

ON THE APPROXIMABILITY OF MAXIMUM AND MINIMUM EDGE CLIQUE PARTITION PROBLEMS

ANDERS DESSMARK, ANDRZEJ LINGAS,
EVA-MARTA LUNDELL and MIA PERSSON*[†]
*Department of Computer Science, Lund University,
Box 118, SE-22100 Lund, Sweden
{andersd, andrzej, emj, mia}@cs.lth.se*

and

JESPER JANSSON[‡]
*Theoretical Computer Science Group (Yamashita Laboratory),
Dept. of Computer Science and Communication Engineering,
Kyushu University, 744 Motoooka, Nishi-ku, Fukuoka 819-0395, Japan
jj@tcslab.csce.kyushu-u.ac.jp*

Received 13 April 2006

Accepted 23 August 2006

Communicated by Joachim Gudmundsson

ABSTRACT

We consider the following clustering problems: given an undirected graph, partition its vertices into disjoint clusters such that each cluster forms a clique and the number of edges within the clusters is maximized (*Max-ECP*), or the number of edges between clusters is minimized (*Min-ECP*). These problems arise naturally in the DNA clone classification. We investigate the hardness of finding such partitions and provide approximation algorithms. Further, we show that greedy strategies yield constant factor approximations for graph classes for which maximum cliques can be found efficiently.

Keywords: Approximation algorithm; clustering; clique partition; inapproximability; DNA clone classification.

1. Introduction

The *correlation clustering* problem has gained a lot of attention recently [1, 2, 3, 5, 7, 15]; given a complete graph with edges labeled “+” (similar) or “-” (dissimilar), find a partition of the vertices into subsets called *clusters* that agrees as much as possible with the edge labels, i.e., that maximizes the *agreements* (the number of “+” edges inside clusters plus the number of “-” edges between clusters) or that minimizes

*Corresponding author.

[†]Also School of Technology and Society, Malmö University, SE-20506 Malmö, Sweden.

[‡]Supported in part by JSPS (Japan Society for the Promotion of Science).

the *disagreements* (the number of “–” edges inside clusters plus the number of “+” edges between clusters).

In this paper, we consider a special variant of the correlation clustering problem in which there are no negative edge labels. Instead, we omit an edge between two vertices of a dissimilar pair. Furthermore, we require an edge between each pair of vertices in a cluster, i.e, every cluster must form a clique. We consider the following two combinatorial optimization problems. The *maximum edge clique partition problem* (*Max-ECP* for short) aims to find a partition of the vertices into cliques such that the total number of edges within all those cliques is maximized. The related minimization version of this problem, the *minimum edge clique partition problem* (*Min-ECP* for short), is defined analogously with the exception that the total number of edges between the cliques is minimized.

The *Max-ECP* and *Min-ECP* problems first came to our attention in the setting of DNA clone classification [9]. In order to characterize cDNA and ribosomal DNA (rDNA) gene libraries, the powerful DNA array based method *oligonucleotide fingerprinting* is commonly used (see, e.g., [6, 10, 16, 17]).

The problem of clustering binarized fingerprint data such that the number of clusters is minimized was first studied and motivated in [8]. In [9], Figueroa *et al.* propose new approaches of partitioning binarized fingerprints into disjoint clusters in order to maximize the number of pairs of similar fingerprints lying inside the clusters (equivalently, minimize the number of pairs of similar fingerprints lying in different clusters). These problems can hence be viewed as the *Max-ECP* and *Min-ECP* problems where the vertices are the binarized fingerprints and the edges between them indicate their similarity.

1.1. Related results

The well studied correlation clustering problem was first introduced for complete graphs by Bansal *et al.* [2]. It has applications in many areas (see, e.g., [2, 5]). As noted in [2], the problem of maximizing agreements and minimizing disagreements are equivalent at optimality but differ from the point of view of approximation. In [2], it was established that these problems are NP-hard for complete graphs, and a PTAS was given in the case of maximizing agreements, whereas a constant factor approximation is given in the case of minimizing disagreements. This constant factor approximation was later improved by Charikar *et al.* [3] where a factor 4 approximation algorithm is given for complete graphs based on linear programming relaxation. The latter problem was also proved to be APX-hard.

The problems of maximizing agreements and minimizing disagreements were later generalized to include non-necessarily complete graphs with edge weights in [3]. A factor 0.7664 approximation algorithm based on the rounding of a semidefinite programming relaxation for the problem of maximizing agreements for general weighted graphs was given in [3], but this factor was later improved to 0.7666 by Swamy [15]. As for the problem of minimizing disagreements, a factor $O(\log n)$ approximation algorithm for general weighted graphs was proposed (independently) in [3], [5], and [7]. Recently, Ailon *et al.* [1] have provided a randomized expected

3-approximation algorithm for minimizing disagreements. In the case of weighted complete graphs, which satisfy probability constraints ($w_{ij}^+ + w_{ij}^- = 1$ for edge (i, j)) and triangle inequality constraints ($w_{ik}^- \leq w_{ij}^- + w_{jk}^-$) on the edges, they have provided a factor 2 approximation algorithm.

The APX-hardness of the unweighted version of *Min-ECP* has been established by Shamir *et al.* [14]. They have also presented results for the case when a solution must contain exactly p clusters; the so restricted problem is solvable in polynomial time for $p = 2$ but NP-complete for $p > 2$.

1.2. Our results

In this paper, we investigate the approximability of *Max-ECP* and *Min-ECP*. Specifically, we prove that *Max-ECP* on general, undirected graphs is hard to approximate within a factor of $n^{1-o(1)}$, unless $\text{NP} \subseteq \text{ZPTIME}(2^{(\log n)^{o(1)}})$, where n denotes the number of vertices in the input graph. On the other hand, we give an n -approximation algorithm running in polynomial time for this problem^a. In the case of *Min-ECP* we provide a polynomial-time $O(\log n)$ -approximation algorithm for undirected graphs with non-negative weights. We also prove that this problem is NP-hard to approximate within $1 + \frac{1}{880} - \epsilon$, for any $\epsilon > 0$. We further consider the greedy heuristic and show that it yields a 2-approximation for both *Max-ECP* and *Min-ECP*, under the assumption that the largest clique can be determined in polynomial time. Thus, the greedy method could be applied in practice only to graph classes for which maximum cliques can be found efficiently, for instance chordal graphs, line graphs and circular-arc graphs (cf. [8]). We also note that these bounds are actually tight. Table 1 summarizes our contributions.

Table 1. Summary of our results on the polynomial-time approximability of *Max-ECP* and *Min-ECP*. The lower bounds listed in the first two rows of the table show inapproximability results for *Max-ECP* and weighted *Min-ECP* whereas the lower bounds shown in the last two rows concern the worst-case behavior of the greedy method.

<i>Problem</i>	<i>Lower Bound</i>	<i>Upper Bound</i>
Max-ECP	$n^{1-o(1)}$	n
weighted Min-ECP	$1 + \frac{1}{880} - \epsilon$	$O(\log n)$
Greedy Max-ECP	2	2
Greedy Min-ECP	2	2

Our paper is structured as follows. We give more formal definitions of *Max-ECP* and *Min-ECP* in Section 2. In Section 3, we provide a factor n approximation algorithm for *Max-ECP*. In Section 4, we give a lower bound on approximability of *Max-ECP*. In Section 5, we provide a polynomial-time $O(\log n)$ -approximation algorithm for the weighted version of *Min-ECP* and in section 6, we derive a lower bound on approximability of *Min-ECP*. Finally, in Section 7, we consider the greedy

^aMore precisely, our algorithm is a k -approximation algorithm, where k is the number of vertices in the largest clique in the input graph.

algorithm for *Max-ECP* and *Min-ECP* and prove that it yields a factor 2 approximation.

2. Preliminaries

The formal definition of *Max-ECP* and *Min-ECP* is as follows.

Definition 1 Let $G = (V, E)$ be an undirected graph and let $n = |V|$. The problem of maximum edge clique partition (*Max-ECP* for short) is to find a partition of V into disjoint subsets V_1, \dots, V_m such that for each $1 \leq i \leq m$, any two vertices in V_i share an edge and the total number of edges within the subsets V_1, \dots, V_m is maximized.

The problem of minimum edge clique partition (*Min-ECP* for short) is defined analogously to *Max-ECP* with the exception that the total number of edges between the subsets V_1, \dots, V_m is minimized.

The subsets V_1, \dots, V_m in the definition above are referred to as *clusters*, and any partition of V into clusters is a *clustering* of V . For *Max-ECP* and *Min-ECP*, the *score* of a clustering is defined as the total number of edges within the clusters and the total number of edges between the clusters, respectively.

Note that an exact solution to *Max-ECP* is an exact solution to *Min-ECP* and vice versa.

The example shown in Figure 1 demonstrates two feasible solutions to *Max-ECP* and *Min-ECP*. As depicted in Figure 1(a), the total number of edges inside the clusters is 18, hence the solution to *Max-ECP* has a total score of 18. On the contrary, the total number of edges outside the clusters in Figure 1(a) is 12, hence the solution to *Min-ECP* has a total score of 12. The optimal clustering is depicted in Figure 1(b), with the total score of 24 for *Max-ECP* and the total score of 6 for *Min-ECP*.

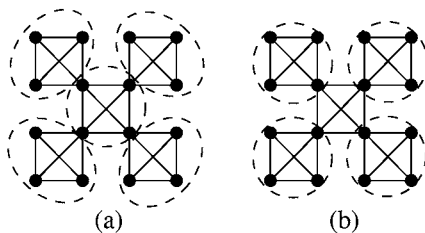


Fig. 1. A feasible solution and the optimal solution to an instance of *Max-ECP* and *Min-ECP*.

3. A Polynomial-Time n -Approximation Algorithm for *Max-ECP*

The *Max-ECP* problem is NP-hard and even hard to approximate within a factor $n^{1-O(1/(\log n)^\gamma)}$, for some constant γ , as proved in the next section. On the positive side, we prove in this section that *Max-ECP* admits a simple polynomial-time, factor k approximation algorithm, where k is the number of vertices in the largest clique in G . The approximation algorithm works as follows: Find a maximum matching

in G and output it and the singletons containing the vertices not covered by the matching as a clique partition.

Theorem 1 *Let k be the number of vertices in the largest clique in G . Max-ECP can be approximated within a factor of k in polynomial time.*

Proof. Denote by $\text{OPT}(G)$ and $\text{APPR}(G)$ the total number of edges within cliques in an optimal solution for *Max-ECP* on G and in the solution returned by the approximation algorithm described above, respectively. Let (V_1, V_2, \dots, V_m) be an optimal solution for *Max-ECP* on G . There is a matching in G which, for $i = 1, \dots, m$, includes at least $\frac{|V_i|-1}{2}$ edges from the clique induced by V_i . Since $k \geq |V_i|$ for all $i = 1, \dots, m$, such a matching includes at least a fraction $\frac{1}{k}$ of the edges from each of the m cliques induced by V_1, V_2, \dots, V_m . Hence, $\text{APPR}(G) \geq \text{OPT}(G) / k$ holds. \square

4. A Lower Bound on the Approximability of Max-ECP

The maximum clique problem is known to not admit an approximation within $n^{1-O(1/(\log n)^\gamma)}$ for some constant γ unless $\text{NP} \subseteq \text{ZPTIME}(2^{(\log n)^{O(1)}})$ [12]. It follows that the aforementioned inapproximability bound holds even if restricted to graphs on n vertices whose largest clique consists of at least n^{1-x} vertices, where $x = O(1/(\log n)^\gamma)$. Below, consider any such graph G . For the sake of contradiction, suppose there exists an n^{1-3x} -approximation algorithm for *Max-ECP*. Let \mathcal{C} be the clustering obtained by running this approximation algorithm on G , and let k be the size of the largest clique in G , i.e., $n^{1-x} \leq k \leq n$ by the assumption above. First observe that an optimal solution to *Max-ECP* for G has at least $\binom{k}{2}$ edges, so \mathcal{C} has at least $k(k-1)/(2n^{1-3x})$ edges. On the other hand, \mathcal{C} contains at most $n(c-1)/2$ edges in total, where c is the size of the largest cluster in \mathcal{C} (to see this, double-count the total number of edges in \mathcal{C} by adding up the number of neighbors in \mathcal{C} of each vertex; clearly, the latter sum is upper-bounded by the number of vertices multiplied by $c-1$). Thus, we get the inequality $k(k-1)/(2n^{1-3x}) \leq n(c-1)/2$, which together with $k \geq n^{1-x}$ yields $c \geq n^x$ for large enough n . This means that by selecting the largest cluster in \mathcal{C} , we would always obtain a clique in G of size at least $1/n^{1-x}$ times n , which is at least $1/n^{1-x}$ times the size of the largest clique in G . Hence, this would yield an n^{1-x} approximation algorithm for the maximum clique problem on G , contradicting [12]. We have proved the following theorem.

Theorem 2 *Unless $\text{NP} \subseteq \text{ZPTIME}(2^{(\log n)^{O(1)}})$, the Max-ECP problem does not admit an $n^{1-O(1/(\log n)^\gamma)}$ approximation, for some constant γ .*

5. A Polynomial-Time $O(\log n)$ -Approximation Algorithm for Weighted Min-ECP

Min-ECP can be approximated within a factor of $O(\log n)$ in polynomial time, even for edge-weighted graphs with arbitrary non-negative weights, as follows.

Let $G = (V, E)$ be a given instance of *Min-ECP* in which each edge e has a non-negative weight $w(e)$. Define $W = \max_{e \in E} w(e)$. Construct an edge-weighted, edge-labeled, complete graph $G' = (V, E')$, where each $e \in E'$ is labeled by $' + '$ and

assigned weight $w(e)$ if $e \in E$, or labeled by $'-'$ and assigned weight $W \cdot n^2 \log^2 n$ if $e \notin E$. Run any one of the polynomial-time $O(\log n)$ -approximation algorithms for Minimum Disagreement Correlation Clustering for weighted graphs [3, 5, 7] on G' to obtain a clustering \mathcal{C}' for V , and return the set \mathcal{S} of subgraphs of G induced by \mathcal{C}' .

Lemma 1 *For any two vertices $u, v \in V$ which are not joined by an edge in G , u and v do not belong to the same cluster in \mathcal{C}' .*

Proof. First note that partitioning the vertex set into singleton clusters would yield a feasible solution whose disagreement score is at most $W \cdot \binom{n}{2}$. Since the disagreement score of an optimal solution is at most this much, the disagreement score for \mathcal{C}' is $O(W \cdot n^2 \log n)$. Now suppose that u and v belong to the same cluster in \mathcal{C}' . Then the disagreement score for \mathcal{C}' must be at least $W \cdot n^2 \log^2 n$, which contradicts the above. \square

By Lemma 1, the vertices from each cluster in \mathcal{C}' form a clique in G . Since the clusters in \mathcal{C}' are disjoint, \mathcal{S} is a partition of G into cliques, which proves the correctness of the method.

Next, we consider the approximation ratio. For any partition M of G into cliques, denote by $ECP(M)$ the ECP score for M , i.e., the sum of all weights of edges whose two endpoints belong to different cliques in M . Similarly, for any clustering M' of G' , let $Disagree(M')$ be the disagreement correlation clustering score for M' . Finally, $MinECP(G)$ and $MinDisagree(G')$ denote the minimum possible scores of ECP for G and $Disagree$ for G' , respectively.

Lemma 2 *$ECP(\mathcal{S})$ is at most $O(\log n)$ times $MinECP(G)$.*

Proof. Let M be a partition of G into cliques which minimizes ECP , and let M' be the clustering of G' induced by the cliques in M . Then, since only edges labeled by $'+'$ contribute to $Disagree(M')$, we obtain $MinECP(G) = ECP(M) = Disagree(M') \geq MinDisagree(G')$.

Next, observe that $ECP(\mathcal{S})$ is equal to $Disagree(\mathcal{C}')$ because only edges labeled by $'+'$ contribute to $Disagree(\mathcal{C}')$ by Lemma 1. Moreover, $Disagree(\mathcal{C}')$ is at most $O(\log n)$ times $MinDisagree(G')$. It follows that $ECP(\mathcal{S})$ is at most $O(\log n)$ times $MinECP(G)$. \square

To summarize:

Theorem 3 *Weighted Min-ECP can be approximated within a factor of $O(\log n)$ in polynomial time.*

6. A Lower Bound for Min-ECP

Shamir *et al.* have established the APX-hardness of unweighted *Min-ECP* by a reduction from a special variant of set cover in [14]. It follows by [14] that the *Min-ECP* problem cannot have a polynomial-time approximation scheme unless $P=NP$. However, no explicit lower bound on the approximation factor for *Min-ECP* achievable in polynomial time is known in the literature.

In this section, we present a new reduction from the so called three way cut problem to the weighted *Min-ECP* problem which yields an explicit lower bound on the approximation factor.

The problem of three way cut (3WC) is to find a minimum number of edges whose removal disconnects three distinguished vertices.

Let A and B be two optimization problems (maximization or minimization). A linearly reduces [4] to B if there are two polynomial time algorithms h and g , and constants $\alpha, \beta > 0$ such that

- For an instance a of A , algorithm h produces an instance $b = h(a)$ of B such that the cost of an optimal solution for b , $opt(b)$, is at most $\alpha \cdot opt(a)$, and
- For $a, b = h(a)$, and any solution y of b , algorithm g produces a solution x of a such that $|cost(x) - opt(a)| \leq \beta|cost(y) - opt(b)|$.

By [13], if A linearly reduces to B and B has a polynomial-time $1 + \epsilon$ approximation algorithm then A has a polynomial-time $(1 + \alpha\beta\epsilon)$ approximation algorithm.

Max-Cut is the problem of finding, for an undirected graph with vertex set V , a partition V_1, V_2 of V such that the number of edges $\{u, v\}$ where $\{u, v\} \cap V_1$ and $\{u, v\} \cap V_2$ are both nonempty is maximized.

In [4], Dahlhaus *et al.* presented a linear reduction of the Max-Cut problem to 3WC in order to prove that 3WC is APX-hard. Since Max-Cut is APX-hard [11], the APX-hardness of 3WC follows. In the aforementioned reduction $\alpha = 56$ and $\beta = 1$ [4]. In fact, α can be decreased to 55 by the proof of Theorem 5 in [4]^b. On the other hand, Håstad has shown that for any $\epsilon > 0$, it is NP-hard to approximate Max-Cut within $1 + \frac{1}{16} - \epsilon$ [11]. Hence, we obtain the following lemma.

Lemma 3 *For any $\epsilon > 0$, it is NP-hard to approximate 3WC within $1 + \frac{1}{880} - \epsilon$.*

To reduce 3WC to weighted *Min-ECP*, fix an arbitrary $\delta > 0$, and transform any given instance of 3WC on n vertices to an instance of *Min-ECP* as follows:

- Assign the weight 1 to each edge in the instance.
- For each non-adjacent pair u, v of vertices in the instance insert an edge of weight δ/n^2 .
- For each distinguished vertex $s_i, i = 1, 2, 3$, add an auxiliary vertex u_i and make it adjacent with each vertex of the instance. Assign the weight n^2 to each of the three edges (s_i, u_i) and the weight δ/n^2 to the remaining edges incident to the vertices $u_i, i = 1, 2, 3$.

Figure 2 demonstrates how the transformation from an instance of 3WC to an instance of *Min-ECP* works.

In this figure, note that a dashed line between a pair of vertices indicates an edge with weight δ/n^2 .

Note that in an optimal *Min-ECP* solution to the transformed instance each of the pairs $s_i, u_i, i = 1, 2, 3$ belongs to a separate clique and the total weight of the edges outside all the cliques in the optimal solution is between cut and $cut + \delta$ where cut stands for the value of an optimal solution to the instance of 3WC.

^bIn the proof of Theorem 5 in [4], observe that $OPT_{3WC}(f(G)) = 56 \cdot \frac{|E|}{2} - K \leq 56 \cdot OPT_{Max-Cut}(G) - OPT_{Max-Cut}(G) = 55 \cdot OPT_{Max-Cut}(G)$.

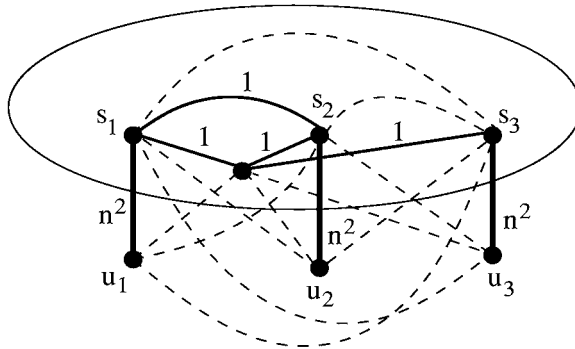


Fig. 2. Transformation from 3WC to *Min-ECP*.

Suppose that for some $\epsilon > 0$, weighted *Min-ECP* could be approximated in polynomial time within a factor of f where $f \leq 1 + \frac{1}{880} - \epsilon$. Then using the set of edges between the three cliques in an approximate solution for weighted *Min-ECP* as an approximate solution for 3WC would yield a three-way cut for the original graph of cardinality at most $f \cdot (cut + \delta) \leq (f + f \cdot \delta) \cdot cut$. By setting $\delta = \frac{\epsilon}{2 \cdot (1 + \frac{1}{880} - \epsilon)}$, we could approximate 3WC in polynomial time within $1 + \frac{1}{880} - \epsilon/2$. We obtain a contradiction with Lemma 3. Hence, we obtain the following theorem.

Theorem 4 *For any $\epsilon > 0$, it is NP-hard to approximate weighted *Min-ECP* within $1 + \frac{1}{880} - \epsilon$.*

7. Greedy Method for Max-ECP and Min-ECP

The greedy strategy applies naturally to the *Max-ECP* and *Min-ECP* problems: iteratively pick the largest clique until all elements have been partitioned into disjoint clusters. However, the problem of finding a maximum clique is itself known to be extremely hard to approximate [12]. Thus, the greedy method could be applied in practice only to graph classes for which maximum cliques can be found efficiently (cf. [8]).

Theorem 5 *The greedy method yields a 2-approximation for *Max-ECP* and *Min-ECP*.*

Proof. Consider an optimal solution to the *Max-ECP* problem (or, the *Min-ECP* problem, respectively) and let us assume that it consists of m cliques. Let E_i be the set of edges in the m cliques, and let E_o be the set of edges of graph G outside these cliques. Let C be the largest clique, say on k vertices, picked by the greedy method. Suppose first that the intersection of C with any clique in the optimal partition is a singleton or empty. Thus, in a way, the at most $k(k - 1)$ edges in E_i are replaced with the $k(k - 1)/2$ edges in C (or, the $k(k - 1)/2$ edges in $C \cap E_o$ with at most $k(k - 1)$ new edges outside the cliques, respectively). In the remaining case, if the intersection of C with any of the cliques in the optimal partition contains more than one vertex, less than $k(k - 1)$ edges in E_i are replaced by the $k(k - 1)/2$ edges in C (or, the $k(k - 1)/2$ edges in $C \cap E_o$ are replaced by

less than $k(k - 1)$ new edges outside the cliques, respectively). By iterating the argument, we obtain the theorem. \square

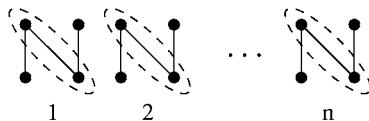


Fig. 3. An example illustrating the worst-case performance of the greedy strategy for *Max-ECP* and *Min-ECP*.

The example shown in Figure 3 demonstrates that our upper bound on the approximation factor of the greedy method for *Max-ECP* is tight. Simply, the greedy method may produce n 2-cliques and $2n$ 1-cliques (singletons) yielding n edges whereas the optimal clique partition consists of $2n$ 2-cliques yielding $2n$ edges.

Figure 3 is also a tight example for greedy *Min-ECP*. Note that the number of edges between cliques will be $2n$ in the approximate solution, whereas the optimum contains n edges between the $2n$ 2-cliques.

8. Final Remarks

By using rather maximum weight matching than maximum cardinality matching we can easily generalize our n -approximation method for *Max-ECP* to include edge weights.

It is an interesting open problem whether or not the gap between the upper and lower bounds on approximability of *Min-ECP* could be tightened.

A careful reader might observe that our approximation hardness result for *Max-ECP* does not hold for the graph classes for which our greedy method could be applied practically. The complexity and approximation status of *Max-ECP* and *Min-ECP* for the aforementioned graph classes are interesting open problems.

References

1. N. Ailon, M. Charikar, and A. Newman, "Aggregating inconsistent information: Ranking and Clustering," *Proc. 37th Annual ACM Symposium on Theory of Computing (STOC 2005)*, Baltimore, MD, May 2005, pp. 684–693.
2. N. Bansal, A. Blum, and S. Chawla, "Correlation Clustering," *Machine Learning* **56**(1–3) (2004) 89–113.
3. M. Charikar, V. Guruswami, and A. Wirth, "Clustering with Qualitative Information," *Proc. 44th Annual Symposium on Foundations of Computer Science (FOCS 2003)*, Cambridge, MA, Oct. 2003, pp. 524–533.
4. E. Dahlhaus, D. S. Johnson, and C. H. Papadimitriou, P.D. Seymour, and M. Yannakakis, "The Complexity of Multiterminal Cuts," *SIAM J. Comput.* **23** (1994) 864–894.
5. E. Demaine and N. Immorlica, "Correlation Clustering with Partial Information," *Proc. 6th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems (APPROX 2003)*, Princeton, NY, Aug. 2003, pp. 1–13.
6. S. Drmanac, N. A. Stavropoulos, I. Labat, J. Vonau, B. Hauser, M. B. Soares, and R. Drmanac, "Gene-representing cDNA clusters defined by hybridization of 57,419

- clones from infant brain libraries with short oligonucleotide probes," *Genomics* **37** (1996) 29–40.
7. D. Emanuel and A. Fiat, "Correlation Clustering – Minimizing Disagreements on Arbitrary Weighted Graphs," *Proc. 11th Annual European Symposium on Algorithms (ESA 2003)*, Budapest, Hungary, Sep. 2003, pp. 208–220.
 8. A. Figueroa, J. Borneman, and T. Jiang, "Clustering binary fingerprint vectors with missing values for DNA array data analysis," *Journal of Computational Biology* **11**(5) (2004) 887–901.
 9. A. Figueroa, A. Goldstein, T. Jiang, M. Kurowski, A. Lingas, and M. Persson, "Approximate Clustering of Fingerprint Vectors with Missing Values," *Proc. 11th Computing: The Australasian Theory Symposium (CATS 2005)*, Newcastle, NSW, Jan. 2005, pp. 57–60.
 10. R. Herwig, A. J. Poustka, C. Müller, C. Bull, H. Lehrach, and J. O'Brien, "Large-scale clustering of cDNA-fingerprinting data," *Genome research* **9** (1999) 1093–1105.
 11. J. Håstad, "Some optimal inapproximability results," *Journal of the ACM* **48**(4) (2001) 798–859.
 12. S. Khot, "Improved inapproximability results for MaxClique, chromatic number, and approximate graph coloring," *Proc. 42th Annual Symposium on Foundations of Computer Science (FOCS 2001)*, Las Vegas, NV, Oct. 2001, pp. 600–609.
 13. C. H. Papadimitriou and M. Yannakakis, "Optimization, Approximation, and Complexity Classes," *Journal of Computer and System Sciences* **43** (1991) 425–440.
 14. R. Shamir, R. Sharan, and D. Tsur, "Cluster Graph Modification Problems," *Proc. 28th International Workshop on Graph-Theoretic Concepts in Computer Science (WG 2002)*, Cesky Krumlov, Czech Republic, Jun. 2002, pp. 379–390.
 15. C. Swamy, "Correlation Clustering: maximizing agreements via semidefinite programming," *Proc. 15th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 2004)*, New Orleans, LA, Jan. 2004, pp. 526–527.
 16. L. Valinsky, G. Della Vedova, T. Jiang, and J. Borneman, "Oligonucleotide fingerprinting of ribosomal RNA genes for analysis of fungal community composition," *Applied and Environmental Microbiology* **68** (2002) 5999–6004.
 17. L. Valinsky, G. Della Vedova, A. Scupham, S. Alvey, A. Figueroa, B. Yin, R. Hartin, M. Chrobak, D. Crowley, T. Jiang, and J. Borneman, "Analysis of bacterial community composition by oligonucleotide fingerprinting of rRNA genes," *Applied and Environmental Microbiology* **68** (2002) 3243–3250.