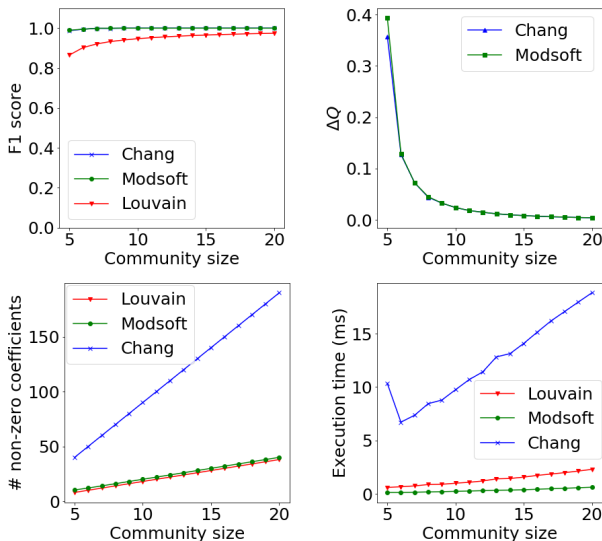


Figure 1: Algorithm performance in function of the cluster size on the overlapping stochastic block model

Proof. We derive this result from the definition of the hard modularity in the supplementary material. \square

Using Proposition 6.1, we see that the update performed by $\text{LouvainUpdate}(i)$ and the update rule (5) are equivalent if (H) is verified and $t > w/\delta$.

7 Experimental results

7.1 Synthetic data

We consider a variant with overlaps of the popular Stochastic Block Model (SBM) [15], defined as follows. We have $V = C_1 \cup \dots \cup C_K$ with $|C_1| = \dots = |C_K| = c$ and $|C_k \cap C_{k+1}| = o$ for all k . For all nodes i and j such that $i \neq j$, the edge (i, j) exists with probability p_{in} if $i, j \in C_k$ for some k , and with probability p_{out} otherwise. In order to numerically evaluate the performance of our algorithm on this model, we first compute the relative difference between the modularity Q_{soft} obtained with our algorithm and the modularity Q_{hard} obtained with the Louvain algorithm, $\Delta Q = (Q_{\text{soft}} - Q_{\text{hard}})/Q_{\text{hard}}$. We also compare the planted clusters of our model C_k with the clusters \hat{C}_k returned by our algorithm. We consider that these clusters are the non-empty sets \hat{C}_k where $\hat{C}_k = \{i \in V : p_{ik} > 0\}$. Note that the clusters \hat{C}_k returned by our algorithm can overlap just as the C_k . In order to compare the clusters C_k and \hat{C}_k , we use the average of the classic F1 score defined as the harmonic mean of the precision and the recall, which takes values between 0 and 1, and is equal to 1 if the estimated clusters \hat{C}_k correspond exactly to the initial clusters C_k . We compute the same two metrics for

Chang’s algorithm [4]. Since the membership matrix returned by this algorithm is not sparse, we define \hat{C}_k as $\hat{C}_k = \{i \in V : p_{ik} > \alpha\}$ for the F1 score evaluation. We present the results for $\alpha = 0.05$, but other values of α leads to similar results.

For different values of the cluster size c , we run our algorithm, Chang’s algorithm and the Louvain algorithm on 100 random instances of this model, with $o = 2$, $K = 2$, $p_{\text{in}} = 0.9$, and $p_{\text{out}} = 0.1$. We display the results in Figure 1. The average F1 scores for our algorithm and Chang’s algorithm are close to 1 (> 0.99 on average), and are higher than those obtained by Louvain (0.94 on average). Besides, the relative difference to the modularity found by Louvain ΔQ for Chang’s algorithm and our algorithm is always positive. However, we note that ΔQ and the differences between the F1 scores become negligible as the cluster size c increases. This can be explained by the fact that, as the cluster size c increases, the proportion of nodes with mixed membership, i.e. nodes that belong to multiple clusters C_k , decreases. If this proportion is low, the misclassification of nodes with mixed membership has a low impact on the modularity score Q and on the F1 score.

We also represent in Figure 1 the sparsity of the solution in term of the number of non-zero coefficients of the matrix \mathbf{p} , and the average execution time in milliseconds for the different algorithms. We observe that our algorithm provides performance comparable to Louvain (but with richer membership information). Besides, we see that MODSOFT runs more than 10 times faster and uses 5 times less memory than Chang’s algorithm.

7.2 Real data

7.2.1 Wikipedia

We first consider a graph built from Wikipedia data, where nodes correspond to Wikipedia articles and edges correspond to hyperlinks between these articles (we consider that the edges are undirected and unweighted). Wikipedia articles can correspond to a wide variety of entities: persons, locations, organizations, concepts etc. In order to restrict ourselves to homogeneous entities, we choose to limit ourselves to Wikipedia pages that correspond to humans. For this purpose, we use Wikidata, a knowledge base that provides structured data about Wikipedia articles. In particular, Wikidata objects have an `instance_of` field that we use to determine if a Wikipedia article represents a human or not. The subgraph induced by the human entities has 731,293 nodes and 3,266,258 edges.

We run the Louvain algorithm, Chang’s algorithm and MODSOFT on this subgraph. As described in Section

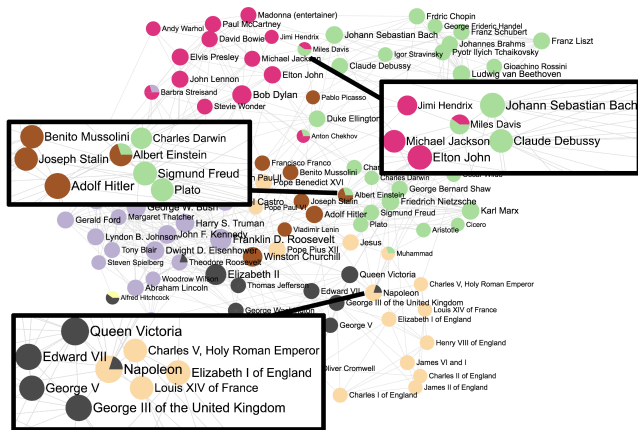


Figure 2: MODSOFT on a small subgraph of Wikipedia

2, Chang’s algorithm requires the maximum number of clusters K to be specified. We run it with $K = n$ (which corresponds to the implicit choice made for our algorithm), K equal to the number of clusters returned by Louvain (i.e. $K = 12,820$) and $K = 100$. Chang’s algorithm is unable to handle the graph (the execution raises a memory error) for all these choices of K , whereas our algorithm runs in 38 seconds on an Intel Xeon Broadwell CPU with 32 GB of RAM. As a baseline for comparison, the Louvain algorithm runs in 20 seconds on this dataset. The difference in memory consumption between our algorithm and Chang’s algorithm can easily be explained by the fact that the membership matrix \mathbf{p} in Chang’s approach is a dense $n \times K$ matrix (see Section 2), i.e. the algorithm needs to store and perform operations on 73, 129, 300 coefficients for $K = 100$. By contrast, the matrix returned by our algorithm is very sparse, since it contains only 760, 546 non-zero coefficients in our experiments, which represents a proportion of approximately 10^{-6} of the coefficients of the $n \times n$ matrix. The average number of positive components of vectors \mathbf{p}_i throughout the execution of the algorithm is 1.21, and the densest vector \mathbf{p}_i has 65 positive components.

Our algorithm leads to a higher modularity than the Louvain algorithm on this dataset. However, this increase represents a modularity increase ΔQ of slightly less than 1%. This can be explained by the fact that the nodes with mixed membership found by our algorithm represent less than 5% of the total number of nodes. Even if the performance increase brought by our algorithm in terms of modularity is marginal, we qualitatively observe that the results of our algorithm, and in particular the nodes with mixed membership identified by our approach, are particularly relevant.

In order to be able to graphically represent the results of our algorithm, we consider the subgraph induced by the top 100 nodes in term of degree. We display

the results of our algorithm on this smaller graph in Figure 2. In particular, we observe that Napoleon, who became Emperor of the French in 1804, appears to belong to the same cluster as European leaders from the old regime, such as Charles V and Louis XIV, and, at the same time, to the same cluster as post-French-revolution leaders such as Queen Victoria and Edward VII. We also see that Miles Davis, a famous American jazz trumpeter, is found to belong to the same cluster as classical music composers, such as Claude Debussy and Johan Sebastian Bach, but also to the same cluster as pop/rock musicians such as Jimi Hendrix and Michael Jackson. Finally, we observe that Albert Einstein who became politically involved during World War II, belongs to the same cluster as thinkers, intellectuals and scientists, such as Sigmond Freud and Plato, and to the same cluster as political leaders during World War II such as Adolf Hitler and Joseph Stalin.

7.2.2 Open Flights

Then, we consider a graph built from the Open Flights database that regroups open-source information on airports and airlines. The graph we consider is built as follows: the nodes correspond to the airports and the edges correspond to the airline routes, i.e. there is an edge between two airports if an airline operates services between these airports.

We apply our algorithm and the Louvain algorithm to this graph. During the execution of the algorithm, the number of positive components of vectors \mathbf{p}_i remains lower than 11 and is on average 1.17. We measure a relative increase in modularity ΔQ of less that 1%, and yet we find that our algorithm brings insightful information about the dataset by identifying bridges between important zones of the world. In Figure 3, we display the results of our algorithm. We use plain colors to code for the cluster membership of *pure* nodes, i.e. nodes that belong to a unique cluster, and bigger

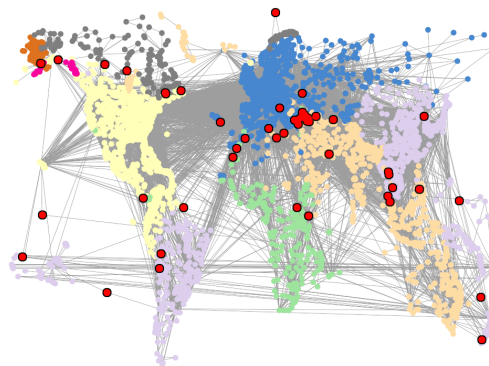


Figure 3: MODSOFT on the OpenFlights graph (red circles represent nodes with mixed-membership)

	Graph size	
	$ V $	$ E $
Amazon	334,863	925,872
DBLP	317,080	1,049,866
YouTube	1,134,890	2,987,624

Table 1: F_1 -scores and execution times in seconds on SNAP datasets

	F_1 -score			Execution time (s)		
	L	O	M	L	O	M
Amazon	0.28	0.47	0.53	15	1038	32
DBLP	0.14	0.35	0.38	17	1717	33
YouTube	0.01	-	0.15	45	-	98

Table 2: F_1 -scores and execution times in seconds on SNAP datasets

red dots to identify nodes with a mixed membership. Note that most of these nodes with mixed membership are located on the borders between large regions of the world. For instance, we remark that numbers of nodes are located on the frontier between Europe, Africa and Middle East, such as the Dubai airport in the United Arab Emirates. We also observe that airports with mixed membership are located all along the border between the United States on the one hand, and Canada and Alaska on the other hand.

7.2.3 SNAP datasets

We present further experiments on datasets from SNAP that come with ground-truth clusters [35] (see Table 1). The Amazon dataset corresponds to a product co-purchasing network. Its nodes represent Amazon products and its edges correspond to frequently co-purchased products. The ground-truth clusters are defined as the product categories. The DBLP dataset corresponds to a scientific collaboration network. The nodes represent the authors and the edges the co-authorship relations. The scientific conferences are used as ground-truth clusters. In the YouTube dataset, nodes represent users and edges connect users who have a friendship relation. User groups are used as ground-truth clusters. We use Louvain and the state-of-the-art algorithm OSLOM [18] as baselines for our experiments. OSLOM is known as one of the most effective overlapping community detection algorithm [33]. Note that Chang’s algorithm does not run on these large graphs. We use a Cython implementation of the Modsoft algorithm that is publicly available². In Table 2, we give the F_1 -scores and the execution times for Louvain (L), OSLOM (O) and Modsoft (M). The experiments were carried out on an AWS instance with

64GB RAM and 16 vCPUs with Intel Xeon. A learning of $t = 0.1$ was used for Modsoft. Note that we did not include the results of OSLOM on the YouTube dataset because its the execution time exceeded 6 hours. We see that Modsoft outperforms OSLOM both in terms of F_1 -score (11% higher on average) and execution time (OSLOM is about 40 times slower). Note that, in our experiments, the objective function of Modsoft is within 1% of its optimal value after 3 to 5 epochs. Moreover, although Modsoft is slower than Louvain (about twice slower on average), it offers F_1 -scores that are at least twice as large as the ones provided by Louvain. This can be explained by the fact that the clusters returned by Louvain are necessary disjoint whereas the ground-truth clusters in all three datasets can overlap. In contrast, Modsoft is able to capture these mixed memberships. Besides, we observe in these experiments that, as expected, the larger t , the sparser the solution. For instance, on the Amazon dataset, the average number of positive components per row is 5.58 for $t = 0.5$, 2.27 for $t = 1$, and 1.87 for $t = 2$.

8 Conclusion

We studied a relaxation of the popular modularity maximization problem, with the aim of performing soft clustering of graph nodes instead of hard partitioning. In order to efficiently solve this relaxation, we introduced a novel algorithm that is based on an alternating projected gradient descent. By diving into the specific form of the gradient descent updates, we were able to guarantee the locality of each algorithm step, and to make good use of the sparsity of the solutions. As a result, unlike existing methods, our approach is both local and memory efficient. Furthermore, we proved that our algorithm includes the main routine of the popular Louvain algorithm for $t > w/\delta$. In this case, Modsoft outputs the sparsest possible solution, with only one component per row. We illustrated on real-life examples, that our algorithm has an execution time and a memory consumption comparable to the Louvain algorithm, but goes further than Louvain by identifying nodes with mixed memberships that represent important bridges between clusters, even if this does not always translate into a large modularity increase.

A generalization of standard modularity introducing a *resolution* parameter has been proposed by different authors [17, 30] to perform multi-level clustering. In future works, we could add this resolution parameter to the soft modularity function in order to perform soft graph clustering at different scale. In the future, we would also like to compare our approach to a non-negative matrix factorization method applied to the matrix \mathbf{W} , which can be performed with a comparable alternating projected gradient descent [23].

²<https://github.com/ahollocou/modsoft>

References

- [1] J. C. Bezdek. Objective function clustering. In *Pattern recognition with fuzzy objective function algorithms*, pages 43–93. Springer, 1981.
- [2] 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*, 2008(10):P10008, 2008.
- [3] U. Brandes, D. Delling, M. Gaertler, R. Görke, M. Hofer, Z. Nikoloski, and D. Wagner. On finding graph clusterings with maximum modularity. In *International Workshop on Graph-Theoretic Concepts in Computer Science*, pages 121–132. Springer, 2007.
- [4] C.-T. Chang and C.-S. Chang. A unified framework for sampling, clustering and embedding data points in semi-metric spaces. *arXiv preprint arXiv:1708.00316*, 2017.
- [5] F. R. Chung. *Spectral graph theory*. Number 92. American Mathematical Soc., 1997.
- [6] J. Duch and A. Arenas. Community detection in complex networks using extremal optimization. *Physical review E*, 72(2):027104, 2005.
- [7] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the l_1 -ball for learning in high dimensions. In *Proceedings of the 25th international conference on Machine learning*, pages 272–279. ACM, 2008.
- [8] J. C. Dunn. A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. 1973.
- [9] S. Fortunato and C. Castellano. Community structure in graphs. In *Computational Complexity*, pages 490–512. Springer, 2012.
- [10] C. Fowlkes, S. Belongie, F. Chung, and J. Malik. Spectral grouping using the nystrom method. *IEEE transactions on pattern analysis and machine intelligence*, 26(2):214–225, 2004.
- [11] A. A. Goldstein. Convex programming in hilbert space. *Bulletin of the American Mathematical Society*, 70(5):709–710, 1964.
- [12] E. Griechisch and A. Pluhár. Community detection by using the extended modularity. *Acta Cybern.*, 20(1):69–85, 2011.
- [13] R. Guimera, M. Sales-Pardo, and L. A. N. Amaral. Modularity from fluctuations in random graphs and complex networks. *Physical Review E*, 70(2):025101, 2004.
- [14] T. C. Havens, J. C. Bezdek, C. Leckie, K. Ramamohanarao, and M. Palaniswami. A soft modularity function for detecting fuzzy communities in social networks. *IEEE Transactions on Fuzzy Systems*, 21(6):1170–1175, 2013.
- [15] P. W. Holland, K. B. Laskey, and S. Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- [16] Y. Kawase, T. Matsui, and A. Miyauchi. Additive approximation algorithms for modularity maximization. *arXiv preprint arXiv:1601.03316*, 2016.
- [17] R. Lambiotte, J.-C. Delvenne, and M. Barahona. Laplacian dynamics and multiscale modular structure in networks. *arXiv preprint arXiv:0812.1770*, 2008.
- [18] A. Lancichinetti, F. Radicchi, J. J. Ramasco, and S. Fortunato. Finding statistically significant communities in networks. *PLoS one*, 6(4):e18961, 2011.
- [19] P. Latouche, E. Birmelé, and C. Ambroise. Overlapping stochastic block models. *arXiv preprint arXiv:0910.2098*, 2009.
- [20] A. Lázár, D. Ábel, and T. Vicsek. Modularity measure of networks with overlapping communities. *EPL (Europhysics Letters)*, 90(1):18001, 2010.
- [21] S. Lehmann and L. K. Hansen. Deterministic modularity optimization. *The European Physical Journal B*, 60(1):83–88, 2007.
- [22] E. S. Levitin and B. T. Polyak. Constrained minimization methods. *USSR Computational mathematics and mathematical physics*, 6(5):1–50, 1966.
- [23] C.-J. Lin. Projected gradient methods for nonnegative matrix factorization. *Neural computation*, 19(10):2756–2779, 2007.
- [24] M. E. Newman. Fast algorithm for detecting community structure in networks. *Physical review E*, 69(6):066133, 2004.
- [25] M. E. Newman. Finding community structure in networks using the eigenvectors of matrices. *Physical review E*, 74(3):036104, 2006.
- [26] M. E. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical review E*, 69(2):026113, 2004.

- [27] V. Nicosia, G. Mangioni, V. Carchiolo, and M. Malgeri. Extending the definition of modularity to directed graphs with overlapping communities. *Journal of Statistical Mechanics: Theory and Experiment*, 2009(03):P03024, 2009.
- [28] B. Perozzi, R. Al-Rfou, and S. Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710. ACM, 2014.
- [29] A. Prat-Pérez, D. Dominguez-Sal, and J.-L. Larriba-Pey. High quality, scalable and parallel community detection for large real graphs. In *Proceedings of the 23rd international conference on World wide web*, pages 225–236. ACM, 2014.
- [30] J. Reichardt and S. Bornholdt. Statistical mechanics of community detection. *Physical Review E*, 74(1):016110, 2006.
- [31] U. Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- [32] J. J. Whang, I. S. Dhillon, and D. F. Gleich. Non-exhaustive, overlapping k-means. In *Proceedings of the 2015 SIAM International Conference on Data Mining*, pages 936–944. SIAM, 2015.
- [33] J. Xie, S. Kelley, and B. K. Szymanski. Overlapping community detection in networks: The state-of-the-art and comparative study. *Acm computing surveys (csur)*, 45(4):43, 2013.
- [34] J. Yang and J. Leskovec. Overlapping community detection at scale: a nonnegative matrix factorization approach. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 587–596. ACM, 2013.
- [35] J. Yang and J. Leskovec. Defining and evaluating network communities based on ground-truth. *Knowledge and Information Systems*, 42(1):181–213, 2015.
- [36] Y. Zhang, E. Levina, and J. Zhu. Detecting overlapping communities in networks using spectral methods. *arXiv preprint arXiv:1412.3432*, 2014.