

Efficient and Privacy-Preserving Multi-User Outsourced K-Means Clustering

Na Li¹, Lianguan Huang², Yanling Li² & Meng Sun²

¹ College of Information Science and Technology, Jinan University, China

² College of Cyber Security, Jinan University, China

Correspondence: Na Li, College of Information Science and Technology, Jinan University, China.

Received: January 22, 2021

Accepted: February 26, 2021

Online Published: March 10, 2021

doi:10.5539/cis.v14n2p26

URL: <https://doi.org/10.5539/cis.v14n2p26>

Abstract

In recent years, with the development of the Internet, the data on the network presents an outbreak trend. Big data mining aims at obtaining useful information through data processing, such as clustering, clarifying and so on. Clustering is an important branch of big data mining and it is popular because of its simplicity. A new trend for clients who lack of storage and computational resources is to outsource the data and clustering task to the public cloud platforms. However, as datasets used for clustering may contain some sensitive information (e.g., identity information, health information), simply outsourcing them to the cloud platforms can't protect the privacy. So clients tend to encrypt their databases before uploading to the cloud for clustering. In this paper, we focus on privacy protection and efficiency promotion with respect to k-means clustering, and we propose a new privacy-preserving multi-user outsourced k-means clustering algorithm which is based on locality sensitive hashing (LSH). In this algorithm, we use a Paillier cryptosystem encrypting databases, and combine LSH to prune off some unnecessary computations during the clustering. That is, we don't need to compute the Euclidean distances between each data record and each clustering center. Finally, the theoretical and experimental results show that our algorithm is more efficient than most existing privacy-preserving k-means clustering.

Keywords: k-means clustering, privacy protection, homomorphic encryption, locality sensitive hashing

1. Introduction

Clustering analysis is one of the most commonly used tasks in data mining area (Kriegel et al., 2009). It is worth noting that our clustering analysis is very different from clustering coefficient (Li, Y. , Shang, Y. , & Yang, Y. ,2017). The former is a method for data analysis, while the latter is an important concept in network topology, describing the degree to which vertices in a graph cluster together. The traditional clustering analysis methods mainly include hierarchical clustering, partitioning clustering, density-based clustering, grid-based clustering and model-based clustering algorithm. Actually, k-means clustering is a kind of partitioning clustering, it is popular owing to its simple characteristics. In addition, k-means algorithm is of great importance in various fields, including image retrieval (Yin & Zhang, 2017), machine learning, pattern recognition, social participatory sensing (Xing, Hu, Yu, Cheng & Zhang, 2017), knowledge discovery and text mining (MUstafi & Sahoo, 2019).

K-means clustering (Jain, Murty & Flynn, 1999) is a typical distance-based clustering algorithm, whose goal is to assign each data record into a cluster which has the shortest distance from it. As a result, similar records will be classified into the same group while different records into different groups. But as the amount of data increases, clients with limited storage and computation resources usually prefer to outsource their datasets and clustering task to the cloud service providers. However, the cloud platform can't be fully trusted, and they may want to gather clients' sensitive information. Thus, we prefer to encrypt the datasets before sending them to the cloud, on which performs clustering without getting any useful information in the end.

The original privacy-preserving k-means clustering schemes (Xia, Hua, Tong & Zhong, 2020; Yu, Luo, Chen & Ding, 2016; Chen, Wang, Zhang, Zhong & Chen, 2018; Su, Cao, Li, Bertiino, Lyu & Jin, 2017) protected the data privacy by adding noises. Then some scholars constructed privacy schemes for k-means clustering (Jiang, Guo, Jin, Lv, Wu, Liu, Fang, Yiu, & Wang, 2020; Rong, Wang, Liu, Hao & Xian, 2017; Yuan & Tian, 2017; Zou, Zhao, Shi, Wang, Peng, Ping & Wang, 2020) with different encryption cryptosystems to improve the accuracy of k-means clustering, but also led to low efficiency. To address this problem, we propose a new

privacy-preserving multi-user outsourced k-means algorithm, which is based on locality sensitive hashing (LSH) (Datar, Immorlica, Indyk & Mirrokni, 2004) for pruning off unnecessary computations in clustering. Our contributions are showed as follows:

- 1) We propose a LSH-based privacy-preserving multi-user outsourced k-means clustering (LSH-PPMOC) algorithm, whose main idea is to prune off the unnecessary computations during clustering task.
- 2) We combine the privacy-preserving techniques with LSH to guarantee the security of users' datasets and clustering results.
- 3) The experiments on real-world datasets show that we improve the clustering efficiency greatly.

The remainder of this paper is organized as follows. In Section 2, we discuss the existing related work. Section 3 presents some definitions and properties related to k-means clustering algorithm, Paillier cryptosystem and Locality Sensitive Hashing as a background. In Section 4, we discuss our proposed LSH-PPMOC solutions in detail. Also, we analyze the security guarantees and complexities of our solution in Section 5. Section 6 presents our experimental results on a real-world dataset under different parameter settings. Finally, we conclude the paper along with the scope for future research in Section 7.

2. Related Work

Privacy-preserving k-means clustering has been widely used in data mining. The private data mining algorithms can be categorized as, 1) data perturbation based, 2) secure multiparty computation (SMC) based and 3) cryptography technologies based. In perturbation, data records are randomized by adding noises, and the work by (Liu, Kargupta & Ryan, 2005) showed that the clustering results based on multiplicative perturbation only had below 5% error rate compared to the results on the original data. However, data perturbation cannot guarantee privacy in any formal sense (Liu, Giannella & Kargupta, 2006; Wong, Cheung, Kao & Mamoulis, 2009). For example, if an adversary gets a few data records in plaintext, he may recover the rest records though they are perturbed (Liu, Giannella & Kargupta, 2006).

Algorithms based on secure multiparty computation (SMC) (Goldreich, 2009) can preserve the security and privacy of users' data. There existed literatures utilizing SMC like (Lindell & Pinkas, 2008; Mohassel, Rosulek & Trieu, 2020; Upmanyu, Namboodiri, Srinathan & Jawahar, 2010), in which algorithms were proposed to perform distributed data mining without revealing private inputs of participants. In addition, Clifton et al. (Clifton, Kantarcioglu, Vaidya, Lin & Zhu, 2002) proposed that a relatively small set of cryptography primitive should be utilized to generate the SMC protocol. However, these solutions suffered from high levels of communication and computing cost. We can conclude that, though the private clustering algorithms based on SMC can provide better security protection, it is too costly to applied in practical applications (Jagannathan & Wright, 2005).

There are also literatures those propose privacy-preserving clustering algorithms using the semantically secure additive or multiplicative homomorphic encryption schemes. Liu et al. (Liu, Bertino & Yi, 2014) first leveraged fully homomorphic encryption technique (Gentry, 2009) to perform clustering outsourcing. However, the encryption scheme adopted in their scheme was not secure according to (Wang, 2015). Besides, this algorithm needed clients to participate and provide some information during each iteration, thus leading to heavy overhead on clients. To reduce the interaction with users, Almutairi et al. (Almutairi, Coenen & Dures, 2017) presented an efficient mechanism by using the concept of an Updatable Distance Matrix (UDM). But, their work revealed partial privacy to the cloud servers, such as the size of each cluster and the distances between data objects and centroids.

Similarly, Samanthulat et al. (Samanthula, Rao, Bertino, Yi & Liu, 2014) proposed a secure and outsourced k-means clustering scheme using the Paillier cryptosystem. Unfortunately, this scheme has a high computation cost because of using the bit array-based comparison. Then Rong et al. (Rong, Wang, Liu, Hao & Xian, 2017) put forward a privacy outsourced k-means clustering under multiple keys based on public key cryptosystem with double decryption (Youn, Park, Kim & Lim, 2005). However, the addition protocol in their scheme was not secure because the assistant server could extract the ratio of messages. To address this problem, (Zou, Zhao, Shi, Wang, Peng, Ping & Wang, 2020) proposed a highly secure outsourced k-means clustering scheme using BCP cryptosystem which had the additive homomorphic property.

To better the clustering performance in cloud computing environment, many scholars proposed MapReduce (Dean & Ghemawat, 2008) based k-means clustering schemes to handle large-scale dataset in parallel (Cui, Zhu, Yang, Li & Ji, 2014; Sardar & Ansari, 2020), but they all didn't consider privacy protection. Yuan et al. (Yuan & Tian, 2017) proposed a privacy-preserving scheme using a lightweight cryptosystem basing on the hardness of

learning with errors (LWE) (Brakerski, Gentry & Halevi, 2013), and incorporated MapReduce to improve the efficiency. However, the scheme was not fully outsourced and didn't support multiple users. Similarly, in (Rong, Wang, Liu, Hao & Xian, 2017), the clustering was executed under the Spark framework.

Besides, there are also other ways to improve efficiency for k-means clustering. Bhaskara et al. (Bhaskara & Wijewardena, 2018) proposed a variant of Locality sensitive hashing (LSH) (Indyk & Motwani, 1998) to speed up the clustering. Similarly, (Li, Wang, Wang, Hu, Li & Li, 2014) utilized the local sensitive property of LSH to prune off unnecessary computations during clustering and carried out the experiment with the help of MapReduce. However, the both schemes did not take privacy into consideration.

Inspired by literature (Li, Wang, Wang, Hu, Li & Li, 2014), we propose a new and efficient LSH-based privacy-preserving outsourced k-means clustering. We consider a scenario, in which there are two cloud service providers and multiple users. We explicitly assume the two cloud servers will never collude during clustering, then our proposed scheme protects all users' data confidentiality.

3. Preliminaries

In this section, we will first introduce the definition of typical k-means clustering. Then we give a brief overview of additive homomorphic Paillier cryptosystem and some basic cryptographic preliminaries for securely performing clustering. At last, we briefly introduce LSH which we use as the basis for our solution.

3.1 K-means Clustering

K-means clustering is an unsupervised clustering algorithm, and it is widely used in various application scenarios. Supposing given d-dimension data records t_1, \dots, t_n , the goal of k-means clustering algorithm is to divide these records into k disjoint groups. That is to say, the objects in the same group are similar, while objects in different groups have low similarity. We use c_1, \dots, c_k to denote k clusters, and μ_1, \dots, μ_k , the corresponding centroids. To cluster one data record $t_i, i \in [1, n]$ into correct cluster, what we should do first is to compute the distance between it and k centroids $\mu_j, j \in [1, k]$. We use squared Euclidean distance here, which is given as

$$Dist(t_i, \mu_j) = \sum_{m=1}^l [(t_i[m] - \mu_j[m])]^2. \quad (1)$$

Then we assign t_i to the cluster c_h , if t_i has the smallest distance with $\mu_h, h \in [1, k]$ among k distances.

A traditional k-means clustering has three stages, (1) Initialization; (2) Clustering; (3) Updating new centroids. For Stage (1), initial k records are selected randomly as cluster centroids μ_1, \dots, μ_k . In Stage (2), we first compute k distances between every one record t_i and k centroids, then cluster them according to the distances. Later, new centroids μ_i' can be derived as the mean values of attributes of records belonging to c_i in Stage (3). If the given maximum iteration number is not reached or the cluster results differ from that getting before, next iteration will restart from Stage (2) to (3) with the new centroids, otherwise, the algorithm terminates.

3.2 The Paillier Cryptosystem

In this paper, we use the Paillier encryption (Paillier, 1999), which is a probabilistic asymmetric encryption scheme. Without loss of generality. Let $Enc(\cdot)$ and $Dec(\cdot)$ denote the encryption and decryption function under the Paillier cryptosystem and N denote the RSA modulus.

We denote the encryption scheme as a triple $\{KeyGen, Enc, Dec\}$.

$$KeyGen(1^\kappa) \rightarrow (pk, sk)$$

Choose two large prime numbers p and q which satisfy that $\gcd(pq, (p-1)(q-1)) = 1$.

Calculate $N = pq$, and $\lambda = \text{lcm}(p-1, q-1)$.

Randomly choose an integer $g \in \mathbb{Z}_{(N^2)}$.

Check whether there exists $u = [(L(g^\lambda \text{ mod } N^2))]^{-1} \text{ mod } N$ where function L is $L(\mu) = (\mu - 1)/N$. Then pk is (N, g) and sk is (λ, μ) .

$$Enc(pk, x) \rightarrow c$$

Select a random value $r \in \mathbb{Z}_N^*$ for the message x and the ciphertext is $c = g^x r^N \text{ mod } N^2$.

$$Dec(sk, c) \rightarrow x$$

Decrypt the message by $x = L(c^\lambda \text{ mod } [N^2]) \mu \text{ mod } n$.

Then, for any $a, b \in Z_N$, the encryption scheme is additive homomorphic:

$$E(a + b) = E(a) \times E(b) \bmod N^2 = E(a + b \bmod N),$$

$$\llbracket E(a) \rrbracket^{b \bmod N^2} = E(a \cdot b \bmod N)$$

Where the symbol " \times " denotes the multiplication on ciphertexts, and " \cdot " means the multiplication on plaintext. Besides, we will omit the term $\bmod N^2$ from homomorphic operations in the rest of the paper.

3.3 Basic Cryptographic Primitives

In this part, we propose some cryptographic primitives as the basis of our methods.

(1) Secure Multiplication (SM) Protocol

Given that Alice holds $\langle E_{\llbracket pk \rrbracket_u}(x), E_{\llbracket pk \rrbracket_u}(y) \rangle$, and Bob holds $\llbracket sk \rrbracket_u$, the goal of Alice and Bob is to securely compute $E_{\llbracket pk \rrbracket_u}(x \cdot y)$. During the execution of SM, no information regarding the contents of x and y is revealed to Alice or Bob. More details are shown in Algorithm 1.

Algorithm 1: $SM(E_{\llbracket pk \rrbracket_u}(x), E_{\llbracket pk \rrbracket_u}(y)) \rightarrow E_{\llbracket pk \rrbracket_u}(x \cdot y)$

Input: Alice has $E_{\llbracket pk \rrbracket_u}(x), E_{\llbracket pk \rrbracket_u}(y)$, Bob has $\llbracket sk \rrbracket_u$

Output: $E_{\llbracket pk \rrbracket_u}(x \cdot y)$

1. Alice: (a) Picks any two different numbers $r_x, r_y \in Z_N$ randomly.

(b) Computes $x' \leftarrow E_{\llbracket pk \rrbracket_u}(x) \times E_{\llbracket pk \rrbracket_u}(r_x)$
 $y' \leftarrow E_{\llbracket pk \rrbracket_u}(y) \times E_{\llbracket pk \rrbracket_u}(r_y)$

(c) Sends x', y' to Bob

2. Bob: (a) Computes $h_x \leftarrow Dec(\llbracket sk \rrbracket_u, x')$, $h_y \leftarrow Dec(\llbracket sk \rrbracket_u, y')$, $h \leftarrow h_x \cdot h_y \bmod N$,

$h' \leftarrow Enc(\llbracket pk \rrbracket_u, h)$

(b) Sends h' to Alice

3. Alice: (a) Computes $s \leftarrow h' \times \llbracket E_{\llbracket pk \rrbracket_u}(x) \rrbracket^{(N - r_y)}$,
 $s' \leftarrow s \times \llbracket E_{\llbracket pk \rrbracket_u}(y) \rrbracket^{(N - r_x)}$

(b) Gets $E_{\llbracket pk \rrbracket_u}(x \cdot y) \leftarrow s' \llbracket E_{\llbracket pk \rrbracket_u}(r_x r_y) \rrbracket^{(N - 1)}$

(2) Secure Squared Euclidean Distance (SSED) Protocol

For k-means algorithm, there are many ways to compute the similarity scores between the data record $\mathbf{t}_i, 1 \leq i \leq n$ and cluster centroid $\boldsymbol{\mu}_j, 1 \leq j \leq k$, such as Euclidean distance, Cosine similarity and Jaccard coefficient. In this paper, we use squared Euclidean distance for its simplicity, denoted by $\|\mathbf{t}_i - \boldsymbol{\mu}_j\|^2$.

Note that $\boldsymbol{\mu}_j$ is a vector, and its components may be rational numbers. However, the ring Z_N doesn't support rational division operation, so we use a new form of expression to represent the cluster center like (Rong, Wang, Liu, Hao & Xian, 2017). Let $\langle \mathbf{s}_j, |c_j| \rangle$ denotes the new form of cluster center, where \mathbf{s}_j and $|c_j|$ represent the sum vector and the total number of records belonging to cluster c_j , respectively. Now we define $\Omega_{\llbracket i, j \rrbracket}$ as the scaled squared Euclidean distance between \mathbf{t}_i and $\boldsymbol{\mu}_j$, which satisfies that $\|\mathbf{t}_i - \boldsymbol{\mu}_j\| = \sqrt{(\Omega_{\llbracket i, j \rrbracket}) / (|c_j|)}$. So $\Omega_{\llbracket i, j \rrbracket}$ can be calculated as follows:

$$\Omega_{\llbracket i, j \rrbracket} = \llbracket (\|\mathbf{t}_i - \boldsymbol{\mu}_j\| \cdot |c_j|) \rrbracket^2 = \sum_{l=1}^d \llbracket (|c_j| \cdot \mathbf{t}_i[l] - \mathbf{s}_j[l]) \rrbracket^2, \quad (2)$$

where $1 \leq i \leq m, 1 \leq j \leq k$, m and d are the total number of objects and dimension size, respectively. Taking $E_{\llbracket pk \rrbracket_u}(\mathbf{t}_i)$ and $\langle E_{\llbracket pk \rrbracket_u}(\mathbf{s}_j), E_{\llbracket pk \rrbracket_u}(|c_j|) \rangle$ as inputs, the output of SSED protocol is $\langle E_{\llbracket pk \rrbracket_u}(\Omega_{\llbracket i, j \rrbracket}), E_{\llbracket pk \rrbracket_u}(|c_j|) \rangle$ which is actually the distance between record \mathbf{t}_i and the centroid of j -th cluster. We present this protocol in Algorithm 2.

Algorithm 2: SSED

Input: Alice has $E_{-}(\llbracket pk \rrbracket_{-u})(t_i)$, $E_{-}(\llbracket pk \rrbracket_{-u})(\mu_j) = (E_{-}(\llbracket pk \rrbracket_{-u})(s_i), E_{-}(\llbracket pk \rrbracket_{-u})(|c_j|))$, $i \in [1, n], j \in [1, k]$,

Bob has $\llbracket sk \rrbracket_{-u}$

$E_{-}(\llbracket pk \rrbracket_{-u})(\Omega(i, j)) = E_{-}(\llbracket pk \rrbracket_{-u})(0)$

for $l = 1$ to d

Alice and Bob jointly compute $\Gamma \leftarrow SM(E_{-}(\llbracket pk \rrbracket_{-u})(|c_j|), E_{-}(\llbracket pk \rrbracket_{-u})(t_i[l]))$,

$\Gamma^{*} \leftarrow \Gamma \times E_{-}(\llbracket pk \rrbracket_{-u})(s_j[m])^{(N-1)}$, $\Gamma^{**} \leftarrow SM(\Gamma^{*}, \Gamma^{*})$

$E_{-}(\llbracket pk \rrbracket_{-u})(\Omega(i, j)) = E_{-}(\llbracket pk \rrbracket_{-u})(\Omega(i, j)) \times \Gamma^{**}$

end for

Alice and Bob jointly compute $E_{-}(\llbracket pk \rrbracket_{-u})(|c_j|) \leftarrow SM(E_{-}(\llbracket pk \rrbracket_{-u})(|c_j|), E_{-}(\llbracket pk \rrbracket_{-u})(|c_j|))$

Output $(E_{-}(\llbracket pk \rrbracket_{-u})(\Omega(i, j)), E_{-}(\llbracket pk \rrbracket_{-u})(|c_j|))$

(3) Secure Comparison (SC) Protocol

Given that Alice has two encrypted values $E_{-}(\llbracket pk \rrbracket_{-u})(x)$, $E_{-}(\llbracket pk \rrbracket_{-u})(y)$, and Bob has sk_{-u} , the goal of SC protocol is to securely compare x and y without knowing both of them. To be specific, Alice generates an arbitrary value r firstly, and then encrypts it as $E_{-}(pk_{-u})(r)$. After that, Alice computes $E_{-}(\llbracket pk \rrbracket_{-u})(t_0) = E_{-}(\llbracket pk \rrbracket_{-u})(x + r)$ and $E_{-}(\llbracket pk \rrbracket_{-u})(t_1) = E_{-}(\llbracket pk \rrbracket_{-u})(y + r)$ utilizing the additive homomorphic property of the Paillier cryptosystem. Alice flips a coin b where $b = 0/1$, then sends two ciphertexts in the order of $E_{-}(\llbracket pk \rrbracket_{-u})(t_b)$, $E_{-}(\llbracket pk \rrbracket_{-u})(t_{(1-b)})$ to Bob. Bob decrypts with sk_{-u} and gets t_b and $t_{(1-b)}$. If $t_b < t_{(1-b)}$, Bob sends $E_{-}(\llbracket pk \rrbracket_{-c})(\gamma) = E_{-}(\llbracket pk \rrbracket_{-c})(1)$ to Alice, else sends $E_{-}(\llbracket pk \rrbracket_{-c})(\gamma) = E_{-}(\llbracket pk \rrbracket_{-c})(0)$. This protocol will return $E_{-}(\llbracket pk \rrbracket_{-c})(\gamma)$ if $b = 0$, otherwise $SBN(E_{-}(\llbracket pk \rrbracket_{-c})(\gamma))$.

(4) Privacy-preserving Minimum (PMIN) Protocol

In this protocol, we aim at comparing two squared Euclidean distances and output the minimum one. To achieve this, we use the privacy-preserving maximum protocol proposed by Liu et al. (Liu, Lu, Ma, Chen & Qin, 2015). However, we will give a slight change to it to transfer the maximum value to the minimum value.

To be specific, Alice has two encrypted squared Euclidean distances between one record $t_i, i \in [1, m]$ and two clusters c_a and c_b ,

$E_{-}(\llbracket pk \rrbracket_{-u})(d_a) = (E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_a), E_{-}(\llbracket pk \rrbracket_{-u})(|c_a|)), E_{-}(\llbracket pk \rrbracket_{-u})(d_b) = (E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_b), E_{-}(\llbracket pk \rrbracket_{-u})(|c_b|))$, their secrets are

$E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_a) = E_{-}(\llbracket pk \rrbracket_{-u})(a)$, $E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_b) = E_{-}(\llbracket pk \rrbracket_{-u})(b)$. Bob has sk_{-u} . After finishing the protocol, Alice will get $E_{-}(\llbracket pk \rrbracket_{-u})(d_{min})$ whose corresponding distance is relatively small, along with $E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_{min})$.

To get d_{min} , Alice first computes

$E_{-}(\llbracket pk \rrbracket_{-u})(\xi_1) = SM(E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_a), E_{-}(\llbracket pk \rrbracket_{-u})(|c_b|^2))$, $E_{-}(\llbracket pk \rrbracket_{-u})(\xi_2) = SM(E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_b), E_{-}(\llbracket pk \rrbracket_{-u})(|c_a|^2))$, then

$E_{-}(\llbracket pk \rrbracket_{-u})(H_1) = [E_{-}(\llbracket pk \rrbracket_{-u})(\xi_1)]^2 \times E_{-}(\llbracket pk \rrbracket_{-u})(1) = E_{-}(\llbracket pk \rrbracket_{-u})(2\Omega_a \cdot |c_b|^2 + 1)$, $E_{-}(\llbracket pk \rrbracket_{-u})(H_2) = E_{-}(\llbracket pk \rrbracket_{-u})(2\Omega_b \cdot |c_a|^2)$. After that, Alice chooses some

random values $R, r_1, r_2, r_3, r_4 \in \mathbb{Z}_N^*$, $\theta \in G$, $2^L < N^2$. If θ is an odd number, Alice computes $C_1 = [E_{-}(\llbracket pk \rrbracket_{-u})(H_1)]^R \times [E_{-}(\llbracket pk \rrbracket_{-u})(H_2)]^{(N-R)} \times g^{(2^L)} [r_4]^N = E_{-}(\llbracket pk \rrbracket_{-u})(R(H_1 - H_2) + 2^L)$,

$C_2 = E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_a - \Omega_b + r_1)$, $C_3 = E_{-}(\llbracket pk \rrbracket_{-u})(|c_a|^2 - |c_b|^2 + r_2)$, $C_4 =$

$E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_a - \lambda_b + r_3)$. else if θ is even, Alice computes
 $C_1 = \llbracket (E_{-}(\llbracket pk \rrbracket_{-u})(H_2)) \rrbracket^R \times \llbracket (E_{-}(\llbracket pk \rrbracket_{-u})(H_1)) \rrbracket^{(N-R) \times g^{(2^L)}} \llbracket r_4 \rrbracket^N =$
 $E_{-}(\llbracket pk \rrbracket_{-u})(R(H_2 - H_1) + 2^L), C_2 = E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_b - \Omega_a + r_1), C_3 =$
 $E_{-}(\llbracket pk \rrbracket_{-u})(|c_b|^2 - \llbracket |c_a| \rrbracket^2 + r_2), C_4 = E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_b - \lambda_a + r_3)$. Then
 C_1, C_2, C_3 and C_4 are sent to Bob.

After receiving C_1, C_2, C_3 and C_4 , Bob generates $r_5, r_6 \in Z_N$ randomly, and decrypts C_1 to get M .
 If $M > 2^L$, assuming $\alpha = 0$, Bob computes
 $E_{-}(\llbracket pk \rrbracket_{-u})(\alpha)$, $A = E_{-}(\llbracket pk \rrbracket_{-u})(0)$, $B = E_{-}(\llbracket pk \rrbracket_{-u})(0)$, $C = E_{-}(\llbracket pk \rrbracket_{-u})(0)$, else if
 $M < 2^L$, assuming $\alpha = 1$, Bob computes $E_{-}(\llbracket pk \rrbracket_{-u})(\alpha)$, $A = C_2 \times r_5^N$, $B = C_3 \times r_6^N$, $C =$
 $C_4 \times r_7^N$. After that, Bob sends $E_{-}(\llbracket pk \rrbracket_{-u})(\alpha)$, A, B, C back to Alice.

For Alice, if θ is odd, she will compute
 $E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_{min}) = E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_b) \times A \times \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\alpha) \rrbracket^{(N-r_1)}$,
 $E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c| \rrbracket_{min}) =$

$E_{-}(\llbracket pk \rrbracket_{-u})(|c_b|^2) \times C \times \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\alpha) \rrbracket^{(N-r_2)}$,
 $E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_{min}) = E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_b) \times B \times \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\alpha) \rrbracket^{(N-r_3)}$. Else if θ is
 even, Alice will compute

$E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_{min}) =$
 $E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_a) \times A \times \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\alpha) \rrbracket^{(N-r_1)}$, $E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c| \rrbracket_{min}) =$
 $E_{-}(\llbracket pk \rrbracket_{-u})(|c_a|^2) \times C \times \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\alpha) \rrbracket^{(N-r_2)}$,

$E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_{min}) = E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_b) \times B \times \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\alpha) \rrbracket^{(N-r_3)}$. Then the output
 of this protocol is $\llbracket E_{-}(\llbracket pk \rrbracket_{-u})(d) \rrbracket_{min} = (E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_{min}), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c| \rrbracket_{min}),$
 $E_{-}(\llbracket pk \rrbracket_{-u})(\lambda_{min}))$.

Algorithm 3: ESDC

Input: Alice has
 $(E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega \rrbracket_a), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_a| \rrbracket)), (E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega \rrbracket_b), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_b| \rrbracket)), \llbracket sk \rrbracket_c$, Bob
 has $\llbracket sk \rrbracket_u, \llbracket pk \rrbracket_c$
 Output: $E_{-}(\llbracket pk \rrbracket_c)(\lambda_2)$

Alice and Bob jointly compute:

$E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_b| \rrbracket^2) \leftarrow SM(E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_b| \rrbracket), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_b| \rrbracket)),$
 $E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_a| \rrbracket^2) \leftarrow SM(E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_a| \rrbracket), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_a| \rrbracket))$
 $E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_a) \leftarrow SM(E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_a), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_b| \rrbracket^2)), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_b) \leftarrow$
 $SM(E_{-}(\llbracket pk \rrbracket_{-u})(\Omega_b), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_a| \rrbracket^2))$
 $E_{-}(\llbracket pk \rrbracket_{-u})(\theta) \leftarrow SM(E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_b| \rrbracket^2), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket |c_a| \rrbracket^2)), E_{-}(\llbracket pk \rrbracket_{-u})(\theta') \leftarrow$
 $\llbracket (SM(E_{-}(\llbracket pk \rrbracket_{-u})(\theta), E_{-}(\llbracket pk \rrbracket_{-u})(\epsilon))) \rrbracket^{(N-2)}$
 $E_{-}(\llbracket pk \rrbracket_c)(\lambda_1) \leftarrow SC(E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_a), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_b))$

Alice decrypts $\lambda_1 \leftarrow Dec(\llbracket sk \rrbracket_c, E_{-}(\llbracket pk \rrbracket_c)(\lambda_1))$

if $\lambda_1 = 1$, computes

$E_{-}(\llbracket pk \rrbracket_c)(\lambda_2) =$
 $SC(E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_b) \times \llbracket (E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_a)) \rrbracket^{(N-1)}, E_{-}(\llbracket pk \rrbracket_{-u})(\theta'))$

else, computes

$E_{-}(\llbracket pk \rrbracket_c)(\lambda_2) =$
 $SC(E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_a) \times \llbracket (E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket \Omega' \rrbracket_b)) \rrbracket^{(N-1)}, E_{-}(\llbracket pk \rrbracket_{-u})(\theta'))$

(5) Enhanced Secure Distance Comparison (ESDC) Protocol

However, the PMIN protocol can't meet the requirement we need for pruning strategies in section 3.5.3, so we design a protocol, enhanced secure distance comparison (ESDC) protocol, the goal of which is to compare the absolute value of the difference between two encrypted Euclidean distance, $E_{pk}(d_a)$, $E_{pk}(d_b)$ and an encrypted constant value ε , where $E_{pk}(d_a) = (E_{pk}(\Omega_a), E_{pk}(c_a))$, $E_{pk}(d_b) = (E_{pk}(\Omega_b), E_{pk}(c_b))$. As a result, we can get λ , where $\lambda = 1$ means $|d_a - d_b| > \varepsilon$, and vice versa. We show it in Algorithm 3.

(6) Privacy-Preserving Minimum Among k Distances ($PMIN_k$) Protocol

The goal of $PMIN_k$ protocol is to gain the encrypted minimum distance and its corresponding secret from k encrypted square Euclidean distances. Assume Alice has k distances in ciphertext along with their secrets $\{E_{pk}(d_1), E_{pk}(1), \dots, E_{pk}(d_k), E_{pk}(k)\}$, where $E_{pk}(d_i) = (E_{pk}(\Omega_i), E_{pk}(c_i))$, $1 \leq i \leq k$ and Bob has the corresponding secret key sk_u . To get the minimum value, Alice and Bob need to execute PMIN with two inputs in a sequential fashion. This process is straightforward, so we omit it here.

3.4 Locality Sensitive Hashing

Locality sensitive hashing (LSH) was introduced in the seminal work (Indyk & Motwani, 1998), which gave an efficient algorithm for nearest neighbor searching in high dimensional space. The main idea of LSH is to use a few special hash functions, called function family to map objects into different buckets, and the objects that are close to each other may be mapped into the same bucket (collision) with higher probability than objects those are far apart. Assume S is the domain of objects, and D is the distance measure between objects.

Definition 1. A function family $\mathcal{H} = \{h: S \rightarrow U\}$ is called (r_1, r_2, p_1, p_2) -sensitive for D if for any $a, b \in S$, it holds

if $\|a - b\| \leq r_1$, then $p[h(a) = h(b)] \geq p_1$

if $\|a - b\| \geq r_2$, then $p[h(a) = h(b)] \leq p_2$

where p_1, p_2 are both probability value, and $0 < r_1 \leq r_2, 0 \leq p_2 < p_1 \leq 1$.

We know different LSH families are adopted for different distance functions. (Datar, Immorlica, Indyk & Mirrokni, 2004) has proposed LSH families for L_p norms based on p -stable distributions. If p is 1, it is Manhattan distance, while p is equal to 2, it is Euclidean distance. As we use Euclidean distance in this paper, we use the hash function in Eq. (3).

$$h(a, b)(v) = \lfloor (a \cdot v + b) / r \rfloor \quad (3)$$

3.5 LSH Based Privacy-Preserving Multi-User Outsourced K-Means Clustering for Massive Datasets

This section discusses a new LSH based privacy-preserving multi-user outsourced k-means clustering algorithm (LSH-PPMOC) which includes two optimization strategies for large-scale high-dimensional data clustering. Our target is to minimize the cost in the process for k-means clustering while we won't reveal any information about users' data records or distances between data records and centroids. Firstly, we use a LSH function family which is determined by all users to get the data skeleton, where the similar points are reduced to the same bucket. Then we can save unnecessary computations or comparisons between k encrypted distances during the work through performing clustering on the data skeleton and utilizing efficient strategies.

3.5.1 Data Skeleton

In this section, we will firstly briefly introduce data skeleton proposed by Li et al. (Li, Wang, Wang, Hu, Li & Li, 2014), and then illustrate what we can do on data skeleton.

Each element of data skeleton is a 2-tuple (p_r, L_p) , where p_r represents a list of points in L_p . To get data skeleton, we first use m hash functions to map the original datasets into different buckets, where the bucket value is an m -dimension value as the result of m hash functions. Then we select the center point of each bucket as a representative data point p_r . Initially, we set L_p for each representative data point to null. The distances

between p_r and all other points in this bucket are computed, and if the distance is smaller than d , the according point will be added into L_p , where d is a user-defined threshold. For the point which further away from p_r is also an element of data skeleton, whose corresponding L_p is null. An example is given as Figure 1. It is worth noting that all the data skeletons in each bucket form a hash table(HT).

It is well known that, as the dataset grows, the cost of distance computing between each point and centroid pair becomes unbearable, especially when the number of clusters is large. And we know that the points in one data skeleton are close in distance, and they may belong to the same cluster with high probability. Then we might save some unnecessary calculations based on this.

We demonstrate that in Figure 2. We say that c_1 and c_2 are two center points, p_1 and p_2 are two different data records, r_1, r_2 are the distances between p_1 and c_1, c_2 , while r_1', r_2' are distances of p_2 , d is the distance between p_1 and p_2 . If we know that r_1 is much smaller than r_2 , and d is near to zero, it is very likely that $r_1' < r_2'$, which means that p_2 shares the same cluster with p_1 . We can generalize it to Theorem 1.

Theorem 1. *Given c_1 and c_2 as two center points, p_1 and p_2 are two points with the distance d . r_1, r_1', r_2 and r_2' are distances between p_1, p_2 and c_1, c_2 respectively. If $r_1 < r_2$ and $r_2 - r_1 > 2 * d$, then it holds that $r_1' < r_2'$.*

The proof was given in (Li, Wang, Wang, Hu, Li & Li, 2014), and we wouldn't go into details here.

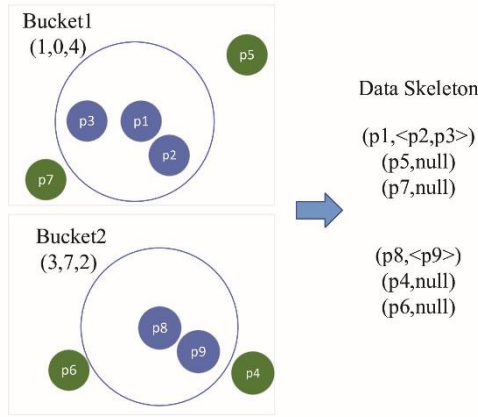


Figure 1. Data Skeleton

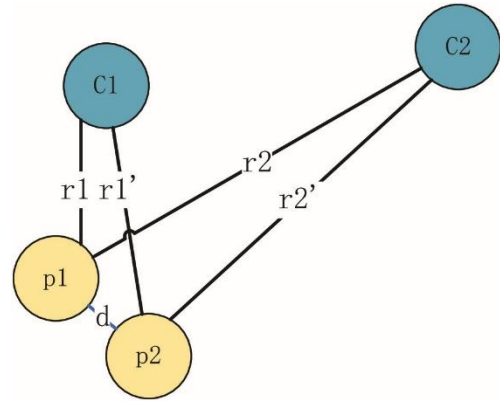


Figure 2. Points and Centers

3.5.2 Privacy-Preserving Pruning Strategies Using LSH

In this section, we explain how we make use of the low bound property of LSH to prune off the unnecessary computations during privacy k-means clustering.

There are two pruning strategies, and both of them are based on Theorem 1. The first one aims at reducing the number of data records that need to find the nearest centers. For each representative point p_r in one bucket, it represents the points in L_p , and the distance between any p_i and p_r is less than d . We first compute k distance between p_r and k centroids $\{c_1, \dots, c_k\}$ using SSED protocol. Let $d(p_r, c_1)$ and $d(p_r, c_2)$ be the smallest and the second smallest distances among k distances. If $d(p_r, c_2) - d(p_r, c_1) \geq 2d$, we can say that all points p_i represented by p_r have a shortest distance with c_1 . In this case, we need not compute the distances between p_i and k centroids in this iteration.

The second strategy is to reduce the number of centroids to be compared for each point p_i in L_p . Let's think about a situation, we have computed k distances for p_r , represented by r_1, \dots, r_k and r_1 is the smallest distance, then p_i is a point represented by p_r . When we tend to assign p_i to a proper cluster, we need to compare $r_i - r_1$ and $2d, i \in [2, k]$. Only the cluster c_i which holds $r_i - r_1 \leq 2d$ will be marked, and added into a set closeSet. That is to say, for p_i , we only need compute distances between it and centroids in the closeSet which is obviously the subset of k centroids.

In Algorithm 4, we give the pseudo-code of privacy-preserving LSH-based multi-user outsourced k-means clustering.

Algorithm 4. Privacy-Prserving LSH-Based Multi-User Outsourced K-Means Clustering

Input: $E_{-}(\llbracket pk \rrbracket_{-u})(D) = \{E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{p}_i)\}, i \in [1, m], C = \{E_{-}(\llbracket pk \rrbracket_{-u})(\mu_j)\} = \{E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{s}_j), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j})\}, j \in [1, k]$

Combined blinded hash table (HT), $E_{-}(\llbracket pk \rrbracket_{-u})(\epsilon)$

for \mathbf{p}_r in each bucket in HT:

for each center μ_j in C:

compute the distance between \mathbf{p}_r and μ_j

$$\begin{aligned} \llbracket dis \rrbracket_{-j} &= (E_{-}(\llbracket pk \rrbracket_{-u})(\Omega(\mathbf{p}_r, j)), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j})) \\ &= SSEDC(E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{p}_r), E_{-}(\llbracket pk \rrbracket_{-u})(\mu_j)), \end{aligned}$$

get closest encrypted center $\mu_{j'}$ for \mathbf{p}_r in C

$$\begin{aligned} \{min, label\} &= \{(E_{-}(\llbracket pk \rrbracket_{-u})(\Omega(\mathbf{p}_r, j')), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j'})), label\} \\ &\leftarrow \llbracket PMIN \rrbracket_{-k}(E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{p}_r), C) \end{aligned}$$

$$\begin{aligned} \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{s}) \rrbracket_{-j'} &= \\ \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{s}) \rrbracket_{-j'} \times E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{p}_r), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j'}) &= \\ E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j'}) \times E_{-}(\llbracket pk \rrbracket_{-u})(1) \end{aligned}$$

Set $closeSet = null$

for $j \in [1, k]$ and $j \neq j'$

CS and AS jointly compute

$$E_{-}(\llbracket pk \rrbracket_{-c})(\lambda_1) \leftarrow ESDC(min, \llbracket dis \rrbracket_{-j}, E_{-}(\llbracket pk \rrbracket_{-u})(d))$$

$$CS \text{ gets } \lambda_1 \leftarrow Dec(\llbracket sk \rrbracket_{-c}, E_{-}(\llbracket pk \rrbracket_{-c})(\lambda_1))$$

if $\lambda_1 = 0$, compute $closeSet = closeSet + \{E_{-}(\llbracket pk \rrbracket_{-u})(\mu_j)\}$

if $closeSet = null$

for each \mathbf{p}_i in L_p do

$$\begin{aligned} \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{s}) \rrbracket_{-j'} &= \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{s}) \rrbracket_{-j'} \times E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{p}_i), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j'}) \\ &= E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j'}) \times E_{-}(\llbracket pk \rrbracket_{-u})(1) \end{aligned}$$

else

for each \mathbf{p}_i in L_p do

get closet center $\mu_{j''}$ for each record \mathbf{p}_i from $closeSet$

$$\begin{aligned} \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{s}) \rrbracket_{-j''} &= \llbracket E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{s}) \rrbracket_{-j''} \times E_{-}(\llbracket pk \rrbracket_{-u})(\mathbf{p}_i), E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j''}) \\ &= E_{-}(\llbracket pk \rrbracket_{-u})(\llbracket c \rrbracket_{-j''}) \times E_{-}(\llbracket pk \rrbracket_{-u})(1) \end{aligned}$$

Then CS and AS will update the bitmap matrix according to $abel$, and they cooperatively judge the termination conditions.

4. The Proposed Solution

4.1 System Model

In our setting, we assume that there are two types of entities: users and cloud service providers as given in Figure 3. There are n users denoted by U_1, \dots, U_n . Suppose that U_i holds a dataset T_i with d-dimension m_i data records, $\sum_{i=1}^n m_i = m$. We use a Computation Server (CS) and an Assistant Server (AS) to excute the k-means clustering task, which are both semi-honest and don't collude. Let AS generate two public-secret key pair $(\llbracket pk \rrbracket_{-u}, \llbracket sk \rrbracket_{-u})$, $(\llbracket pk \rrbracket_{-c}, \llbracket sk \rrbracket_{-c})$ based on the Paillier cryptosystem and the public key $\llbracket pk \rrbracket_{-u}$ is sent to all users and CS, while AS remains $\llbracket sk \rrbracket_{-u}$, while $(\llbracket pk \rrbracket_{-c}, \llbracket sk \rrbracket_{-c})$ is sent to CS.

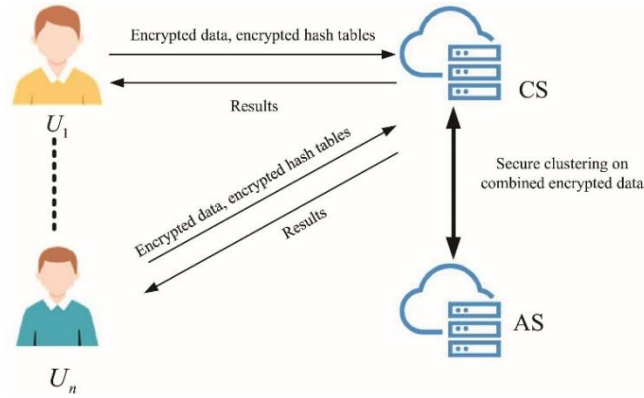


Figure 3. System architecture

When clustering starts, all users first use the same LSH function family to get their own hash table (Data Skeleton), and then blind the bucket values through MD5 to preserve the privacy of buckets. After that, each user U_i encrypts his/her own dataset attribute-wise using $\llbracket pk \rrbracket_u$, and sends hash table and encrypted dataset to CS. Then CS will first aggregate all received hash tables into one hash table (HT), in this process, buckets with the same value will be merged. Next, CS performs k-means clustering on ciphertexts with the help of AS. When the clustering results do not change anymore or a predefined number of the iteration is reached, the clustering is finished. The results will be sent to users.

4.2 The Proposed LSH-PPMOC Protocol

In this section, we will discuss our proposed LSH-PPMOC algorithm. There are three steps in our proposed algorithm. Specifically, our algorithm is composed of the following: 1) Data uploading, 2) LSH-based pruning clustering, 3) Termination.

Data uploading. At the beginning, AS generates two public/private key pairs $(\llbracket pk \rrbracket_c, \llbracket sk \rrbracket_c)$, $(\llbracket pk \rrbracket_u, \llbracket sk \rrbracket_u)$ based on the Paillier cryptosystem, then $\llbracket pk \rrbracket_u$ is sent to all users and CS, $\llbracket sk \rrbracket_u$ remains secret, while $(\llbracket pk \rrbracket_c, \llbracket sk \rrbracket_c)$ are sent to CS for decrypting intermediate results. Then each user generates their own key pair $(\llbracket pk \rrbracket_v, \llbracket sk \rrbracket_v)$ for decrypting the final clustering centroids. Before clustering, all users agree on a hash function family, and generate their own hash tables. The table is made up of different buckets. Unfortunately the bucket value is gotten from users' original data records, it needs to be blinded by a one-way hash function (SHA5, MD5 and so on) to protect information about data records from leaking. The hash table will be sent to CS. Besides, each user encrypts their dataset using $\llbracket pk \rrbracket_u$ before sending to CS.

LSH-based pruning clustering. After receiving encrypted hash tables, and encrypted original datasets, CS will aggregate tables firstly and then perform the clustering work cooperating with AS. This process can be divided into four steps.

(1) CS chooses initial k centroids randomly from all encrypted objects for k clusters, denoted as $E_{\llbracket pk \rrbracket_u}(\mu_i) = \langle E_{\llbracket pk \rrbracket_u}(s_i), E_{\llbracket pk \rrbracket_u}(|c_i|) \rangle, (i \in [1, k])$, where s_i equals $p_{j,j} \in [1, m]$ and $E_{\llbracket pk \rrbracket_u}(|c_i|) = E_{\llbracket pk \rrbracket_u}(1)$.

(2) For every representative point p_r in each bucket, CS and AS perform original privacy-preserving k-means clustering on them. In other words, they compute k squared Euclidean distances as r_1, \dots, r_k for each point using SSED, where we always assume r_1 is the smallest distance, then we assign this point to proper cluster through running PMIN_k protocol on the k distances. After clustering, we update the bitmap vector.

(3) For any point p_i in L_p in each bucket, whose distance to representative point p_r in this bucket is smaller than ε , we add cluster j into a set closeSet if $r_j - r_1 \leq 2d, j \in [2, k]$. The next step we should do is to compute distances between p_i and centroids in closeSet, cluster p_i accurately, and update according bitmap vector.

(4) Update the new cluster centroids.

Termination. The Step (2)-(4) in LSH-based pruning clustering process will be executed repeatedly until the clustering results do not change any more or the given maximum iteration number is reached.

5. Security Analysis

In this section, we will give a security proof of our LSH-PPMOC protocol under the semi-honest model. Since the proof of preliminaries we utilize in this paper are similar, we only take the SM protocol as an example and give a formal proof under "Real-versus-Ideal" framework (Goldreich, 2009).

Theorem 2. The SM Protocol securely computes the multiplication on ciphertexts using the Paillier cryptosystem under two semi-honest but passive cloud servers.

Proof. Our SM protocol is performed by two semi-honest parties, Alice and Bob. We need to prove that SM is secure against both of semi-honest attacker Alice \mathcal{A}_A and Bob \mathcal{A}_B .

(1) Security against \mathcal{A}_A : In Step (1), the real world view of \mathcal{A}_A includes inputs $\{E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(x), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(y)\}$, random values $\{r_x, r_y\}$ and outputs $E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(x + r_x), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(y + r_y)$. In Step (2), the real world view includes $E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})((x + r_x)(y + r_y))$ and $E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(x \cdot y)$. However, \mathcal{A}_A can't decrypt these ciphertexts without $\llbracket sk \rrbracket_{\mathcal{U}}$. Then we can construct a simulator \mathcal{S}_A in the ideal world. \mathcal{S}_A generates ciphertexts $\{E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(m_1), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(m_2), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})((m_1)^\wedge),$

$E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})((m_2)^\wedge)\}$ and random values r_1, r_2 by randomly selecting from Z_N . Then \mathcal{S}_A executes SM protocol. Considering the Paillier cryptosystem is semantic secure, it is computationally difficult for \mathcal{S}_A to distinguish ideal world and the real world, which means that

$$\llbracket \llbracket Ideal \rrbracket_{\mathcal{U}}(\llbracket f, \mathcal{S} \rrbracket_A) \rrbracket_{\mathcal{U}} \approx^c \llbracket \llbracket Real \rrbracket_{\mathcal{U}}(SM(\llbracket \mathcal{S} \rrbracket_A)) \rrbracket_{\mathcal{U}} \quad (4)$$

where the symbol f means multiplication function on

$E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(m_1), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(m_2), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})((m_1)^\wedge), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})((m_2)^\wedge), r_1, r_2$ and

\approx^c means computationally indistinguishable.

(2) Security against \mathcal{A}_B : In Step (2), the real world view of \mathcal{A}_B comprises of input $E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(x + r_x), E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})(y + r_y)$, and output $E_{\mathcal{U}}(\llbracket pk \rrbracket_{\mathcal{U}})((x + r_x)(y + r_y))$. Though \mathcal{A}_B can decrypt the ciphertexts and get the messages $x + r_x, y + r_y$ with $\llbracket sk \rrbracket_{\mathcal{U}}$, it is still hard for \mathcal{A}_B to get any useful information about the original data records x and y . Since r_x, r_y are random values which are chosen by Alice, they are random values in the point of view of \mathcal{A}_B . Then we can build a simulator \mathcal{S}_B in ideal world, where we just choose random messages as ciphertexts. There is no doubt that it is computationally hard for \mathcal{A}_B to distinguish the ideal world and the real world. That is to say

$$\llbracket \llbracket Ideal \rrbracket_{\mathcal{U}}(\llbracket f, \mathcal{S} \rrbracket_B) \rrbracket_{\mathcal{U}} \approx^c \llbracket \llbracket Real \rrbracket_{\mathcal{U}}(SM(\llbracket \mathcal{S} \rrbracket_B)) \rrbracket_{\mathcal{U}} \quad (5)$$

Combine the above two analyses, we prove the correctness of Theorem 2.

In the data uploading process, even though CS holds the data records in ciphertexts form, it can't get any information about original datasets due to the semantic security of Paillier cryptosystem. As for blinded hash tables, though it may disclose the relationship among user's data records, it is a compromise for efficiency. During the clustering process, AS with decryption key will assist CS to perform various computation operations. It is worth noting that the plaintexts obtained by AS are randomized. Considering Paillier cryptosystem is semantic secure, and blinding factors are all randomly chosen, no additional information regarding users' data or clusters is revealed to the cloud servers. Besides, (Liu, Lu, Ma, Chen & Qin, 2015) has given security proof of PMIN and PMIN_k. Then preliminaries utilized in our paper are proven to be privacy-protected, according to Composition Theorem (Goldreich, 2009), we conclude the sequential composition in proposed LSH-PPMOC protocol is secure against the semi-honest cloud servers.

6. Performance Analysis

In this section, we give the performance analysis of LSH-PPMOC protocol from the view of theory and experiment.

6.1 Theoretical Analysis

To better illustrate the computational cost and communicational cost of our LSH-PPMOC protocol, we use the

symbol exp , mul being the modular exponentiation and multiplication operations. We assume $g = N + 1$ in the Paillier cryptosystem which is a common setting like (Samanthula, Rao, Bertino, Yi & liu, 2014), then

$$E_{pk}(a) = (N + 1)^a \cdot r^N \bmod N^2 = (a \cdot N + 1) \cdot r^N \bmod N^2 \quad (6)$$

From Equation (6), we can conclude that the computational cost for one encryption is $1exp + 2mul$, also the same cost for decryption. The computational overheads and communicational cost in one iteration for our main preliminaries are given in Table 1. Here d is the dimension of a data record, k is the number of predefined clusters and m is number of all users' records. Besides, we use N to denote the size (in bits) of the Paillier encryption key. It is worth noting that the Stage 1 of our proposed algorithm is run only once, while the Stage 2 and Stage 3 are run in an iterative fashion until the clustering is finished.

Table 1. Computational and communication costs of primary algorithms

Algorithm	Computational cost	Communicational cost (in bits)
SM	$7exp + 6mul$	$3 N $
SSED	$15l exp + 14l mul$	$4 N $
SC	$3exp + 6mul$	$6 N $
ESDC	$57exp + 111mul$	$55 N $
PMIN	$19exp + 38mul$	$14 N $

6.2 Experimental Analysis

In this part, we will conduct our experiments on local terminals, the server running Windows10 has Intel(R) Core(TM) i7-8700 CPU @ 3.20GHz 3.19GHz and 8.00 GB RAM. We compare our proposed LSH-PPOMC with PPOCM (Rong, Wang, Liu, Hao & Xian, 2017), because our models are similar and both under public key cryptosystem, what's more, we have the same security goal for outsourced k-means clustering. We choose the key size $|N|$ as 1024-bit.

For authenticity of our experiments, we choose a real dataset, KEGG Metabolic Reaction Network dataset, which consists of 65,554 data records and 28 attributes, to perform our algorithm. But we find some data records are missing attribute values or repeating, we delete these corresponding records. And because some of the attribute values are decimal, we normalize all records so that each attribute value is scaled into integer which maintains $[0,1000]$.

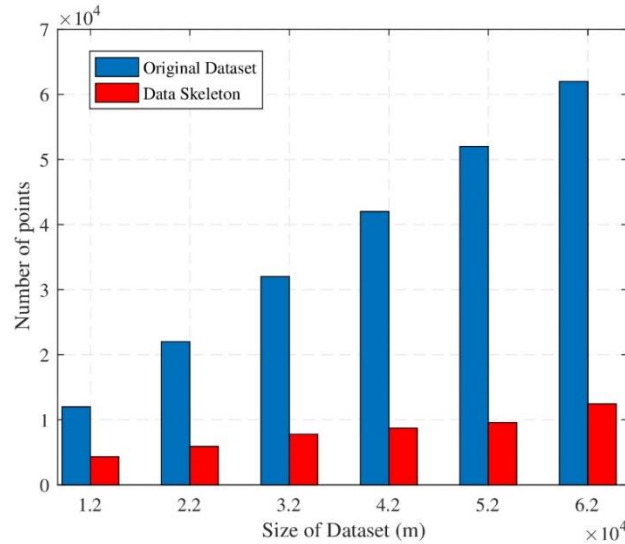


Figure 4. Original Dataset and Data Skeleton for KEGG

Before we perform k-means clustering work on the dataset, we want to illustrate the significant decrease between data points in original dataset and data skeleton in the form of a chart. The results are shown in Figure 4. From

this chart, we can see that, after users' preprocessing, the number of data points falls off drastically.

As for the secure k-means clustering on KEGG dataset, there are few factors affecting the final effect: the number of clusters (k), the size of aggregated dataset (m), and the number of attributes (d). We first asset the performance of PPOCM and LSH-PPMOC on the datasets of different sizes with dimension $d=10$, while $k=10$ and 20. The results was given in Figure 5. We can see that, the cloud running time grows almost linearly with k for both algorithms. Besides, we can also see that our LSH-PPMOC outperforms PPOCM with an increase in the number of clusters k owing to our pruning strategies.

Then we compare two k-means algorithms, with or without pruning strategies, while keeping everything else the same. Notably, the algorithm with pruning technique is exactly our LSH-PPMOC. More specifically, we change the size of aggregated dataset and the size of attribute values d respectively, and we assume $k=10$. The result is given in Figure 6. From the former diagram, we can learn that the pruning strategies can save about half the cloud running time compared with the algorithm without pruning strategies, while from the latter, we observe that our LSH-PPMOC can save nearly 70% time in terms of distance calculation. For example, when $m=10000$, $d=10$ and $k=10$, the time of distance calculation is 92min after pruning, while 275min without pruning operation. The difference can be attributed to the pruning strategies, whose goal is to reduce unnecessary distance calculation during the clustering. Besides, we can also learn that, as the amount of data increasing, the pruning is more effective.

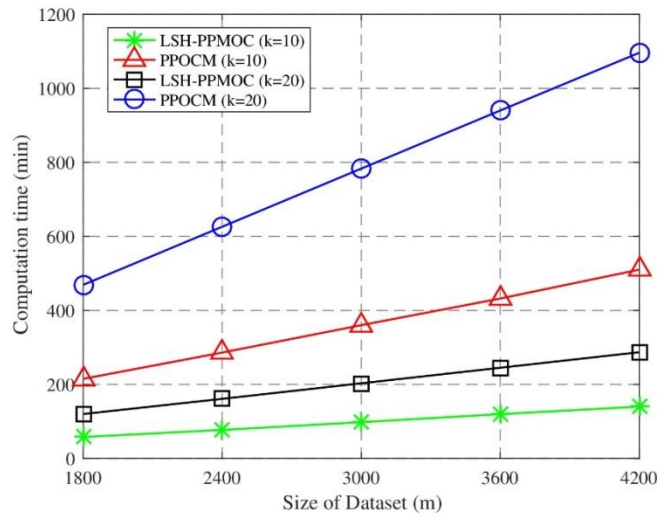
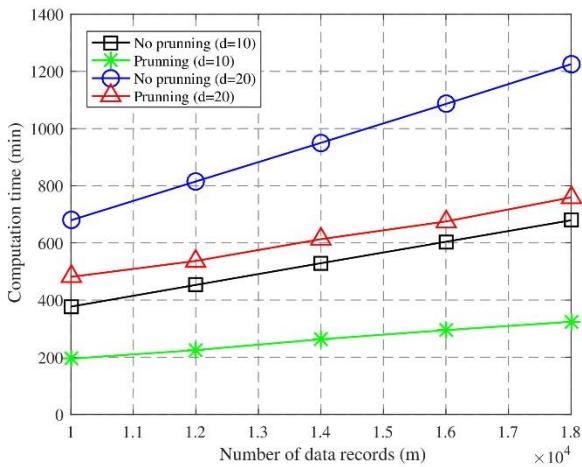
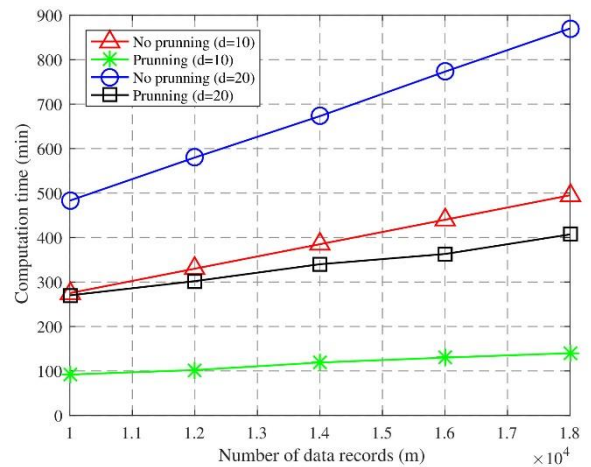


Figure 5. Cloud Running Time



(a) Time for Complete Algorithm



(b) Time for Distance Computing

Figure 6. Cloud Running Time for Complete Algorithm and Distance Computing

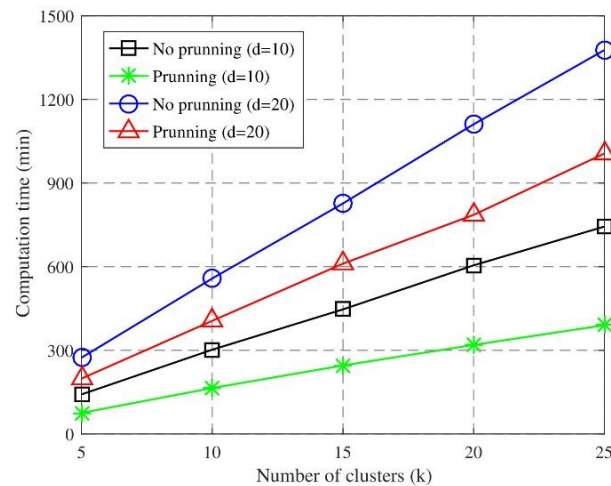


Figure 7. Cloud Running Time

Similarly, we give Figure 7 to show the effect of various k on performing clustering with two methods. It is obvious that the cloud running time saved is proportion to the number of clusters k , that is to say, when k is larger, the pruning effect is more obvious. So our LSH-PPMOC is suitable for k -means clustering on big dataset even though k is large.

7. Conclusion

In this paper, we proposed an efficient and privacy-preserving outsourced k -means clustering algorithm (LSH-PPMOC) for big data mining. For that purpose, we gave a series of building blocks to help achieve ciphertext multiplication, squared Euclidean distances computation, comparison and so on, which would never leak any useful information. The main idea in our algorithm is utilizing the local sensitivity of locally sensitive hash to prune some unnecessary computations when clustering. Besides, the experiments on KEGG dataset show that our algorithm is more efficient than some existing job. In the future, we are more willing to focus on the data integrity verification during k -means clustering and achieve strong secure under malicious model.

References

- Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). *Data clustering: a review*. ACM Comput. Surv. 31, 3 (Sept. 1999), 264-323. <https://doi.org/10.1145/331499.331504>
- Almutairi, N., Coenen, F., & Dures, K. (2017, August). *K-means clustering using homomorphic encryption and an updatable distance matrix: secure third party data clustering with limited data owner interaction*. In International Conference on Big Data Analytics and Knowledge Discovery (pp. 274-285). Springer, Cham. https://doi.org/10.1007/978-3-319-64283-3_20
- Bhaskara, A., & Wijewardena, M. (2018, July). *Distributed clustering via lsh based data partitioning*. In International Conference on Machine Learning (pp. 570-579).
- Brakerski, Z., Gentry, C., & Halevi, S. (2013, February). *Packed ciphertexts in LWE-based homomorphic encryption*. In International Workshop on Public Key Cryptography (pp. 1-13). Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-36362-7_1
- Chen, Z., Wang, Y., Zhang, S., Zhong, H., & Chen, L. (2018). *Differentially private user-based collaborative filtering recommendation based on K-means clustering*. arXiv preprint arXiv:1812.01782.
- Clifton, C., Kantarcioglu, M., Vaidya, J., Lin, X., & Zhu, M. Y. (2002). Tools for privacy preserving distributed data mining. *ACM Sigkdd Explorations Newsletter*, 4(2), 28-34. <https://doi.org/10.1145/772862.772867>
- Cui, X., Zhu, P., Yang, X., Li, K., & Ji, C. (2014). Optimized big data K-means clustering using MapReduce. *The Journal of Supercomputing*, 70(3), 1249-1259. <https://doi.org/10.1007/s11227-014-1225-7>
- Datar, M., Immorlica, N., Indyk, P., & Mirrokni, V. S. (2004). *Locality-sensitive hashing scheme based on p-stable distributions*. In Proceedings of the twentieth annual symposium on Computational geometry (pp. 253-262). <https://doi.org/10.1145/997817.997857>

- Dean, J., & Ghemawat, S. (2008). MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1), 107-113. <https://doi.org/10.1145/1327452.1327492>
- Gentry, C. (2009, May). *Fully homomorphic encryption using ideal lattices*. In Proceedings of the forty-first annual ACM symposium on Theory of computing (pp. 169-178). <https://doi.org/10.1145/1536414.1536440>
- Goldreich, O. (2009). *Foundations of cryptography: volume 2, basic applications*. Cambridge university press.
- Indyk, P., & Motwani, R. (1998, May). *Approximate nearest neighbors: towards removing the curse of dimensionality*. In Proceedings of the thirtieth annual ACM symposium on Theory of computing. <https://doi.org/10.1145/276698.276876>
- Jagannathan, G., & Wright, R. N. (2005, August). *Privacy-preserving distributed k-means clustering over arbitrarily partitioned data*. In Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining (pp. 593-599). <https://doi.org/10.1145/1081870.1081942>
- Jiang, Z. L., Guo, N., Jin, Y., Lv, J., Wu, Y., Liu, Z., ... Wang, X. (2020). Efficient two-party privacy-preserving collaborative k-means clustering protocol supporting both storage and computation outsourcing. *Information Sciences*, 518, 168-180. <https://doi.org/10.1016/j.ins.2019.12.051>
- Kriegel, H. P., Kröger, P., & Zimek, A. (2009). Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 3(1), 1-58. <https://doi.org/10.1145/1497577.1497578>
- Li, Q., Wang, P., Wang, W., Hu, H., Li, Z., & Li, J. (2014, April). *An efficient K-means clustering algorithm on MapReduce*. In International Conference on Database Systems for Advanced Applications (pp. 357-371). Springer, Cham. https://doi.org/10.1007/978-3-319-05810-8_24
- Li, Y., Shang, Y., & Yang, Y. (2017). Clustering coefficients of large networks. *Information Sciences*, 382-383, 350-358. <https://doi.org/10.1016/j.ins.2016.12.027>
- Lindell, Y., & Pinkas, B. (2008). *Secure Multiparty Computation for Privacy-Preserving Data Mining*. IACR Cryptol. ePrint Arch., 2008, 197. <https://doi.org/10.29012/jpc.v1i1.566>
- Liu, D., Bertino, E., & Yi, X. (2014, June). *Privacy of outsourced k-means clustering*. In Proceedings of the 9th ACM symposium on Information, computer and communications security (pp. 123-134). <https://doi.org/10.1145/2590296.2590332>
- Liu, K., Giannella, C., & Kargupta, H. (2006). *An attacker's view of distance preserving maps for privacy preserving data mining*. In European Conference on Principles of Data Mining and Knowledge Discovery (pp. 297-308). Springer, Berlin, Heidelberg. https://doi.org/10.1007/11871637_30
- Liu, K., Kargupta, H., & Ryan, J. (2005). Random projection-based multiplicative data perturbation for privacy preserving distributed data mining. *IEEE Transactions on knowledge and Data Engineering*, 18(1), 92-106. <https://doi.org/10.1109/TKDE.2006.14>
- Liu, X., Lu, R., Ma, J., Chen, L., & Qin, B. (2015). Privacy-preserving patient-centric clinical decision support system on naive Bayesian classification. *IEEE journal of biomedical and health informatics*, 20(2), 655-668. <https://doi.org/10.1109/JBHI.2015.2407157>
- Mohassel, P., Rosulek, M., & Trieu, N. (2020). Practical privacy-preserving k-means clustering. *Proceedings on Privacy Enhancing Technologies*, 2020(4), 414-433. <https://doi.org/10.2478/popets-2020-0080>
- Mustafi, D., & Sahoo, G. (2019). A hybrid approach using genetic algorithm and the differential evolution heuristic for enhanced initialization of the k-means algorithm with applications in text clustering. *Soft Computing*, 23(15), 6361-6378. <https://doi.org/10.1007/s00500-018-3289-4>
- Paillier, P. (1999, May). *Public-key cryptosystems based on composite degree residuosity classes*. In International conference on the theory and applications of cryptographic techniques (pp. 223-238). Springer, Berlin, Heidelberg. https://doi.org/10.1007/3-540-48910-X_16
- Rong, H., Wang, H., Liu, J., Hao, J., & Xian, M. (2017). *Outsourced k-Means clustering over encrypted data under multiple keys in spark framework*. In International Conference on Security and Privacy in Communication Systems (pp. 67-87). Springer, Cham. https://doi.org/10.1007/978-3-319-78813-5_4
- Samanthula, B. K., Rao, F. Y., Bertino, E., Yi, X., & Liu, D. (2014). *Privacy-preserving and outsourced multi-user k-means clustering*. arXiv preprint arXiv:1412.4378.
- Sardar, T. H., & Ansari, Z. (2020). An Analysis of Distributed Document Clustering Using MapReduce Based

- K-Means Algorithm. *Journal of The Institution of Engineers (India): Series B*, 101(6), 641-650. <https://doi.org/10.1007/s40031-020-00485-2>
- Su, D., Cao, J., Li, N., Bertino, E., Lyu, M., & Jin, H. (2017). Differentially private k-means clustering and a hybrid approach to private optimization. *ACM Transactions on Privacy and Security (TOPS)*, 20(4), 1-33. <https://doi.org/10.1145/3133201>
- Upmanyu, M., Namboodiri, A. M., Srinathan, K., & Jawahar, C. V. (2010, June). Efficient privacy preserving k-means clustering. In *Pacific-Asia Workshop on Intelligence and Security Informatics* (pp. 154-166). Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-13601-6_17
- Wang, Y. (2015). *Notes on Two Fully Homomorphic Encryption Schemes Without Bootstrapping*. IACR Cryptol. ePrint Arch., 2015, 519.
- Wong, W. K., Cheung, D. W. L., Kao, B., & Mamoulis, N. (2009). *Secure kNN computation on encrypted databases*. In Proceedings of the 2009 ACM SIGMOD International Conference on Management of data (pp. 139-152). <https://doi.org/10.1145/1559845.1559862>
- Xia, C., Hua, J., Tong, W., & Zhong, S. (2020). Distributed K-Means clustering guaranteeing local differential privacy. *Computers & Security*, 90, 101699. <https://doi.org/10.1016/j.cose.2019.101699>
- Xing, K., Hu, C., Yu, J., Cheng, X., & Zhang, F. (2017). Mutual privacy preserving \$k\$-means clustering in social participatory sensing. *IEEE Transactions on Industrial Informatics*, 13(4), 2066-2076. <https://doi.org/10.1109/TII.2017.2695487>
- Yin, C., & Zhang, S. (2017). Parallel implementing improved k-means applied for image retrieval and anomaly detection. *Multimedia Tools and Applications*, 76(16), 16911-16927. <https://doi.org/10.1007/s11042-016-3638-1>
- Youn, T. Y., Park, Y. H., Kim, C. H., & Lim, J. (2005, August). An efficient public key cryptosystem with a privacy enhanced double decryption mechanism. In *International Workshop on Selected Areas in Cryptography* (pp. 144-158). Springer, Berlin, Heidelberg. https://doi.org/10.1007/11693383_10
- Yu, Q., Luo, Y., Chen, C., & Ding, X. (2016). Outlier-eliminated k-means clustering algorithm based on differential privacy preservation. *Applied Intelligence*, 45(4), 1179-1191. <https://doi.org/10.1007/s10489-016-0813-z>
- Yuan, J., & Tian, Y. (2017). *Practical privacy-preserving mapreduce based k-means clustering over large-scale dataset*. IEEE Transactions on Cloud Computing.
- Zou, Y., Zhao, Z., Shi, S., Wang, L., Peng, Y., Ping, Y., & Wang, B. (2020). *Highly Secure Privacy-Preserving Outsourced k-Means Clustering under Multiple Keys in Cloud Computing*. Security and Communication Networks, 2020. <https://doi.org/10.1155/2020/1238505>

Copyrights

Copyright for this article is retained by the author(s), with first publication rights granted to the journal.

This is an open-access article distributed under the terms and conditions of the Creative Commons Attribution license (<http://creativecommons.org/licenses/by/4.0/>).