

Utility-efficient Differentially Private K-means Clustering based on Cluster Merging

Tianjiao Ni, Minghao Qiao, Zhili Chen^{*}, Shun Zhang, Hong Zhong
School of Computer Science and Technology, Anhui University, Hefei, China

Abstract

Differential privacy is widely used in data analysis. State-of-the-art k -means clustering algorithms with differential privacy typically add an equal amount of noise to centroids for each iterative computation. In this paper, we propose a novel differentially private k -means clustering algorithm, DP-KCCM, that significantly improves the utility of clustering by adding adaptive noise and merging clusters. Specifically, to obtain k clusters with differential privacy, the algorithm first generates $n \times k$ initial centroids, adds adaptive noise for each iteration to get $n \times k$ clusters, and finally merges these clusters into k ones. We theoretically prove the differential privacy of the proposed algorithm. Surprisingly, extensive experimental results show that: 1) cluster merging with equal amounts of noise improves the utility somewhat; 2) although adding adaptive noise only does not improve the utility, combining both cluster merging and adaptive noise further improves the utility significantly.

Keywords: K-means, Cluster, Differential Privacy

1. Introduction

With the rapid development of Internet technology, third-party applications have produced a large amount of user data. The correct use of these data is able to create incalculable value for governments, companies and individuals. How to extract useful information from user data is currently a hot research

^{*}Zhili Chen is the corresponding author. Email: zlchen@ahu.edu.cn

direction [1, 2]. Clustering algorithms are widely used to complete this task in the field of data analysis [3, 4]. The goal of clustering is to divide elements of a dataset into different groups so that the elements in the same group have high similarity. There are a large number of clustering algorithms [5, 6, 7, 8, 9, 10]. Among them, the k -means clustering is one of the most popular methods for numeric data.

Recently, many applications adopted the k -means clustering algorithm. Javadi et al. [11] classified aquifer vulnerability using k -means cluster analysis. Han et al. [12] factored out k -means cluster-based location privacy protection scheme for Internet of Things. Shakeel et al. [13] used k -means clustering to diagnosis of diabetes mellitus. Wu et al. [14] used k -means in the compressing convolutions of convolutional neural network. Reza et al. [15] adopted k -means clustering with graph-cut segmentation to estimate rice yield. Omrani et al. [16] presented the artificial neural network-based land transformation model, which uses the k -means clustering algorithm implemented within the Spark high-performance compute environment. However, most of these applications fail to consider disclosure of sensitive information, which might bring immeasurable threats to users [17, 18].

To solve the privacy problem, differential privacy [19, 20] is proposed as a powerful privacy protection technique and has been extensively used [21, 22, 23]. Recently, several state-of-the-art k -means clustering algorithms with differential privacy have been proposed. For example, Yu et al. [24] presented a differentially private k -means clustering scheme and improved its utility by selecting initial centroids with the distribution density of elements. Su et al. [25] analyzed several existing differentially private k -means clustering algorithms and improved one of them by selecting initial centroids based on the concept of sphere packing. However, these algorithms still suffer the issue of lacking high utility due to adding large amounts of noise. Thus, how to improve the utility of differentially private k -means clustering remains as a key question.

In this paper, to address the above-mentioned utility issue, we propose a novel differentially private k -means clustering algorithm based on cluster merg-

ing (DP-KCCM). DP-KCCM first partitions the data into $n \times k$ clusters with differential privacy, and then merges these clusters into required k ones. The main idea is that the Laplace noises added to cluster centroids are random, and cluster merging would cancel the noises each other, and thus improve the utility. More interestingly, we find that combining cluster merging with adaptive noise is able to further improve the cluster utility.

The main contributions of this paper are as follows: 1) We propose a utility-efficient, differentially private k -means clustering algorithm based on cluster merging. 2) We design a privacy budget (i.e., the privacy parameter of differential privacy, cf. Definition 1) allocation to work with cluster merging to further improve the cluster utility. 3) Extensive experimental results show that our algorithm is superior to the state-of-the-art ones.

The rest of this paper is organized as follows. In Section 2, we introduce the background knowledge of this paper. Section 3 describes our algorithm in detail and establishes that the algorithm satisfies differential privacy. In Section 4, we carry out extensive experiments, and compare the utility of our algorithm with those of the existing algorithms. We conclude in Section 5.

2. Background

In this section, we first introduce the notion of differential privacy and the algorithms of K-means clustering, we then present the problem statement. Some notations used in this paper are described in Table 1.

2.1. Differential Privacy

The notion of differential privacy requires that the outputs of a data analysis mechanism should be similar over any two adjacent datasets. The formal definition is as follows.

Definition 1. (*ϵ -Differential Privacy [19]*). *A randomized algorithm \mathcal{M} satisfies ϵ -differential privacy (ϵ -DP), if and only if for any pair of neighboring*

Table 1: Notations and descriptions

Notations	Descriptions
D	dataset
N	the number of data points in the dataset
d	the dimension of the dataset
x_i	i -th data point in the dataset and ranging from $[-r, r]^d$
k	the number of the clusters
C	the set of the centroids
C_j	j -th centroid
C_j^i	i -th dimension of j -th centroid
c_j	j -th noisy centroid
c_j^i	i -th dimension of j -th noisy centroid
C_j^*	j -th cluster
ϵ	the privacy budget
$sum(\cdot)$	the sum of data point values in a cluster
$sum'(\cdot)$	$sum'(\cdot)$ is $sum(\cdot)$ with noise
$num(\cdot)$	the number of data points in a cluster
$num'(\cdot)$	$num'(\cdot)$ is $num(\cdot)$ with noise
$dist(x, y)$	the distance between data points x and y
Δ	the global sensitivity of an iteration
$iter$	the number of iterations

datasets D and D' , and any $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$, we have

$$\Pr[\mathcal{M}(D) = \mathcal{S}] \leq e^\epsilon \cdot \Pr[\mathcal{M}(D') = \mathcal{S}] \quad (1)$$

In this definition, D 's neighboring dataset D' can be obtained by adding an element to or removing an element from D , and they can be denoted by $D \simeq D'$. The $\text{Range}(\mathcal{M})$ represents the set of all possible outputs of the algorithm \mathcal{M} . It is worth noting that the parameter ϵ is called *privacy budget*, which indicates the privacy level. A smaller ϵ value means more similar outputs resulted

from neighboring datasets due to Eq. (1), and thus represents stronger privacy achieved. On the other hand, the greater the ϵ value is, the weaker the privacy is preserved. In the extreme, when there is no differential privacy protection, it is equivalent that the ϵ is infinitely large, and thus the privacy may be easily disclosed.

Differential privacy has good properties in composition of multiple algorithms, which are described in Lemmas 1, 2 and 3.

Lemma 1 (Parallel Composition [26]). *If there are algorithms $\mathcal{M}_1, \dots, \mathcal{M}_k$ satisfying $\epsilon_1, \dots, \epsilon_k$ -DP, respectively, then for disjoint datasets D_1, D_2, \dots, D_k , composition algorithm $\mathcal{M}(\mathcal{M}_1(D_1), \dots, \mathcal{M}_k(D_k))$ provides $(\max_{i \in \{1, \dots, k\}} \epsilon_i)$ -DP.*

Lemma 1 shows that if the input datasets are disjoint, the privacy level provided by a parallel composition depends on the algorithm with the lowest privacy level, namely the one with the largest privacy budget.

Lemma 2 (Sequential Composition [27, 26]). *If there are algorithms \mathcal{M}_1 satisfying ϵ_1 -DP, and \mathcal{M}_2 satisfying ϵ_2 -DP, then $\mathcal{M}(D) = \mathcal{M}_1(\mathcal{M}_2(D), D)$ satisfies $(\epsilon_1 + \epsilon_2)$ -DP.*

Lemma 2 implies that if multiple algorithms are applied to the same data set sequentially, the resulted privacy budget is the sum of these algorithms' privacy budgets.

Lemma 3 (Post-processing [27]). *If there is an algorithm $\mathcal{M}_1(\cdot)$ satisfying ϵ -DP, then for any algorithm $\mathcal{M}_2(\cdot)$, $\mathcal{M}_2(\mathcal{M}_1(\cdot))$ satisfies ϵ -DP.*

Lemma 3 shows the post-processing property of differential privacy, namely, if an algorithm takes as input the output of another algorithm that satisfies ϵ -differential privacy, then the resulted algorithm still satisfies ϵ -differential privacy.

In this paper, we use the Laplace mechanism [20] to design algorithms. The Laplace mechanism preserves differential privacy by adding random noise satisfying the Laplace distribution to any query function f (e.g., count query) over

D . The magnitude of noise depends on the sensitivity of f , Δf , which represents the maximum deviation of the query result on any adjacent datasets. For instance, $\Delta f = 1$ for count query. The idea of Laplace mechanism is to add smallest but sufficient noise to any query result of f , such that the query results of any neighboring datasets (only differing in an element) are indistinguishable, and thus the personal privacy (i.e., any single element information) cannot be inferred from query results. The *Laplace mechanism* \mathcal{M}_f is given in Definition 2, and its differential privacy is guaranteed by Lemma 4

Definition 2. (*Laplace Mechanism [20]*) *The Laplace mechanism is defined as:*

$$\mathcal{M}_L(D, f, \epsilon) = f(D) + (Y_1, Y_2, \dots, Y_d) \quad (2)$$

where $f(D)$ is a given query function $f(D) : D \rightarrow \mathbb{R}^d$ with sensitivity

$$\Delta f = \max_{(D, D') : D \simeq D'} \|f(D) - f(D')\|_1 \quad (3)$$

and $Y_i (1 \leq i \leq d)$ are i.i.d. random variables drawn from $\text{Lap}(\frac{\Delta f}{\epsilon})$.

The probability density function of Laplace distribution is as follows:

$$\text{Lap}(b) = \text{Lap}(x|b) = \frac{1}{2b} e^{-|x|/b} \quad (4)$$

where for the Laplace mechanism, $b = \Delta f / \epsilon$.

Lemma 4. [27, 26] *The Laplace mechanism preserves ϵ -DP.*

2.2. K-means Clustering Algorithms

The k -means clustering divides a set of data points into different subsets. In the clustering, there are typically two following steps.

2.2.1. Initial Centroid Selection

Given a value $a > 0$, paper [25] randomly selects k initial centroids one by one. Each selection of a centroid must follow two principles: 1) the distance between any centroid and the boundary of the domain is at least a . 2) the distance between any two centroids is at least $2a$. If a randomly selected centroid

does not satisfy the above two conditions, it is discarded and another centroid is reselected until k initial centroids are obtained. When it fails to get k centroids several times, it may be that the given value a is too large, and then a smaller a is tried again. The initial value of a can be determined according to the domain size of data points. In our context, since the data is normalized into $[-1, 1]$, the initial value of a is set to 0.5. During the experiment, the optimal value of a can be obtained by the binary search or it is simply set by experience. The above process of selecting the initial centroids only depends on the domain of data points rather than the data points themselves, so this process does not impact the privacy, and it can be applied directly in the privacy-preserving algorithms.

2.2.2. *K-means Clustering*

For a dataset $D = \{x_1, x_2, \dots, x_N\}$, $x_i \in R^d$, the standard k -means clustering aims at partitioning data into k disjoint subsets $(C_1^*, C_2^*, \dots, C_k^*)$. The evaluation metric of clustering results called the Normalized Intra-Cluster Variance (NICV)[25] is as follows.

$$\frac{1}{N} \sum_{j=1}^k \sum_{x_i \in C_j^*} \|x_i - C_j\|_2 \quad (5)$$

where C_j is the centroid of the cluster C_j^* , and the smaller the NICV value, the better the clustering result.

Specifically, the algorithm selects k data points as the initial centroids by initial centroids selection algorithm, then the quality of the centroids is improved iteratively until the centroids do not change. In each iteration, the algorithm traverses all the data points of the dataset and assigns the data points to the nearest cluster, then updates the centroid of each cluster.

$$C_j^t = \frac{\sum_{x_i \in C_j^*} x_i^t}{|C_j^*|}, \forall t \in \{1, \dots, d\} \quad (6)$$

where C_j^t is the t -th dimension of j -th centroid and x_i^t is the t -th dimension of x_i .

2.3. Problem Statement

We focus on the privacy problem for data analysis as follows. A trustable data holder (e.g., a government agency) collects personal records from a great number of users, analyzes these data with certain machine learning algorithms (e.g., in our context K-means clustering algorithm), and publishes the analytic results to dishonest third parties, which make use of these results for some statistical purpose. However, the privacy problem is that the dishonest third parties may also infer the personal privacy through exploiting the analytic results, beyond the legitimate use, if there is no appropriate privacy protection measures, as shown in Figure 1. This results in the disclosure of individual privacy.

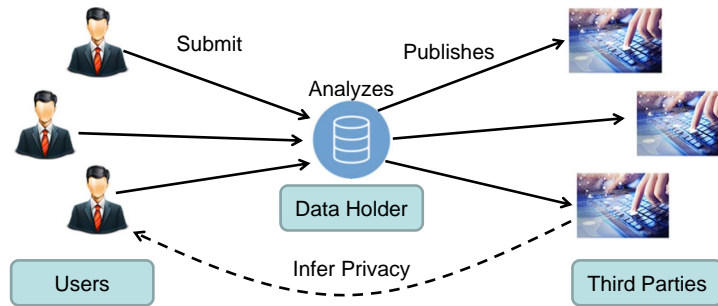


Figure 1: Problem Model

We use the notion of differential privacy to solve the privacy problem described above. As described in Section 2.1, differential privacy ensures that the analytic results of any two neighboring datasets are similar by adding appropriate noise to the results. This means that the differentially private analytic result of a dataset remains roughly the same when any personal record opts into or out of the dataset. Conversely, the inference of any personal record (i.e., individual privacy) from the differentially private analytic results is thus hard.

In our context, we design a K-means clustering algorithm with ϵ -differential privacy to protect individual privacy (i.e., any data point) from being inferred from the clustering results by dishonest third parties. A proper ϵ value can be chosen to ensure a certain level of privacy, e.g., $\epsilon = 1$ (Note that the smaller

ϵ , the stronger the privacy). Furthermore, since ϵ -differential privacy requires that noise be added to the clustering process, we aim to add as little noise as possible while preserving ϵ -differential privacy, to improve the clustering utility.

3. DP-KCCM Algorithm

In this section, we propose our differentially private k -means clustering algorithm (DP-KCCM) and prove its privacy in detail.

Considering previous works on k -means clustering with differential privacy, such as DPLloyd [28, 29] and DPLloyd-Impr [25], they all added equal amounts of noise to the centroids of each iteration in the clustering process. Moreover, it is implicitly suggested that privacy budget should be divided equally across the iterations of the clustering [25]. It seems hard to improve the cluster utility merely through the privacy budget allocation. Our idea is that we may reduce the amounts of noise added by merging some noisy clusters and canceling the amounts of noise. Additionally, we may combine cluster merging with privacy budget allocation to further improve the utility.

Therefore, we are going into two questions: 1) Can we add noise adaptively in the process of iteration to improve the utility? 2) Can we merge adjacent clusters to reduce the noise added, and hence improve the utility? We describe these two aspects in detail.

3.1. Idea 1

In our differentially private k -means clustering algorithm, initial centroids are first selected using the initial centroids selection algorithm [25] (cf. Section 2.2.1), which ensures that the initial centroids are separated as much as possible. Then, these centroids are iteratively updated. Intuitively, in the first several iterations of the algorithm, the centroids change greatly, and we could inject relatively more noise. As the number of iterations increases, the changes of cluster centroids become less, and we could add a small amount of noise to ensure better clustering results. Since the noise volume is controlled by ϵ , we introduce the partition of ϵ as follows.

According to the previous analysis, we know that as the number of iterations increases, the clustering centroids tends to be stable, and hence the noise added should become smaller, and ϵ should become bigger. Thus, our privacy budget allocation policy is to increase ϵ share gradually as the cluster centroids are updated iteratively.

We do a lot of experiments based on this allocation policy, and find that the increase of the ϵ share should be relatively slow. Otherwise, the first several ϵ shares would become too small that the cluster result deteriorates severely. Finally we settled on the following division. The clustering algorithm in this paper has carried out 12 iterations, and we set ϵ share for each iteration as follows: the values of ϵ shares for 1st to 4th iterations are $\frac{1}{24}\epsilon$, those for 5th to 8th iterations are $\frac{1}{12}\epsilon$, and those for 9th to 12th iterations are $\frac{1}{8}\epsilon$.

3.2. Idea 2

We know that the noise added to each centroid is random. Can we reduce the influence of noise on the centroid by merging adjacent clusters? For k -means clustering algorithm, we can divide dataset into $n \times k$ clusters. After clustering, we merge $n \times k$ clusters into k ones. By merging multiple clusters, the noises added to clusters were empirically proved to cancel each other out.

We describe this idea with the example as shown in Figure 2. Suppose there are four points A , B , C and D , which are probably clustered into the same category in a differentially private clustering. There are two ways to do this. The first way is to cluster these points into a category one shot without cluster merging, and get its noisy clustering centroid C_{ABCD}^* . The second way is to first cluster the points into two clusters, with A, B in one and C, D in the other, compute their respective cluster centroids C_{AB}^* and D_{CD}^* , and then merge the two clusters to get the final cluster centroid C_{ABCD}^{**} . As long as the noises added to centroids C_{AB}^* , C_{CD}^* and C_{ABCD}^* are roughly the same, namely the distances from the noiseless centroids C_{AB} , C_{CD} and C_{ABCD} to the noisy counterparts are approximately equal, the merged cluster centroid C_{ABCD}^{**} is probably less noisy than the centroid C_{ABCD}^* , and thus this likely leads to a better clustering

utility.

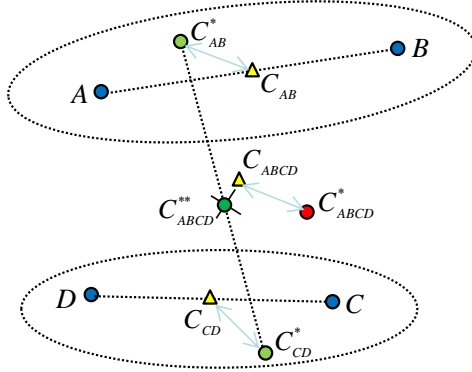


Figure 2: An example for cluster merging.

3.3. Algorithm DP-KCCM

We combine idea 1 with idea 2 to design the algorithm DP-KCCM, which can be described with the following steps.

- (1) Obtain $n \times k$ initial centroids by initial centroids selection algorithm.
- (2) Divide data points into $n \times k$ clusters.
- (3) Recalculate the centroids.
- (4) Add adaptive noise according to the number of iterations (cf. Section 3.1).
- (5) Repeat Steps (2), (3) and (4) until the maximum number of iterations is reached.
- (6) Merge $n \times k$ clusters into k clusters.

The formal description of the algorithm is shown in Algorithm 1. We also detail the key steps in the following.

3.3.1. Divide Data Points into $n \times k$ Clusters

For a data point, we set two variables min_dist and $cent$. The former represents the distance from this data point to its centroid, and the latter is the index of the cluster (i.e., $cent \in [1, n \times k]$) containing the data point. Firstly, we traverse all the data points in the dataset and calculate the distances of

Algorithm 1 DP-KCCM

Input: dataset D , cluster number k , clustering number $max_clustering = 12$,

global sensitivity $\Delta = d \cdot r + 1$, privacy budget $\epsilon = \sum_{i=1}^{max_clustering} \epsilon_i$.

Output: The k centroids.

initialize $C = n \times k$ centroids by initial centroid selection

for $iter \leftarrow 1$ to $max_clustering$ **do**

 Get $n \times k$ clusters through the standard k -means algorithm

 Recalculate the centroid of each cluster

for $j \leftarrow 1$ to $n \times k$ **do**

for $i \leftarrow 1$ to d **do**

$$sum'(C_j^*)[i] = sum(C_j^*)[i] + Lap(\frac{\Delta}{\epsilon_{iter}})$$

end for

$$num'(C_j^*) = num(C_j^*) + Lap(\frac{\Delta}{\epsilon_{iter}})$$

$$c_j = \frac{sum'(C_j^*)}{num'(C_j^*)}$$

end for

end for

$$C_num = |C|$$

while $C_num > k$ **do**

 Find the two nearest cluster C_p^* , C_q^* and combine them into cluster C_o^* :

for $i \leftarrow 1$ to d **do**

$$c_o^i = \min(c_p^i, c_q^i) + |c_p^i - c_q^i| \cdot \frac{num'(C_m^*)}{num'(C_p^*) + num'(C_q^*)}$$

end for

$$C_num = C_num - 1$$

end while

 Get k centroids

each data point to all centroids, then each data point is assigned to the cluster determined by the corresponding minimum distance. Furthermore, we save the minimum distance value of each data point into min_dist and record the index of the corresponding centroid into $cent$. The distance from x_i to the j -th

centroid C_j is computed as follow:

$$dist(x_i, C_j) = \|x_i - C_j\|_2. \quad (7)$$

where each dimension of data points x_i is normalized to $[-r, r]$, and we choose $r = 1$ in our context for simplicity.

3.3.2. Add Noise to Centroids

When all data points are divided, we obtain $n \times k$ clusters. Then we recalculate the centroid of each cluster to get $n \times k$ new centroids. The calculation of the centroid is as follows:

$$C_j = \frac{sum(C_j^*)}{num(C_j^*)}, \forall j \in \{1, \dots, n \times k\} \quad (8)$$

where $sum(C_j^*) = \sum_{x_i \in C_j^*} x_i$, $num(C_j^*) = |C_j^*|$.

In order to protect the information of data points, we need to add a certain amount of Laplace noise to both d -dimension sum of data points and the number of data points during the above calculation. And then we get $n \times k$ noisy centroids after noise addition. The function for calculating noisy centroid j is

$$c_j = \frac{sum'(C_j^*)}{num'(C_j^*)}. \quad (9)$$

with

$$sum'(C_j^*) = sum(C_j^*) + (Y_1, \dots, Y_d) \quad (10)$$

$$num'(C_j^*) = num(C_j^*) + Y_{d+1} \quad (11)$$

where Y_i (for $1 \leq i \leq d+1$) are i.i.d. random variables drawn from $Lap(\Delta/\epsilon_{iter})$, $\Delta = d \cdot r + 1$, r is the maximum absolute value of each dimension, and ϵ_{iter} is the privacy budget for the current iteration $iter$.

Note that the global sensitivity Δ can be computed as follows. For each iteration, each data point is involved in answering d sum queries and one count query. Moreover, each dimension of data points is normalized to $[-r, r]$. Thus, the global sensitivity $\Delta = d \cdot r + 1$.

3.3.3. Merge $n \times k$ Clusters into k Clusters

After the number of iterations reaches the specified value, we can obtain $n \times k$ clusters. And then we merge the two nearest clusters iteratively. We merge two clusters using their noisy centroids and noisy cluster sizes. Note that, to achieve differential privacy, cluster sizes are needed to add Laplace noise, and the noisy sizes does not represent real counts of elements in clusters, but are only used for computing merged centroids. After merging of cluster C_p^* and cluster C_q^* , the i -th dimension of the centroid of the new cluster C_o^* is

$$c_o^i = \min(c_p^i, c_q^i) + |c_p^i - c_q^i| \cdot \frac{\text{num}'(C_m^*)}{\text{num}'(C_p^*) + \text{num}'(C_q^*)} \quad (12)$$

where $m = \arg \max(c_p^i, c_q^i)$, namely, $m = p$ if $c_p^i \geq c_q^i$, and $m = q$ otherwise.

3.4. Privacy Analysis

Theorem 1 states that the DP-KCCM algorithm achieves ϵ -differential privacy. We mainly prove the theorem below.

Theorem 1. *The DP-KCCM algorithm preserves ϵ -differential privacy.*

Proof. We prove the theorem in the following three parts.

First, the initial centroid selection is independent of the data, so the privacy is not impacted in this step.

Second, we show that the algorithm achieves ϵ_{iter} -differential privacy for iteration $iter$, and achieves ϵ -differential privacy for all iterations due to the sequential composition property (cf. Lemma 2).

For iteration $iter$, each data point is involved in d sum queries and one count query. Conversely, each iteration queries the function $f : D^d \rightarrow D^d \times N$ over each cluster, with the global sensitivity $\Delta f = \Delta = d \cdot r + 1$. Specifically, let D and $D' = D - \{x\}$, for any point $x \in D$, be neighboring datasets. Let D be divided into disjoint clusters $C_1^*, C_2^*, \dots, C_{n \cdot k}^*$. Then, the neighboring dataset D' is correspondingly divided into disjoint clusters $C_1^{*'}, C_2^{*'}, \dots, C_{n \cdot k}'$, satisfying that there exists an J , such that $C_J^{*'} = C_J^* - \{x\}$, and $C_j^{*'} = C_j^*$ for $j \neq J$.

Therefore, each iteration can be regarded as the parallel composition of mechanisms querying function $f(\cdot)$ over $n \times k$ disjoint clusters (cf. Lemma 1), and differential privacy achievement is determined by the mechanism over cluster C_J^* (since mechanisms over other clusters C_j^* achieve 0-differential privacy for $C_j^* = C_j^{*'}$).

Let $p(\cdot)$ and $p'(\cdot)$ denote the probability density functions of the mechanism over cluster C_J^* . For any point $v \in D^d \times N$, the probability density ratio between the cases of C_J^* and $C_J^{*'}$ are as follows.

$$\begin{aligned} \frac{p(v)}{p'(v)} &= \frac{\exp(-\frac{\epsilon_{iter}\|f(C_J^*)-v\|_1}{\Delta f})}{\exp(-\frac{\epsilon_{iter}\|f(C_J^{*'})-v\|_1}{\Delta f})} \\ &= \exp\left(\frac{\epsilon_{iter}(\|f(C_J^{*'})-v\|_1 - \|f(C_J^*)-v\|_1)}{\Delta f}\right) \\ &\leq \exp\left(\frac{\epsilon_{iter} \cdot \|f(C_J^*) - f(C_J^{*'})\|_1}{\Delta f}\right) \\ &\leq \exp(\epsilon_{iter}) \end{aligned}$$

Symmetrically, we have $\frac{p(v)}{p'(v)} \geq \exp(-\epsilon_{iter})$. The mechanism over C_J^* achieves ϵ_{iter} -differential privacy, and thus the iteration mechanism $\mathcal{DP} - \mathcal{ITER}$ achieves ϵ_{iter} -differential privacy according to the parallel composition property (cf. Lemma 1).

Since the computation of noisy centroids c_j is a post-process of iteration mechanism $\mathcal{DP} - \mathcal{ITER}$, iteration $iter$ achieves ϵ_{iter} -differential privacy due to the post-processing property (cf. Lemma 3).

For all iterations, the sequential composition property (cf. Lemma 2) is applied, and the resulted mechanism satisfies the ϵ -differential privacy, where $\epsilon = \epsilon_1 + \epsilon_2 + \dots + \epsilon_{max_clustering}$.

Finally, merging $n \times k$ clusters into k ones involves only the noisy cluster centroids and noisy cluster sizes, and it is actually a post-process of the composition of iterations, impacting nothing on the differential privacy achieved.

Therefore, we conclude that the DP-KCCM algorithm preserves ϵ -differential privacy. \square

4. Performance Evaluation

4.1. Methodology

We implement the proposed differentially private k -means clustering algorithm, DP-KCCM, and do experiments to evaluate the algorithm based on six datasets. The detailed description of these six datasets is as Table 2. The data attributes contained in these six datasets are all of numerical type. We normalize the domain of each attribute to the range of $[-1, 1]$.

Table 2: Description of the Datasets.

Dataset	tuples	dims	cluster	description
Blood	748	5	4	The dataset records individual blood donations, and is taken from the Blood Transfusion Service Center.
Adult	32561	6	5	This is a census dataset that records personal information.
Tripadvisor-review	980	10	4	The dataset is the reviews on destinations in 10 categories mentioned across East Asia.
Electrical	10000	13	5	The dataset is the simulated data for the local stability analysis of the 4-node star system (electricity producer is in the center) implementing Decentral Smart Grid Control concept.
Review-ratings	5454	24	4	The dataset contains google reviews on attractions from 24 categories across Europe.
Credit-card	30000	24	5	The dataset contains customer default payments in Taiwan.

All datasets are downloaded from website <http://archive.ics.uci.edu/ml/datasets.php>

We mainly focus on comparisons of algorithm performances from the following two aspects:

- Comparing the effect of different algorithms with a fix k value under different ϵ values.
- Comparing the effect of different algorithms with a fix ϵ value under different k values.

We compare the following four differentially private k -means clustering algorithms.

- *average_k* : k initial centroids are generated by the initial centroid selection algorithm, and the average noise is added to all centroids during each iteration.
- *allocation_k* : k initial centroids are generated by the initial centroid selection algorithm, and adaptive noise is added to all centroids during each iteration.
- *average_nk* : $n \times k$ initial centroids are generated by the initial centroid selection algorithm, and average noise is added to all centroids in the process of each iteration. When the clustering is stable, $n \times k$ clusters are combined into k clusters.
- *allocation_nk* : $n \times k$ initial centroids are generated by the initial centroid selection algorithm, and adaptive noise is added to all centroids during each iteration. When the clustering is stable, $n \times k$ clusters are combined into k clusters.

The *average_k* algorithm is in fact the state-of-the-art algorithm, DPLloyd-Impr, which is reported to be the best one overall among several existing differentially private k -means clustering algorithms [25]. Thus, in our experiments, we use it a benchmark algorithm for the performance comparisons. Then, on the basis of *average_k* algorithm, the two ideas are introduced separately, generating algorithms *allocation_k* and *average_nk*, respectively. Finally, the two

ideas are combined into the algorithm *allocation_nk*, which is the proposed DP-KCCM algorithm. We compare these four algorithms in the experiments to demonstrate the effectiveness of the two ideas.

The above algorithms all output k centroids $\mathbf{C} = \{C_1, C_2, \dots, C_k\}$. The quality of clustering are assessed by Normalized Intra-Cluster Variance (NICV). For all algorithms, we apply the initial centroid selection algorithm to get the initial centroids. In our experiments, we first adopt the initial centroid selection algorithm to generate 20 sets of initial centroids. Then, we run 50 times on each set of initial centroids. So we take the average of NICV in 1000 experiments. Through a lot of experiments, we make the following settings for some parameters in the process of cluster with each dataset. We found that the clustering tends to be stable after the number of iterations reaches 10, so we set *max_clustering* = 12. The experimental results obtained by merging $n \times k$ clustering into k clusters when n is set to 3 are relatively ideal.

Note that in this paper, we measure the clustering utility with NICV through extensive experiments. The underlying reason is that NICV is the objective function. NICV is directly effected by the differentially private clustering. To some extent, we can regard NICV value as the utility of the differentially private algorithm, and NICV can reflect straightforwardly how the noise addition for achieving differential privacy impacts the clustering result. By using multiple datasets and averaging a great number of independently random runs, we expect NICV to measure the clustering utility reasonably. In the future, other reasonable measures can be investigated for fully evaluating the clustering utility.

4.2. Experimental Results.

We now show two groups of experimental results, and make corresponding discussions.

(1) Performance Comparison in term of ϵ

Figures 3 to 8 show the influence of different ϵ values on the clustering results for different datasets. The k values are fixed at 4 or 5. The six figures can be

divided into three groups: Figures 3 and 6, Figures 4 and 7, Figures 5 and 8. The dataset dimensionality values for the first group are in the range [1..10], those of the second group are in the range [11..20], and those of the third group are in the range [21..30].

From these 6 figures, we make the following observations. 1) For all cases, *allocation.k* performs worse than *average.k*, which illustrates that the average allocation of privacy budget seems to be the best choice for differentially private k -means clustering algorithms without cluster merging. 2) Nearly for all cases, *average.nk* performs better than *average.k*, which indicates that cluster merging indeed improves the clustering utility. 3) For all cases, *allocation.nk* performs significantly better than *average.nk*, demonstrating the surprising result that combining both cluster merging and adaptive privacy budget allocation is able to further improve the clustering utility. 4) The performances of all algorithms become better in a similar way as the ϵ value increases, which shows the stability of the improvement over clustering utility. 5) It seems that the dimensionality values have no obvious impact on the utility improvement, and some datasets seems have more effect than others of utility improvement maybe due to their data characteristics.

(2) Performance Comparison in term of k

Figures 9 to 14 show the influence of different k values on the clustering results for different datasets when the ϵ value is fixed. From these figures, we make observations as follows. 1) The result of performance comparisons between these four algorithms is similar to that of Figures 3 to 8, and the ranking of the algorithms in decreasing order of the performance is *allocation.nk*, *average.nk*, *average.k*, and *allocation.k*. 2) For each figure, there is a point of certain k value that has the most effect in utility improvement, which seems to show that this k is the most appropriate number of clusters for the corresponding dataset.

Based on the experimental analysis above, we can conclude that the proposed algorithm combining both cluster merging and adaptive privacy budget allocation improves the clustering utility significantly, and it is superior to the state-of-the-art algorithms.

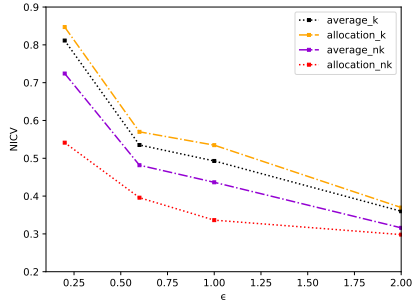


Figure 3: Blood ($k = 4$)

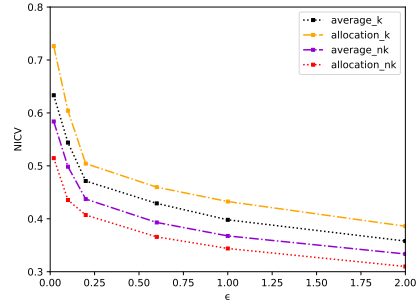


Figure 4: Tripadvisor-review ($k = 4$)

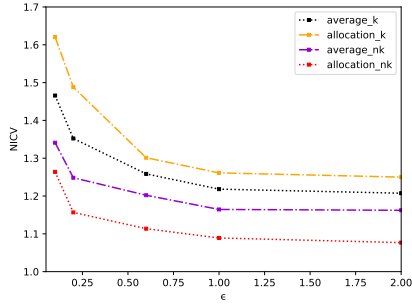


Figure 5: Travel-review ($k = 4$)

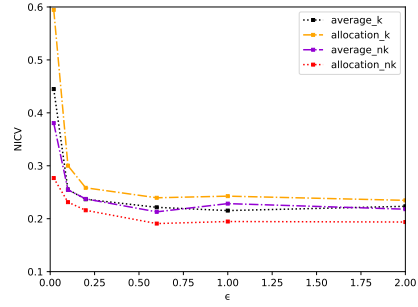


Figure 6: Adult ($k = 5$)

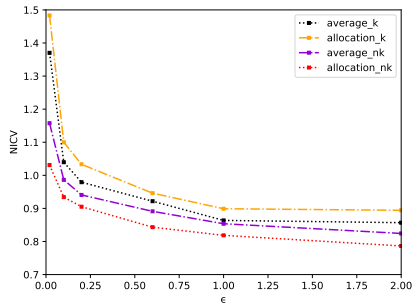


Figure 7: Electrical ($k = 5$)

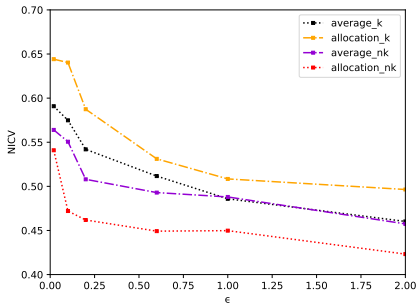


Figure 8: Credit-card ($k = 5$)

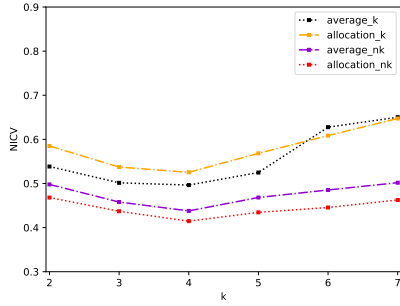


Figure 9: Blood ($\epsilon = 0.6$)

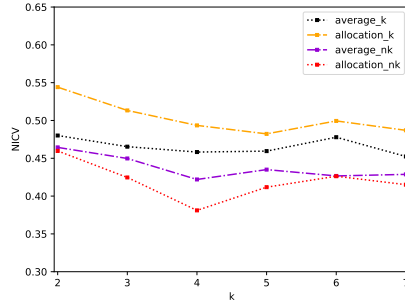


Figure 10: Tripadvisor-review ($\epsilon = 0.6$)

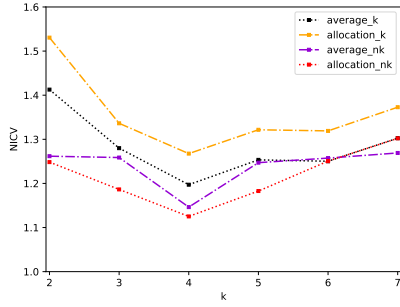


Figure 11: Travel-review ($\epsilon = 0.2$)

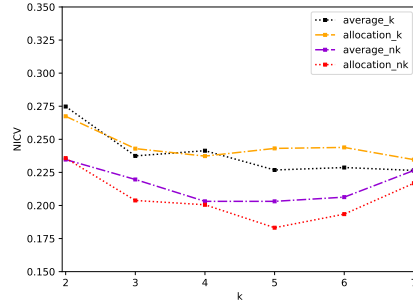


Figure 12: Adult ($\epsilon = 0.2$)

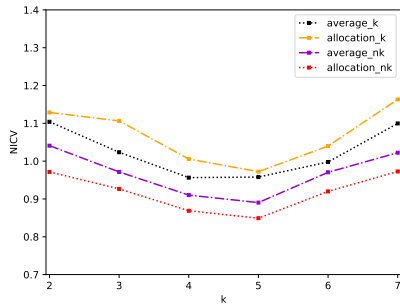


Figure 13: Electrical ($\epsilon = 0.2$)

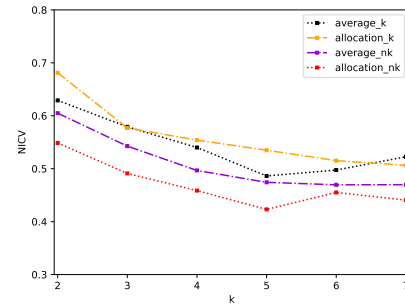


Figure 14: Credit-card ($\epsilon = 0.2$)

We can explain our experimental results with the stability theory by David et al. [30]. According to the definition, given a clustering algorithm, the stability is mainly determined by the probability distribution of the data. Specifically, if the data distribution is of some non-trivial symmetry structure, the algorithm would be unstable; otherwise, it would be stable. When adding differential privacy to the clustering algorithm, the data distribution is flattened. This makes the distribution is more symmetry, and the clustering becomes more unstable. In the extreme, when $\epsilon = 0$, the distribution is uniformly symmetry, and the clustering is completely unstable. Thus, adding differential privacy deteriorates the clustering. When merging the nearest clusters, the clustering algorithm makes the clusters more asymmetry by canceling out noises added, and the clustering becomes more stable. So, we can see that merging clusters improves the clustering utility. Furthermore, in the context of merging clusters, the privacy budget allocation works well. The underlying reason may be that the cluster merging allows more instability in the first several iterations, while allows less in the last ones.

5. CONCLUSION

In this paper, we have proposed a differentially private k -means clustering algorithm based on cluster merging. This algorithm improves the utility of k -means clustering by first partitioning the data into more clusters than required, and then merging the clusters into required number of clusters. It is shown that this cluster merging improves the clustering utility, and when combined with privacy budget allocation, it can further improves the utility. Extensive experiments show that our algorithm outperforms the state-of-the-art algorithms significantly. Besides, we only consider numerical data, without taking into account the non-numerical and mixed data. The future work is to design differentially private k -means clustering algorithms for more types of data.

References

References

- [1] M. S. Chen, J. Han, P. S. Yu, Data mining: an overview from a database perspective, *IEEE Transactions on Knowledge & Data Engineering* 8 (6) (2002) 866–883. doi:<http://dx.doi.org/10.1109/69.553155>.
- [2] S. Tambe, B. Thakur, A. Vishwakarma, Effective data mining using neural networks, *IEEE Transactions on Knowledge & Data Engineering* 8 (6) (2016) 957–961.
- [3] J. Grabmeier, A. Rudolph, Techniques of cluster algorithms in data mining, *Data Mining & Knowledge Discovery* 6 (4) (2002) 303–360.
- [4] K. M. A. Patel, P. Thakral, The best clustering algorithms in data mining, in: *International Conference on Communication & Signal Processing*, IEEE, 2016. doi:<http://dx.doi.org/10.1109/ICCSP.2016.7754534>.
- [5] J. G. Sun, J. Liu, L. Y. Zhao, Clustering algorithms research, *Journal of Software* 19 (19). doi:<http://dx.doi.org/10.3724/SP.J.1001.2008.00048>.
- [6] P. D. Antonenko, S. Toy, D. S. Niederhauser, Using cluster analysis for data mining in educational technology research, *Educational Technology Research & Development* 60 (3) (2012) 383–398. doi:<http://dx.doi.org/10.1007/s11423-012-9235-8>.
- [7] S. Ding, F. Wu, J. Qian, H. Jia, F. Jin, Research on data stream clustering algorithms, *Artificial Intelligence Review* 43 (4) (2015) 593–600. doi:<http://dx.doi.org/10.1007/s10462-013-9398-7>.
- [8] A. Mohebbi, Iterative big data clustering algorithms: a review, *Software Practice & Experience* 46 (1) (2016) 107–129. doi:<http://dx.doi.org/10.1002/spe.2341>.

- [9] K. Nissim, U. Stemmer, Clustering algorithms for the centralized and local models, arXiv preprint arXiv:1707.04766.
- [10] U. Stemmer, Locally private k-means clustering, arXiv preprint arXiv:1907.02513.
- [11] S. Javadi, S. Hashemy, K. Mohammadi, K. Howard, A. Neshat, Classification of aquifer vulnerability using k-means cluster analysis, *Journal of hydrology* 549 (2017) 27–37. doi:<http://dx.doi.org/10.1016/j.jhydro.2017.03.060>.
- [12] G. Han, H. Wang, M. Guizani, S. Chan, W. Zhang, Kclp: A k-means cluster-based location privacy protection scheme in wsns for iot, *IEEE Wireless Communications* 25 (6) (2018) 84–90. doi:<http://dx.doi.org/10.1109/MWC.2017.1800061>.
- [13] P. M. Shakeel, S. Baskar, V. S. Dhulipala, M. M. Jaber, Cloud based framework for diagnosis of diabetes mellitus using k-means clustering, *Health information science and systems* 6 (1) (2018) 16. doi:<http://dx.doi.org/10.1007/s13755-018-0054-0>.
- [14] J. Wu, Y. Wang, Z. Wu, Z. Wang, A. Veeraraghavan, Y. Lin, Deep k -means: Re-training and parameter sharing with harder cluster assignments for compressing deep convolutions, arXiv preprint arXiv:1806.09228.
- [15] M. N. Reza, I. S. Na, S. W. Baek, K.-H. Lee, Rice yield estimation based on k-means clustering with graph-cut segmentation using low-altitude uav images, *Biosystems engineering* 177 (2019) 109–121. doi:<http://dx.doi.org/10.1016/j.biosystemseng.2018.09.014>.
- [16] H. Omrani, B. Parmentier, M. Helbich, B. Pijanowski, The land transformation model-cluster framework: Applying k-means and the spark computing environment for large scale land change analytics, *Environmental modelling & software* 111 (2019) 182–191. doi:<http://dx.doi.org/10.1016/j.envsoft.2018.10.004>.

- [17] J. Xiong, Z. Yao, J. Ma, F. Li, X. Liu, A secure self-destruction scheme with ibe for the internet content privacy, *Chinese Journal of Computers* 37 (1) (2014) 139–150.
- [18] F. Gao, J. He, X. Wu, An approach for tracking privacy disclosure, in: *The 6th International Conference on Networked Computing and Advanced Information Management*, IEEE, 2010, pp. 294–299.
- [19] C. Dwork, Differential privacy, in: *Proceedings of the 33rd international conference on Automata, Languages and Programming-Volume Part II*, Springer-Verlag, 2006, pp. 1–12. doi:http://dx.doi.org/10.1007/11787006_1.
- [20] C. Dwork, F. McSherry, K. Nissim, A. Smith, Calibrating noise to sensitivity in private data analysis, in: *Theory of cryptography conference*, Springer, 2006, pp. 265–284. doi:http://dx.doi.org/10.1007/11681878_14.
- [21] C. Dwork, Differential privacy: A survey of results, *Theory and Applications of Models of Computation* (2008) 1–19doi:http://dx.doi.org/10.1007/978-3-540-79228-4_1.
- [22] F. K. Dankar, K. El Emam, The application of differential privacy to health data, in: *Proceedings of the 2012 Joint EDBT/ICDT Workshops*, ACM, 2012, pp. 158–166. doi:<http://dx.doi.org/10.1145/2320765.2320816>.
- [23] Y. Xiao, L. Xiong, Protecting locations with differential privacy under temporal correlations, in: *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, ACM, 2015, pp. 1298–1309. doi:<http://dx.doi.org/10.1145/2810103.2813640>.
- [24] Q. Yu, Y. Luo, C. Chen, X. Ding, Outlier-eliminated k-means clustering algorithm based on differential privacy preservation, *Applied Intelligence* 45 (4) (2016) 1179–1191. doi:<http://dx.doi.org/10.1007/s10489-016-0813-z>.

- [25] D. Su, J. Cao, N. Li, E. Bertino, H. Jin, Differentially private k-means clustering, in: Proceedings of the sixth ACM conference on data and application security and privacy, ACM, 2016, pp. 26–37. doi:<http://dx.doi.org/10.1145/2857705.2857708>.
- [26] N. Li, L. Min, S. Dong, W. Yang, Differential privacy: From theory to practice, Synthesis Lectures on Information Security Privacy & Trust 8 (4) (2016) 1–138. doi:<http://dx.doi.org/10.2200/S00735ED1V01Y201609SPT018>.
- [27] C. Dwork, A. Roth, et al., The algorithmic foundations of differential privacy, Foundations and Trends® in Theoretical Computer Science 9 (3–4) (2014) 211–407. doi:<http://dx.doi.org/10.1561/04000000042>.
- [28] A. Blum, C. Dwork, F. McSherry, K. Nissim, Practical privacy: The sulq framework, in: Twenty-Fourth ACM Sigmod-Sigact-Sigart Symposium on Principles of Database Systems, ACM, 2005, pp. 128–138. doi:<http://dx.doi.org/10.1145/1065167.1065184>.
- [29] McSherry, Frank, Privacy integrated queries: An extensible platform for privacy-preserving data analysis, Communications of the Acm 53 (9) (2010) 89–97. doi:<http://dx.doi.org/10.1145/1810891.1810916>.
- [30] S. Ben-David, U. Von Luxburg, D. Pál, A sober look at clustering stability, in: International Conference on Computational Learning Theory, Springer, 2006, pp. 5–19.