

Clustering Aggregation

Aristides Gionis, Heikki Mannila, and Panayiotis Tsaparas
Helsinki Institute for Information Technology, BRU
Department of Computer Science
University of Helsinki, Finland
first.lastname@cs.helsinki.fi

Abstract

We consider the following problem: given a set of clusterings, find a clustering that agrees as much as possible with the given clusterings. This problem, clustering aggregation, appears naturally in various contexts. For example, clustering categorical data is an instance of the problem: each categorical variable can be viewed as a clustering of the input rows. Moreover, clustering aggregation can be used as a meta-clustering method to improve the robustness of clusterings. The problem formulation does not require a-priori information about the number of clusters, and it gives a natural way for handling missing values. We give a formal statement of the clustering-aggregation problem, we discuss related work, and we suggest a number of algorithms. For several of the methods we provide theoretical guarantees on the quality of the solutions. We also show how sampling can be used to scale the algorithms for large data sets. We give an extensive empirical evaluation demonstrating the usefulness of the problem and of the solutions.

1 Introduction

Clustering is an important step in the process of data analysis with applications to numerous fields. Informally, clustering is defined as the problem of partitioning data objects into groups (clusters), such that objects in the same group are similar, while objects in different groups are dissimilar. This definition assumes that there is some well defined *quality measure* that captures intra-cluster similarity and/or inter-cluster dissimilarity, and then clustering becomes the problem of grouping together data objects so that the quality measure is optimized.

In this paper we propose a novel approach to clustering that is based on the concept of *aggregation*. We assume that given the data set we can obtain some information on how these points should be clustered. This information comes in the form of m clusterings $\mathcal{C}_1, \dots, \mathcal{C}_m$. The objective is to

	\mathcal{C}_1	\mathcal{C}_2	\mathcal{C}_3	\mathcal{C}
v_1	1	1	1	1
v_2	1	2	2	2
v_3	2	1	1	1
v_4	2	2	2	2
v_5	3	3	3	3
v_6	3	4	3	3

Figure 1. An example of clustering aggregation

produce a single clustering \mathcal{C} that agrees as much as possible with the m clusterings. We define a disagreement between two clusterings \mathcal{C} and \mathcal{C}' as a pair of objects (v, u) such that \mathcal{C} places them in the same cluster, while \mathcal{C}' places them in a different cluster, or vice versa. If $d(\mathcal{C}, \mathcal{C}')$ denotes the number of disagreements between \mathcal{C} and \mathcal{C}' , then the task is to find a clustering \mathcal{C} that minimizes $\sum_{i=1}^m d(\mathcal{C}_i, \mathcal{C})$.

As an example, consider the dataset $V = \{v_1, v_2, v_3, v_4, v_5, v_6\}$ that consists of six objects, and let $\mathcal{C}_1 = \{\{v_1, v_2\}, \{v_3, v_4\}, \{v_5, v_6\}\}$, $\mathcal{C}_2 = \{\{v_1, v_3\}, \{v_2, v_4\}, \{v_5\}, \{v_6\}\}$, and $\mathcal{C}_3 = \{\{v_1, v_3\}, \{v_2, v_4\}, \{v_5, v_6\}\}$ be three clusterings of V . Figure 1 shows the three clusterings, where each column corresponds to a clustering, and a value i denotes that the tuple in that row belongs in the i -th cluster of the clustering in that column. The rightmost column is the clustering $\mathcal{C} = \{\{v_1, v_3\}, \{v_2, v_4\}, \{v_5, v_6\}\}$ that minimizes the total number of disagreements with the clusterings $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$. In this example the total number of disagreements is 5: one with the clustering \mathcal{C}_2 for the pair (v_5, v_6) , and four with the clustering \mathcal{C}_1 for the pairs (v_1, v_2) , (v_1, v_3) , (v_2, v_4) , (v_3, v_4) . It is not hard to see that this is the minimum number of disagreements possible for any partition of the dataset V .

We define *clustering aggregation* as the optimization problem where, given a set of m clusterings, we want to find the clustering that minimizes the total number of disagreements with the m clusterings. Clustering aggregation provides a general framework for dealing with a variety of problems related to clustering: (i) it gives a natural cluster-

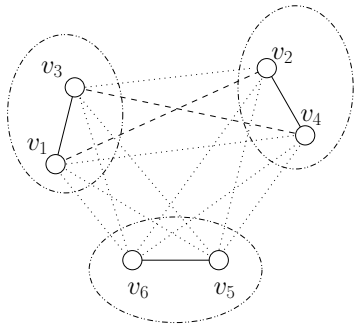


Figure 2. Correlation clustering instance for the dataset in Figure 1. Solid edges indicate distances of $1/3$, dashed edges indicate distances of $2/3$, and dotted edges indicate distances of 1 .

ing algorithm for categorical data, which allows for a simple treatment of missing values, (ii) it handles heterogeneous data, where tuples are defined over incomparable attributes, (iii) it determines the appropriate number of clusters and it detects outliers, (iv) it provides a method for improving the clustering robustness, by combining the results of many clustering algorithms, (v) it allows for clustering of data that is vertically partitioned in order to preserve privacy. We elaborate on the properties and the applications of clustering aggregation in Section 2.

The algorithms we propose for the problem of clustering aggregation take advantage of a related formulation, which is known as *correlation clustering* [2]. We map clustering aggregation to correlation clustering by considering the tuples of the dataset as vertices of a graph, and summarizing the information provided by the m input clusterings by weights on the edges of the graph. The weight of the edge (u, v) is the fraction of clusterings that place u and v in different clusters. For example, the correlation clustering instance for the dataset in Figure 1 is shown in Figure 2. Note that if the weight of the edge (u, v) is less than $1/2$ then the majority of the clusterings place u and v together, while if the weight is greater than $1/2$, the majority places u and v in different clusters. Ideally, we would like to *cut* all edges with weight more than $1/2$, and *not cut* all edges with weight less than $1/2$. The goal in correlation clustering is to find a partition of the vertices of the graph that it cuts as few as possible of the edges with low weight (less than $1/2$), and as many as possible of the edges with high weight (more than $1/2$). In Figure 2, clustering $\mathcal{C} = \{\{v_1, v_3\}, \{v_2, v_4\}, \{v_5, v_6\}\}$ is the optimal clustering.

Clustering aggregation has been previously considered under a variety of names (consensus clustering, clustering ensemble, clustering combination) in a variety of different areas: machine learning [19, 12], pattern recognition [14], bio-informatics [13], and data mining [21, 5]. The problem of correlation clustering is interesting in its own

right, and it has recently attracted a lot of attention in the theoretical computer-science community [2, 6, 8, 10]. We review some of the related literature on both clustering aggregation, and correlation clustering in Section 6.

Our contributions can be summarized as follows.

- We formally define the problem of clustering aggregation, and we demonstrate the connection between clustering aggregation and correlation clustering.
- We present a number of algorithms for clustering aggregation and correlation clustering. We also propose a sampling mechanism that allows our algorithms to handle large datasets. The problems we consider are NP-hard, yet we are still able to provide approximation guarantees for many of the algorithms we propose. For the formulation of correlation clustering we consider we give a combinatorial 3-approximation algorithm, which is an improvement over the best known 9-approximation algorithm.
- We present an extensive experimental study, where we demonstrate the benefits of our approach. Furthermore, we show that our sampling technique reduces the running time of the algorithms, without sacrificing the quality of the clustering.

The rest of this paper is structured as follows. In Section 2 we discuss the various applications of the clustering-aggregation framework, which is formally defined in Section 3. In Section 4 we describe in detail the proposed algorithms for clustering aggregation and correlation clustering, and the sampling-based algorithm that allows us to handle large datasets. Our experiments on synthetic and real datasets are presented in Section 5. Finally, Section 6 contains a review of the related work, and Section 7 is a short conclusion.

2 Applications of clustering aggregation

Clustering aggregation can be applied in various settings. We will now present some of the main applications and features of our framework.

Clustering categorical data: An important application of clustering aggregation is that it provides a very natural method for clustering categorical data. Consider a dataset with tuples t_1, \dots, t_n over a set of categorical attributes A_1, \dots, A_m . The idea is to view each attribute A_j as a way of producing a simple clustering of the data: if A_j contains k_j distinct values, then A_j partitions the data in k_j clusters – one cluster for each value. Then, clustering aggregation considers all those m clusterings produced by the m attributes and tries to find a clustering that agrees as much as possible with all of them.

For example, consider a `Movie` database. Each tuple in the database corresponds to a movie that is defined over a set of attributes such as `Director`, `Actor`, `Actress`, `Genre`, `Year`, etc, some of which take categorical values. Note that each of the categorical attributes defines naturally a clustering. For example, the `Movie.Genre` attribute groups the movies according to their genre, while the `Movie.Director` according to who has directed the movie. The objective is to combine all these clusterings into a single clustering.

Clustering heterogeneous data: The clustering aggregation method can be particularly effective in cases where the data are defined over heterogeneous attributes that contain incomparable values. Consider for example the case that there are many numerical attributes whose units are incomparable (say, `Movie.Budget` and `Movie.Year`) and so it does not make sense to compare numerical vectors directly using an L_p -type distance measure. A similar situation arises in the case where the data contains a mix of categorical and numerical values. In such cases the data can be partitioned vertically into sets of homogeneous attributes, obtain a clustering for each of these sets by applying the appropriate clustering algorithm, and then aggregate the individual clusterings into a single clustering.

Missing values: One of the major problems in clustering, and data analysis in general, is the issue of missing values. These are entries that for some reason (mistakes, omissions, lack of information) are incomplete. The clustering aggregation framework provides several ways for dealing with missing values in categorical data. One approach is to average them out: an attribute that contains a missing value in some tuple does not have any information about how this tuple should be clustered, so we should let the remaining attributes decide. When computing the fraction of clusterings that disagree over a pair of tuples, we only consider the attributes that actually have a value on these tuples.

Another approach, which we adopt in this paper, is to assume that given a pair of tuples for which an attribute contains at least one missing value, the attribute tosses a random coin and with probability p it reports the tuples as being clustered together, while with probability $1 - p$ it reports them as being in separate clusters. Each pair of tuples is treated independently. We are then interested in minimizing the *expected* number of disagreements between the clusterings.

Identifying the correct number of clusters: One of the most important features of the formulation of clustering aggregation is that there is no need to specify the number of clusters in the result. The automatic identification of the appropriate number of clusters is a deep research problem that has attracted lots of attention [16, 18]. For most clustering approaches the quality (likelihood, sum of distances

to cluster centers, etc.) of the solution improves as the number of clusters is increased. Thus, the trivial solution of all singleton clusters is the optimal. There are two ways of handling the problem. The first is to have a hard constraint on the number of clusters, or on their quality. For example, in agglomerative algorithms one can either fix in advance the number of clusters in the final clustering, or impose a bound on the distance beyond which no pair of clusters will be merged. The second approach is to use model selection methods, e.g., Bayesian information criterion (BIC) or cross-validated likelihood [18] to compare models with different numbers of clusters.

The formulation of clustering aggregation takes automatically care of the selection of the number of clusters. If many input clusterings place two objects in the same cluster, then it will not be beneficial for a clustering-aggregation solution to split these two objects. Thus, the solution of all singleton clusters is not a trivial solution for our objective function. Furthermore, if there are k subsets of data objects in the dataset, such that the majority of the input clusterings places them together, and separates them from the rest, then the clustering aggregation algorithm will correctly identify the k clusters, without any prior knowledge of k . Note that in the example in Figures 1 and 2 the optimal solution \mathcal{C} discovers naturally a set of 3 clusters in the dataset.

Indeed, the structure of the objective function ensures that the clustering aggregation algorithms will make their own decisions, and settle naturally to the appropriate number of clusters. As we will show in our experimental section, our algorithms take advantage of this feature and for all our datasets they generate clusterings with very reasonable number of clusters. On the other hand, if the user insists on a predefined number of clusters, most of our algorithms can be easily modified to return that specific number of clusters. For example, the agglomerative algorithm described in Section 4 can be made to continue merging clusters until the predefined number is reached.

Detecting outliers: The ability to detect outliers is closely related with the ability to identify the correct number of clusters. If a node is not close to any other nodes, then, from the point of view of the objective function, it would be beneficial to assign that node in a singleton cluster. In the case of categorical data clustering, the scenarios for detecting outliers are very intuitive: if a tuple contains many uncommon values, it does not participate in clusters with other tuples, and it will be identified as an outlier. Another scenario where it pays off to consider a tuple as an outlier is when the tuple contains common values (and therefore it participates in big clusters in the individual input clusterings) but there is no consensus to a common cluster (for example, a horror movie featuring actress `Julia.Roberts` and directed by the “independent” director `Lars.vonTrier`).

Improving clustering robustness: Different clustering algorithms have different qualities and different shortcomings. Some algorithms might perform well in specific datasets but not in others, or they might be very sensitive to parameter settings. For example, the single-linkage algorithm is good in identifying elongated regions but it is sensitive to clusters being connected with narrow strips of points. The k -means method is a widely-used algorithm, but it is sensitive to clusters of uneven size, and it can get stuck in local optima.

We suggest that by aggregating the results of different clustering we can improve significantly the robustness and the quality of the final clustering. The idea is that different algorithms make different “mistakes” that can be “canceled out” in the final aggregation. Furthermore, for objects that are outliers or noise, it is most likely that there will be no consensus on how they should be clustered, and thus they will be singled out by the aggregation algorithm. The intuition is similar to performing *rank aggregation* for improving the results of web searches [9]. Our experiments indicate that clustering aggregation can improve significantly the results of individual algorithms.

Privacy-preserving clustering: Consider a situation where a database table is vertically split and different attributes are maintained in different sites. Such a situation might arise in cases where different companies or governmental administrations maintain various sets of data about a common population of individuals. For such cases, our method offers a natural model for clustering the data maintained in all sites as a whole in a privacy-preserving manner, that is, without the need for the different sites to reveal their data to each other, and without the need to rely on a trusted authority. Each site clusters its own data independently and then all resulting clusterings are aggregated. The only information revealed is which tuples are clustered together; no information is revealed about data values of any individual tuples.

3 Description of the framework

We begin our discussion of the clustering aggregation framework by introducing our notation. Consider a set of n objects $V = \{v_1, \dots, v_n\}$. A clustering \mathcal{C} of V is a *partition* of V into k disjoint sets C_1, \dots, C_k , that is, $\bigcup_i C_i = V$ and $C_i \cap C_j = \emptyset$, for all $i \neq j$. The k sets C_1, \dots, C_k are the clusters of \mathcal{C} . For each $v \in V$ we use $\mathcal{C}(v)$ to denote the label of the cluster to which the object v belongs. It is $\mathcal{C}(v) = j$ if and only if $v \in C_j$. In the sequel we consider m clusterings: we write \mathcal{C}_i to denote i -th clustering, and k_i for the number of clusters of \mathcal{C}_i .

In the clustering aggregation problem the task is to find a clustering that agrees as much as possible with a number of already-existing clusterings. To make the notion more

precise, we need to define a measure of disagreement between clusterings. Consider first two objects u and v in V . The following simple 0/1 distance function checks if two clusterings \mathcal{C}_1 and \mathcal{C}_2 place u and v in the same clusters.

$$d_{u,v}(\mathcal{C}_1, \mathcal{C}_2) = \begin{cases} 1 & \text{if } \mathcal{C}_1(u) = \mathcal{C}_1(v) \text{ and } \mathcal{C}_2(u) \neq \mathcal{C}_2(v), \\ & \text{or } \mathcal{C}_1(u) \neq \mathcal{C}_1(v) \text{ and } \mathcal{C}_2(u) = \mathcal{C}_2(v), \\ 0 & \text{otherwise.} \end{cases}$$

The distance between two clusterings \mathcal{C}_1 and \mathcal{C}_2 is defined as the number of pairs of objects on which the two clusterings disagree, that is,

$$d_V(\mathcal{C}_1, \mathcal{C}_2) = \sum_{(u,v) \in V \times V} d_{u,v}(\mathcal{C}_1, \mathcal{C}_2).$$

For the distance measure $d_V(\cdot, \cdot)$ one can easily prove is that it satisfies the *triangle inequality* on the space of clusterings. We omit the proof due to space constraints.

Observation 1 *Given a set of objects V and clusterings $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$ on V , we have*

$$d_V(\mathcal{C}_1, \mathcal{C}_3) \leq d_V(\mathcal{C}_1, \mathcal{C}_2) + d_V(\mathcal{C}_2, \mathcal{C}_3)$$

The clustering aggregation problem can now be formalized as follows.

Problem 1 (Clustering aggregation) *Given a set of objects V and m clusterings $\mathcal{C}_1, \dots, \mathcal{C}_m$ on V , compute a new clustering \mathcal{C} that minimizes the total number of disagreements with all the given clusterings, i.e., it minimizes*

$$D(\mathcal{C}) = \sum_{i=1}^m d_V(\mathcal{C}_i, \mathcal{C}).$$

The clustering aggregation problem is also defined in [13], where it is shown to be NP-complete using the results of Barthelemy and Leclerc [3].

The algorithms we propose for the problem of clustering aggregation take advantage of a related formulation, which is known as *correlation clustering* [2]. Formally, correlation clustering is defined as follows.

Problem 2 (Correlation clustering) *Given a set of objects V , and distances $X_{uv} \in [0, 1]$ for all pairs $u, v \in V$, find a partition \mathcal{C} for the objects in V that minimizes the score function*

$$d(\mathcal{C}) = \sum_{\substack{(u,v) \\ \mathcal{C}(u)=\mathcal{C}(v)}} X_{uv} + \sum_{\substack{(u,v) \\ \mathcal{C}(u) \neq \mathcal{C}(v)}} (1 - X_{uv}).$$

Correlation clustering is a generalization of clustering aggregation. Given the m clusterings $\mathcal{C}_1, \dots, \mathcal{C}_m$ as input one can construct an instance of the correlation clustering problem by defining the distances X_{uv} appropriately. In particular, let $X_{uv} = \frac{1}{m} \cdot |\{i \mid 1 \leq i \leq m \text{ and } \mathcal{C}_i(u) \neq \mathcal{C}_i(v)\}|$ be the *fraction* of clusterings that assign the pair (u, v) into *different* clusters. For a candidate solution \mathcal{C} of correlation clustering, if \mathcal{C} places u, v in the same cluster it will disagree with mX_{uv} of the original clusterings, while if \mathcal{C} places u, v in different clusters it will disagree with the remaining $m(1 - X_{uv})$ clusterings. Thus, clustering aggregation can be reduced to correlation clustering, and as a result correlation clustering is also an NP-complete problem. We note that an instance of correlation clustering produced by an instance of clustering aggregation is a restricted version of the correlation clustering problem. It is not hard to show that the values X_{uv} obey the triangle inequality, that is, $X_{uw} \leq X_{uv} + X_{vw}$ for all u, v and w in V .

Since both problems we consider are NP-complete it is natural to seek algorithms with provable approximation guarantees. For the clustering aggregation problem, it is trivial to obtain a 2-approximation solution. The idea is to take advantage of the triangle inequality property of the distance measure $d_V(\cdot, \cdot)$. Assume that we are given m objects in a metric space and we want to find a new point that minimizes the sum of distances from the given objects. Then it is a well known fact that selecting the best among the m original objects yields a factor $2(1 - 1/m)$ approximate solution. For our problem, this method suggests taking as the solution to clustering aggregation the clustering \mathcal{C}_i that minimizes $D(\mathcal{C}_i)$. Despite the small approximation factor, this solution is non-intuitive, and we observed that it does not work well in practice.

The above algorithm cannot be used for the problem of correlation clustering – there are no input clusterings to choose from. In general, the correlation clustering problem we consider is not equivalent to the clustering aggregation. There is an extensive literature in the theoretical computer science community on many different variants of the correlation clustering problem. We review some of these results in Section 6. Our problem corresponds to the weighted correlation clustering problem with linear cost functions [2], where the weights on the edges obey the triangle inequality. The best known approximation algorithm for this case can be obtained by combining the results of Bansal et al. [2] and Charikar et al. [6], giving an approximation ratio of 9. The algorithm does not take into account the fact that the weights obey the triangle inequality.

4 Algorithms

In this section we present the suggested algorithms for clustering aggregation. Most of our algorithms approach

the problem through the correlation-clustering problem, and most of the algorithms are parameter-free.

The BESTCLUSTERING algorithm: This is the trivial algorithm that was already mentioned in the previous section. Given m clusterings $\mathcal{C}_1, \dots, \mathcal{C}_m$, BESTCLUSTERING finds the clustering \mathcal{C}_i from them that minimizes the total number of disagreements $D(\mathcal{C}_i)$. Using the data structures described in [3] the best clustering can be found in time $O(mn)$. As discussed, this algorithm yields a solution with approximation ratio of at most $2(1 - 1/m)$. We can show that this bound is tight, that is, there exists an instance of the clustering aggregation problem, where the algorithm BESTCLUSTERING produces a solution that is exactly $2(1 - 1/m)$ worse than the optimal. The proof of the lower bound appears in the full version of the paper.

The algorithm is specific to clustering aggregation—it cannot be used for correlation clustering. An interesting approach to the correlation clustering problem that might lead to a better approximation algorithm is to reconstruct m clusterings from the input distance matrix $[X_{uv}]$. Unfortunately, we are not aware of a polynomial-time algorithm for this problem of decomposition into clusterings.

The BALLS algorithm: The BALLS algorithm works on the correlation clustering problem, and thus, it takes as input the matrix of pairwise distances X_{uv} . Equivalently, we view the input as a graph whose vertices are the tuples of a dataset, and the edges are weighted by the distances X_{uv} .

The algorithm is defined with an input parameter α , and it is the only algorithm that requires an input parameter. By the theoretical analysis we sketch below, we can set the parameter α to a constant value. However, different values of α can lead to better solutions in practice.

The intuition of the algorithm is to find a set of vertices that are close to each other and far from other vertices. Given such a set, we consider it to be a cluster, we remove it from the graph, and we proceed with the rest of the vertices. The difficulty lies in finding such a set, since in principle any subset of the vertices can be a candidate. We overcome the difficulty by resorting again to the triangle inequality—this time for the distances X_{uv} . In order to find a good cluster we take all vertices that are close (within a “ball”) to a vertex u . The triangle inequality guarantees that if two vertices are close to u , then they are also relatively close to each other. We also note that it is intuitive that good clusters should be ball-shaped: since our cost function penalizes for long edges that are not cut, we do not expect to have elongated clusters in the optimal solution.

More formally the algorithm is described as follows. It first sorts the vertices in increasing order of the total weight of the edges incident on each vertex, a heuristic that we observed to work well in practice. At every step, the algorithm picks the first unclustered node u in that ordering. It

then finds the set of nodes S that are at a distance of at most $1/2$ from the node u , and it computes the average distance $d(u, S)$ of the nodes in S to node u . If $d(u, S) \leq \alpha$ then the nodes in $S \cup \{u\}$ are considered to form a cluster, otherwise, node u forms a singleton cluster.

The cost of BALLS algorithm is guaranteed to be at most 3 times the cost of the optimal clustering. This result improves the previously best-known 9 approximation. The proof is deferred to the full version of the paper, and here we give only a short sketch.

Theorem 1 For $\alpha = \frac{1}{4}$, the approximation ratio of the BALLS algorithm is at most 3.

Sketch of the proof: The proof proceeds by bounding the cost that the algorithm pays for each edge (u, v) in terms of the cost that the optimal algorithm pays for the same edge. We make use of the fact that for an edge of distance c , if the algorithm takes the edge, then it pays at most $c/(1-c)$ times the cost of the optimal, while if it does not take the edge then it pays at most $(1-c)/c$ times the optimal cost. It follows that taking an edge of cost less than $1/2$, or splitting an edge of cost more than $1/2$ is an optimal decision.

Given a node u , if the BALLS algorithm decides to place u into a singleton cluster, then this means that $\sum_{i \in S} X_{ui} \geq \frac{1}{4}|S|$. Therefore, the cost paid by BALLS is at most $\frac{3}{4}|S|$. At the same time, since the weight of all edges (u, i) , for $i \in S$ is less than $1/2$, the optimal algorithm takes all edges and pays at least $\frac{1}{4}|S|$, giving approximation ratio 3.

In the case that the algorithm creates the cluster $C = S \cup \{u\}$ the decision to merge edges (u, i) for $i \in S$, and split the edges (u, j) for $j \notin S$ is optimal. We need to consider the edges X_{ij} where $i, j \in S$, or $i \in S$ and $j \notin S$. We consider the case that $i, j \in S$; the other case is handled symmetrically. We sort the nodes in order of increasing distance from node u , and for each node j we consider the cost of all edges X_{ij} with $i < j$. First, we note that if $X_{uj} \leq \frac{1}{4}$ then $X_{ij} \leq X_{ui} + X_{uj} \leq \frac{1}{2}$, so the decision of the algorithm to take edge (i, j) is optimal. For $X_{uj} > \frac{1}{4}$, let C_j denote the set of nodes i , where $i < j$. A key observation is that $\sum_{i \in C_j} X_{ui} \leq \frac{1}{4}|C_j|$: since we omit only edges with high cost, the average cost cannot be more than $\frac{1}{4}$. Using the triangle inequality we can prove that $\sum_{i \in C_j} X_{ij} \leq \frac{3}{4}|C_j|$, that is, the sum of X_{ij} for $i \in C_j$ is also small. This implies that the optimal algorithm cannot gain a lot by splitting some of the X_{ij} edges. The cost of the optimal can be shown to be at least $\frac{1}{4}|C_j|$, giving again approximation ratio 3. \square

There are special cases where it is possible to prove a better approximation ratio. For the case that $m = 3$ it is easy to show that the cost of the BALLS algorithm is at most 2 times that of the optimal solution. The complexity of the algorithm is $O(mn^2)$ for generating the table and $O(n^2)$ for running the algorithm.

In our experiments we have observed that the value $\frac{1}{4}$ tends to be small as it creates many singleton clusters. For many of our real datasets we have found that $\alpha = \frac{2}{5}$ leads to better solutions.

The AGGLOMERATIVE algorithm: The AGGLOMERATIVE algorithm is a standard bottom-up algorithm for the correlation clustering problem. It starts by placing every node into a singleton cluster. It then proceeds by considering the pair of clusters with the smallest average distance. The average distance between two clusters is defined as the average weight of the edges between the two clusters. If the average distance of the closest pair of clusters is less than $1/2$ then the two clusters are merged into a single cluster. If there are no two clusters with average distance smaller than $1/2$, then no merging of current clusters can lead to a solution with improved cost $d(\mathcal{C})$. Thus, the algorithm stops, and it outputs the clusters it has created so far.

The AGGLOMERATIVE algorithm has the desirable feature that it creates clusters where the average distance of any pair of nodes is at most $1/2$. The intuition is that the opinion of the majority is respected on average. Using this property we are able to prove that when $m = 3$, the AGGLOMERATIVE algorithm produces a solution with cost at most 2 times that of the optimal solution. The complexity of the algorithm is $O(mn^2)$ for creating the matrix plus $O(n^2 \log n)$ for running the algorithm.

The FURTHEST algorithm: The FURTHEST algorithm is a top-down algorithm that works on the clustering correlation problem. It is inspired by the *furthest-first traversal* algorithm, for which Hochbaum and Shmoys [17] showed that it achieves a 2-approximation for the clustering formulation of p -centers. As the BALLS algorithm uses a notion of a “center” to find clusters and repeatedly remove them from the graph, the FURTHEST algorithm uses “centers” to partition the graph in a top-down fashion.

The algorithm starts by placing all nodes into a single cluster. Then it finds the pair of nodes that are furthest apart, and places them into different clusters. These two nodes become the centers of the clusters. The remaining nodes are assigned to the center that incurs the least cost. This procedure is repeated iteratively: at each step a new center is generated that is the furthest from the existing centers, and the nodes are assigned to the center that incurs the least cost. At the end of each step, the cost of the new solution is computed. If it is lower than that of the previous step then the algorithm continues. Otherwise, the algorithm outputs the solution computed in the previous step. The complexity of the algorithm is $O(mn^2)$ for creating the matrix and $O(k^2n)$ for running the algorithm, where k is the number of clusters created.

The LOCALSEARCH algorithm: The LOCALSEARCH algorithm is an application of a local-search heuristic to the

problem of correlation clustering. The algorithm starts with some clustering of the nodes. This clustering could be a random partition of the data, or it could be obtained by running one of the algorithms we have already described. The algorithm then goes through the nodes and it considers placing them into a different cluster, or creating a new singleton cluster with this node. The node is placed in the cluster that yields the minimum cost. The algorithm iterates until there is no move that can improve the cost. The LOCALSEARCH can be used as a clustering algorithm, but also as a post-processing step, to improve upon an existing solution.

When considering a node v , the cost $d(v, C_i)$ of assigning a node v to a cluster C_i is computed as follows.

$$d(v, C_i) = \sum_{u \in C_i} X_{vu} + \sum_{u \in S \cap \overline{C_i}} (1 - X_{vu})$$

The first term is the cost of merging v in C_i , while the second term is the cost of *not* merging node v with the nodes not in C_i . We compute $d(v, C_i)$ efficiently as follows. For every cluster C_i we compute and store the cost $M(v, C_i) = \sum_{u \in C_i} X_{vu}$ and the size of the cluster $|C_i|$. Then the distance of v to C_i is

$$d(v, C_i) = M(v, C_i) + \sum_{j \neq i} (|C_j| - M(v, C_j))$$

The cost of assigning node v to a singleton cluster is $\sum_j (|C_j| - M(v, C_j))$.

The running time of the LOCALSEARCH algorithm, given the distance matrix X_{uv} , is $O(In^2)$, where I is the number of local search iterations until the algorithm converges to a solution for which no better move can be found. Our experiments showed that the LOCALSEARCH algorithm is quite effective, and it improves significantly the solutions found by the previous algorithms. Unfortunately, the number of iterations tends to be large, and thus the algorithm is not scalable to large datasets.

4.1 Handling large datasets

The algorithms we described in Section 4 take as input the distance matrix, so their complexity is quadratic in the number of data objects in the dataset. The quadratic complexity is inherent in the correlation clustering problem, since the input to the problem is a complete graph. Given a node, the decision of placing the node to a cluster has to take into account not only the cost of merging the node to the cluster, but also the cost of *not* placing the node to the other clusters. Furthermore, the definition of the cost function does not allow for an easy summarization of the clusters, a technique that is commonly used in many clustering algorithms. However, the quadratic complexity makes the

algorithms inapplicable to large datasets. We will now describe the algorithm SAMPLING, which uses sampling to reduce the running time of the algorithms.

The SAMPLING algorithm is run on top of the algorithms we described Section 4. The algorithm performs a pre-processing and post-processing step that are linear on the size of the dataset. In the pre-processing step the algorithm samples a set of nodes, S , uniformly at random from the dataset. These nodes are given as input to one of the clustering aggregation algorithms. The output is a set of ℓ clusters $\{C_1, \dots, C_\ell\}$ of the nodes in S . In the post-processing step the algorithm goes through the nodes in the dataset not in S . For every node it decides whether or to place it in one of the existing clusters, or to create a singleton cluster. In order to perform this step efficiently, we use the same technique as for the LOCALSEARCH algorithm. We observed experimentally that at the end of the assignment phase there are too many singleton clusters. Therefore, we collect all singleton clusters and we run the clustering aggregation again on this subset of nodes.

The size of the sample S is determined so that if there is a large cluster in the dataset the sample will contain at least one node from the cluster. Large cluster means a cluster that contains a constant fraction of the nodes in the dataset. Using the Chernoff bounds we can prove that sampling $O(\log n)$ nodes is sufficient to ensure that we will select at least one of the nodes in a large cluster with high probability. Note that although nodes in small clusters may not be selected, these will be assigned in singleton clusters in the post-processing step. When clustering the singletons, they will be clustered together. Since the size of these clusters is small this does not incur a significant overhead on the algorithm.

5 Experimental evaluation

We have conducted extensive experiments to test the quality of the clusterings produced by our algorithms on a varied collection of synthetic and real datasets. Furthermore, for our SAMPLING algorithm, we have experimented with the quality vs. running-time trade off.

5.1 Improving clustering robustness

The goal in this set of experiments is to show how clustering aggregation can be used to improve the quality and robustness of widely used vanilla-clustering algorithms. For the two experiments we are describing we used synthetic datasets of two-dimensional points.

The first dataset is shown in Figure 3. An intuitively good clustering for this dataset consists of the seven perceptually distinct groups of points. We ran five different

clustering algorithms implemented in Matlab: single linkage, complete linkage, average linkage, Ward’s clustering, and k -means. For all of the clusterings we set the number of clusters to be 7, and for the rest parameters, if any, we used Matlab’s defaults. The results for the five clusterings are shown in the first five panels of Figure 3. One sees that all clusterings are imperfect. In fact, the dataset contains features that are known to create difficulties for the selected algorithms, e.g., narrow “bridges” between clusters, uneven-sized clusters, etc. The last panel in Figure 3 shows the results of aggregating the five previous clusterings. The aggregated clustering is better than any of the input clusterings (although average-linkage comes very close), and it shows our intuition of how mistakes in the input clusterings can be “canceled out”.

In our second experiment the goal is to show how clustering aggregation can be used to identify the “correct” clusters, as well as outliers. Three datasets were created as follows: k^* cluster centers were selected uniformly at random in the unit square, and 100 points were generated from the normal distribution around each cluster center. For the three datasets we used $k^* = 3, 5, \text{ and } 7$, respectively. An additional 20% of the total number of points were generated uniformly from the unit square and they were added in the datasets. For each of the three datasets we ran k -means with $k = 2, 3, \dots, 10$, and we aggregated the resulting clusterings, that is, in each dataset we performed clustering aggregation on 9 input clusterings. For lack of space, the input clusterings are not shown; however, most are imperfect. Obviously, when k is too small some clusters get merged, and when k is too large some clusters get split. On the other hand, the aggregated clusterings, for the three datasets, are shown in Figure 4. We see that the main clusters identified are precisely the correct clusters. Small additional clusters that are found contain only points from the background “noise”, and they can be clearly characterized as outliers.

5.2 Clustering categorical data

In this section we use the ideas we discussed in Section 2 for performing clustering of categorical datasets. We used three datasets from the UCI Repository of machine learning databases [4]. The first dataset, **Votes**, contains voting information for 435 persons. For each person there are votes on 16 issues (yes/no vote viewed as categorical values), and a class label classifying a person as republican or democrat. There are a total of 288 missing values. The second dataset, **Mushrooms**, contains information on physical characteristics of mushrooms. There are 8,124 instances of mushrooms, each described by 22 categorical attributes, such as shape, color, odor, etc. There is a class label describing if a mushroom is poisonous or edible, and there are 2,480 missing values in total. Finally,

the third dataset, **Census**, has been extracted from the census bureau database, and it contains demographic information on 32,561 people in the US. There are 8 categorical attributes (such as education, occupation, marital status, etc.) and 6 numerical attributes (such as age, capital gain, etc.). Each person is classified according to whether they receive an annual salary of more than \$50K or less.

For each of the datasets we perform clustering based on the categorical attributes and we compare the clusters using the class labels of the datasets. The intuition is that clusterings with “pure” clusters, i.e., clusters in which all objects have the same class label, are preferable. Thus, if a clustering contains k clusters with sizes s_1, \dots, s_k , and the sizes of the majority class in each cluster are m_1, \dots, m_k , respectively, then we measure the quality of the clustering by the *classification error*, defined as

$$E_C = \frac{\sum_{i=1}^k (s_i - m_i)}{\sum_{i=1}^k s_i} = \frac{\sum_{i=1}^k (s_i - m_i)}{n}.$$

If a clustering has E_C value equal to 0 it means that it contains only pure clusters. Notice that clusterings with many clusters tend to have smaller E_C values—in the extreme case if $k = n$ then $E_C = 0$ since singleton clusters are pure. We remark that the use of the term “classification error” is somehow abusing since no classification is performed in the formal sense, (i.e., using training and test data, cross validation, etc); the classification-error measure is only indicative of the cluster quality, since clustering and classification are two different problems. It is not clear that the best clusters in the dataset correspond to the existing classes. Depending on the application one may be interested in discovering different clusters.

We also run comparative experiments with the categorical clustering algorithm ROCK [15], and the much more recent algorithm LIMBO [1]. ROCK uses the *Jaccard coefficient* to measure tuple similarity, and places a *link* between two tuples whose similarity exceeds a threshold θ . For our experiments, we used values of θ suggested by Guha et al. [15] in the original ROCK paper. LIMBO uses information theoretic concepts to define clustering quality. It clusters together tuples so that the conditional entropy of the attribute values within a cluster is low. For the parameter ϕ of LIMBO we used again values suggested in [1].

The results for the **Votes** and **Mushrooms** datasets are shown in Tables 2 and 3, respectively. Except for the classification error (E_C), we also show the number of clusters of each clustering (k), and the disagreement error (E_D), which is the measure explicitly optimized by our algorithms. Since the clustering aggregation algorithms make their own decisions for the resulting number of clusters, we have run the other two algorithms for the same values of k so that we ensure fairness. Overall the classification errors are comparable with the exception of LIMBO’s impressive performance

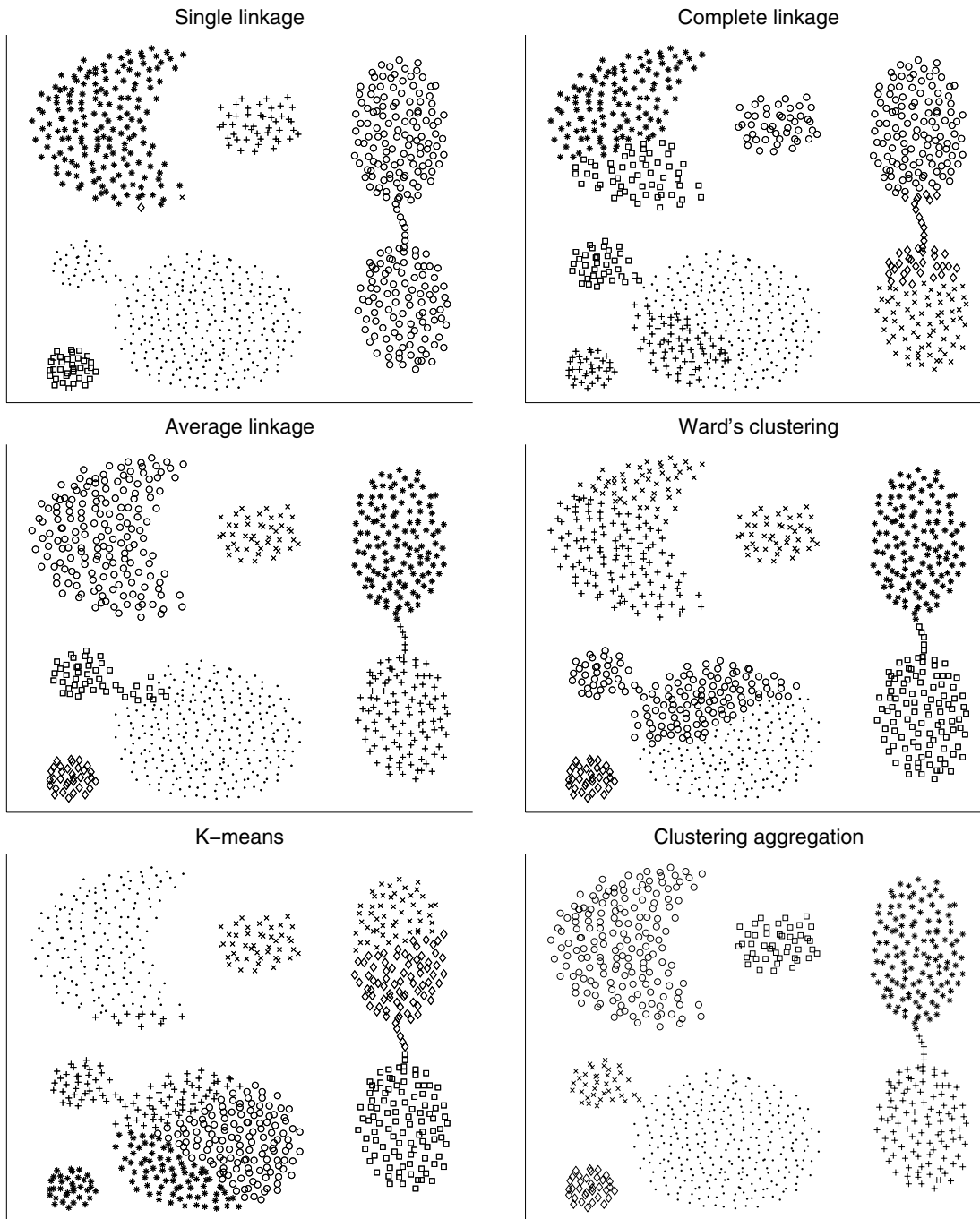


Figure 3. Clustering aggregation on five different input clusterings. To obtain the last plot, which is the result of aggregating the previous five plots, the `AGGLOMERATIVE` algorithm was used.

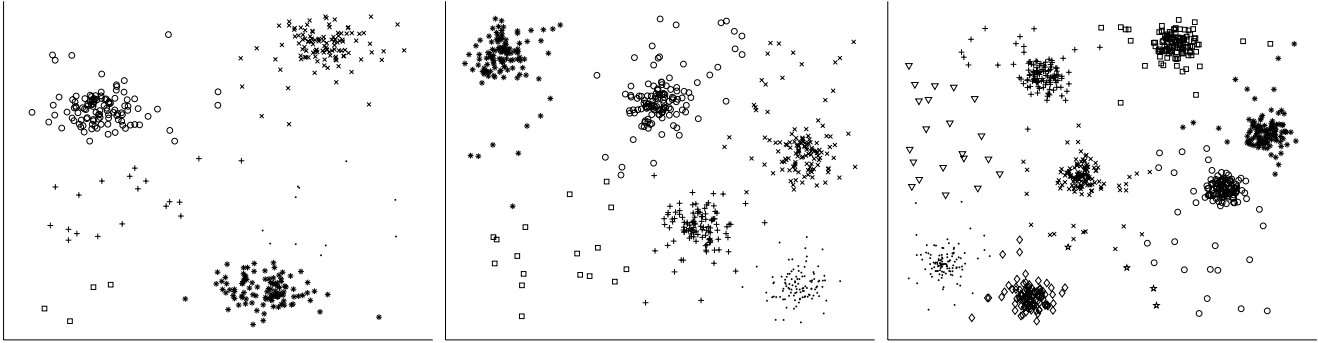


Figure 4. Finding the correct clusters and outliers.

	c_1	c_2	c_3	c_4	c_5	c_6	c_7
Poisonous	808	0	1296	1768	0	36	8
Edible	2864	1056	0	96	192	0	0

Table 1. Confusion matrix for class labels and clusters found by the AGGLOMERATIVE algorithm on Mushrooms dataset.

on **Mushrooms** for $k = 7$ and $k = 9$. Our algorithms achieve the lowest distance error, and classification error that is close to that of the optimal. Furthermore, the attractiveness of the algorithms AGGLOMERATIVE, FURTHEST, and LOCALSEARCH lies in the fact that they are completely parameter-free! Neither a threshold nor the number of clusters need to be specified. The number of clusters discovered by our algorithms seem to be very reasonable choices: for the **Votes** dataset, most people vote according to the official position of their political parties, so having two clusters is natural; for the **Mushrooms** dataset, notice that both ROCK and LIMBO achieve much better quality for the suggested values $k = 7$ and $k = 9$, so it is quite likely that the correct number of clusters is around these values. Indicatively, in Table 1 we present the confusion matrix for the clustering produced by the AGGLOMERATIVE algorithm on the **Mushrooms** dataset.

For the **Census** dataset, clustering aggregation algorithms report about 50-60 clusters. To run clustering aggregation on the **Census** dataset we need to resort to the SAMPLING algorithm. As an indicative result, when the SAMPLING uses the FURTHEST algorithm to cluster a sample of 4,000 persons, we obtain 54 clusters and the classification error is 24%. ROCK does not scale for a dataset of this size, while LIMBO with parameters $k = 2$ and $\phi = 1.0$ gives classification error 27.6%. For contrasting these numbers, we mention that supervised classification algorithms (like decision trees and Bayes classifiers) yield classification error between 14 and 21%—but again, clustering is

	k	$E_C(\%)$	E_D
Class labels	2	0	34,184
Lower bound			28,805
BESTCLUSTERING	3	15.1	31,211
AGGLOMERATIVE	2	14.7	30,408
FURTHEST	2	13.3	30,259
BALLS $_{\alpha=0.4}$	2	13.3	30,181
LOCALSEARCH	2	11.9	29,967
ROCK $_{k=2, \theta=0.73}$	2	11	32,486
LIMBO $_{k=2, \phi=0.0}$	2	11	30,147

Table 2. Results on Votes dataset.

a conceptually different task than classification. We visually inspected the smallest of the 54 different clusters, and many corresponded to distinct social groups, for example, male Eskimos occupied with farming-fishing, married Asian-Pacific islander females, unmarried executive-manager females with high-education degrees, etc. We omit a more detailed report due to lack of space.

5.3 Handling large datasets

In this section we describe our experiments with the SAMPLING algorithm that allows us to apply clustering aggregation to large datasets. First we use the **Mushrooms** dataset to experiment with the behavior of our algorithms as a function of the sample size. As we saw in Table 3, the number of clusters found with the non-sampling algorithms is around 10. When sampling is used, the number of clusters found in the sample remains close to 10. For small sample size, clustering the sample is relatively fast compared to the post-processing phase of assigning the non-sampled points to the best cluster, and the overall running time of the SAMPLING algorithm is linear. In Figure 5 (left), we plot the running time of the SAMPLING algorithm, as a fraction of the running time of the non-sampling algorithm, and we show how it changes as we increase the sample size. For a sam-

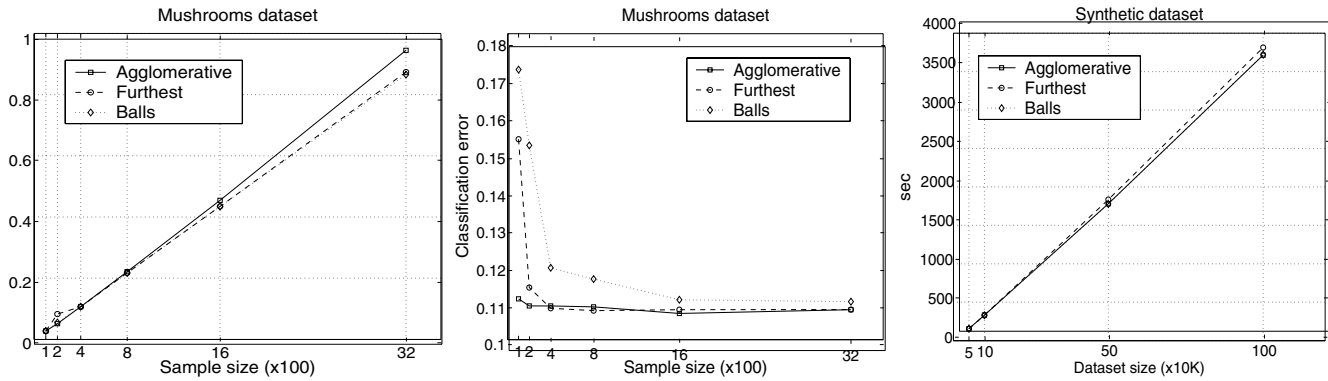


Figure 5. Scalability Experiments

	k	$E_C(\%)$	E_D
Class labels	2	0	13.537 M
Lower bound			8.388 M
BESTCLUSTERING	5	35.4	8.542 M
AGGLOMERATIVE	7	11.1	9.990 M
FURTHEST	9	10.4	10.169 M
BALLS $_{\alpha=0.4}$	10	14.2	11.448 M
LOCALSEARCH	10	10.7	9.929 M
ROCK $_{k=2, \theta=0.8}$	2	48.2	16.777 M
ROCK $_{k=7, \theta=0.8}$	7	25.9	10.568 M
ROCK $_{k=9, \theta=0.8}$	9	9.9	10.312 M
LIMBO $_{k=2, \phi=0.3}$	2	10.9	13.011 M
LIMBO $_{k=7, \phi=0.3}$	7	4.2	10.505 M
LIMBO $_{k=9, \phi=0.3}$	9	4.2	10.360 M

Table 3. Results on Mushrooms dataset.

ple of size 1600 we achieve more than 50% reduction of the running time. At the same time, the classification error of the algorithm converges very fast to the error of the non-sampling algorithms. This is shown in Figure 5 (middle). For sample size 1600 we have almost the same classification error, with only half of the running time.

We also measured the running time of the SAMPLING algorithm for large synthetic datasets. We repeated the configuration of the experiments shown in Figure 4 but on a larger scale. Each dataset consists of points generated from clusters normally distributed around five centers plus an additional 20% of uniformly distributed points. We generate datasets of sizes 50K, 100K, 500K, and 1M points. We then cluster the points using Matlab’s k -means for $k = 2, \dots, 10$, and we run SAMPLING clustering aggregation on the resulting 9 clusterings. The results are shown in Figure 5 (right). These results are for sample size equal to 1000. Once again, the five correct clusters were identified in the sample, and the running time is dominated by the time to assign the non-sampled points in the clusters of the

sample, resulting to the linear behavior shown in the figure.

6 Related Work

A source of motivation for our work is the literature on comparing and merging multiple rankings [9, 11]. Dwork et al. [9] demonstrated that combining multiple rankings in a meta-search engine for the Web yields improved results and removes noise (spam). The intuition behind our work is similar. By combining multiple clusterings we improve the clustering quality, and remove noise (outliers).

The problem of clustering aggregation has been previously considered in the machine learning community, under the name *Clustering Ensemble* and *Consensus Clustering*. Strehl and Ghosh [19] consider various formulations for the problem, most of which reduce the problem to a hyper-graph partitioning problem. In one of their formulations they consider the same graph as in the correlation clustering problem. The solution they propose is to compute the best k -partition of the graph, which does not take into account the penalty for merging two nodes that are far apart. All of their formulations assume that the correct number of clusters is given as a parameter to the algorithm.

Fern and Brodley [12] apply the clustering aggregation idea to a collection of soft clusterings they obtain by random projections. They use an agglomerative algorithm similar to ours, but again they do not penalize for merging dissimilar nodes. Fred and Jain [14] propose to use a single linkage algorithm to combine multiple runs of the k -means algorithm. Dana Cristofor and Dan Simovici [7] observe the connection between clustering aggregation and clustering of categorical data. They propose information theoretic distance measures, and they propose genetic algorithms for finding the best aggregation solution. Boulis and Ostendorf [5] use Linear Programming to discover a correspondence between the labels of the individual clusterings and those of an “optimal” meta-clustering. Topchy et al [21] define clustering

aggregation as a maximum likelihood estimation problem, and they propose an EM algorithm for finding the consensus clustering. Filkov and Skiena [13] consider the same distance measure between clusterings as ours. They propose a simulating annealing algorithm for finding an aggregate solution, and a local search algorithm similar to ours. They consider the application of clustering aggregation to the analysis of microarray data.

There is an extensive literature in the field of theoretical computer science for the problem of correlation clustering. The problem was first defined by Bansal et al. [2]. In their definition, the input is a complete graph with +1 and -1 weights on the edges. The objective is to partition the nodes of the graph so as to minimize the number of positive edges that are cut, and the number of negative edges that are not cut. The best known approximation algorithm for this problem is by Charikar et al. [6] who give an LP-based algorithm that achieves a 4 approximation factor. When the edge weights are arbitrary, the problem is equivalent to the multi-way cut, and thus there is a tight $O(\log n)$ -approximation algorithm [8, 10]. If one considers the corresponding maximization problem, that is, maximize the agreements rather than minimize disagreements, then the situation is much better. Even in the case of graphs with arbitrary edge weights there is a 0.76-approximation algorithm using semi-definite programming [6, 20].

7 Concluding remarks

In this paper we proposed a novel approach to clustering based on the concept of aggregation. Simply stated, the idea is to cluster a set of objects by trying to find a clustering that agrees as much as possible with a number of already-existing clusterings. We motivated the problem by describing in detail various applications of clustering aggregation including clustering categorical data, dealing with heterogeneous data, improving clustering robustness, and detecting outliers. We formally defined the problem and we showed its connection with the problem of correlation clustering. For the problems of clustering aggregation and correlation clustering we gave a number of algorithms, including a sampling algorithm that allows us to handle large datasets with no significance loss in the quality of the solutions. Finally, we verified the intuitive appeal of the proposed approach and we studied the behavior of our algorithms with experiments on real and synthetic datasets.

Acknowledgments We thank Alexander Hinneburg for many useful comments, and Periklis Andritsos for pointers to code, data, and his help with the experiments.

References

- [1] P. Andritsos, P. Tsaparas, R. J. Miller, and K. C. Sevcik. LIMBO: Scalable clustering of categorical data. In *EDBT*, 2004.
- [2] N. Bansal, A. Blum, and S. Chawla. Correlation clustering. In *FOCS*, 2002.
- [3] J.-P. Barthelemy and B. Leclerc. The median procedure for partitions. *DIMACS Series in Discrete Mathematics*, 1995.
- [4] C. L. Blake and C. J. Merz. UCI repository of machine learning databases, 1998.
- [5] C. Bouslis and M. Ostendorf. Combining multiple clustering systems. In *PKDD*, 2004.
- [6] M. Charikar, V. Guruswami, and A. Wirth. Clustering with qualitative information. In *FOCS*, 2003.
- [7] D. Cristofor and D. A. Simovici. An information-theoretical approach to genetic algorithms for clustering. Technical Report TR-01-02, UMass/Boston, 2001.
- [8] E. D. Demaine and N. Immerlica. Correlation clustering with partial information. In *APPROX*, 2003.
- [9] C. Dwork, R. Kumar, M. Naor, and D. Sivakumar. Rank aggregation methods for the Web. In *WWW*, 2001.
- [10] D. Emanuel and A. Fiat. Correlation clustering: Minimizing disagreements on arbitrary weighted graphs. In *ESA*, 2003.
- [11] R. Fagin, R. Kumar, and D. Sivakumar. Comparing top k lists. In *SODA*, 2003.
- [12] X. Z. Fern and C. E. Brodley. Random projection for high dimensional data clustering: A cluster ensemble approach. In *ICML*, 2003.
- [13] V. Filkov and S. Skiena. Integrating microarray data by consensus clustering. In *International Conference on Tools with Artificial Intelligence*, 2003.
- [14] A. Fred and A. K. Jain. Data clustering using evidence accumulation. In *ICPR*, 2002.
- [15] S. Guha, R. Rastogi, and K. Shim. ROCK: A robust clustering algorithm for categorical attributes. *Information Systems*, 25(5):345–366, 2000.
- [16] G. Hamerly and C. Elkan. Learning the k in k -means. In *NIPS*. 2003.
- [17] D. Hochbaum and D. Shmoys. A best possible heuristic for the k -center problem. *Mathematics of Operations Research*, pages 180–184, 1985.
- [18] P. Smyth. Model selection for probabilistic clustering using cross-validated likelihood. *Statistics and Computing*, 10(1):63–72, 2000.
- [19] A. Strehl and J. Ghosh. Cluster ensembles —A knowledge reuse framework for combining multiple partitions. *Journal of Machine Learning Research*, 2002.
- [20] C. Swamy. Correlation clustering: maximizing agreements via semidefinite programming. In *SODA*, 2004.
- [21] A. Topchy, A. K. Jain, and W. Punch. A mixture model of clustering ensembles. In *SDM*, 2004.