

Scalable and Accurate Algorithm for Graph Clustering*

Hristo N. Djidjev¹ and Melih Onus²

¹ Los Alamos National Laboratory
Los Alamos, NM 87545
djidjev@lanl.gov

² Department of Computer Engineering
Cankaya University
Ankara, Turkiye
melih@asu.edu

Abstract. One of the most useful measures of quality for graph clustering is the modularity of the partition, which measures the difference between the number of the edges with endpoints in the same cluster and the expected number of such edges in a random graph. In this paper we show that the problem of finding a partition maximizing the modularity of a given graph G can be reduced to a minimum weighted cut problem on a complete graph with the same vertices as G . We then show that the resulting minimum cut problem can be efficiently solved by adapting existing graph partitioning tools. Our algorithm is accurate and finds a graph clusterings much faster than alternative algorithms that produce a comparable clustering quality.

Keywords: Graph clustering, graph partitioning.

1 Introduction

One way to analyze and understand the information contained in the huge amount of data available on the WWW and the relationships between the individual items is to organize them into "communities," maximal groups of related items. Determining the communities is of great theoretical and practical importance since they correspond to entities such as collaboration networks, online social networks, scientific publications or news stories on a given topic, related commercial items, etc. Communities also arise in other types of networks such as computer and communication networks (the Internet, ad-hoc networks) and biological networks (protein interaction networks, genetic networks).

The problem of identifying communities in a network is usually modeled as a *graph clustering* (GC) problem, where vertices correspond to individual items and edges describe relationships. Then the communities correspond to subgraphs with a lot of edges between vertices belonging to the same subgraph (called *in-cluster* edges) and fewer edges between vertices from different subgraphs (called *between-cluster* edges). The GC problem has been intensively studied in the recent years in relation to its applications in the analysis of networks. Girvan and Newman propose in [22], [42] algorithms based on the *betweenness* of the edges of a graph, a characteristic that measures the number of the shortest paths in a graph that use any given edge. In [32] Newman describes an algorithm based on a characteristic of clustering quality called *modularity*, a measure that takes into account the number of in-cluster edges and the expected number of such edges. (We formally define and discuss modularity in more detail in

* An early version of this paper was presented at the workshop *Algorithms and Models for the Web-Graph (WAW 2006)*, Lecture Notes in Computer Science, vol. 4936, pp. 117–128.

the next section.) A faster version of the algorithm from [32] was described by Clauset *et al.* in [15]. Several algorithms have been proposed based on other techniques such as computing eigenvectors of the graph Laplacian, e.g., [50], [39], [40], simulated annealing [45], [24], belief propagation [25] and greedy methods [8]. In all previous cases the algorithms reported in the literature are either not fast enough, or are inaccurate. The problem of finding a partition that maximizes the modularity was shown to be NP-hard [13].

In this paper we will describe a new approach for GC that uses our newly discovered relationship between the GC and the minimum weighted cut problems. The *minimum weighted cut* (MWC) problem is, given a graph $G = (V, E)$ with real weights on its edges, find a partition of V such that the set of all edges of G that join vertices from different sets of the partition, called a *cut* of the partition, is of minimum weight. GC looks related to the MWC problems since, in a good quality clustering, the weight of the edges between different sets of the partition (the cut) should be small compared to the weight of the edges inside the sets. But the MWC problem can not be directly applied to solve the GC problem since it does not take into account the sizes of the subgraphs induced by the cut (e.g., it is likely that the minimum cut will consist of the edges incident to a single vertex). There are some minimum cut based clustering algorithms, e.g., [20], that use maximum flow computations combined with heuristics, but they are typically slower than modularity based algorithms, e.g. [15], and, moreover, they cannot determine the optimal number of clusters and, instead, construct a hierarchical decomposition of the set of all vertices of the graph.

In this paper we prove that the problem of finding a partition of a graph G that maximizes the modularity can be reduced to the problem of finding a MWC of a weighted complete graph on the same set of vertices as G . We then show that the resulting minimum cut problem can be solved by modifying existing fast algorithms for graph partitioning. We demonstrate by experiments that our algorithm has generally a better quality and is much faster than the best existing GC algorithms.

2 Our clustering algorithm

2.1 Preliminaries

A *graph* G is an ordered pair $(V(G), E(G))$ of sets, where $V(G)$ is the set of the *vertices* and $E(G)$ is the set of the *edges* of G and each edge is an unordered pair (v, w) of vertices. If $E' \subseteq E(G)$, then by $G - E'$ we denote the graph $(V(G), E(G) \setminus E')$. A graph is bipartite, if $E(G) \subseteq \{(v_1, v_2) \mid v_1 \in V_1, v_2 \in V_2\}$, where $V_1 \cup V_2 = V$ and $V_1 \cap V_2 = \emptyset$. A *path* p in G is a sequence (v_1, \dots, v_k) of vertices such that $(v_i, v_{i+1}) \in E(G)$ for $1 \leq i < k$. If $v_1 = v_k$, then p is a *cycle*. G is *connected* if there is a path between any pair of vertices of G . The *components* of G are its maximal connected subgraphs. A *partition* \mathcal{P} of G is a division of $V(G)$ into subsets V_1, \dots, V_s such that $V_i \cap V_j = \emptyset$ for $i \neq j$ and $\bigcup_{i=1}^s V_i = V(G)$. If $s = 2$, then \mathcal{P} is a *bisection*. Note that, in contrast to other works, our definitions of partition and bisection does not require the parts to be balanced in size. A set $C \subseteq E(G)$ is a *cut* of G if there exists a partition \mathcal{P} of G such that C is the set of the edges of G joining vertices from different sets of \mathcal{P} . We will use the notation $C = \text{cut}(\mathcal{P})$ and $\mathcal{P} = \text{part}(C)$ and define $\text{cut}(V_i, V_j) = \{(t, u) \in E(G) \mid t \in V_i, u \in V_j\}$. Two partitions \mathcal{P}_1 and \mathcal{P}_2 are *equivalent*, if $\text{cut}(\mathcal{P}_1) = \text{cut}(\mathcal{P}_2)$. If there are weights $\text{wt}(\cdot)$ associated with the edges of G , then by $\text{cutWt}(\mathcal{P}) = \text{wt}(C)$ we denote the sums of the weights of all edges in C . If M is a finite set, by $|M|$ we denote the number of the elements of M .

The following property follows directly from the definitions.

Lemma 1. *An edge set C is a cut of G if and only if, for each $(v, w) \in C$, each path between v and w contains an edge from C .*

Lemma 2. *A partition \mathcal{P} of G is equivalent to a bisection of G if and only if any cycle in G contains an even number of edges from $\text{cut}(\mathcal{P})$.*

Proof. Let V_1, \dots, V_k be the connected components of $G - \text{cut}(\mathcal{P})$. Construct a graph G' with a vertex set $V(G') = \{v_1, \dots, v_k\}$ and edge between v_i and v_j , if there is an edge in $\text{cut}(\mathcal{P})$ joining a vertex from V_i with a vertex from V_j . Clearly, \mathcal{P} is equivalent to a bisection if and only if G' is bipartite. Moreover, for each cycle in G with k edges from $\text{cut}(\mathcal{P})$, there is a cycle in G' with k edges. The claim follows from the facts that a graph is bipartite if and only if each of its cycles has an even number of edges. \square

2.2 Modularity optimization as a minimum cut problem

As there is no formal definition of clustering and what the clusters of a given graph are, in general it is not possible to determine if a certain partition represents the "correct" clustering or which of two alternative partitions of a graph corresponds to a better clustering. For that reason, researchers have used their intuition to define measures for cluster quality that can be used for comparing different partitions of the same graph. One such measure, introduced in [42, 37], which has received considerable attention recently, is the *modularity* of a graph. Given an n -vertex m -edge graph $G = (V(G), E(G))$ and a partition \mathcal{P} of $V(G)$ into k subsets (clusters) V_1, \dots, V_k , the modularity $Q(\mathcal{P}, G, \mathcal{G})$ of \mathcal{P} with respect to \mathcal{G} (or $Q(\mathcal{P})$ for short if G and \mathcal{G} are clear from the context) is a number between -1 and 1 defined as

$$Q(\mathcal{P}) = Q(\mathcal{P}, G, \mathcal{G}) = \frac{1}{m} \sum_{i=1}^k (|E(V_i)| - \text{Ex}(V_i, \mathcal{G})),$$

where $E(V_i)$ is the set of all edges of G with endpoints in V_i and $\text{Ex}(V_i, \mathcal{G})$ is the expected number of such edges in a random graph with a vertex set V_i from a given random graph distribution \mathcal{G} on $V(G)$. $Q(\mathcal{P})$ measures the difference between the number of in-cluster edges and the expected value of that number for \mathcal{P} in a random (e.g., without cluster structure) graph on the same vertex set. Larger values of $Q(\mathcal{P})$ correspond to better clusterings.

Having the definition of $Q(\mathcal{P})$, we can formulate the clustering problem as finding a partition $\mathcal{P} = \{V_1 \cup \dots \cup V_k\}$ of $V(G)$ such that

$$\sum_{i=1}^k (|E(V_i)| - \text{Ex}(V_i, \mathcal{G})) \rightarrow \max. \quad (1)$$

Clearly

$$\begin{aligned} \max_{\mathcal{P}} \left\{ \sum_{i=1}^k (|E(V_i)| - \text{Ex}(V_i, \mathcal{G})) \right\} &= - \min_{\mathcal{P}} \left\{ - \sum_{i=1}^k (|E(V_i)| - \text{Ex}(V_i, \mathcal{G})) \right\} \\ &= - \min_{\mathcal{P}} \left\{ (|E(G)| - \sum_{i=1}^k |E(V_i)|) - (|E(G)| - \sum_{i=1}^k \text{Ex}(V_i, \mathcal{G})) \right\}. \end{aligned}$$

Denote

$$\text{ExCut}(\mathcal{P}, G, \mathcal{G}) = |E(G)| - \sum_{i=1}^k \text{Ex}(V_i, \mathcal{G}).$$

Intuitively, $\text{ExCut}(\mathcal{P}, G, \mathcal{G})$ is the expected value of $|\text{cut}(\mathcal{P})|$ with respect to the random graph class \mathcal{G} , assuming the expected number of edges for \mathcal{G} is $|E(G)|$. Then

$$\max_{\mathcal{P}} \left\{ \sum_{i=1}^k (|E(V_i)| - \text{Ex}(V_i, \mathcal{G})) \right\} =$$

$$- \min_{\mathcal{P}} \{ |\text{cut}(\mathcal{P})| - \text{ExCut}(\mathcal{P}, G, \mathcal{G}) \}.$$

Hence, instead of problem (1), one can address the problem of computing

$$\operatorname{argmin}_{\mathcal{P}} \{ |\text{cut}(\mathcal{P})| - \text{ExCut}(\mathcal{P}, G, \mathcal{G}) \}. \quad (2)$$

The last expression shows that we can solve (1) as a problem of finding a MWC in a complete graph G' with a vertex set $V(G)$ and weight $\text{weight}(i, j)$ on any edge $(i, j) \in E(G')$ defined by

$$\text{weight}(i, j) = \begin{cases} 1 - p_{ij}, & \text{if } (i, j) \in E(G) \\ -p_{ij}, & \text{if } (i, j) \notin E(G), \end{cases} \quad (3)$$

where p_{ij} is the probability that there is an edge between vertices i and j in a random graph from the class \mathcal{G} . Then, problem (1) is equivalent to the problem of computing

$$\operatorname{argmin}_{\mathcal{P}'} \{ \text{cutWt}(\mathcal{P}') \}, \quad (4)$$

where $\text{cutWt}(\mathcal{P}')$ denotes the weight of the cut of \mathcal{P}' .

We summarize these observations in the following theorem.

Theorem 1. *The problem of finding a partition of a graph $G = (V, E)$ that minimizes the modularity can be reduced in $O(|V| + |E|)$ time to the problem of finding a minimum weighted cut in a complete graph $G' = (V, E')$ with edge weights given by (3).*

For the reduction time bound in Theorem 1 we assume that the edges of $E' \setminus E$ are defined implicitly. There are several choices for \mathcal{G} that have been favored by various researchers. The random graph model $G(n, p)$ of Erdős-Renyi [18] defines n vertices and puts an edge between each pair with probability p . Clearly, the expected number of edges of $G(n, p)$ is $\binom{n}{2}p$. Hence, for a graph with expected number of edges m

$$p_{ij} = p = \frac{m}{\binom{n}{2}}. \quad (5)$$

One disadvantage of the $G(n, p)$ model is that it fails to capture important features of the real-world networks, in particular, the degree distribution. As has been recently observed [6], many important types of networks like technological networks (the Internet, the WWW), social networks (collaboration networks, online social networks), biological networks (protein interactions) have degree distributions that follow a *power law*, e.g., the fraction of the vertices that have degree $k > 0$ is roughly proportional to $\alpha k^{-\lambda}$ for some constants α and $\lambda > 0$. Such networks are called *scale-free*. In comparison, the degrees of a random graph from the $G(n, p)$ model follow a Poisson distribution, i.e., the probability that a given vertex has degree k is $\binom{n}{k} p^k (1-p)^{n-k}$ and the expected degree of each vertex is pn . Hence, the Erdős-Renyi model may not be suitable as a choice for \mathcal{G} when used for determining the community structure of graphs of the above type.

One model that takes into account the degrees of the vertices is studied by Chung and Lu in [14]. In that model, the probability that there is an edge between a vertex i and a vertex j is

$$p_{ij} = \frac{d_i d_j}{\sum_{k=1}^n d_k}, \quad (6)$$

where d_1, \dots, d_n are positive reals corresponding to the degrees of the vertices such that $\max_{1 \leq i \leq n} d_i^2 < \sum_{i=1}^n d_i$. (The last condition guarantees that such a graph exists if all numbers d_i are integers and will be always satisfied if numbers d_i are chosen to be the degrees of G .)

We will refer to that model as the Chung-Lu (CL) model. Clearly, in the CL model, the expected degree of vertex i is d_i , compared with pn (i.e., independent on i) in the $G(n, p)$ model.

Note that for both of the above choices of \mathcal{G} the expected number of edges for a graph in \mathcal{G} is $|E(G)|$.

In the next section we will describe an efficient method for finding a MWC of a complete graph G' with weights on the edges satisfying (3) and p_{ij} defined by (5) or (6).

2.3 Finding a MWC using multilevel graph partitioning

Above we established an important relationship between the modularity optimization and the MWC problems, i.e., that the problem of finding a partition of a given graph that maximizes the modularity can be reduced to the problem of finding a minimum weight cut. Most existing work on the MWC problem considers the case where all weights are non-negative. The MWC problem in the case of non-negative weights is known to be polynomially solvable, e.g., by using algorithms for computing maximum flows [2]. In contrast, the MWC problem in case of real-value weights is NP-hard and algorithmic aspects of the problem are much less studied. Here we show that available heuristics for another related problem, graph partitioning, can be adapted to solve this version of the MWC problem.

Overview of the multilevel graph partitioning. Formally, the *graph partitioning* (GP) problem is, given a graph $G = (V, E)$, to find a partition (V_1, V_2) of V such that $||V_1| - |V_2|| \leq 1$ (i.e., the partition is *balanced*) and $\text{cut}(V_1, V_2)$ is minimum. (Some versions of the problem consider partitions into an arbitrary number of parts.) Hence, in comparison with the minimum cut problem, there is the additional requirement for a balanced partition. Because of its important applications, e.g., in high performance computing and VLSI design, GP is a well-researched problem for which very efficient methods have been developed. One such approach is the multilevel GP. This method is both fast and accurate for a wide class of graphs that appear in practical applications. Inspired by the multigrid method from computational mathematics, it has been used in the works of Barnard and Simon [7], Hendrickson and Leland [26], Karypis and Kumar [27, 28], and others. The method for bisecting a graph consists of the following three phases (Figure 1):

Coarsening phase. The original graph G is coarsened by partitioning it into connected subgraphs, replacing each of the subgraphs by a single vertex, and replacing the set of the edges between any pair of shrunk subgraphs by a single edge. Moreover, a weight of each new vertex (respectively edge) is assigned equal to the sum of the weights of the vertices (respectively edges) that it represents. Weights on the original vertices of G are defined 1, in the case of the $G(n, p)$ model, or their degrees, in the case of the the CL model, as justified in Corollary 1 below. (The coarsening procedure, including alternative methods for determining the set of the shrunk subgraphs and analysis of their effect on the quality of the final partition, is described in much detail in [28].) The resulting graph is coarsened repeatedly by the same procedure until one gets a graph of a sufficiently small size. Let $G_0 = G, G_1, \dots, G_l$ be the resulting graph sequence.

Partitioning phase. The graph G_l is partitioned into two parts using any available partitioning method (e.g., spectral partitioning or the Kernighan-Lin (KL) algorithm [29]).

Uncoarsening and refinement phase. The partition of G_l is projected on G_{l-1} . Since the weight of each vertex of G_l is a sum of the weights of the corresponding vertices of G_{l-1} , then

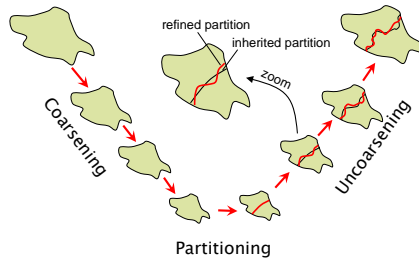


Fig. 1. The stages of multilevel partitioning.

the partition of G_{l-1} will be balanced if the partition of G_l is and the cut of both partitions will have the same weight. However, since G_{l-1} has more vertices than G_l , it has more degrees of freedom and, therefore, it is possible to refine the partition of G_{l-1} in order to reduce its cut size. To this end, the projection of the partition of G_l is followed by a refinement phase, which is usually based on the KL algorithm. In the same way, the resulted partition of G_{l-1} is converted into a partition of G_{l-2} and refined, and so on until a partition of G_0 is found.

Kernighan-Lin refinement. Since the refinement step is the most involved part of the algorithm, which ultimately determines its accuracy and efficiency, we will describe it in more detail. It has been shown [28] that the KL algorithm can be a good choice for performing the refinement.

The KL algorithm involves several iterations, each consisting of moving a vertex from one set of the partition to the other. Let $\mathcal{P} = \{P_1, P_2\}$ be the current partition. For each vertex u of the graph a *gain* for u is defined as

$$gain(u) = \sum_{v \in N(u) \setminus P(u)} weight(u, v) - \sum_{v \in N(u) \cap P(u)} weight(u, v), \quad (7)$$

where $N(u)$ is the set of all neighbors of u and $P(u)$ is that set of \mathcal{P} that contains u . $gain(u)$ measures how the weight of the cut will be affected if u is moved from $P(u)$ to the other set of \mathcal{P} . The KL algorithm then selects a vertex w from the larger set of the partition with a maximum gain, moves it to the other set, and updates the gains of the vertices adjacent to w . Moreover, w is marked so that it will not be moved again during that refinement step. The process is continued until either all vertices have been moved, or the S most recent moves have not led to a better partition. (S is a user chosen parameter that is set to 50 in the current implementation.) At the end of the refinement step, the last $s \leq S$ moves that have not improved the partition are reversed.

Implementation The implementation of our algorithm for clustering is based on the version of multilevel partitioning implemented by Karypis and Kumar [27, 28], which has been made freely available as a software package under the name METIS. Note that graph partitioning, minimum cut, and clustering are related, but with important differences, problems, as illustrated in Table 1. We have already shown how the clustering problem can be reduced to a minimum cut problem and here we will show how the resulting minimum cut problem

can be solved by a graph partitioning algorithm based on METIS. Because of the differences between graph partitioning and MWC, we have to make some evident changes. For instance, since graph partitioning requires balanced partitions, we have to drop the requirement for balance of the partition. We have also to determine the cardinality of the partition that minimizes the cut size. But the main implementation difficulty is related to the size of G' . Although the original graph, G , is often sparse, i.e., it has n vertices and $O(n)$ edges, the transformed one, G' , is always dense, as it has $\binom{n}{2} = \Omega(n^2)$ edges. The main challenge will be to construct an algorithm whose complexity is close to linear on the size of the original graph, rather than on the size of the transformed one. Next we show that it is possible to simulate an execution of a KL refinement step on G' by explicitly maintaining information only about the edges from the original graph G and implicitly taking into account the remaining edges by modifying the formulae for computing weights and gains.

In order to give intuition about why this works, assume that the edges of G' belong to two types that we call *visible* and *invisible*. The visible edges correspond to the edges of the original graph G and are therefore few (assuming G is sparse). These edges carry weight 1 and are maintained explicitly. The invisible edges are between any two vertices of G' . (Note that for each visible edge there is also an invisible one parallel to it, i.e., joining the same endpoints.) The weight of invisible edge (i, j) is $-p_{ij}$. Although the number of invisible edges is $\Omega(n^2)$, because of their uniform distribution, the contribution of these edges to the cut is easy to compute by maintaining additional information of size $O(1)$ only. The next two lemmas formalize this notion.

Problem	Modularity optimization	Minimum Cut	Graph Partitioning
Objective	Maximize modularity	Minimize cut size	Minimize cut size
Balance of partition	Sizes may differ	Sizes may differ	Equal sizes
Cardinality of partition	To be computed	To be computed	An input parameter

Table 1. Comparison between modularity optimization, minimum cut, and graph partitioning problems.

Lemma 3. *Let $\mathcal{P} = \{V_1, V_2\}$ be a partition of G and let G' be the transformed weighted graph with respect to the $G(n, p)$ random graph model. Let \mathcal{P}' be the cut in G' corresponding to \mathcal{P} . Then*

$$\text{cutWt}(\mathcal{P}') = \text{cutWt}(\mathcal{P}) - |V_1||V_2|p,$$

where $p = m/\binom{n}{2}$.

Proof. Follows from formulae (3) and (5). There is an edge in G' joining any vertex from V_1 with any vertex in V_2 . For an edge from G the corresponding weight is $1 - p$, and an edge in G' not in G the corresponding weight is $-p$. \square

The lemma shows that if one maintains the values of $|V_1|$ and $|V_2|$ during a KL refinement, one can work with the original graph G rather than with the modified G' , updating at each step the value of the cut in $O(1)$ time using Lemma 3.

A similar formula holds for the case of the CL model.

Lemma 4. *Let $\mathcal{P} = \{V_1, V_2\}$ be a partition of G and let G' be the corresponding weighted graph with respect to the CL random graph model. Assign a weight $\text{wt}(v)$ to each vertex v equal to its degree. Let \mathcal{P}' be the cut in G' corresponding to \mathcal{P} . Then*

$$\text{cutWt}(\mathcal{P}') = \text{cutWt}(\mathcal{P}) - \text{wt}(V_1)\text{wt}(V_2)p, \tag{8}$$

where $\text{wt}(V_i) = \sum_{v \in V_i} \text{wt}(v)$ and $p = \left(\sum_{v \in V(G)} \text{wt}(v) \right)^{-1}$.

Proof. Follows from formulae (3) and (6) and the equality

$$\sum_{v \in V_1} \sum_{w \in V_2} \frac{\text{wt}(v)\text{wt}(w)}{p} = \left(\sum_{v \in V_1} \text{wt}(v) \right) \left(\sum_{w \in V_2} \text{wt}(w) \right) / p. \quad \square$$

According to the lemma, the cut weight of \mathcal{P}' can be computed in $O(1)$ time given the cut weight of \mathcal{P} , if one maintains the values of the weights of V_1 and V_2 during the KL refinement.

In the case of both the $G(n, p)$ and the CL random graph models, for moving a vertex v from one partition to another during a KL refinement we need only to update the gains of the neighbors of v in G . Having those gains, one can maintain $\text{cutWt}(\mathcal{P})$ in total time proportional to the size of G , excluding the time for priority queue operations needed to extract vertices with maximum gains, which is $O(n \log n)$ in total. By Lemma 3 or Lemma 4, one can at any time compute $\text{cutWt}(\mathcal{P}')$ from $\text{cutWt}(\mathcal{P})$ and the weights of the partitions in $O(1)$ additional time.

From Lemma 3 and Lemma 4 it follows that in the case of both models the same KL refinement algorithm can be used, if the vertex weights are appropriately defined.

Corollary 1. *Let $\mathcal{P} = \{V_1, V_2\}$ be a partition of G and let G' be the corresponding weighted graph with respect to either the $G(n, p)$ or the CL random graph model. Define the weight of any vertex v to be 1, in the case of the $G(n, p)$ model, or the degree of v , in the case of the CL model. Then $\text{cutWt}(\mathcal{P}')$ can be computed by formula (8).*

Time analysis. By using the analysis of Fiduccia and Mattheyses of the KL algorithm from [19], it follows that clustering any network of n vertices and m edges into two communities by our algorithm takes $O(n \log n + m)$ time, where n and m are the numbers of the nodes and links of the network, respectively. Finding a clustering in optimal number of k parts, our algorithm first divides to 2 parts, then to $2^2 = 4$ parts, then to $2^3 = 8$ parts, and so on. Finding a clustering in optimal number of k parts takes $O((n \log n + m)d)$ time, where d is the depth of the dendrogram describing the clustering hierarchy. Since the dendrogram is represented by a binary tree, $\log_2 k \leq d \leq \log_2 k + 1$.

3 Experiments and performance evaluation

We present two type of experimental results. The first type, described in Section 3.1, compares our algorithm with other well known algorithms, testing its speed and accuracy. The second type, described in Section 3.2, presents the results of our algorithm when applied to graphs from the DIMACS Challenge testbed.

3.1 Comparison against other clustering algorithms

We performed a number of experiments on randomly generated graphs, in order to measure the accuracy of our algorithm and its efficiency as well as to compare it with previous algorithms. First we present the results of an experiment measuring the algorithm accuracy, the so called Newman-Girvan test. We include this test as an example of non-modularity based accuracy test and because of its popularity. Its disadvantages are that it uses graphs of very special structure and of relatively small sizes. That is why we concentrate most of our effort and describe in most detail the results of another type of experiments, included in the third subsection, that use graphs of different size and structures, and on which we are able to test both the speed and the accuracy of our algorithm versus several others. In all experiments the CL version of our algorithm was used.

Newman-Girvan accuracy test Following the experimental setting of [42], we generated random graphs with 128 vertices and 4 communities of size 32 each. The expected degree of any vertex is 16, but the expected *outdegree* (the expected number of neighbors of a vertex that belong to a different community) is set to i in the i -th experiment ($i \leq 16$). Hence, higher values of i correspond to graphs with weaker cluster structures. The experiment is intended to measure the sensitivity of the algorithm to the strength of the communities.

In order to decide whether to include an edge (v, w) in the graph in the i -th experiment, a random number r in the interval $[0, 1]$ is generated and (v, w) is accepted if $r \geq i/31$ and v and w belong to the same community or $r \geq (16 - i)/96$ and v and w belong to different communities, and is rejected otherwise.

Outdegree	Degree	Newman-Girvan	Ours
1	16	1.00	1.00
2	16	1.00	1.00
3	16	0.98	1.00
4	16	0.97	1.00
5	16	0.95	1.00
6	16	0.85	0.99
7	16	0.60	0.95
8	16	0.30	0.79

Table 2. Comparing the quality of the clustering of our algorithm and Newman-Girvan’s algorithm.

Table 2 compares the quality of the clusterings produced by Newman-Girvan’s algorithm and ours. A clustering produced by any of the algorithms is considered ”correct” if it matches the original partition of communities from the graph generation phase. (Note that, due to the probabilistic nature of the graphs, the clustering that maximizes the modularity might be different from the original partition, especially if the modularity is low.)

Our algorithm classifies correctly more than 99% of the edges for outdegrees 0, 1, 2, 3, 4, 5, 6 and in all cases it is better than Newman-Girvan’s (more than twice better for the case outdegree=8).

Testing both speed and accuracy Table 3 compares the performance of our algorithm with four other algorithms that are considered among the best with respect to their speeds and/or accuracies. Clauset, Newman, and Moore’s algorithm [15] is an *agglomerative* algorithm that is an improvement of a previous algorithm [32] in terms of the speed and is claimed to have the same quality of the partition. Agglomerative algorithms start with a community partition, where each single vertex represents a community. At each iteration a pair of communities are merged into a single one such that a measure of cluster quality, in this case the modularity, is improved. The second algorithm is Newman’s algorithm described in [40], which is a spectral algorithm based on eigenvector computations. The other two algorithms, of Guimera and Amaral [24] and Reichardt and Bornholdt [45], are based on simulated annealing optimization.

Most of the algorithms tested, notably Guimera-Amaral and Reichardt-Bornholdt algorithms, have parameters that can be played with in order to improve the accuracy of the algorithms on particular graphs. It is possible that by varying the parameters from experiment to experiment and from graph to graph, the quality of some partitions would have

improved. Our algorithm also has parameters that allow trading off speed for accuracy. However, such type of optimization and fine-tuning of the algorithms is beyond the scope of this paper. In all experiments, we have used the recommended or default values of all parameters.

The test graphs in our experiments are random graphs with varying numbers of clusters, sizes, densities, and modularities. The graphs are generated by initially assigning a set of isolated vertices into a number of clusters with preset sizes. Then, for each pair of vertices v and w , an edge (v, w) is generated with probability p_{in} , if v and w belong to the same cluster, and with probability p_{out} , otherwise, where p_{in} and p_{out} are input parameters. Experiment 1–10 have been run 100 times on different random graphs and experiments 11–13 have been run 10 times. All experiments have been run on an Intel Xeon CPU 1.60GHz processor desktop computer with 4G of memory.

For each experiment, the table shows the number of the vertices and the average number of edges of the test graph, the number of the clusters in the original partition during generation, and the average modularity of that partition. Then, for each of the algorithms, the average running time and modularity of the partition are listed.

Experiments 1–4 study how the performance of the algorithms depends on the number of clusters, which vary from 2 to 9. The results indicate that the qualities of the clusterings are comparable, while Newman’s (N) and Guimera-Amaral’s (GA) algorithms time performance is more sensitive to the number of the clusters.

In experiments 5–7, the test graphs have the same numbers of vertices, numbers of cluster, and modularities, but different densities. All algorithms were quite accurate and showed little variance in their performance when sparsity changes.

In experiments 8–10, we compare the algorithms when the modularity (the quality of the original clustering) is low. In these experiments, the Clauset, Newman, and Moore’s (CNM) algorithm considerably underperformed the other four with respect to the quality of the partition.

Finally, in experiments 11-13, we compared the scalability of the algorithms. Because of the low scalability of some algorithms and the long time it takes to run a single experiment, those experiments were run only 10 times. As such, small differences in the modularity should be taken with caution, and attention should be paid to the running times, which vary significantly from algorithm to algorithm. The experiments show that the GA algorithm is the slowest, followed by the other simulated annealing based Reichardt and Bornholdt’s (RB) algorithm, which is about 4 times faster. Neither of these two algorithms can be used in reasonable time for graphs containing more than a few hundred thousand edges. Algorithm N is much faster than those two and can be used for graphs of size several million edges. The only algorithm that can scale to graphs of sizes up to tens of millions of edges is the CNM algorithm, but it is also the least accurate of all, as seen in the sensitivity tests (experiments 8–10). Yet, our algorithm is about 30 times faster than the CNM algorithm.

In the Appendix we present some more details of the experiment. Figures 2 and 3 show the distribution of the degrees of the vertices for each of the experiments. The *in-degree* of a vertex in those tables is defined here as the number of the adjacent vertices from the same cluster as defined during the generation process and the *out-degree* as the number of adjacent vertices from a different partition. One would expect that the support interval (the interval where the density function is positive) for the in-degrees will always be greater than (to the right of) the one for the out-degrees, in order to have well defined community structures, but this is not always the case. When the number of the communities is large (as in experiment 4), it is possible for the out-degrees of vertices to exceed their in-degrees, while the average number of neighbors to any *fixed* neighboring community to be still lower than the in-degree. In experiments 8, 9, and 10, this effect is further amplified by the low modularity of the partitions, which translates into weaker community structures.

Since Table 3 shows only averages, we give on Figures 4 and 5 the distribution of the modularities for each algorithm and each experiment, represented as differences between the modularities of our algorithm and those of the other algorithms. Those figures show that in experiments 1 through 8 our algorithm not only produces equal or better quality clusterings on average, but virtually on any single graph in those tests. The only experiments where the quality is worse in some instances, in spite of the good quality of our algorithm on average, are experiments 9 and 10, where the modularity is very low – 0.123 and 0.081, respectively. Algorithm N performs the best in those experiments, which shows that it can be a good alternative to our algorithm in low to moderate size graphs (up to 4-5 million edges).

In summary, in those experiments our algorithm produced partitions of quality comparable to the most accurate existing algorithms, in times orders of magnitude smaller. Ours is the only one of the tested algorithms that can produce high quality clusterings on graph of sizes exceeding several million edges.

3.2 Testing on DIMACS Testbed graphs

We ran our algorithm on the Co-author and Citation Networks and the Clustering Instances datasets of the DIMACS Challenge testbed. Each file was preprocessed in order to convert it from the Metis format to our Metis-like format. The difference is in the weights of the nodes. In order to use the CL model modularity, which is the type of modularity chosen for this Challenge, each node of the network should have a weight equal to the sum of the weights of its adjacent nodes. In case the graph is unweighted, i.e., all weight are one, the weight becomes the degree of the node. This has been done in order to have the same algorithm handle without change both cases of the CL and the Erdős-Renyi models, the difference being only in the definition of node weights in the input file. The two graphs that we have not processed are uk-2002.graph and uk-2007-05.graph of the Clustering Instances dataset. Because of their large size, they could not fit in the memory of the computer that we ran the experiments on.

The graphs from the Citation Networks dataset were originally used in [46] and the graphs from the Clustering Instances dataset were used in [47], [3], [5], [10], [12], [11], [49], [1], [38], [30], [35], [21], [31], [34], [17], [36], [48], [33], [43], [41], [9], [23], [16], [44], and [4].

4 Conclusion

This paper proposes a new approach for modularity optimization by reducing it to a minimum cut problem and then solving the latter problem by applying methods for graph partitioning. Our proof-of-concept implementation, based on the METIS partitioning package, demonstrated the practicality of the approach. The changes we made to METIS were relatively small and various improvements and refinements that take into account the specifics of the clustering problem, use alternative minimum cut or graph partitioning algorithms, or apply heuristics and parameter adjustments in order to improve the accuracy are possible and will be topics of further research.

References

1. Adamic, L.A., Glance, N.: The political blogosphere and the 2004 us election. WWW-2005 Workshop on the Weblogging Ecosystem (2005)

2. Ahuja, R.K., Magnanti, T.L., Orlin, J.B.: Network flows : theory, algorithms, and applications. Prentice Hall (1993)
3. Albert, R., Barabási, A.L.: Emergence of scaling in random networks. *Science* (1999)
4. Arenas, A.: <http://deim.urv.cat/~aarenas/data/welcome.htm>
5. Baird, D., Ulanowicz, R.: The seasonal dynamics of the chesapeake bay ecosystem. *Ecol. Monogr.* 59, 329–364 (1989)
6. Barabasi, A.L., Albert, R.: Emergence of scaling in random networks. *Science* 286, 509 (1999), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/9910332>
7. Barnard, S.T., Simon, H.D.: A fast multilevel implementation of recursive spectral bisection for partitioning unstructured problems. *Concurrency: Practice and Experience* 6, 101–107 (1994), citeseer.ist.psu.edu/barnard94fast.html
8. Blondel, V.D., Guillaume, J.L., Lambiotte, R., Lefebvre, E.: Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment* 2008(10), P10008 (2008), <http://stacks.iop.org/1742-5468/2008/i=10/a=P10008>
9. Bogua, M., Pastor-Satorras, R., Diaz-Guilera, A., Arenas, A.: Pgp network. *Physical Review E* 70 (September 2004)
10. Boldi, P., Codenotti, B., Santini, M., Vigna, S.: Ubicrawler: A scalable fully distributed web crawler. *Software: Practice & Experience* 34(8), 711–726 (2004)
11. Boldi, P., Rosa, M., Santini, M., Vigna, S.: Layered label propagation: A multiresolution coordinate-free ordering for compressing social networks. In: *Proceedings of the 20th international conference on World Wide Web*. ACM Press (2011)
12. Boldi, P., Vigna, S.: The WebGraph framework I: Compression techniques. In: *Proc. of the Thirteenth International World Wide Web Conference (WWW 2004)*. pp. 595–601. ACM Press, Manhattan, USA (2004)
13. Brandes, U., Dellinger, D., Gaertler, M., Görke, R., Hoefer, M., Nikoloski, Z., Wagner, D.: On modularity clustering. *IEEE Trans. Knowl. Data Eng.* 20(2), 172–188 (2008)
14. Chung, F., Lu, L.: Connected components in random graphs with given degree sequences. *Annals of Combinatorics* 6, 125–145 (2002)
15. Clauset, A., Newman, M.E.J., Moore, C.: Finding community structure in very large networks. *Physical Review E* 70, 066111 (2004), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0408187>
16. Duch, J., Arenas, A.: C. elegans metabolic network. *Physical Review E* 72 (2005)
17. Duch, J., Arenas, A.: Condensed matter collaborations 2003. *Phys. Rev. E* 72 (2005)
18. Erdos, P., Renyi, A.: On random graphs. *Publicationes Mathematicae* 6, 290–297 (1959)
19. Fiduccia, C.M., Mattheyses, R.M.: A linear time heuristic for improving network partitions. In: *IEEE Design Automation Conference*. pp. 175–181 (1982)
20. Flake, G., Tarjan, R., Tsioutsoulklis, K.: Graph clustering and minimum cut trees. *Internet Mathematics* 1, 385–408 (2004)
21. Girvan, M., Newman, M.E.J.: American college football. *Proc. Natl. Acad. Sci. USA* 99, 7821–7826 (2002)
22. Girvan, M., Newman, M.E.J.: Community structure in social and biological networks. *Proc. Natl. Acad. Sci. USA* 99, 7821–7826 (2002), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0112110>
23. Guimera, R., Danon, L., Diaz-Guilera, A., Giralt, F., Arenas, A.: E-mail network urv. *Physical Review E* 68 (2003)
24. Guimera, R., Amaral, L.A.N.: Functional cartography of complex metabolic networks. *Nature* 433, 895 (2005), <http://www.citebase.org/abstract?id=oai:arXiv.org:q-bio/0502035>
25. Hastings, M.B.: Community detection as an inference problem. *Phys.Rev.E* 74, 035102 (2006), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0604429>
26. Hendrickson, B., Leland, R.W.: A multi-level algorithm for partitioning graphs. In: *ACM/IEEE Conference on Supercomputing* (1995), citeseer.ist.psu.edu/hendrickson93multilevel.html
27. Karypis, G., Kumar, V.: Multilevel graph partitioning schemes. In: *International Conference on Parallel Processing*. pp. 113–122 (1995)
28. Karypis, G., Kumar, V.: A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.* 20(1), 359–392 (1998)

29. Kernighan, B.W., Lin, S.: An efficient heuristic procedure for partitioning graphs. *Bell Sys. Tech. J.* 49(2), 291–308 (1970)
30. Knuth, D.E.: *Les miserables: coappearance network of characters in the novel les miserables*. The Stanford GraphBase: A Platform for Combinatorial Computing (1993)
31. Lusseau, D., Schneider, K., Boisseau, O.J., Haase, P., Slooten, E., Dawson, S.M.: Dolphin social network. *Behavioral Ecology and Sociobiology* 54, 396–405 (2003)
32. Newman, M.: Fast algorithm for detecting community structure in networks. *Physical Review E* 69, 066133 (2004), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0309508>
33. Newman, M.E.J.: Astrophysics collaborations. *Proc. Natl. Acad. Sci.* 98, 404–409 (September 2001)
34. Newman, M.E.J.: Condensed matter collaborations 2005. *Proc. Natl. Acad. Sci. USA* 98, 404–409 (2001)
35. Newman, M.E.J.: High-energy theory collaborations. *Proc. Natl. Acad. Sci. USA* 98, 404–409 (2001)
36. Newman, M.E.J.: The structure of scientific collaboration networks. *Proc. Natl. Acad. Sci.* 98, 404–409 (2001)
37. Newman, M.E.J.: Mixing patterns in networks. *Physical Review E* 67, 026126 (2003), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0209450>
38. Newman, M.E.J.: Coauthorships in network science. *Phys. Rev. E* 74 (May 2006)
39. Newman, M.E.J.: Finding community structure in networks using the eigenvectors of matrices. *Physical Review E* 74, 036104 (2006), doi:10.1103/PhysRevE.74.036104
40. Newman, M.E.J.: Modularity and community structure in networks. *Proc. Natl. Acad. Sci. USA* 103, 8577 (2006), doi:10.1073/pnas.0601602103
41. Newman, M.E.J.: Word adjacencies. *Phys. Rev. E* 74 (2006)
42. Newman, M.E.J., Girvan, M.: Finding and evaluating community structure in networks. *Physical Review E* 69, 026113 (2004), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0308217>
43. Newman, M.: Internet: a symmetrized snapshot of the structure of the internet at the level of autonomous systems. The University of Oregon Route Views Project (July 2006)
44. P.Gleiser, Danon, L.: Jazz musicians network. *Adv. Complex Syst.* 565 (2003)
45. Reichardt, J., Bornholdt, S.: Detecting fuzzy community structures in complex networks with a potts model. *Physical Review Letters* 93, 218701 (2004), <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0402349>
46. Sanders, R.G.P., Schultes, D.: Better approximation of betweenness centrality. 10th Workshop on Algorithm Engineering and Experimentation pp. 90–108 (2008)
47. Watts, Strogatz: Collective dynamics of small-world networks. *Nature* (1998)
48. Watts, D., Strogatz, S.: Neural network. *Nature* 393, 440–442 (1998)
49. Watts, D.J., Strogatz, S.H.: Power grid. *Nature* 393, 440–442 (1998)
50. White, S., Smyth, P.: A spectral clustering approach to finding communities in graph. In: Proceedings of the SIAM International Conference on Data Mining (2005), citeseer.ist.psu.edu/734075.html

	Exp	# vert.	# edges	# clust	Q_{orig}	Q_{CNM}	Q_{N}	Q_{GA}	Q_{RB}	Q_{here}
						t_{CNM}	t_{N}	t_{GA}	t_{RB}	t_{here}
# communities	1	200	8934	2	0.388	0.387	0.388	0.387	0.386	0.388
						1.15	0.70	88.45	35.55	0.07
	2	400	21811	4	0.476	0.474	0.476	0.472	0.473	0.476
						2.45	3.35	335.50	102.40	0.15
	3	600	38743	6	0.447	0.445	0.447	0.445	0.445	0.447
						4.15	9.95	928.20	189.95	0.30
	4	900	71654	9	0.386	0.370	0.386	0.385	0.384	0.386
						7.85	23.05	2539.15	388.25	0.50
sparsity	5	200	9919	2	0.298	0.296	0.298	0.296	0.296	0.298
						1.05	0.65	98.60	38.70	0.10
	6	200	4958	2	0.299	0.297	0.299	0.297	0.297	0.299
						0.95	0.30	37.85	21.25	0.05
	7	200	2483	2	0.300	0.299	0.300	0.300	0.299	0.300
						0.95	0.40	27.50	22.40	0.05
sensitivity	8	400	38783	4	0.209	0.206	0.209	0.208	0.208	0.209
						3.00	3.40	716.65	184.80	0.10
	9	400	47775	4	0.123	0.113	0.123	0.122	0.122	0.122
						3.45	3.30	819.90	229.85	0.05
	10	400	53864	4	0.081	0.060	0.081	0.081	0.080	0.081
						3.50	3.80	1242.90	248.15	0.35
scalability	11	1000	174990	2	0.357	0.357	0.357	0.356	0.358	0.357
						10.33	17.00	15808.67	1333.67	0.47
	12	5000	3749007	2	0.333	0.332	0.333	–	0.333	0.333
						329.50	2973.00	–	53119.50	8.00
	13	20000	24995617	2	0.300	0.297	0.300	–	–	0.300
						2199.33	18234.67	–	–	76.33

Table 3. Comparing the scalability of our algorithm with the algorithms of Clauset, Newman, and Moore (CNM) [15], Newman (N) [40], Guimera and Amaral (GA) [24], and Reichardt and Bornholdt (RB) [45]. Time is measured in seconds and Q_X and t_X denote the average modularity and the average time for algorithm X .

Experiment	# clusters	Modularity	# vertices	# edges	Run Time (Sec)
coPapersDBLP	124	0.833225	540486	15245729	70.158
coPapersCiteseer	127	0.897382	434102	16036720	61.033
coAuthorsCiteseer	118	0.884720	227320	814134	7.875
CitationCiteseer	55	0.793215	268495	1156647	16.156
CoAuthorsDBLP	104	0.808873	299067	977676	12.187

Table 4. Citation Networks

Experiment	# clusters	Modularity	# vertices	# edges	Run Time (Sec)
adjnoun	6	0.293686	112	425	0.015
as-22july06	33	0.644198	22963	48436	1.546
astro-ph	31	0.716611	16706	121251	0.515
celegans-metabolic	6	0.423146	453	2025	0.000
cnr-2000	25	0.894582	325557	2738969	104.487
cond-mat	59	0.831337	16726	47594	0.390
cond-mat-2003	31	0.750067	31163	120029	0.750
cond-mat-2005	52	0.718460	40421	175691	1.312
dolphins	4	0.526799	62	159	0.000
email	8	0.568339	1133	5451	0.031
football	10	0.600912	115	613	0.015
hep-th	59	0.835649	8361	15751	0.203
jazz	3	0.444469	198	2742	0.015
karate	3	0.402038	34	78	0.000
PGPgiantcompo	56	0.875764	10680	24316	0.234
polblogs	4	0.425972	1490	16715	0.046
polbooks	5	0.523920	105	441	0.015
power	31	0.932475	4941	6594	0.093
caidaRouterLevel	56	0.847032	192244	609066	6.047
celegansneural	7	0.480297	297	2148	0.015
chesapeake	3	0.254654	39	170	0.000
eu-2005	78	0.928244	862664	16138468	130.769
G-n-pin-pout	4	0.380990	100000	501198	3.093
in2004	62	0.968524	1382908	13591473	260.945
lesmis	6	0.565822	77	254	0.000
preferentialAttachment	4	0.280632	100000	499985	2.578
smallWorld	478	0.747160	100000	499998	4.765

Table 5. Clustering Instances

Appendix

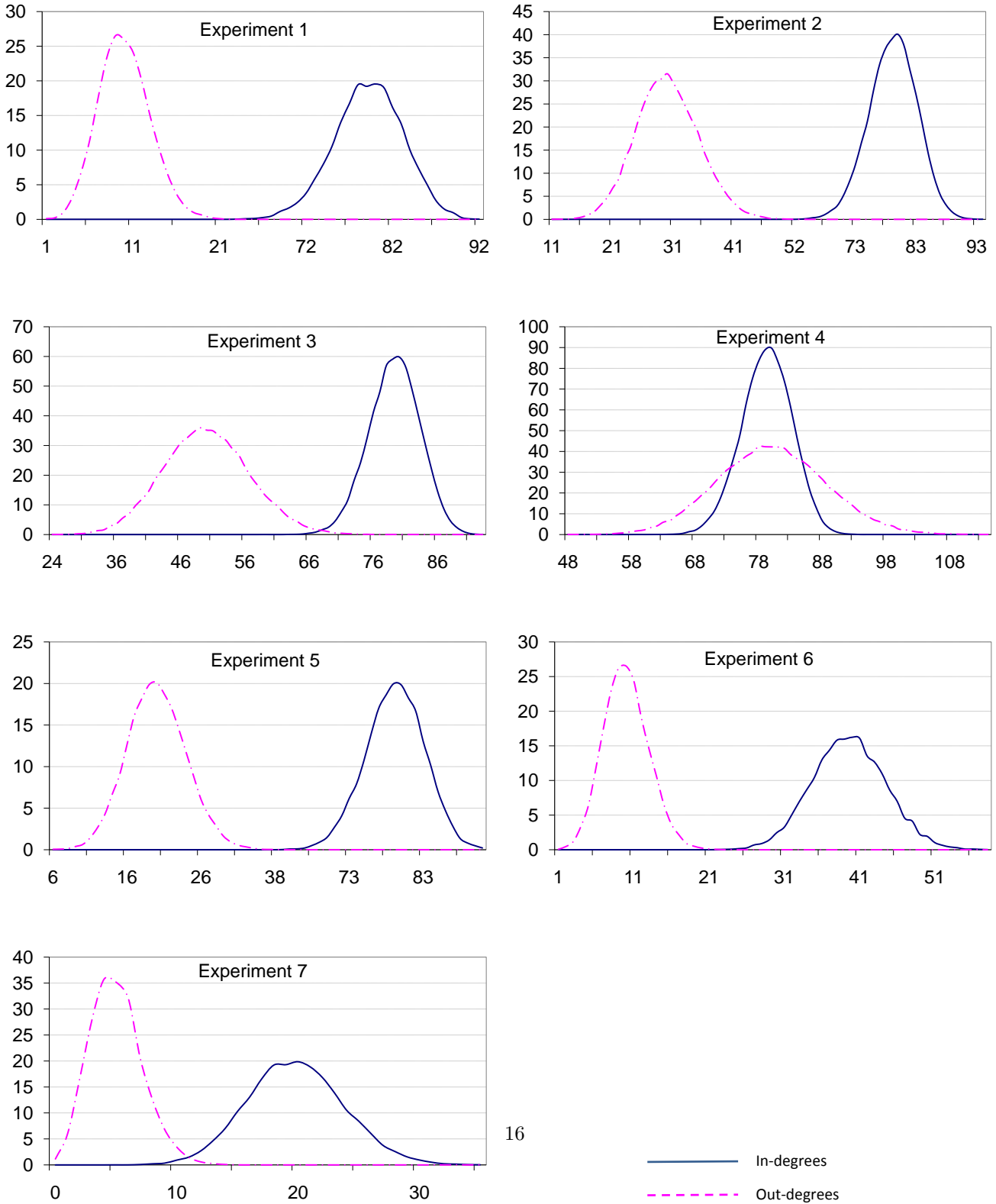


Fig. 2. Distribution of the degrees for experiments 1 through 7. The horizontal axis gives the degree and the vertical the number of the vertices with that degree.

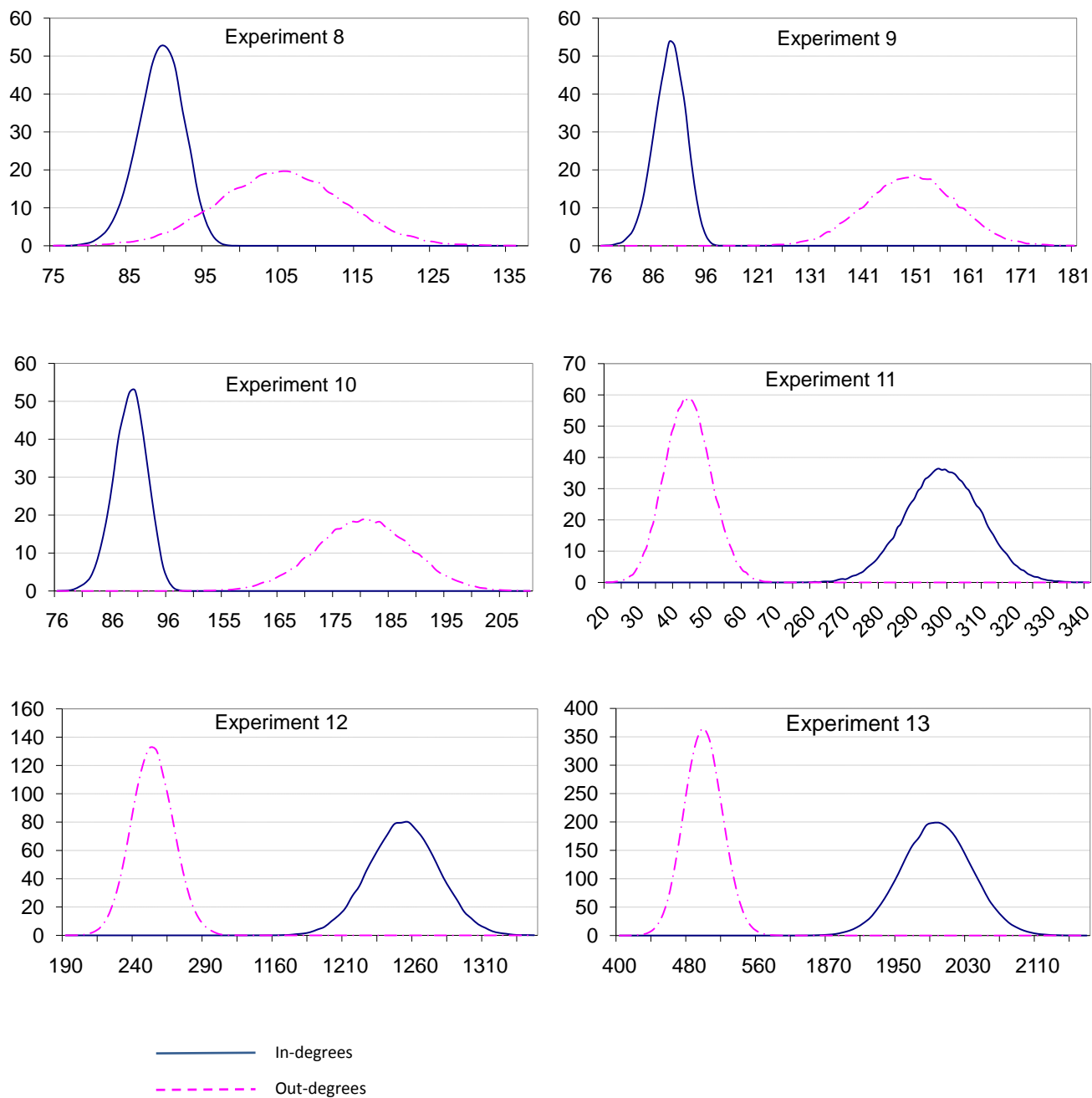


Fig. 3. Distribution of the degrees for experiments 8 through 13. The horizontal axis gives the degree and the vertical the number of the vertices with that degree.

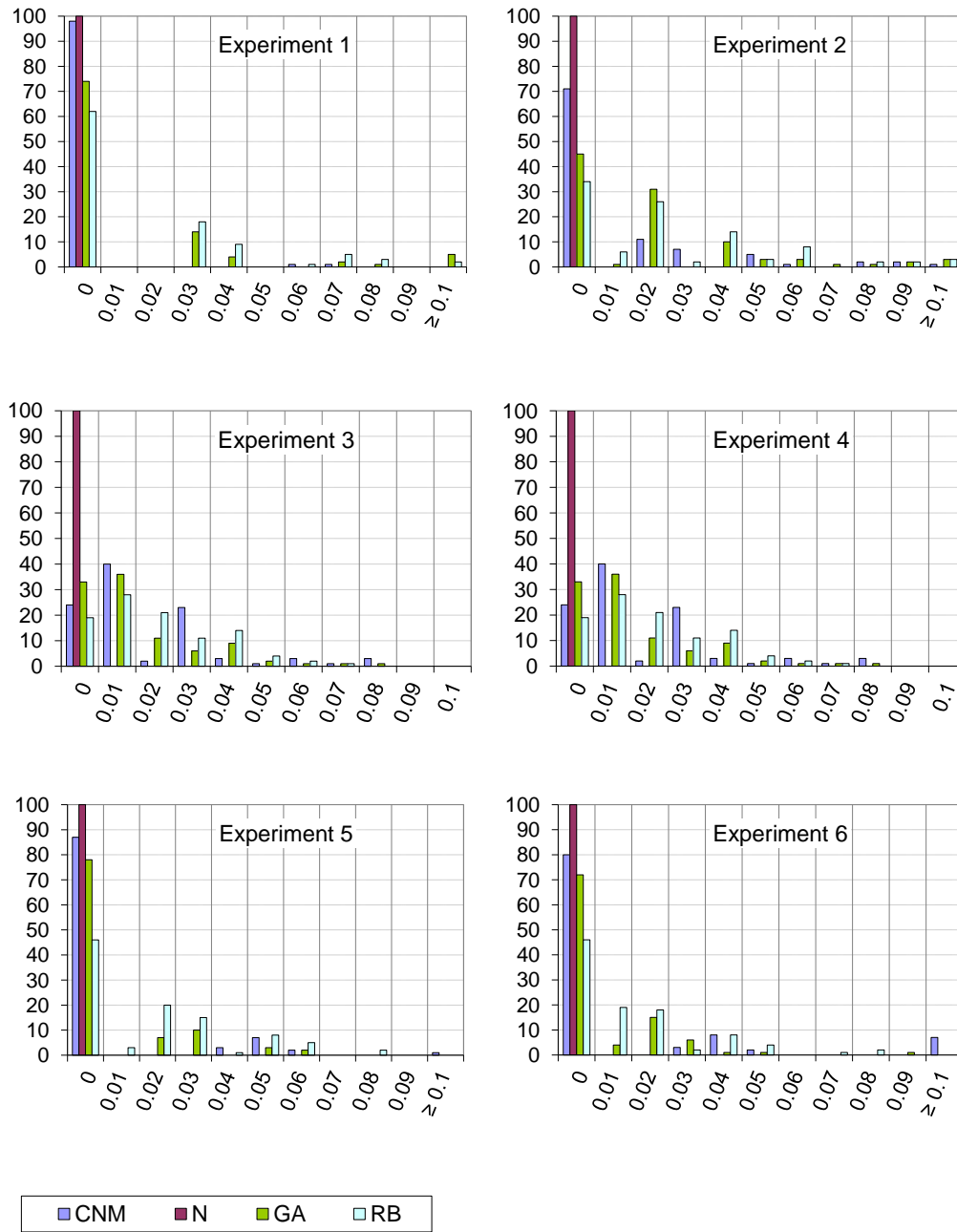


Fig. 4. Distribution of the modularities produced by each algorithm in experiments 1–6. The x -axis gives the difference between the modularities of our algorithm and each of the other algorithms and the y -axis gives the frequency as a percentage of the total number of runs.

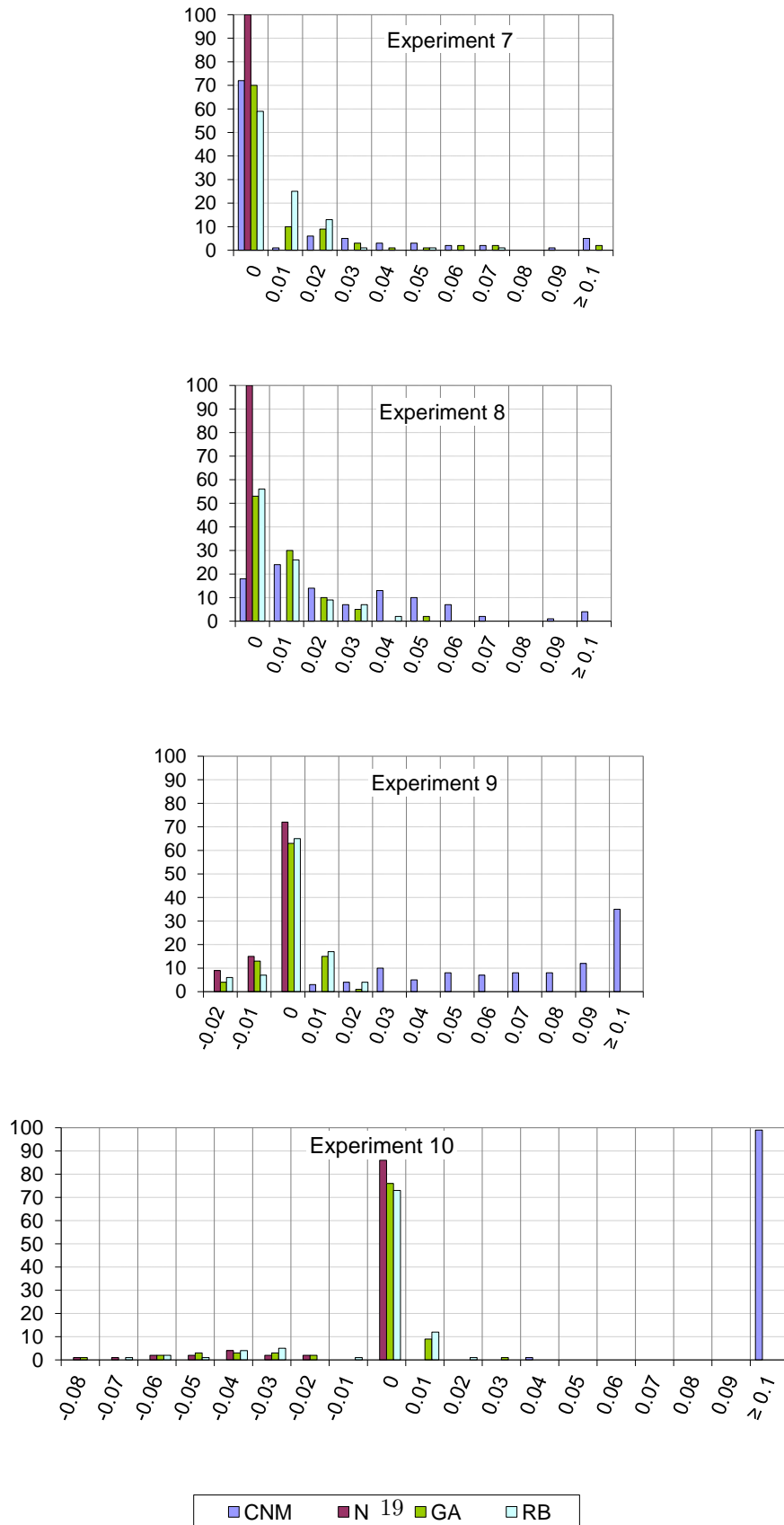


Fig. 5. Distribution of the modularities produced by each algorithm in experiments 7–10. The notation is the same as in Figure 4.