

Achieving Anonymity via Clustering

Gagan Aggarwal¹

Tomás Feder²

Krishnaram Kenthapadi²

Samir Khuller³

Rina Panigrahy^{2,4}

Dilys Thomas²

An Zhu¹

ABSTRACT

Publishing data for analysis from a table containing personal records, while maintaining individual privacy, is a problem of increasing importance today. The traditional approach of de-identifying records is to remove identifying fields such as social security number, name etc. However, recent research has shown that a large fraction of the US population can be identified using non-key attributes (called quasi-identifiers) such as date of birth, gender, and zip code [15]. Sweeney [16] proposed the k -anonymity model for privacy where non-key attributes that leak information are suppressed or generalized so that, for every record in the modified table, there are at least $k-1$ other records having exactly the same values for quasi-identifiers. We propose a new method for anonymizing data records, where quasi-identifiers of data records are first clustered and then cluster centers are published. To ensure privacy of the data records, we impose the constraint that each cluster must contain no fewer than a pre-specified number of data records. This technique is more general since we have a much larger choice for cluster centers than k -Anonymity. In many cases, it lets us release a lot more information without compromising privacy. We also provide constant-factor approximation algorithms to come up

with such a clustering. This is the first set of algorithms for the anonymization problem where the performance is independent of the anonymity parameter k . We further observe that a few outlier points can significantly increase the cost of anonymization. Hence, we extend our algorithms to allow an ϵ fraction of points to remain unclustered, i.e., deleted from the anonymized publication. Thus, by not releasing a small fraction of the database records, we can ensure that the data published for analysis has less distortion and hence is more useful. Our approximation algorithms for new clustering objectives are of independent interest and could be applicable in other clustering scenarios as well.

Categories and Subject Descriptors

H.2.8 [Database Applications]: Data mining;

H.3.3 [Information Search and Retrieval]: Clustering

General Terms

Algorithms

Keywords

Privacy, Anonymity, Clustering, Approximation Algorithms

¹Google Inc., Mountain View, CA 94043. This work was done when the authors were Computer Science PhD students at Stanford University. {gagan,anzhu}@cs.stanford.edu

²Department of Computer Science, Stanford University, Stanford, CA 94305. Supported in part by NSF Grant ITR-0331640. This work was also supported in part by TRUST (The Team for Research in Ubiquitous Secure Technology), which receives support from the National Science Foundation (NSF award number CCF-0424422) and the following organizations: Cisco, ESCHER, HP, IBM, Intel, Microsoft, ORNL, Qualcomm, Pirelli, Sun and Symantec. {tomas,kngk,rinap,dilys}@cs.stanford.edu

³Department of Computer Science, University of Maryland, College Park, MD 20742. Supported by NSF Award CCF-0430650. samir@cs.umd.edu

⁴Supported in part by Stanford Graduate Fellowship.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PODS'06, June 26–28, 2006, Chicago, Illinois, USA.

Copyright 2006 ACM 1-59593-318-2/06/0003 ...\$5.00.

1. INTRODUCTION

With the rapid growth in database, networking, and computing technologies, a large amount of personal data can be integrated and analyzed digitally, leading to an increased use of data-mining tools to infer trends and patterns. This has raised universal concerns about protecting the privacy of individuals [17].

Combining data tables from multiple data sources allows us to draw inferences which are not possible from a single source. For example, combining patient data from multiple hospitals is useful to predict the outbreak of an epidemic. The traditional approach of releasing the data tables without breaching the privacy of individuals in the table is to de-identify records by removing the identifying fields such as name, address, and social security number. However, joining this de-identified table with a publicly available database (like the voters database) on columns like race, age, and zip code can be used to identify individuals. Recent research [15] has shown that for 87% of the population in the United States, the combination of non-key fields like date of birth, gender, and zip code corresponds to a unique person. Such non-key fields are called *quasi-identifiers*. In what follows we assume that the identifying fields have been removed

Age	Location	Disease
α	β	Flu
$\alpha + 2$	β	Flu
δ	$\gamma + 3$	Hypertension
δ	γ	Flu
δ	$\gamma - 3$	Cold

(a) Original table

Age	Location	Disease
*	β	Flu
*	β	Flu
δ	*	Hypertension
δ	*	Flu
δ	*	Cold

(b) 2-anonymized version

Age	Location	NumPoints	Disease
$\alpha + 1$	β	2	Flu Flu
δ	γ	3	Hypertension Flu Cold

(c) 2-gather clustering, with maximum radius 3

Age	Location	NumPoints	Radius	Disease
$\alpha + 1$	β	2	1	Flu Flu
δ	γ	3	3	Hypertension Flu Cold

(d) 2-cellular clustering, with total cost 11

Figure 1: Original table and three different ways of achieving anonymity

and that the table has two types of attributes: (1) the quasi-identifying attributes explained above and (2) the sensitive attributes (such as disease) that need to be protected.

In order to protect privacy, Sweeney [16] proposed the k -Anonymity model, where some of the quasi-identifier fields are suppressed or generalized so that, for each record in the modified table, there are at least $k - 1$ other records in the modified table that are identical to it along the quasi-identifying attributes. For the table in Figure 1(a), Figure 1(b) shows a 2-anonymized table corresponding to it. The columns corresponding to sensitive attributes, like disease in this example, are retained without change. The aim is to provide a k -anonymized version of the table with the minimum amount of suppression or generalization of the table entries. There has been a lot of recent work on k -anonymizing a given database table [3, 12]. An $O(k \log k)$ approximation algorithm was first proposed for the problem of k -Anonymity with suppressions only [14]. This was recently improved to an $O(k)$ approximation for the general version of the problem [1].

In this paper, instead of generalization and suppression, we propose a new technique for anonymizing tables before their release. We first use the quasi-identifying attributes to define a metric space (i.e., pairwise distances satisfying the triangle inequality) over the database records, which are then viewed as points in this space. This is similar to the approach taken in [5], except that we do not restrict ourselves to points in \mathcal{R}^d ; instead, we allow our points to be in an arbitrary metric space. We then cluster the points and publish only the final cluster centers along with some cluster size and radius information. Our privacy requirement is similar to the k -Anonymity framework – we require each cluster to have at least r points¹. Publishing the cluster centers instead of the individual records, where each cluster represents at least r records, gives privacy to individual records, but at the same time allows data-mining tools to infer macro trends from the database.

¹We use r instead of k , as k is traditionally used in clustering to denote the number of clusters.

In the rest of the paper we will assume that a metric space has been defined over the records, using the quasi-identifying attributes. For this, the quasi-identifying attributes may need to be remapped. For example, zip codes could first be converted to longitude and latitude coordinates to give a meaningful distance between locations. A categorical attribute, i.e., an attribute that takes n discrete values, can be represented by n equidistant points in a metric space. Furthermore, since the values of different quasi-identifying attributes may differ by orders of magnitude, we need to weigh the attributes appropriately while defining the distance metric. For example, the attribute location may have values that differ in orders of 10 miles with a maximum of 1000 miles, while the attribute age may differ by a single year with a maximum of 100 years. In this case we assume that the attribute location is divided by 10 and the attribute age retained without change if both attributes are needed to have the same relative importance in the distance metric. For the example we provide in Figure 1, we assume that the quasi-identifying attributes have already been scaled. As we see above, it is quite complicated to algorithmically derive a metric space over quasi-identifying attributes of records; we do not pursue it any further in this paper and leave it for future work.

To publish the clustered database, we publish three types of features for each cluster: (1) the quasi-identifying attribute values for the cluster center (age and location in our example), (2) the number of points within the cluster, and (3) a set of values taken by the sensitive attributes (disease in our example). We’ll also publish a measure of the quality of the clusters. This will give a bound on the error introduced by the clustering.

In this paper we consider two cluster-quality measures. The first one is the maximum cluster radius. For this we define the r -GATHER problem, which aims to minimize the maximum radius among the clusters, while ensuring that each cluster has at least r members. As an example, r -GATHER clustering with minimum cluster size $r = 2$, applied to the table in Figure 1(a) gives the table in Figure 1(c). In this

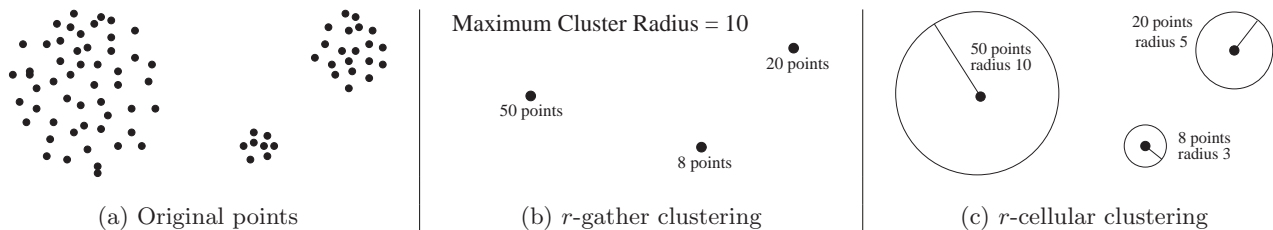


Figure 2: Publishing anonymized data

example, the maximum radius over all clusters is 3. As another example, Figure 2(b) gives the output of the r -GATHER algorithm applied to the quasi-identifiers, shown as points in a metric space in Figure 2(a). Our formulation of the r -GATHER problem is related to, but not to be confused with, the classic k -CENTER problem [8]. The k -CENTER problem has the same objective of minimizing the maximum radius among the clusters, however, the constraint is that we can have no more than k clusters in total. The r -GATHER problem is different from k -CENTER problem in that instead of specifying an upper bound on the number of clusters, we specify a lower bound on the number of points per cluster as part of the input. It's also worth noting that the constraint of at least r points per cluster implies that we can have no more than n/r number of clusters, where n is the total number of points in our data set.

We also consider a second (more verbose) candidate for indicating cluster-quality, whereby we publish the radius of each cluster, rather than just the maximum radius among all clusters. For each point within a cluster, the radius of the cluster gives an upper bound on the distortion error introduced. Minimizing this distortion error over all points leads to the *cellular* clustering measurement that we introduce in this paper. More formally, the *cellular* clustering measurement over a set of clusters, is the sum, over all clusters, of the products of the number of points in the cluster and the radius of the cluster. Using this as a measurement for anonymizing tables, we define the r -CELLULAR CLUSTERING problem as follows: Given points in a metric space, the goal is to partition the points into cells, a.k.a. clusters, each of size at least r , and the cellular clustering measurement is minimized. Consider again the data in Figure 1(a). Figure 1(d) shows a r -cellular cluster solution with minimum cluster size $r = 2$. The total cost is $2 \times 1 + 3 \times 3 = 11$. Also, Figure 2(c) gives the output of the r -CELLULAR CLUSTERING algorithm applied to the quasi-identifiers shown as points in a metric space in Figure 2(a). The total cost of the solution in Figure 2(c) is: $50 \times 10 + 20 \times 5 + 8 \times 3 = 624$. As this cellular clustering objective could be relevant even in contexts other than anonymity, we study a slightly different version of the problem: similar to the FACILITY LOCATION problem [9], we add an additional setup cost for each potential cluster center, associated with opening a cluster centered at that point, but we don't have the lower bound on number of points per cluster. We call this the CELLULAR CLUSTERING problem. In fact, we will use the setup costs in the CELLULAR CLUSTERING problem formulation to help us devise an algorithm that solves r -CELLULAR CLUSTERING.

Comparison with k -Anonymity. While k -Anonymity

forces one to suppress or generalize an attribute value even if all but one of the records in a cluster have the same value, the above clustering-based anonymization technique allows us to pick a cluster center whose value along this attribute dimension is the same as the common value, thus enabling us to release more information without losing privacy. For example, consider the table in Figure 3 with the Hamming distance metric on the row vectors. If we wanted to achieve 5-Anonymity, we will have to hide all the entries in the table, resulting in a total distortion of 20. On the other hand, a 5-CELLULAR CLUSTERING solution could use (1, 1, 1, 1) as the cluster center with a cluster radius of 1. This will give a total distortion bound of 5 (the actual distortion is only 4).

	Attr1	Attr2	Attr3	Attr4
Record 0	1	1	1	1
Record 1	0	1	1	1
Record 2	1	0	1	1
Record 3	1	1	0	1
Record 4	1	1	1	0

Figure 3: A sample table where there is no common attribute among all entries.

Just like k -Anonymity, r -GATHER and r -CELLULAR CLUSTERING is sensitive to outlier points, with just a few outliers capable of increasing the cost of the clustering significantly. To deal with this problem, we generalize the above algorithms to allow an ϵ fraction of the points to be deleted before publication. By not releasing a small fraction of the database records, we can ensure that the data published for analysis has less distortion and hence is more useful. This can be done as long as our aim is to infer macro trends from the published data. On the other hand, if the goal is to find out anomalies, then we should not ignore the outlier points. There has been no previous work for k -Anonymity with this generalization.

We note that, as in k -Anonymity, the objective function is oblivious to the sensitive attribute labels. Extensions to the k -Anonymity model, like the notion of l -diversity [13], can be applied independently to our clustering formulation.

We provide constant-factor approximation algorithms for both the r -GATHER and r -CELLULAR CLUSTERING problems. In particular, we first show that it is NP-hard to approximate the r -GATHER problem better than 2 and provide a matching upper bound. We then provide extensions of both these algorithms to allow for an ϵ fraction of unclustered points, which we call the (r, ϵ) -GATHER and (r, ϵ) -CELLULAR CLUSTERING, respectively. These are the

first constant-factor approximation algorithms for publishing an anonymized database. The best known algorithms [1, 14] for previous problem formulations had an approximation ratio linear in the anonymity parameter r .

The rest of the paper is organized as follows. First, in Section 2, we present a tight 2-approximation algorithm for the r -GATHER problem and its extension to the (r, ϵ) -GATHER problem. In Section 3, motivated by the desire to reduce the sum of the distortions experienced by the points, we introduce the problem of CELLULAR CLUSTERING. We present a primal-dual algorithm for the problem without any cluster-size constraints that achieves an approximation ratio of 4. We then study the additional constraint of having a minimum cluster size of r . Finally, we relax the problem by allowing the solution to leave at most an ϵ fraction of the points unclustered. We conclude in Section 4.

2. r -GATHER CLUSTERING

To publish the clustered database, we publish three types of features for each cluster: (1) the quasi-identifying attribute values for the cluster center, (2) the number of points within the cluster, and (3) a set of values taken by the sensitive attributes. The maximum cluster radius is also published to give a bound on the error introduced by clustering. This is similar to the traditionally studied k -CENTER clustering. In order to ensure r -Anonymity, we don't restrict the total number of clusters, instead, we pose the alternative restriction that each cluster should have at least r records assigned to it. We call this problem r -GATHER, which we formally define below.

DEFINITION 2.1. *The r -GATHER problem is to cluster n points in a metric space into a set of clusters, such that each cluster has at least r points. The objective is to minimize the maximum radius among the clusters.*

We note that the minimum cluster size constraint has been considered earlier in the context of facility location [10].

We first show the reduction for NP-completeness and hardness proofs.

2.1 Lower Bound

We show that this problem is NP-complete by a reduction from the 3-Satisfiability problem, where each literal belongs to at most 3 clauses [6].

Suppose that we have a boolean formula \mathcal{F} in 3-CNF form with m clauses and n variables. Let $\mathcal{F} = C_1 \wedge \dots \wedge C_m$, be a formula composed of variables $x_i, i = 1 \dots n$ and their complements \bar{x}_i .

From the boolean formula, we create a graph $G = (V, E)$ with the following property: There is a solution to the r -GATHER problem with a cluster radius of 1, with respect to the shortest distance metric on the graph G , if and only if \mathcal{F} has a satisfying assignment.

We create the graph as follows: For each variable x_i , create two vertices v_i^T and v_i^F , and create an edge (v_i^T, v_i^F) between

the two vertices; in addition create a set S_i of $(r - 2)$ nodes and add edges from each node in S_i to both v_i^T and v_i^F . Picking v_i^T (v_i^F) as a center corresponds to setting $x_i = T$ (F). (Note that we cannot choose both v_i^T and v_i^F since there are not enough nodes in S_i .) For each clause C_j , create a new node u_j that is adjacent to the nodes corresponding to the literals in the clause. For example, if $C_1 = (x_1 \vee \bar{x}_2)$ then we add edges from u_1 to v_1^T and v_2^F .

If the formula is indeed satisfiable, then there is a clustering by picking v_i^T as a center if $x_i = T$ and picking v_i^F otherwise. Each clause is true, and must have a neighbor chosen as a center. Moreover by assigning S_i to the chosen center, we ensure that each center has at least r nodes in its cluster.

Now suppose there is an r -gather clustering. If $r > 6$ then both v_i^T and v_i^F cannot be chosen as centers. In addition, the clause nodes u_j have degree at most 3 and cannot be chosen as centers. If exactly one of v_i^T or v_i^F is chosen as a center, then we can use this to find the satisfying assignment. The assignment is satisfying as each clause node has some neighbor at distance 1 that is a chosen center, and makes the clause true.

This completes the NP-completeness proof. Note that this reduction also gives us a hardness of 2. We just showed that there is a solution to the r -GATHER problem with a cluster radius of 1 if and only if \mathcal{F} had a satisfying assignment. The next available cluster radius is 2 in the metric defined by the graph G .

2.2 Upper Bound

We first use the threshold method used for k -CENTER clustering to guess R , the optimal radius for r -GATHER. The choices for R are defined as follows. We will try all values $\frac{1}{2}d_{ij}$ where d_{ij} is the distance between points i and j . Note that this defines a set of $O(n^2)$ distance values. We find the smallest R for which the following two conditions hold:

Condition (1) Each point p in the database should have at least $r - 1$ other points within distance $2R$ of p .

Condition (2) Let all nodes be unmarked initially. Consider the following procedure: Select an arbitrary point p as a center. If there are at least r unmarked points within distance $2R$ of p (including p), then form a cluster $C(p)$ with center at p , and mark all points within distance $2R$ of p (including p). This set of points is now covered. Repeat this as long as possible, until all points are marked. Now we try to reassign points to clusters to meet the requirement that each cluster have size at least r . This is done as follows. Create a flow network as follows. Create a source s and sink t . Let C be the set of centers that were chosen. Add edges with capacity r from s to each node in C . Add an edge of unit capacity from a node $c \in C$ to a node $v \in V$ if their distance is at most $2R$. Add edges of unit capacity from nodes in V to t and check to see if a flow of value $r|C|$ can be found (saturating all the edges out of s). If so, then we can obtain the clusters by choosing the nodes to which r units of flow are sent by a node $c \in C$. All remaining nodes of V can be

assigned to any node of C that is within distance $2R$. If no such flow exists, we exit with failure.

The following lemma guarantees that the smallest R that satisfies these conditions is a lower bound on the value of the optimal solution for r -GATHER. Suppose we have an optimal clustering S_1, \dots, S_ℓ with ℓ clusters. Let the maximum diameter of any of these clusters be d^* (defined as the maximum distance between any pair of points in the same cluster).

LEMMA 2.1. *When we try $R = \frac{d^*}{2}$, then the above two conditions are met.*

PROOF. By the definition of r -GATHER, every point has at least $r - 1$ other points within the optimal diameter, and hence within distance $2R$. Consider an optimal r -GATHER clustering. For each point i , all points belonging to the same optimal cluster c as the point i are within a distance $2R$ of i . Thus, in the procedure of Condition (2), as soon as any point in c is selected to open a new cluster, all remaining points belonging to c get assigned to this new cluster. So at most one point from each optimal cluster is chosen as a center and forms a new cluster. We would now like to argue that the assignment phase works correctly as well. Let S be the set of chosen centers. Now consider an optimal solution with clusters, each of size at least r . We can assign each point of a cluster to the center that belongs to that cluster, if a center was chosen in the cluster. Otherwise, since the point was marked by the algorithm, some center was chosen that is within distance $2R$. We can assign it to the center that marked it covered. Each chosen center will have at least r points assigned to it (including itself). \square

Since we find the smallest R , we will ensure that $R \leq d^*/2 \leq R^*$ where R^* is the radius of the optimal clustering. In addition, our solution has radius $2R$. This gives us a 2-approximation.

THEOREM 2.2. *There exists a polynomial time algorithm that produces a 2-approximation to the r -GATHER problem.*

2.3 (r, ϵ) -Gather Clustering

A few outlier points can significantly increase the clustering cost under the minimum cluster size constraint. We consider a relaxation whereby the clustering solution is allowed to leave an ϵ fraction of the points unclustered, i.e., to delete an ϵ fraction of points from the published k -anonymized table. Charikar et al. [4] studied various facility location problems with this relaxation and gave constant-factor approximation algorithms for them.

For the (r, ϵ) -GATHER problem, where each cluster is constrained to have at least r points and an ϵ fraction of the points are allowed to remain unclustered, we modify our r -GATHER algorithm to achieve a 4-approximation. We re-define the condition to find R . We find the smallest R that satisfies the following condition: There should be a subset S of points containing at least $1 - \epsilon$ fraction of the points,

such that each point in S has at least $r - 1$ neighbors within distance $2R$ in S .

This condition can be checked in $O(n^2)$ time by repeatedly removing any point in S that has fewer than $r - 1$ other points in S within distance $2R$ of itself, with S initially being the entire vertex set. It is clear that the smallest R we found is no more than R^* , the optimal radius.

Let R be the value that we found. Let $N(v)$ denote the set of points in G within distance $2R$ of v , including v itself. We know then $N(v) \geq r$. We then consider the following procedure: Select an arbitrary point v from G . If there are at least $r - 1$ other points within distance $2R$ of p , then form a new cluster and assign p and all points within distance $2R$ of p to this cluster. Remove all these points from further consideration and repeat this process until all remaining points have fewer than $r - 1$ other points within distance $2R$ of them. Let U be the set of points left unclustered at the end of this process. For each $u \in U$, there exists a point $p \in N(u)$ such that p is assigned to some cluster c in the procedure of forming clusters. We can see this as follows. Since u was left unassigned at the end of the procedure, there are fewer than r unassigned points remaining in $N(u)$. This implies that there is at least one point $p \in N(u)$ which is already assigned to some cluster c . We assign u to c , which already has at least r points.

Thus, we have assigned all points to clusters, such that each cluster has at least r points. Note that the radius of each cluster is no more than $4R$. This gives us the following theorem.

THEOREM 2.3. *There exists a polynomial time algorithm that produces a 4-approximation to the (r, ϵ) -GATHER problem.*

We note that in the problem formulation of (r, ϵ) -GATHER, if we require the cluster centers to be input points, instead of arbitrary points in the metric, then we can improve the approximation factor to 3. The details are deferred to the journal version.

2.4 Combining r -Gather with k -Center

We can combine the r -GATHER problem with the k -CENTER problem and have the two constraints present at the same time. That is, we minimize the maximum radius, with the constraint that we have no more than k clusters, each must have at least r members. We call this the (k, r) -CENTER problem.

It is worth mentioning that a similar problem has been studied before in the k -CENTER literature. That is, instead of having a lower bound r on the cluster size as an additional constraint to the original k -CENTER formulation, an upper bound on the cluster size is specified. This is called the CAPTITATED k -CENTER problem [11]. Bar-Ilan, Kortsarz, and Peleg [2] gave the first constant approximation factor of 10 for this problem. The bound was improved subsequently to 5 by Khuller and Sussmann [11]. In this subsection though we only concentrate on the (k, r) -CENTER problem defined above.

We note here that the algorithm developed for r -GATHER in Subsection 2.2 can be extended to provide a 2-approximation for the (k, r) -CENTER problem. We just have to add to Condition (2) the extra criteria that if the number of centers chosen exceeds k then exit with failure, i.e., try a different value for R . We can show that Lemma 2.1 holds for the modified conditions, hence an approximation factor of 2.

We also consider the outlier version of this problem, namely, the (k, r, ϵ) -CENTER problem. Combining the techniques presented in this paper and the techniques for the (k, ϵ) -CENTER problem by Charikar et. al [4], one can devise a 6-approximation algorithm. We defer the details to the journal version.

3. CELLULAR CLUSTERING

As mentioned in the introduction, a second approach is to publish the radius of each cluster in addition to its center and the number of points within it. In this case, for each point within a cluster, the radius of the cluster gives an upper bound on the distortion error introduced. The CELLULAR CLUSTERING problem aims to minimize the overall distortion error, i.e., it partitions the points in a metric space into cells, each having a cell center, such that the sum, over all cells, of the products of the number of points in the cell and the radius of the cell is minimized. We even allow each potential cluster center to have a facility (setup) cost $f(v)$ associated with opening a cluster centered at it. This will later allow us to solve the problem in the case when each cluster is required to have at least r points within it.

DEFINITION 3.1. *A cluster consists of a center along with a set of points assigned to it. The radius of the cluster is the maximum distance between a point assigned to the cluster and the cluster center. To open a cluster with cluster center v and radius r incurs a facility cost $f(v)$. In addition, each open cluster incurs a service cost equal to the number of points in the cluster times the cluster radius. The sum of these two costs is called the cellular cost of the cluster. The CELLULAR CLUSTERING problem is to partition n points in a metric space into clusters with the minimum total cellular cost.*

The CELLULAR CLUSTERING problem is NP-complete via reduction from dominating set. We present a primal-dual algorithm for the CELLULAR CLUSTERING problem that achieves an approximation factor of 4.

Let $c = (v_c, d_c)$ denote a cluster c whose cluster center is the node v_c and whose radius is d_c . By definition, the setup cost $f(c)$ for a cluster $c = (v_c, d_c)$ depends only on its center v_c ; thus $f(c) = f(v_c)$. For each possible choice of cluster center and radius $c = (v_c, d_c)$, define a variable y_c , a 0/1 indicator of whether or not the cluster c is open. There are $O(n^2)$ such variables. For a cluster $c = (v_c, d_c)$, any point p_i within a distance of d_c of its center v_c is said to be a *potential member* of the cluster c . For all potential members p_i of a cluster c , let x_{ic} be a 0/1 indicator of whether or not point p_i joins cluster c . Note that the pair (i, c) uniquely identifies an edge between p_i and the center of cluster c . We relax the integer program formulation to get the following linear program:

$$\begin{aligned} \text{Minimize:} & \quad \sum_c (\sum_i x_{ic} d_c + f_c y_c) \\ \text{Subject to:} & \quad \sum_c x_{ic} \geq 1 & \forall i, \\ & \quad x_{ic} \leq y_c & \forall i, c \\ & \quad 0 \leq x_{ic} \leq 1 & \forall i, c \\ & \quad 0 \leq y_c \leq 1 & \forall c \end{aligned}$$

And the dual program is:

$$\begin{aligned} \text{Maximize:} & \quad \sum_i \alpha_i \\ \text{Subject to:} & \quad \sum_i \beta_{ic} \leq f_c & \forall c \\ & \quad \alpha_i - \beta_{ic} \leq d_c & \forall i, c \\ & \quad \alpha_i \geq 0 & \forall i \\ & \quad \beta_{ic} \geq 0 & \forall i, c \end{aligned}$$

The above formulation is similar to the primal-dual formulation of facility location [9]. However, since the assignment of additional points to clusters increases the service cost incurred by existing members of the cluster, we need a different approach to assign points to clusters.

Procedure 1 describes the details of the growth of dual variables and the assignment of points to clusters. We say an edge (i, c) is *tight* if $\alpha_i \geq d_c$. When an edge (i, c) becomes tight, the corresponding cluster c becomes partially open and p_i contributes an amount of $(\alpha_i - d_c)$ to the fixed facility cost of $f(c)$. At any step of the procedure, a point is labeled *unassigned*, *idle* or *dead*. Initially, all points are *unassigned*. As some cluster becomes tight, all *unassigned* or *idle* points having tight edges to it become *dead*. In addition, some of the *unassigned* points become *idle* as described in the procedure.

Procedure 1 A PRIMAL DUAL METHOD

- 1: **repeat**
 - 2: Grow the unfrozen dual variables α_i uniformly.
 - 3: **if** $\alpha_i \geq d_c$ for some cluster c and its potential member p_i , i.e., edge (i, c) is tight, and c has not been shut down **then**
 - 4: Open the cluster c partially, and grow the dual variable β_{ic} at the same rate as α_i .
 - 5: **end if**
 - 6: **if** $\sum_i \beta_{ic} = f_c$ for some cluster c **then**
 - 7: Freeze all variables α_i for which the edge (i, c) is tight.
 - 8: All *unassigned* points with a tight edge to c are assigned to c . Call this set V_c^U .
 - 9: Let V_c^I be the set of all *idle* points that have a tight edge to c .
 - 10: Permanently shut down any cluster $c' \neq c$ for which a point p_i in $V_c^U \cup V_c^I$ has a tight edge (i, c') . Assign to c all *unassigned* points p_j with a tight edge to c' . Call this newly-assigned set of points V_c^{IU} .
 - 11: All points in V_c^{IU} are labeled *idle* and their dual variables are frozen.
 - 12: All points in V_c^U and V_c^I are labeled *dead*.
 - 13: **end if**
 - 14: **until** All points become *dead* or *idle*.
-

We now show that the primal solution constructed has a cost of at most 4 times the value of the dual solution found using Procedure 1. For this, we note the following properties:

- (1) At any instant, the value of α_i for all *unassigned* points i is the same. Moreover, this value is no less than the value of α_j for any *dead* or *idle* point j .
- (2) Once a point has a tight edge to a particular cluster c (i.e., a cluster is partially open), all *unassigned* potential members of that cluster (i.e. points within a distance d_c of the cluster center v_c) have tight edges to it.
- (3) When a cluster opens, all its *unassigned* potential members are assigned to it and become *dead*.
- (4) When a point p_i becomes *dead*, all but one facility partially supported by p_i is shut down.
- (5) When a cluster shuts down, all its *unassigned* potential members are assigned to some open cluster and become *idle*.

Property (1) follows from the definition of our procedure. Property (2) follows from property (1) and the fact that the edge (i, c) becomes tight when the dual variable α_i equals d_c . Property (3) then follows from (2). Property (4) again follows from the definition of the procedure. Property (5) can be seen as follows: we shut down a cluster c only when one of its *unassigned* or *idle* members has a tight edge to the cluster c' currently being opened, and also has a tight edge to c . By property (2), all *unassigned* members of c have tight edges to c . Hence in Steps 10 and 11 of the procedure, these members will be assigned to c' and become *idle*.

LEMMA 3.1. *The service cost for each point, $\sum_c x_{ic}d_c$, is no more than $3\alpha_i$.*

PROOF. Consider the cluster c to which point i is assigned. When cluster c opens, points in V_c^U and V_c^{IU} are assigned to c . We need to bound the radius of the cluster consisting of $V_c^U \cup V_c^{IU}$. By property (1), all points in V_c^U and V_c^{IU} have the same dual variable value, say α . Let p be the cluster center of c . Clearly, for a point $q \in V_c^U$, $d(q, p) \leq d_c \leq \alpha$. For a point $r \in V_c^{IU}$, let c' be its cluster that was shut down (in Step 10) when r was assigned to c . Let p' be the cluster center of c' , and let $q' \in V_c^U$ be the point that was partially supporting c' . Clearly, $\alpha \geq d_{c'}$ since q' is partially supporting c' . Combined with the fact that r and q' are potential members of c' , we get that $d(r, p) \leq d(r, p') + d(p', q') + d(q', p) \leq 2d_{c'} + d_c \leq 3\alpha$. Thus, the cluster made of V_c^U and V_c^{IU} has overall radius no more than $3\alpha = 3\alpha_i$. \square

LEMMA 3.2. *The cost of opening the clusters, $\sum_c y_c f_c$, is no more than $\sum_i \alpha_i$.*

PROOF. A cluster c is opened when $\sum_i \beta_{ic}$ equals f_c . Thus, for each open cluster c , we need to find $V_c \subseteq V$, s.t. $\sum_i \beta_{ic}$ can be charged to $\sum_{i \in V_c} \alpha_i$. To avoid charging any point i more than once, we need to make sure that the V_c 's are disjoint. We begin by noting that when a cluster c opens, only points i with a tight edge to c can contribute to $\sum_i \beta_{ic}$. When a point is labeled *dead*, by Property 4, all

the clusters to which it has a tight edge are shut down and are not opened in future. This implies that clusters which are opened do not have tight edges to *dead* points. Thus, when a cluster c is opened, V_c^U and V_c^I are the only points which have tight edges to c . If we let $V_c = V_c^U \cup V_c^I$, then $\sum_{i \in V_c} \alpha_i \geq \sum_i \beta_{ic}$. Also, since the points in $V_c^U \cup V_c^I$ are labeled *dead* in this iteration, they will not appear in $V_{c'}^U \cup V_{c'}^I$ for any other cluster c' . \square

We thus obtain the following theorem.

THEOREM 3.3. *The primal-dual method in Procedure 1 produces a 4-approximation solution to the CELLULAR CLUSTERING problem.*

3.1 r -Cellular Clustering

We now extend the above primal-dual algorithm to get an approximation algorithm for the r -CELLULAR CLUSTERING problem which has the additional constraint that each cluster is required to have at least r members. The notation (r, C) is used to denote a solution having a total cost of C , and having at least r members in each cluster.

Comparison with prior clustering work. Since our algorithm can be viewed as an extension of facility location, we briefly discuss related results. The facility location (and k -median) problems have been studied with the minimum cluster size constraint [10], as well as in the context of leaving an ϵ fraction of the points unclustered [4]. Let OPT_r be the optimal facility location cost with minimum cluster size r . If as stated before (r, C) denotes a solution with minimum cluster size r and solution cost C , bi-criteria approximation for the facility location problem of $(r/2, 5.184OPT_r)$ was achieved independently by Guha, Meyerson and Munagala and by Karger and Minkoff [7, 10]. It is not known whether it is possible to achieve a one-sided approximation on facility location cost alone. In contrast, for the r -CELLULAR CLUSTERING problem, we provide an one-sided approximation algorithm, specifically we obtain a $(r, 80OPT_r)$ solution, where OPT_r is the cost of the optimal solution with cluster size at least r ,

To achieve this, we first study a *sharing* variant of this problem, where a point is allowed to belong to multiple clusters, thus making it easier to satisfy the minimum cluster size constraint. Interestingly, allowing sharing changes the value of the optimal solution by at most a constant factor. We note that this observation does not hold for facility location, where a shared solution might be arbitrarily better than an unshared one. The algorithm consists of three main steps:

1. Augmenting with Setup Costs. Given an instance of r -CELLULAR CLUSTERING, we first construct an instance of CELLULAR CLUSTERING as follows: augment the cluster cost f_c of a cluster c by $r \times d_c$. In addition, if a cluster $c = (v_c, d_c)$ has fewer than r points within distance d_c of its center v_c , this cluster is eliminated from the instance. If the original r -CELLULAR CLUSTERING instance has an optimal solution with cost OPT_r , it is not hard to see that the same solution works for the CELLULAR CLUSTERING instance constructed above with a total cost of at most $2OPT_r$. We

invoke the 4-approximation algorithm for CELLULAR CLUSTERING on this new instance to find a solution with cost at most $8OPT_r$.

2. Sharing Points between Clusters. We now describe the notion of a *shared* solution for r -CELLULAR CLUSTERING. In a shared solution, points are allowed to be assigned to multiple clusters, as long as they pay the service cost for each cluster they are assigned to. A shared solution is *feasible* if all clusters have at least r (potentially shared) members. We modify the solution obtained above to get a feasible shared solution for r -CELLULAR CLUSTERING as follows: for each open cluster c with center P , assign the r closest neighbors of P to c as well, regardless of where they are initially assigned. The extra service cost of at most $r \times d_c$ for these r points can be accounted for by the extra facility cost of $r \times d_c$ being paid by the open cluster c in the CELLULAR CLUSTERING solution. Thus, we have obtained an $(r, 8OPT_r)$ shared solution for the r -CELLULAR CLUSTERING instance.

3. Making the Clusters Disjoint. Finally we show how to convert a shared solution to a valid solution where each point is assigned to only one cluster, with only a constant blowup in cost. We note that for the corresponding facility location problem, it is not feasible to do this “unsharing” without a large blowup in cost in the worst case.

Initially, all points are labeled *unassigned*. We consider the clusters in order of increasing cluster radius d_c . If a cluster c has at least r *unassigned* members, then it is opened and all its *unassigned* members are assigned to c and labeled *assigned*. We stop this process when all the remaining clusters have fewer than r *unassigned* members each. The remaining clusters are called *leftover* clusters. We temporarily assign each of the *unassigned* points arbitrarily to one of the leftover clusters it belongs to. Since each cluster had at least r members in the shared solution, each leftover cluster c' must have a member in the shared solution, which is now *assigned* to an open cluster o , s.t. $d_{c'} \geq d_o$. We thus have the situation illustrated in Figure 4.

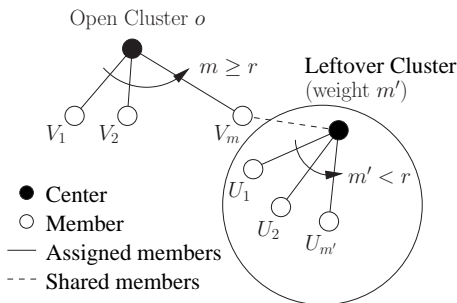


Figure 4: Structures of open and leftover clusters

The points are organized in a forest structure, where each tree has two “levels”. We can regroup points into clusters, on a per tree basis. It is obvious that each tree has at least r points, since it contains at least one open cluster o . We further simplify the structure into a true two-level structure as in Figure 4, by collapsing each leftover cluster into a single

node with weight equal to the number of points temporarily assigned to it. Nodes in the first level of the tree have weight 1. We apply the following greedy grouping procedure: first consider only the nodes at the second level of the tree and collect nodes until the total weight exceeds r for the first time. We group these nodes (belonging to leftover clusters) into a cluster, and repeat the process. Notice that since we did not touch the first-level nodes, the total weight of remaining nodes in the tree is at least r . If the total weight of remaining nodes in the second level, W_s , is less than r , then we extend the grouping into the first level nodes. Let m denote the total weight of nodes in the first level. If $W_s + m \geq 2r$, then we group the nodes in the second level with $r - W_s$ first level nodes together into a cluster; the remaining nodes in the first level form a cluster. Otherwise, all the remaining nodes (both the first and second level) are grouped into a cluster. If we break up the tree using the procedure above, each resulting cluster has size at least r .

LEMMA 3.4. *For a cluster that contains any second-level nodes, the total number of points in the cluster is no more than $2r - 1$.*

PROOF. Since a single second-level node has weight less than r , a cluster containing only second-level nodes has at most $2r - 1$ members. If the cluster contains both the first and second-level nodes, then we must have reached the case where the total weight of remaining nodes in the second level is less than r . In that case, by definition, the cluster formed containing these second-level nodes has size either r or less than $2r - 1$. \square

There could be a cluster that only contains the first level nodes, and its entire cost (both the service and cluster cost) can be accounted for by its cost in the original $(r, 8OPT_r)$ shared solution. We now bound the cost of clusters containing the second-level nodes.

LEMMA 3.5. *For each cluster c formed that contains second level nodes, there exists a leftover cluster c' unique to c , such that the following holds: let p be the center of c' , if we center the cluster c at p , then the radius of cluster c , $\text{radius}(c) \leq 5d_{c'}$.*

PROOF. Among all the leftover clusters that contributed to c , let c' be the one with the maximum radius. By definition, all nodes assigned to a leftover cluster get assigned to a single cluster, guaranteeing the uniqueness of c' . Let d_o be the radius of the open cluster at level 1 of this tree. Consider a point $q \in c$. If q is a first-level node, then $d(q, p) \leq 2d_o + d_{c'} \leq 3d_{c'}$. If q is a second-level node, then let c'' be the leftover cluster that q was assigned to, then $d(q, p) \leq 2d_{c''} + 2d_o + d_{c'} \leq 5d_{c'}$. \square

The above lemma implies that by choosing p as the cluster center, the service cost of each point in c is no more than $5d_{c'}$ and the total facility cost incurred within our solution is no more than that of the shared solution. Together with Lemma 3.4, we conclude that the service cost of points in c

is no more than $10r \times d_{c'}$. Notice that in the shared solution, points in cluster c' are paying a total service cost of at least $r \times d_{c'}$. We thus have the following theorem.

THEOREM 3.6. *The above procedure produces a solution with minimum cluster size r and total cost no more than $80OPT_r$, i.e., a $(r, 80OPT_r)$ solution, where OPT_r is the value of the optimal solution with a minimum cluster size of r .*

We note that the above algorithm and analysis can be combined with the technique developed in [4] to give an constant approximation to the (r, ϵ) -CELLULAR CLUSTERING problem. The above algorithm can also be adapted to provide a constant-factor approximation for the problem where the diameter of any cluster is not allowed to exceed a certain pre-specified threshold. Details are deferred to the full version of the paper.

4. CONCLUSIONS

Publishing data about individuals without revealing sensitive information is an important problem. The notion of privacy called k -Anonymity has attracted a lot of research attention recently. In a k -anonymized database, values of quasi-identifying attributes are suppressed or generalized so that for each record there are at least $k - 1$ records in the modified table that have exactly the same values for the quasi-identifiers. However, the performance of the best known approximation algorithms for k -Anonymity depends linearly on the anonymity parameter k . In this paper, we introduced clustering as a technique to anonymize quasi-identifiers before publishing them. We studied r -GATHER as well as a newly introduced clustering metric called r -CELLULAR CLUSTERING and provided the first constant-factor approximation algorithms for publishing an anonymized database table. Moreover, we generalized these algorithms to allow an ϵ fraction of points to remain unclustered.

5. REFERENCES

- [1] G. Aggarwal, T. Feder, K. Kenthapadi, R. Motwani, R. Panigrahy, D. Thomas, and A. Zhu. Approximation Algorithms for k -Anonymity. *Journal of Privacy Technology*, Paper number: 20051120001, 2005.
- [2] J. Bar-Ilan, G. Kortsarz, and D. Peleg. How to allocate network centers. *Journal of Algorithms*, 15, pages 385–415, 1993.
- [3] R. J. Bayardo and R. Agrawal. Data Privacy through Optimal k -Anonymization. In *Proceedings of the International Conference on Data Engineering*, pages 217–228, 2005.
- [4] M. Charikar, S. Khuller, D. Mount and G. Narasimhan. Algorithms for Facility Location with Outliers. In *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms*, pages 642–651, 2001.
- [5] S. Chawla, C. Dwork, F. McSherry, A. Smith, and H. Wee. Toward Privacy in Public Databases. In *Proceedings of the Theory of Cryptography Conference*, pages 363–385, 2005.
- [6] M. R. Garey and D. S. Johnson. *Computers and intractability, a guide to the theory of NP-completeness*. W. H. Freeman and Company, New York, Nov. 1990.
- [7] S. Guha, A. Meyerson, and K. Munagala. Hierarchical Placement and Network Design Problems. In *Proceedings of the IEEE Symposium on Foundations of Computer Science*, pages 603–612, 2000.
- [8] D. Hochbaum and D. Shmoys. A best possible approximation algorithm for the k -center problem. *Mathematics of Operations Research*, 10(2), pages 180–184, 1985.
- [9] K. Jain and V. V. Vazirani. Primal-Dual Approximation Algorithms for Metric Facility Location and k -Median Problems. In *Proceedings of the IEEE Symposium on Foundations of Computer Science*, pages 2–13, 1999.
- [10] D. Karger and M. Minkoff. Building Steiner Trees with Incomplete Global Knowledge. In *Proceedings of the IEEE Symposium on Foundations of Computer Science*, pages 613–623, 2000.
- [11] S. Khuller and Y. Sussmann. The Capacitated K -Center Problem. *SIAM Journal on Discrete Mathematics*, 13(3), pages 403–418, 2000.
- [12] K. LeFevre, D. DeWitt, and R. Ramakrishnan. Incognito: Efficient Full-Domain K -Anonymity. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 49–60, 2005.
- [13] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam. l -Diversity: Privacy Beyond k -Anonymity. In *Proceedings of the International Conference on Data Engineering*, 2006.
- [14] A. Meyerson and R. Williams. On the Complexity of Optimal k -Anonymity. In *Proceedings of the Symposium on Principles of Database Systems*, pages 223–228, 2004.
- [15] L. Sweeney. Uniqueness of Simple Demographics in the U.S. Population. LIDAP-WP4. *Carnegie Mellon University, Laboratory for International Data Privacy, Pittsburgh, PA*, 2000.
- [16] L. Sweeney. k -Anonymity: A Model for Protecting Privacy. *International Journal on Uncertainty Fuzziness Knowledge-based Systems*, 10(5), pages 557–570, 2002.
- [17] Time. *The Death of Privacy*, August 1997.