

Centerless Clustering: An Efficient Variant of K -means Based on K-NN Graph

Shenfei Pei, Huimin Chen, Feiping Nie, Rong Wang, and Xuelong Li, *Fellow, IEEE*

Abstract—Although lots of clustering models have been proposed recently, k -means and the family of spectral clustering methods are both still drawing a lot of attention due to their simplicity and efficacy. We first reviewed the unified framework of k -means and graph cut models, and then proposed a clustering method called k -sums where a k -nearest neighbor (k -NN) graph is adopted. The main idea of k -sums is to minimize directly the sum of the distances between points in the same cluster. To deal with the situation where the graph is unavailable, we proposed k -sums-x that takes features as input. The computational and memory overhead of k -sums are both $O(nk)$, indicating that it can scale linearly w.r.t. the number of objects to group. Moreover, the costs of computational and memory are irrelevant to the product of the number of points and clusters. The computational and memory complexity of k -sums-x are both linear w.r.t. the number of points. To validate the advantage of k -sums and k -sums-x on facial datasets, extensive experiments have been conducted on 10 synthetic datasets and 17 benchmark datasets. While having a low time complexity, the performance of k -sums is comparable with several state-of-the-art clustering methods.

Index Terms—Efficient clustering, fast, k -nearest graph, spectral clustering, k -means.

1 INTRODUCTION

WITH the popularity of the Internet, large-scale datasets can be found anywhere in our lives, such as shopping records generated by e-commerce sites and photos (videos) uploaded by users on social media platforms. Take Facebook as an example, there are over 2.5 billion monthly active users in the United States and Canada, about 1 million links are shared, and 3 million messages are sent, every 20 minutes, 100 million hours of video are watched on Facebook every day, as of April 2020, reported by statista¹. Although these large-scale datasets are relatively easy to obtain, further analysis is essential to get useful information.

As one of the fundamental unsupervised technologies in the fields of machine learning, data mining, pattern recognition, and others, clustering is drawing more and more attention over the last few years. Up to now, a series of models have been presented for cluster analysis and put into practice successfully in various fields, such as face recognition [1], [2], document clustering [3], [4], social networks [5], [6], and image segmentation [7], [8].

Among various clustering algorithms, the family of spectral clustering (SC) and k -means are still drawing lots of attention, due to their simplicity and efficacy. The former

takes graphs as input and the latter takes features as input.

The models proposed in this article are closely related to these traditional clustering algorithms. Therefore, a brief review of them is presented first. For convenience, we use $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$, $\mathbf{W} \in \mathbb{R}^{n \times n}$, and c to denote the dataset with n samples, the similarity matrix, and the number of clusters to construct, respectively.

Taking the spectral clustering method introduced in [9] as an example, it can be seen as a relaxation of the normalized-cut problem. The objective function of normalized-cut can be formulated as

$$\min_{\mathcal{A}_1, \dots, \mathcal{A}_c} \sum_{k=1}^c \frac{\sum_{\mathbf{x}_i \in \mathcal{A}_k, \mathbf{x}_j \in \bar{\mathcal{A}}_k} w_{ij}}{vol(\mathcal{A}_k)}, \quad (1)$$

where $vol(\mathcal{A}_k) = \sum_{\mathbf{x}_i \in \mathcal{A}_k} \sum_{j=1}^n w_{ij}$. The spectral clustering method [9] transforms the problem in Eq. (1) into an eigen-decomposition problem by relaxing the discrete condition and often yields more superior experimental performance. However, the computational cost of the eigendecomposition involved increases quadratically with the number of samples, which is very time-consuming.

With these notations mentioned above, the objective function of k -means can be formulated as

$$\min_{\mathcal{A}_1, \dots, \mathcal{A}_c} \sum_{k=1}^c \sum_{\mathbf{x}_i \in \mathcal{A}_k} \|\mathbf{x}_i - \mathbf{m}_k\|_2^2, \quad (2)$$

where $\mathbf{m}_k = \frac{1}{n_k} \sum_{\mathbf{x}_i \in \mathcal{A}_k} \mathbf{x}_i$ denotes the centroid of cluster \mathcal{A}_k containing n_k samples. The goal of k -means is to minimize the sum of distances between the samples and the center in each cluster. In general, the computational overhead of k -means is linear with respect to the number of samples, but it cannot separate clusters non-linearly separable in input space, which is a major disadvantage.

In this paper, we are aiming at designing a fast clustering algorithm whose computational complexity increases

• S. Pei, H. Chen, F. Nie (corresponding author) are with the School of Computer Science, School of Artificial Intelligence, Optics and Electronics (iOPEN), and the Key Laboratory of Intelligent Interaction and Applications (Ministry of Industry and Information Technology), Northwestern Polytechnical University, Xi'an 710072, Shaanxi, P. R. China. E-mail: shenfeipei@gmail.com, chenhuimin@mail.nwpu.edu.cn, feipingnie@gmail.com.

• R. Wang and X. Li are with the School of Artificial Intelligence, Optics and Electronics (iOPEN), and the Key Laboratory of Intelligent Interaction and Applications (Ministry of Industry and Information Technology), Northwestern Polytechnical University, Xi'an 710072, Shaanxi, P. R. China. E-mail: wangrong07@tsinghua.org.cn, li@nwpu.edu.cn

Code release: <https://github.com/ShenfeiPei/KSUMS2>

1. <https://www.statista.com/>

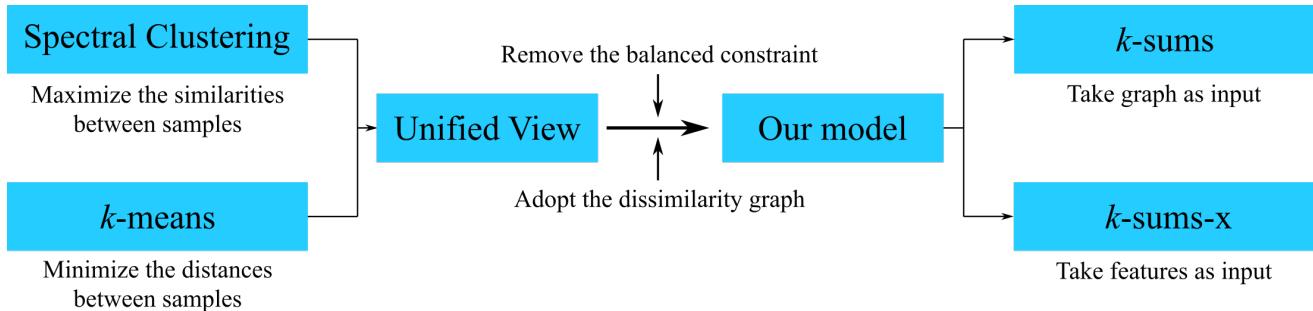


Fig. 1: The motivation of the proposed algorithms (k -sums and k -sums-x). (1) The unified view of k -means and spectral clustering is firstly revisited. (2) Then, our model is proposed by making some modifications to this unified view. (3) Two versions of the proposed method have been developed, called k -sums and k -sums-x. The former (k -sums) takes a k -NN graph as input and the latter takes the features as input.

linearly with the number of samples. Moreover, its performance is comparable to traditional spectral clustering.

The main idea of the proposed model comes from the unified view of k -means and spectral clustering. Specifically, we found that the unified view has a relatively simple form if we focus on the situation where the number of samples in each cluster is strictly equal, and the model (without the constraint) naturally tends to produce more balanced clustering results, if a dissimilarity matrix is used. The motivation of the algorithm is shown in Fig. 1.

Based on these findings, a novel clustering method that is dedicated to directly minimizing the sum of the distances between points in the same cluster is proposed and called k -sums. It is worth mentioning that it takes the nearest neighbor graph as input and only cares about the distance between points that are neighbors of each other. The standard coordinate descent (CD) algorithm is adopted to optimize our model, and an extremely efficient optimization algorithm is designed for the sub-problems involved in the CD algorithm. Therefore, the computational complexity of k -sums is greatly reduced.

In addition, in order to apply the proposed clustering algorithm in the situation where the nearest neighbour graph is not available, a version that takes features as input is also discussed and called k -sums-x. Again, the coordinate descent algorithm is used to optimize the problem of k -sums-x, and the property of squared Euclidean distance is used to reduce the computational overhead.

It is worthwhile to highlight the main contributions of this article as follows:

- Based on the unified view of k -means and spectral clustering, a novel clustering method that is dedicated to directly minimizing the sum of the distances between points in the same cluster is proposed. It cares about the squared Euclidean distance between points only, so, the trouble of measuring the similarity between points can be avoided.
- The optimization algorithm employed, i.e., coordinate descent method, guarantees that no empty cluster will appear in the clustering result, with any initialization. In particular, our model naturally tends to produce a more balanced partition, and a detailed explanation is provided.

- For k -sums, the version of the proposed model that takes a k -nearest neighbor graph as input, its computational and memory overhead are both linear with respect to the number of points. More importantly, the time complexity is independent of the product of the number of points and clusters. These properties mean that it is easily scalable and applicable to large-scale problems.
- To apply our model in the situation where a k -NN graph is not available, k -sums-x, the version of the proposed model that takes features as input, is also provided. With some pre-processing, its computational complexity is still linear with respect to the number of samples.
- In order to validate the superiority of our algorithms, extensive experiments have been conducted on 10 synthetic datasets, 13 middle-scale real-world datasets, and 4 large-scale datasets, and our algorithms have shown satisfactory results compared to the state-of-the-art algorithms.

The rest of this paper is organized as follows: In Section 2, we briefly introduce some algorithms related to spectral clustering or k -means. In Section 3, the unified view of k -means and spectral clustering is briefly revisited, and our model is presented based on the unified view. In Section 4, a variant of our model is proposed, and the computational and memory overhead are analyzed. The experimental results are reported in Section 5 and the conclusions are given in Section 6.

Notations: All matrices are denoted by boldface uppercase letters, and all vectors are represented by boldface lowercase letters. For matrix $\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_n]^T \in \mathbb{R}^{n \times d}$, m_{ij} and \mathbf{m}_i are used to represent the (i, j) -th entry and transpose of the i -th row of matrix \mathbf{M} , respectively. The squared l_2 norm of the vector $\mathbf{v} \in \mathbb{R}^n$ is defined as $\|\mathbf{v}\|_2^2 = \sum_{i=1}^n v_i^2$. $\mathbf{X} \in \mathbb{R}^{n \times d}$ is adopted to denote the dataset with n samples, where \mathbf{x}_i denotes the i -the point, and c is used the number of clusters to construct. For matrix \mathbf{Y} , if there is only one element equal to 1 in $\mathbf{y}_i, \forall i = 1, \dots, n$, and the rest elements are all 0, then \mathbf{Y} is referred to cluster indicator matrix. We use $\Phi^{n \times c}$ to denote the set of all cluster indicator matrices. $\mathbf{1}$ is the column vector of all ones. $\mathbf{x}_j^{(i)}$ denotes the j -th sample in cluster i .

2 RELATED WORK

Our model comes from the unified view of k -means and spectral clustering, so, we mainly introduce the related algorithms of the two.

2.1 Algorithms Related to Spectral Clustering

The workflow of SC can be summarized as follows: 1. Construct an undirected graph, i.e., $\mathbf{G} = \langle \mathbf{V}, \mathbf{W} \rangle$, where \mathbf{V} and \mathbf{W} denote the collection of input patterns and the adjacency matrix, respectively. In general, w_{ij} represents the similarity between samples \mathbf{x}_i and \mathbf{x}_j , and $\mathbf{W} \in \mathbb{R}^{n \times n}$ can be determined by the heat kernel as

$$w_{ij} = \begin{cases} e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{t}} & \text{if } \mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_j) \text{ or } \mathbf{x}_j \in \mathcal{N}_k(\mathbf{x}_i) \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where $\mathcal{N}_k(\mathbf{x}_i)$ denotes the set of k nearest neighbors of \mathbf{x}_i . $t > 0$ is a parameter. Obviously, \mathbf{W} is a symmetric matrix. 2. Compute the Laplacian matrix $\mathbf{L} = \Delta - \mathbf{W}$, where Δ is a diagonal matrix and defined as $\Delta_{jj} = \sum_{i=1}^n w_{ij}$. 3. Compute the first c eigenvectors of the normalized Laplacian matrix \mathbf{L}_{norm} , with

$$\mathbf{L}_{norm} = \begin{cases} \Delta^{-1/2} \mathbf{L} \Delta^{-1/2} & \text{Ng et al [9]} \\ \Delta^{-1} \mathbf{L} & \text{Shi and Malik [10]} \end{cases} \quad (4)$$

Despite its promising performance, it has a high computational overhead, which can be seen from the two stages of graph construction and eigendecomposition. To this end, A lot of effort has been made. Among these studies, Nyström-based and anchor-based methods gain the most popularity.

The Nystrom method can be seen as an algorithm for approximately solving eigendecomposition problems. The approximation works by first solving the eigendecomposition problem for a small random subset of points and then extrapolating this solution to the full set of data points. [11], [12], [13] applied this algorithm to accelerate the calculation of spectral clustering. [14] proposed a variational Nystrom method. Using fewer sampled points, it achieved a lower approximation error.

The anchor-based strategy is also a very popular method for accelerating the spectral clustering algorithm. These methods learn a sparse representation for each data point firstly. The similarity matrix is then designed as

$$\mathbf{W} = \mathbf{Z}^T \Delta^{-1} \mathbf{Z}, \quad (5)$$

where $\mathbf{Z} \in \mathbb{R}^{m \times n}$ is the sparse representation, m denotes the number of anchors, Δ is a diagonal matrix with $\Delta_{kk} = \sum_{j=1}^n z_{jk}$. Generally, the computational overhead of these methods is $O(nm^2)$.

[15] proposed a fast implementation of spectral clustering, where random projection and sampling methods are used for reducing the dimensionality and cardinality of data. [16] proposed a fast method, where the spectral clustering is performed on a hierarchical bipartite graph that much smaller than the similarity graph. In addition, the method is able to deal with the out-of-sample problem. For coping with large-scale clustering problems in the case of limited memory, [17], [18] proposed sequential SC algorithms. The time complexity of them is nearly linear with respect to the number of points, and the memory overhead is independent of the number of points.

2.2 Algorithms Related to k -means

k -means, as one of the most popular clustering algorithms, is dedicated to minimizing within-cluster variances (squared Euclidean distances). The main idea of k -means is to assign each point to the cluster whose center is closest to the point. With the help of the aforementioned notations, the objective function of k -means can be expressed as

$$\min_{\mathbf{Y} \in \Phi^{n \times c}, \mathbf{M}} \|\mathbf{X} - \mathbf{YM}\|_F^2, \quad (6)$$

where $\mathbf{Y} \in \Phi^{n \times c}$ is a cluster indicator matrix, $\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_c]^T \in \mathbb{R}^{c \times d}$, \mathbf{m}_i denotes the center (mean) of i -th cluster. The problem (6) is NP-hard, but many efficient optimization methods [27, 37, 13] are proposed.

Although problem (6) can be optimized efficiently, k -means has the following disadvantages:

- It cannot separate clusters that are not separable linearly in input space, which is a major disadvantage of it.
- The quality of the clustering result will be greatly affected by the initialization of the clustering partition.

To solve the first problem, kernel-based k -means was proposed [19], [20], where a feature function $\phi : \mathbb{R}^d \rightarrow \mathcal{F}$ is used to map the input features non-linearly into output vectors, and the traditional k -means is then performed in these vectors instead of the original features. However, most kernel-based k -means algorithms cannot be applied directly to large-scale problems for their high computational complexity. To this end, by employing a randomized approach, Chitta et al. [21] proposed an approximation scheme, termed approximate kernel k -means, where the computational and memory overhead are both reduced. [47, 17, 7] adopted the incomplete Cholesky factorization and Nystrom methods to solve approximately kernel learning problems. In order to overcome the influence of the initialization on performance, k -means++, a variant of k -means [22] which is widely used, is proposed. It seeded the initial centers by a specific rule and the initialization is provably competitive with the optimal solution. Despite its promising performance, k -means++ requires c full passes over the whole dataset. To this end, Bachem et al. [23] adopted Markov Chain Monte Carlo (MCMC) method to approximate k -means++ with a lower time complexity. In addition, Bachem et al. [24] proposed a fast seeding method, where the quality of the clustering result is provably competitive with the optimal solution.

From the above discussion, it can be seen that the k -means based and spectral-based clustering algorithms are loosely related. In order to establish the connection spectral clustering and k -means, a lot of effort has been done, which narrowed the gap between the two to some extent. An explicit theoretical connection between kernel k -means and spectral clustering has been proposed in [25], [26]. In addition, Ding et al. [27] has shown that Nonnegative Matrix Factorization (NMF) of \mathbf{W} , i.e.,

$$\min_{\mathbf{H} \geq 0} \|\mathbf{W} - \mathbf{HH}^T\|_F^2 \quad (7)$$

is equivalent mathematically to kernel k -means with the strict orthogonality relation relaxed.

2.3 The Unified View

As mentioned before, the graph partitioning problem from which spectral clustering is derived can be formulated as

$$\min_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr} \left((\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{L} \mathbf{Y} \right). \quad (8)$$

where $\mathbf{L} = \Delta - \mathbf{W}$, Δ is the degree matrix, $\Delta_{jj} = \sum_{i=1}^n w_{ij}$. In practice, a normalized Laplacian matrix \mathbf{L}_{norm} is usually used, and the meaning of \mathbf{L}_{norm} is different in different versions, as mentioned above. If the similarity matrix \mathbf{W} is doubly stochastic, that is to say

$$\sum_{i=1}^n w_{ij} = 1, \sum_{j=1}^n w_{ij} = 1, \forall i, j = 1, 2, \dots, n, \quad (9)$$

then, we have $\Delta = \mathbf{I}$. Therefore the objective function in Eq. (8) can be simplified as

$$\max_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr} \left((\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{W} \mathbf{Y} \right) \quad (10)$$

It's worth noting that many studies [28], [29] have shown that the performance of SC can be enhanced by employing the doubly stochastic matrix. Thus, the simplified model (10) is meaningful. Next, we pay attention to k -means and transform its objective function into a form similar to Eq. (10).

For convenience, we focus on the matrix form of k -means. It is not difficult to find that the following equations are hold.

$$\min_{\mathbf{Y} \in \Phi^{n \times c}, \mathbf{M}} \|\mathbf{X} - \mathbf{Y}\mathbf{M}\|_F^2 \quad (11)$$

$$\Leftrightarrow \min_{\mathbf{Y} \in \Phi^{n \times c}, \mathbf{M}} \text{Tr} \left(-2\mathbf{X}^T \mathbf{Y}\mathbf{M} + \mathbf{M}^T \mathbf{Y}^T \mathbf{Y}\mathbf{M} \right) \quad (12)$$

Taking the derivation of Eq. (12) with respect to \mathbf{M} , we have

$$-2\mathbf{Y}^T \mathbf{X} + 2\mathbf{Y}^T \mathbf{Y}\mathbf{M} = 0 \quad (13)$$

$$\mathbf{M} = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \quad (14)$$

We call the cluster "empty cluster" if it contains no samples. Obviously, if no empty cluster appears in the clustering result, $\mathbf{Y}^T \mathbf{Y}$ is reversible.

Here, we only focus on the situation where $\mathbf{Y}^T \mathbf{Y}$ is reversible since the optimization algorithm adopted in this paper guarantees that no empty cluster will appear in the clustering result, in the proposed model.

Replacing \mathbf{M} with $(\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X}$, we arrive at

$$\min_{\mathbf{Y} \in \Phi^{n \times c}, \mathbf{M}} \text{Tr} \left(-2\mathbf{X}^T \mathbf{Y}\mathbf{M} + \mathbf{M}^T \mathbf{Y}^T \mathbf{Y}\mathbf{M} \right) \quad (15)$$

$$\Leftrightarrow \min_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr} \left((-2\mathbf{X}^T \mathbf{Y} + \mathbf{X}^T \mathbf{Y}) \mathbf{M} \right) \quad (16)$$

$$\Leftrightarrow \max_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr} \left(\mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \right) \quad (17)$$

$$\Leftrightarrow \max_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr} \left((\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \mathbf{X}^T \mathbf{Y} \right) \quad (18)$$

It is not difficult to find that both the spectral clustering and k -means can be unified into the framework below.

$$\max \text{Tr} \left((\mathbf{Y}^T \mathbf{Y})^{-\frac{1}{2}} \mathbf{Y}^T \mathbf{G} \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-\frac{1}{2}} \right). \quad (19)$$

Equation (19) denotes the objective function of spectral clustering if $\mathbf{G} = \mathbf{W}$, and the objective function of k -means if $\mathbf{G} = \mathbf{X} \mathbf{X}^T$.

The unified view discussed here is different from that discussed in [25], [26]. Specifically, they mainly explored the relationship between kernel k -means and spectral clustering and proved the equivalence of the two. In our article, we are committed to discussing the essential difference between traditional k -means and spectral clustering with doubly stochastic matrix.

2.4 The Difference and Connection Between k -means and Spectral Clustering

Here, we write the two into another unified mathematical form. Based on it, we found that the essential difference between the two is that the measurements they employed are different.

According to [30], k -means (2) can be rewritten as

$$\min_{A_1, \dots, A_c} \sum_{k=1}^c \frac{1}{|A_k|} \sum_{\mathbf{x}_i, \mathbf{x}_j \in A_k} d_{ij}^2 \quad (20)$$

where $d_{ij}^2 = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$.

It is not difficult to find that the objective function of spectral clustering can be written similarly as

$$\max_{A_1, \dots, A_c} \sum_{k=1}^c \frac{1}{|A_k|} \sum_{\mathbf{x}_i, \mathbf{x}_j \in A_k} w_{ij} \quad (21)$$

Based on the discussion above, we have the following observations: 1. Spectral clustering adopted doubly stochastic matrix shares the same framework as the k -means, as shown in Eq. (19). 2. For these two algorithms, the only difference is that different metrics are used to measure the similarity between samples. 3. From another view, the goal of SC is to **maximize** the mean of the similarities between samples in the same cluster, while the goal of k -means is to **minimize** the mean of distances between samples in the same cluster, as shown in Eq. (20) and Eq. (21).

3 AN EFFICIENT CLUSTERING MODEL

Taking the unified framework as starting point, we developed our model. Interestingly, our model has two forms: k -sums that take graph as input, and k -sums-x that take features as input. k -sums can be regard as a competitor proposed for the family of spectral clustering and k -sums-x can be seen as a competitor proposed for k -means-like algorithms, as shown in Figure 1.

3.1 Our Model

Although those clustering algorithms where a k -nearest neighbor graph is adopted usually yield better experimental performance, in terms of commonly used metrics, such as ACC, NMI, compared to other methods, their optimization algorithms are time-consuming usually.

Fortunately, we found that a simple model can be derived from the unified view shown in Eq. (19), if we focus on the situation where the number of samples in each cluster is strictly equal. In other words, we have the following equations

$$\max_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr} \left((\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{G} \mathbf{Y} \right) \quad s.t. \mathbf{Y}^T \mathbf{Y} = \bar{n} \mathbf{I} \quad (22)$$

$$\max_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr} \left(\mathbf{Y}^T \mathbf{G} \mathbf{Y} \right) \quad s.t. \mathbf{Y}^T \mathbf{Y} = \bar{n} \mathbf{I} \quad (23)$$

We remove the balance constraint to make the problem easy to solve and adopted a distance matrix instead of a similarity matrix to avoid the trivial solution. Therefore the final objective function of our model is as follows

$$\min_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr}(\mathbf{Y}^T \tilde{\mathbf{D}} \mathbf{Y}) \quad (24)$$

with

$$\tilde{d}_{ij} = \begin{cases} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2, & x_i \leftrightarrow x_j \\ \gamma, & \text{otherwise} \end{cases} \quad (25)$$

where $\mathbf{x}_i \leftrightarrow \mathbf{x}_j$ means $\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_j)$ and $\mathbf{x}_j \in \mathcal{N}_k(\mathbf{x}_i)$, γ means the maximum of the set \mathcal{A} consisting of defined values, that is to say

$$\mathcal{A} = \{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \mid \mathbf{x}_i \leftrightarrow \mathbf{x}_j, i = 1, \dots, n, j = 1, \dots, n\}. \quad (26)$$

The proposed model aims to minimize the sum of the squared Euclidean distances between points in the same cluster, so we term our model as k -sums.

3.2 Analysis

To begin with, some notations are presented here. Note that these notations only work in this subsection. Let $\mathcal{A}(\mathbf{x}_i)$ be the set of $\{\mathbf{x}_l \mid \mathbf{x}_l \in \mathcal{N}_k(\mathbf{x}_i) \text{ and } \mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_l)\}$, \mathcal{F}_i be the set of samples in cluster i , and $\mathbf{x}_j^{(i)}$ be the j -th sample in cluster i , we use the following notations,

$$\begin{cases} k_i^* = |\mathcal{A}(\mathbf{x}_i)|, \\ k_{ij}^* = |\mathcal{A}(\mathbf{x}_j^{(i)})|, \\ k_{ij} = |\mathcal{A}(\mathbf{x}_j^{(i)}) \cap \mathcal{F}_i|, \\ \bar{k}_{ij} = |\mathcal{A}(\mathbf{x}_j^{(i)}) - \mathcal{F}_i|, \end{cases}$$

to represent the number of samples in corresponding sets. Then, we have

$$k_{ij}^* = k_{ij} + \bar{k}_{ij}. \quad (27)$$

It means that the samples in $\mathcal{A}(\mathbf{x}_i)$ can be divided into two categories (in or not in the i -th cluster, \mathcal{F}_i).

Similarly, using the following notations

$$\begin{cases} s_i^* = \sum_{\mathbf{x}_j \in \mathcal{B}} d(\mathbf{x}_i, \mathbf{x}_j), \mathcal{B} = \mathcal{A}(\mathbf{x}_i), \\ s_{ij}^* = \sum_{\mathbf{x}_l \in \mathcal{B}} d(\mathbf{x}_j^{(i)}, \mathbf{x}_l), \mathcal{B} = \mathcal{A}(\mathbf{x}_j^{(i)}), \\ s_{ij} = \sum_{\mathbf{x}_l \in \mathcal{B}} d(\mathbf{x}_j^{(i)}, \mathbf{x}_l), \mathcal{B} = \mathcal{A}(\mathbf{x}_j^{(i)}) \cap \mathcal{F}_i, \\ \bar{s}_{ij} = \sum_{\mathbf{x}_l \in \mathcal{B}} d(\mathbf{x}_j^{(i)}, \mathbf{x}_l), \mathcal{B} = \mathcal{A}(\mathbf{x}_j^{(i)}) - \mathcal{F}_i, \end{cases}$$

with

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2,$$

to represent the sum of distances in these sets, we have

$$s_{ij}^* = s_{ij} + \bar{s}_{ij}. \quad (28)$$

That is to say, for distance $d(\mathbf{x}_j^{(i)}, \mathbf{x}_l)$, there are two cases: (1) \mathbf{x}_l belongs to \mathcal{F}_i ; (2) \mathbf{x}_l does not belong to \mathcal{F}_i , which is similar to Equation (27).

With these notations, we can rewrite our model Eq. (24) into a more detailed form, from which we can see how our model works.

$$\sum_{i=1}^c \sum_{j=1}^{n_i} (s_{ij} + (n_i - k_{ij})\gamma) \quad (29)$$

$$= \gamma \sum_{i=1}^c n_i^2 + \sum_{i=1}^c \sum_{j=1}^{n_i} (s_{ij}^* - \bar{s}_{ij}) - \sum_{i=1}^c \sum_{j=1}^{n_i} (k_{ij}^* - \bar{k}_{ij}) \gamma \quad (30)$$

$$= \gamma \sum_{i=1}^c n_i^2 + \sum_{i=1}^c (s_i^* - \gamma k_i^*) - \sum_{i=1}^c \sum_{j=1}^{n_i} \bar{s}_{ij} + \gamma \sum_{i=1}^c \sum_{j=1}^{n_i} \bar{k}_{ij} \quad (31)$$

Because s_i^* and k_i^* are constants, the problem (24) can be written equivalently as

$$\min_{\mathcal{F}_1, \dots, \mathcal{F}_c} \gamma \sum_{i=1}^c n_i^2 - \sum_{i=1}^c \sum_{j=1}^{n_i} \bar{s}_{ij} + \gamma \sum_{i=1}^c \sum_{j=1}^{n_i} \bar{k}_{ij} \quad (32)$$

Assuming $d(\mathbf{x}_i, \mathbf{x}_j) = e, \forall \mathbf{x}_i \in \mathbf{X}, \mathbf{x}_j \in \mathcal{A}(\mathbf{x}_i)$, the problem (32) can be simplified further as

$$\min_{\mathcal{F}_1, \dots, \mathcal{F}_c} \gamma \sum_{i=1}^c n_i^2 + (\gamma - e) \sum_{i=1}^c \sum_{j=1}^{n_i} \bar{k}_{ij} \quad (33)$$

To make the clustering results meaningful, $\gamma > e$ is required.

It can be seen from Eq. (32) that (i). our model tends to produce more balanced clustering results (the first term), (ii). the proposed model tends to group \mathbf{x}_i and its neighbors into the same cluster (the third term), (iii). if grouping \mathbf{x}_i and its neighbors into different clusters will result in a smaller objective function value, then the model will give priority to the neighbors that are far away from the sample.

The simplified model (33) retains the properties (i) and (ii) of the model (32), which means the proposed model is also applicable to unweighted graphs.

3.3 Graph Construction

In practical applications, a weighted graph containing similarities between samples is more common. In other words, the squared Euclidean distance between two points that are connected in the graph can not be obtained directly. In this situation, it is recommended to use the following formula to convert the similarities into the distances.

$$\tilde{d}_{ij} = -\log(s_{ij}), \quad (34)$$

where s_{ij} represents the similarity between points \mathbf{x}_i and \mathbf{x}_j . In some situations, only features are provided, at this time, it is recommended to construct a k -NN graph using fast approximate nearest neighbor algorithms. In this article, the k-d tree and EFANNA are employed to construct the approximate nearest neighbor graph required by the model for synthetic and real-world datasets, respectively.

3.4 Initialization

An intuitive way is to group the closer samples (the distance between the two is less than γ , for example) into one cluster. However, it is difficult for us to find an appropriate value for γ . In this paper, we treat the nearest neighbours as connected points and establish a path for each sample by breadth-first search method. Then group the first n/c points of the path into one cluster. If the number of clusters obtained is more than c , the clusters with the smaller number of samples are randomly assigned to other clusters.

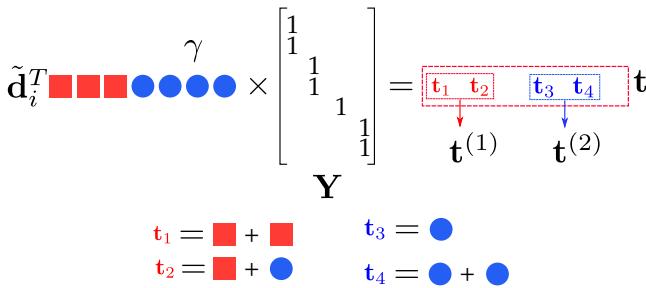


Fig. 2: Schematic diagram of calculating the product of $\tilde{\mathbf{d}}_i^T$ and \mathbf{Y} . Red squares represent the squared Euclidean distances between samples, and blue circles denote the elements that are equal to γ .

3.5 Optimization

To solve the i -th row of \mathbf{Y} , we fixed the other rows as constants. To begin with, let $\mathbf{U}_{(i)}$ be the permutation matrix to swap the i -th and first rows and \mathbf{Y}_r^T be the matrix $\mathbf{U}_{(i)}\mathbf{Y}$ with the first row removed, we have

$$\min_{\mathbf{y}_i} \operatorname{Tr}(\mathbf{Y}^T \tilde{\mathbf{D}} \mathbf{Y}) \quad (35)$$

$$\Leftrightarrow \min_{\mathbf{y}_i} \operatorname{Tr}(\mathbf{Y}^T \mathbf{U}_{(i)}^T \mathbf{U}_{(i)} \tilde{\mathbf{D}} \mathbf{U}_{(i)}^T \mathbf{U}_{(i)} \mathbf{Y}) \quad (36)$$

$$\Leftrightarrow \min_{\mathbf{y}_i} \operatorname{Tr}\left([\mathbf{y}_i \quad \mathbf{Y}_r] \begin{bmatrix} \tilde{d}_{ii} & \mathbf{v}^T \\ \mathbf{v} & \tilde{\mathbf{D}}_r \end{bmatrix} [\mathbf{y}_i^T \quad \mathbf{Y}_r^T]\right) \quad (37)$$

$$\Leftrightarrow \min_{\mathbf{y}_i} \operatorname{Tr}(\mathbf{y}_i \tilde{d}_{ii} \mathbf{y}_i^T) + 2\operatorname{Tr}(\mathbf{Y}_r \mathbf{v} \mathbf{y}_i^T) \quad (38)$$

where $\mathbf{v}^T = [\tilde{d}_{i,2} \dots \tilde{d}_{i,i-1} \tilde{d}_{i,1} \tilde{d}_{i,i+1} \dots \tilde{d}_{i,n}]$.

Since $\tilde{d}_{ii} = 0$, the model can be further simplified to

$$\min_{\mathbf{y}_i} \operatorname{Tr}(\mathbf{Y}_r \mathbf{v} \mathbf{y}_i^T) \quad (39)$$

$$\Leftrightarrow \min_{\mathbf{y}_i} \mathbf{y}_i^T \mathbf{Y}_r \mathbf{v} \quad (40)$$

Obviously $\mathbf{Y}_r \mathbf{v} = \mathbf{Y}^T \tilde{\mathbf{d}}_i$, so problem (40) is equivalent to

$$\min_{\mathbf{y}_i} \operatorname{Tr}(\mathbf{y}_i^T \mathbf{Y}^T \tilde{\mathbf{d}}_i) \quad (41)$$

From Eq. (35) to Eq. (41), we know that our model (24) can be optimized by solving Eq. (41) ($\forall i = 1, \dots, n$) iteratively until converging.

It is not difficult to find that the optimal solution of the i -th row of \mathbf{Y} , \mathbf{y}_i can be obtained by

$$y_{il} = \begin{cases} 1 & l = \arg \min_j (\tilde{\mathbf{d}}_i^T \mathbf{Y})_j \\ 0 & \text{otherwise} \end{cases} \quad (42)$$

3.6 Acceleration

Without any acceleration technique, it takes $O(nc)$ time to calculate the product of $\tilde{\mathbf{d}}_i^T$ and \mathbf{Y} . Note that \mathbf{Y} is an indicator matrix, that is, there are only n non-zero elements in \mathbf{Y} . Therefore, calculating the product of $\tilde{\mathbf{d}}_i^T$ and \mathbf{Y} can be completed in $O(n)$ time. However, the computational cost of the entire algorithm is $O(n^2)$, which means that the algorithm is difficultly scalable to large-scale problems.

We denote the elements in $\tilde{\mathbf{d}}_i^T$ that are smaller than γ by red squares, and denote the elements that are equal to γ by

Algorithm 1: Accelerated algorithm for solving problem (24).

Data: $\tilde{\mathbf{D}} \in \mathbb{R}^{n \times n}$, the number of cluster, c .

Result: Clustering result, \mathbf{g} .

Initialize \mathbf{g} randomly;

Compute \mathbf{n} (n_i is the number of points in cluster i);

Initialize vectors $\mathbf{u} \in \mathbb{R}^c$, $\mathbf{t} \in \mathbb{R}^c$, and $\mathbf{v} \in \mathbb{R}^c$ with 0;

while not converge **do**

for $i = 1, \dots, n$ **do**

$\tilde{c} = g_i$;

$\mathcal{A}(\mathbf{x}_i) = \{\mathbf{x}_j | \mathbf{x}_j \in \mathcal{N}_k(\mathbf{x}_i) \text{ and } \mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_j)\}$;

$\mathbf{t}, \mathbf{v} = f_1(\tilde{\mathbf{d}}_i, \mathbf{g}, \mathbf{n}, \mathbf{u}, \mathbf{v}, \mathbf{t}, \mathcal{A}(\mathbf{x}_i))$;

 // Compute the elements in $\mathbf{t}^{(2)}$;

for $l = 1 \dots c$ **do**

if $v_l == 0$ **then** $t_l = n_l * \gamma$;

for $\mathbf{x}_j \in \mathcal{A}(\mathbf{x}_i)$ **do**

$l = g_j$; $v_l = 0$;

 Find the minimum of \mathbf{t} , denote it by $t_{\hat{c}}$;

$g_i = \hat{c}$; $n_{\hat{c}-} = 1$; $n_{\hat{c}+} = 1$;

blue circles, as shown in Fig. 2. Let \mathbf{t} denote the product of $\tilde{\mathbf{d}}_i^T$ and \mathbf{Y} , i.e., $\mathbf{t} = \tilde{\mathbf{d}}_i^T \mathbf{Y}$, t_j can be regarded as the sum of some elements in $\tilde{\mathbf{d}}_i^T$. Therefore, there are k elements at most of \mathbf{t} , called $\mathbf{t}^{(1)}$, that is related to the red squares, and the remaining elements, called $\mathbf{t}^{(2)}$, are only related to the blue circles. We denote the j -th element in $\mathbf{t}^{(2)}$ by the i -th element in \mathbf{t} , i.e., $t_j^{(2)} = t_i^{(1)}$. Then we have

$$t_j^{(2)} = n_i \gamma. \quad (43)$$

If we record the number of samples of each cluster at the beginning of the optimization, then the elements in $\mathbf{t}^{(2)}$ can be directly obtained by Eq. (43) when solving \mathbf{y}_i . In addition, since the coordinate descent algorithm is adopted, only $O(1)$ time is needed to update the number of samples in each cluster if \mathbf{y}_i is updated. So far, the computational overhead of our algorithm is $O(nk + nc)$. The detailed algorithm to solve the problem (24) is summarized in Algorithm 1.

3.6.1 Speed Up Further

Although the computational overhead of Algorithm 1 is linear with respect to the number of samples, it involves the product of the number of points and clusters. Therefore, it is very time-consuming, if the number of clusters to construct is very large.

We find the minimum value of \mathbf{t} by

$$t^* = \begin{cases} \min(\mathbf{t}^{(1)}), & \text{if } \min(\mathbf{t}^{(1)}) \leq \min(\mathbf{t}^{(2)}) \\ \min(\mathbf{t}^{(2)}), & \text{otherwise} \end{cases} \quad (44)$$

where t^* denotes the minimum of \mathbf{t} . Obviously, it takes $O(k)$ and $O(c - k)$ time to find the minimum of $\mathbf{t}^{(1)}$ and $\mathbf{t}^{(2)}$, respectively. Fortunately, the elements in $\mathbf{t}^{(2)}$ are all multiples of γ , thus, the cluster with the smallest number of samples in $\mathbf{t}^{(2)}$ corresponds to its minimum value, so we can find its minimum in a faster way.

Before introducing the acceleration algorithm, the following lemma is introduced.

Algorithm 2: The function used to calculate $\mathbf{t}^{(1)}$, called $f_1(\tilde{\mathbf{d}}_i, \mathbf{g}, \mathbf{n}, \mathbf{u}, \mathbf{v}, \mathbf{t}, \mathcal{A}(\mathbf{x}_i))$.

Data: $\tilde{\mathbf{d}}_i \in \mathbb{R}^n$, $\mathbf{g} \in \mathbb{R}^n$, $\mathbf{n}, \mathbf{u}, \mathbf{v}, \mathbf{t} \in \mathbb{R}^c$, $\mathcal{A}(\mathbf{x}_i)$.
Result: Vector \mathbf{t} (t_j has been calculated if $t_j \in \mathbf{t}^{(1)}$).
 Vector \mathbf{v} that is modified.

```

for  $\mathbf{x}_j \in \mathcal{A}(\mathbf{x}_i)$  do
     $l = g_j; t_l = 0; u_l = 0;$ 
for  $\mathbf{x}_j \in \mathcal{A}(\mathbf{x}_i)$  do
     $l = g_j; t_l + = \tilde{d}_{ij}; u_l + = 1;$ 
for  $\mathbf{x}_j \in \mathcal{A}(\mathbf{x}_i)$  do
     $l = g_j;$ 
    if  $v_l == 0$  then
         $v_l = 1;$ 
         $t_l + = \gamma(n_l - u_l);$ 
    
```

Lemma 1. Let $\mathbf{n} = [n_1, \dots, n_c]^T \in \mathbb{R}^c$ denote an ordered vector, where $0 \leq n_i \leq N$, and $N < \infty$ is a positive integer. It takes only $O(1)$ time to maintain its order if one element of it is increased or reduced by 1.

The vector \mathbf{n} representing the number of samples of each cluster just satisfies the conditions described Lemma 2, thus, it takes only $O(1)$ time to maintain its order if y_i is updated. In other words, if we sort \mathbf{n} at the beginning of optimization, it takes $O(k)$ time at most to get the minimum value of $\mathbf{t}^{(2)}$.

So far, the computational overhead of our algorithm has been reduced to $O(nk)$. However, in practice, we find the minimum value of t through the simpler formula (45) instead of the formula (44).

$$\min(\mathbf{t}) = \begin{cases} \min(\mathbf{t}^{(1)}), & \text{if } \min(\mathbf{t}^{(1)}) \leq \gamma q_1 \\ \gamma q_1, & \text{otherwise} \end{cases} \quad (45)$$

where q_1 is the number of samples in p_1 which is the cluster with the smallest number of samples in \mathbf{t} .

Equation (45) holds because of the following facts:

- γq_1 is an element of $\mathbf{t}^{(2)}$. In this case, $\gamma q_1 = \min(\mathbf{t}^{(2)})$, and Eq. (45) is then equivalent to Eq. (44).
- γq_1 does not belong to $\mathbf{t}^{(2)}$. In this case, we have

$$\min(\mathbf{t}^{(1)}) \leq \gamma q_1 \leq \min(\mathbf{t}^{(2)}).$$

Thus, $\min(\mathbf{t}^{(1)})$ is the minimum of \mathbf{t} , Eq. (45) holds.

The finally adopted algorithm for solving the problem (24) is summarized in Algorithm 3.

In short, we have carefully designed the optimization algorithm for problem (24), which reduces the computational complexity from $O(n^2c)$ to $O(nk)$, as shown in Table 1, where k is the number of neighborhoods in graph construction, not the number of clusters.

TABLE 1: The time complexity of algorithms

Cases	Time complexity
Original	$O(n^2c)$
Consider the sparsity of \mathbf{Y}	$O(n^2)$
Consider the property of $\tilde{\mathbf{d}}_i^T$	$O(n(k+c))$
The final optimization algorithm	$O(nk)$

Algorithm 3: The finally adopted algorithm for solving problem (24).

Data: $\tilde{\mathbf{D}} \in \mathbb{R}^{n \times n}$, the number of cluster, c .
Result: Clustering result, \mathbf{g} .
 Initialize \mathbf{g} randomly;
 Initialize vectors $\mathbf{u} \in \mathbb{R}^c$, $\mathbf{t} \in \mathbb{R}^c$, and $\mathbf{v} \in \mathbb{R}^c$ with 0;
 Compute \mathbf{n} (n_i is the number of points in cluster i);
 // Sort $\mathbf{n} \in \mathbb{R}^c$ by bucket sort algorithm, $n[p_i] = q_i$;
 $\mathbf{p}, \mathbf{q} = \text{sort}(\mathbf{n})$;
while not converge **do**
for $i = 1, \dots, n$ **do**
 $\tilde{c} = g_i$;
 $\mathcal{A}(\mathbf{x}_i) = \{\mathbf{x}_j | \mathbf{x}_j \in \mathcal{N}_k(\mathbf{x}_i) \text{ and } \mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_j)\}$;
 $\mathbf{t}, \mathbf{v} = f_1(\tilde{\mathbf{d}}_i, \mathbf{g}, \mathbf{n}, \mathbf{u}, \mathbf{v}, \mathbf{t}, \mathcal{A}(\mathbf{x}_i))$;
for $\mathbf{x}_j \in \mathcal{A}(\mathbf{x}_i)$ **do**
 $l = g_j; v_l = 0$;
 Get the minimum of $\mathbf{t}^{(1)}$;
 Obtain the minimum of \mathbf{t} by judging whether $\min(\mathbf{t}^{(1)}) < \gamma q_1$ holds, and denote it by $t_{\hat{c}}$;
 $g_i = \hat{c}$;
 $n_{\hat{c}-} = 1; n_{\hat{c}+} = 1$;
 Update \mathbf{p} and \mathbf{q} based on Lemma 1.;

3.7 Time and space complexity

- Space complexity: In finally adopted algorithm, the memory is mainly occupied by matrix $\tilde{\mathbf{D}}$, which is a dense matrix. But with the property shown in Eq. (32), we can find that there are at most nk elements in $\tilde{\mathbf{D}}$ that need to be stored. Thus, the space complexity of the algorithm is $O(nk)$.
- Time complexity: $O(n)$ time is needed to sort \mathbf{n} if the bucket sort algorithm is adopted. $O(k)$ time is needed to find the minimum of $\mathbf{t}^{(1)}$, as shown in Algorithm 2. According to Lemma 1, it takes only $O(1)$ time to update \mathbf{p} and \mathbf{q} . Hence, the computational overhead of Algorithm 3 is $O(nk)$.

3.8 Advantages

Our model has the following advantages:

- **Efficient.** An extremely efficient optimization algorithm is designed for the problem (24), and the time and space complexity is greatly reduced. Specifically, the computational and memory overhead are both $O(nk)$, which has been discussed in subsection 3.7;
- **Effective.** Since a k-NN graph is adopted in k -sums, it can separate clusters that are not separable linearly in input space. Thus, it should be yield better experimental performance, which is verified in Section 5;
- **Robust.** The distance between \mathbf{x}_i and \mathbf{x}_j not calculated if $\mathbf{x}_i \notin \mathcal{N}_k(\mathbf{x}_j)$ or $\mathbf{x}_j \notin \mathcal{N}_k(\mathbf{x}_i)$, but a parameter γ is used instead, which means that our model is robust to outliers, and this point is verified by the experiment performed on “Outlier”, in Section 5.

In addition, our model only involves the squared Euclidean distance, so it has no trouble measuring the similarities between samples.

4 A VARIANT OF K-SUMS

If the dimensionality of the data is high, even if the value of k is appropriate, the quality of the graph will be poor, which will lead to poor performance of the clustering algorithm. In addition, it is non-trivial to choose an appropriate k value for all datasets. Thus, it is desirable to get an algorithm that takes features as input. We consider the following model and refer to it k -sums-x.

$$\min_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr}(\mathbf{Y}^T \mathbf{D} \mathbf{Y}), \quad (46)$$

where $\mathbf{D} \in \mathbb{R}^{n \times n}$ denotes the distance matrix, $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$. Again, the standard CD algorithm is used to optimize the problem (46). Because the matrix \mathbf{D} involved in the model is a dense matrix of size n by n , the following two issues that users care about are naturally raised.

- The dense matrix \mathbf{D} is included in the proposed model. Although \mathbf{D} is symmetric, there are still n elements with different values. Does the algorithm need to store the values of these elements? In that case, the algorithm will have a high memory overhead.
- If the standard CD algorithm is used to optimize the problem (46), the product of \mathbf{d} and \mathbf{Y} will still be involved in the sub-problem. However, all elements of \mathbf{d} may not be equal, so the acceleration technique described earlier does obviously not work, which may cause the computational overhead to increase quadratically with the number of samples.

Considering the problems discussed above, we designed a new acceleration technology for the sub-problems involved in the coordinate descent algorithm. This method relies on the property of the squared Euclidean distance, which is one of the reasons why the squared Euclidean distance is used instead of other distances.

It is worth noting that by adopting the proposed acceleration algorithm, the computational and memory overhead are both reduced largely reduced. Specifically, The time and space complexity of the proposed optimization algorithm are both linear with respect to the number of points.

4.1 Analysis

Similar to the previously proposed model, model (46) also tends to produce a more balanced partition. Specifically, the following equations hold.

$$\min_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr}(\mathbf{Y}^T \mathbf{D} \mathbf{Y}) \quad (47)$$

$$\Leftrightarrow \min_{\mathbf{Y} \in \Phi^{n \times c}} \sum_{k=1}^c \sum_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{A}_k} d_{ij} \quad (48)$$

$$\Leftrightarrow \min_{\mathbf{Y} \in \Phi^{n \times c}} \sum_{k=1}^c n_i^2 \bar{d}_i \quad (49)$$

where \bar{d}_i denotes the average of the distances between samples in the i -th cluster, $\bar{d}_i = \frac{1}{n_i^2} \sum_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{A}_k} d_{ij}$. We consider the case where the average distance of each cluster is approximately equal, i.e., $\bar{d}_1 \approx \dots \approx \bar{d}_c$. Then the problem (49) can be written equivalently as

$$\min_{\mathbf{Y} \in \Phi^{n \times c}} \sum_{k=1}^c n_i^2. \quad (50)$$

Obviously, the $n_i = n/c, \forall i = 1, \dots, c$ is the optimal solution of the problem (50).

4.2 Connection with k -means

Theorem 2. *The problem (46) is equivalent to the problem of k -means 2, if we constrain the partitions must be balanced ($\mathbf{Y}^T \mathbf{Y} = \beta I$, where β is a constant).*

Proof. With the constraint of $\mathbf{Y}^T \mathbf{Y} = \beta I$, problem (46) can be rewritten as follows.

$$\min_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr}(\mathbf{Y}^T \mathbf{D} \mathbf{Y}) \quad (51)$$

$$\Leftrightarrow \min_{\mathbf{Y} \in \Phi^{n \times c}} \text{Tr}\left(\left(\mathbf{Y}^T \mathbf{Y}\right)^{-1} \mathbf{Y}^T \mathbf{D} \mathbf{Y}\right) \quad (52)$$

$$\Leftrightarrow \min_{\mathbf{Y} \in \Phi^{n \times c}} \text{diag}\left(\left(\mathbf{Y}^T \mathbf{Y}\right)^{-1}\right)^T \text{diag}(\mathbf{Y}^T \mathbf{D} \mathbf{Y}) \quad (53)$$

$$\Leftrightarrow \min_{\mathbf{Y} \in \Phi^{n \times c}} \sum_{k=1}^c \frac{1}{n_k} \sum_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{A}_k} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \quad (54)$$

where, $\text{diag}(\mathbf{A})$ is the function to extract the diagonal vector from the matrix, \mathbf{A} , i.e., $\text{diag}(\mathbf{A}) = [a_{1,1}, \dots, a_{n,n}]^T$. The Eq. (54) is exactly the problem of k -means as shown in Section 2.4, completing the proof. \square

4.3 Optimization

Similar to k -sums, the coordinate descent algorithm is used again to solve problem (46). Therefore, the optimal solution of \mathbf{y}_i can be obtained by

$$y_{il} = \begin{cases} 1 & l = \arg \min_j (\mathbf{d}_i^T \mathbf{Y})_j \\ 0 & \text{otherwise} \end{cases} \quad (55)$$

Let \mathbf{t} represent the product of \mathbf{d}_i^T and \mathbf{Y} , i.e., $\mathbf{t} = \mathbf{d}_i^T \mathbf{Y}$, then t_j can be regarded as the sum of some elements of \mathbf{d}_i . Fortunately, we have a simpler formula for summing these squared Euclidean distances as shown in Figure ??.

$$\sum_{i=1}^n \|\mathbf{z} - \mathbf{x}_i\|_2^2 = n\|\mathbf{z}\|_2^2 + \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 - 2\mathbf{z}^T \sum_{i=1}^n \mathbf{x}_i \quad (56)$$

With the help of Formula (56), the k -th element of \mathbf{t} can be expressed as

$$t_k = \sum_{\mathbf{x}_j \in \mathcal{A}_k} \|\mathbf{x}_j\|_2^2 + \sum_{\mathbf{x}_j \in \mathcal{A}_k} \|\mathbf{x}_j\|_2^2 - 2\mathbf{x}_k^T \sum_{\mathbf{x}_j \in \mathcal{A}_k} \mathbf{x}_j \quad (57)$$

For convenience, let

$$\mathbf{S} \in \mathbb{R}^{c \times d}, \mathbf{s}_k = \sum_{\mathbf{x}_j \in \mathcal{A}_k} \mathbf{x}_j, \quad (58)$$

$$\mathbf{u} \in \mathbb{R}^n, u_j = \|\mathbf{x}_j\|_2^2, \quad (59)$$

$$\mathbf{n} \in \mathbb{R}^c, n_k = \sum_{\mathbf{x}_j \in \mathcal{A}_k} 1 \quad (60)$$

$$\mathbf{v} \in \mathbb{R}^c, v_k = \sum_{\mathbf{x}_j \in \mathcal{A}_k} \|\mathbf{x}_j\|_2^2, \quad (61)$$

the vector \mathbf{t} can be expressed as

$$t_k = n_k u_i + v_k - 2\mathbf{x}_k^T \mathbf{s}_k, \quad (62)$$

and can be calculated in $O(dc)$ time, if the variables (n_k , u_k , v_k , and \mathbf{s}_k) are computed in advance. In addition, it is not difficult to find that we can update variables \mathbf{s}_k , v_k , and n_k in only $O(d)$ time, if \mathbf{y}_i is updated. Therefore, the computational complexity of the whole algorithm is $O(ndc)$.

Algorithm 4: The finally adopted algorithm for solving problem (46).

Data: $\mathbf{X} \in \mathbb{R}^{n \times d}$, the number of cluster, c .
Result: Clustering result, \mathbf{g} .
 Initialize \mathbf{g} randomly;
 Calculate the l_2 norm of each point, i.e.,
 $\|\mathbf{x}_i\|, \forall i = 1, \dots, n$;
 Compute \mathbf{n} (n_i is the number of points in cluster i);
 Compute \mathbf{v} and \mathbf{s} according to Eq. (61) and Eq. (58);
while not converge **do**
for $i = 1, \dots, n$ **do**
 $\hat{c} = g_i$;
 Compute \mathbf{t} according to Eq. (57);
 Obtain the minimum of \mathbf{t} and denote it by $t_{\hat{c}}$;
 // Update \mathbf{g} ;
 $g_i = \hat{c}$;
 // Update \mathbf{n} ;
 $n_{\hat{c}-} = 1; n_{\hat{c}+} = 1$;
 // Update \mathbf{v} ;
 $v_{\hat{c}-} = \|\mathbf{x}_i\|; v_{\hat{c}+} = \|\mathbf{x}_i\|$;
 // Update \mathbf{S} ;
 $s_{\hat{c}-} = \mathbf{x}_i; s_{\hat{c}+} = \mathbf{x}_i$;

4.4 Time and space complexity

- Space complexity: Variables $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{S} \in \mathbb{R}^{c \times d}$, $\mathbf{u} \in \mathbb{R}^n$, $\mathbf{n} \in \mathbb{R}^c$, and $\mathbf{v} \in \mathbb{R}^c$ need to be stored, and the temporary variables involved in the algorithm occupy much less memory than these variables, so the space complexity is $O(nd)$, which is linear with respect to the number of points.
- Time complexity: It takes $O(n)$, $O(nd)$, $O(n)$, and $O(nd)$ time to calculate variables \mathbf{n} , \mathbf{u} , \mathbf{v} , and \mathbf{S} in advance, and it takes $O(1)$, $O(1)$, and $O(d)$ time to update variables \mathbf{n} , \mathbf{v} , and \mathbf{S} , if y_i is updated. In addition, the vector \mathbf{t} can be computed in $O(dc)$ time if these variables are computed in advance. $O(c)$ time is needed to find the minimum value of \mathbf{t} . Therefore, the computational overhead of the proposed algorithm is $O(ndc)$, which is the same as traditional k -means.

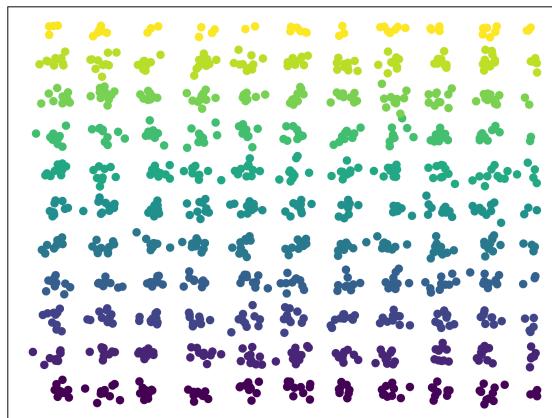


Fig. 3: Distribution of D_1 .

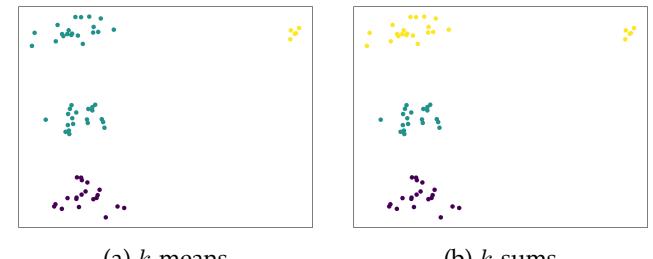


Fig. 4: Performance on Outlier

TABLE 2: The performance on Grid-like datasets. (n : the number of samples, c : the number of clusters).

Data info n	c	FIN		k -means		k -sums		
		Time	F_1	Time	F_1	Time	F_1	
D_1	100K	5K	4.07	0.895	14.40	0.892	0.49	0.992
D_2	100K	10K	4.31	0.961	24.55	0.882	0.36	0.991
D_3	100K	20K	3.57	0.872	36.87	0.867	0.31	0.986
D_4	200K	5K	7.60	0.951	34.96	0.898	1.73	0.992
D_5	200K	10K	7.44	0.893	59.71	0.891	1.03	0.992
D_6	200K	20K	7.74	0.961	96.80	0.881	1.01	0.992
D_7	300K	5K	11.73	0.954	55.16	0.900	3.66	0.992
D_8	300K	10K	11.86	0.874	98.36	0.896	2.05	0.992
D_9	300K	20K	11.72	0.940	168.03	0.888	1.83	0.991

5 EXPERIMENTS

Extensive experiments on synthetic and benchmark datasets have been conducted, to validate the superiority of our algorithms, i.e., k -sums and k -sums-x. In this section, these experiments are introduced in detail.

5.1 Experiments conducted on synthetic datasets

5.1.1 Datasets

Ten synthetic datasets are used, including D_1 , D_2 , D_3 , D_4 , D_5 , D_6 , D_7 , D_8 , D_9 , and Outlier. The first nine are grid-like datasets, specifically, their distribution is roughly as shown in Figure 3, and the statistics of these datasets are summarized in Table 2.

5.1.2 Experimental results

Experiments based on synthetic datasets are mainly used to verify the advantages mentioned in Section 3.8.

1. From the results in Table 2, it can be seen that the running time of k -sums (The time of constructing the k-NN graph by KD-Tree is included) is approximately linear with respect to the number of samples in the data, and is irrelevant to the number of clusters. So, for the datasets with more clusters, the advantage of our algorithm is more obvious, which verifies the first advantage.

2. It can be seen from Table 2 while having a low time complexity, k -sums shows high-quality clustering results, which verifies the second advantage.

3. As shown in Figure 4, the cluster composed of a few abnormal points heavily affected the performance of k -means but nearly has no effect on our algorithm. Therefore, our algorithm is robust to abnormal points.

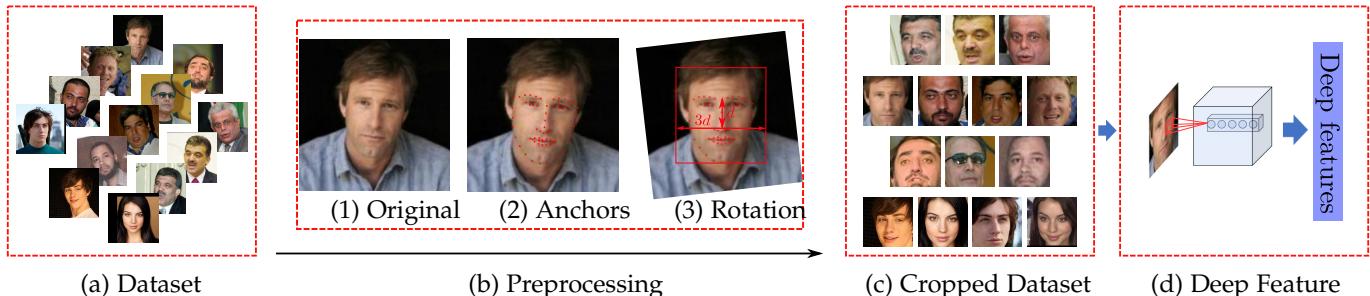


Fig. 5: Schematic diagram of face representation.

5.2 Experiments conducted on benchmark datasets

5.2.1 Datasets

Eighteen real-world datasets are involved for the verification of the clustering performance of the proposed methods (k -sums and k -sums-x). The description of them are as follows:

FaceV5² is a widely used dataset in face recognition, it contains 2,500 images of 500 subjects under various illumination, pose, expression, and imaging distance.

CFPW [31] is a widely used face verification dataset, containing 7,000 images of 500 celebrities in frontal and profile views. Each image contains 16×16 pixels.

FEI [32] dataset is a Brazilian face database. There are 14 images for each of 200 individuals (students or staff), adding up to a total number of 2,800 images. We use a subset (700 images of 50 individuals) in this paper.

Grimace [33] dataset contains 360 face images with 20 individuals. Images of individuals from various racial origins have different facial expressions on their beards and glasses.

IMM [34] dataset consisting of 240 annotated face images of 40 (7 females and 33 males) different individuals, each with 6 images. In this data, no one wears glasses.

MSRA25 [35] dataset consists of 1,799 face gray images in 12 classes represented by 256 features, each class has 113—186 samples with different poses.

MUCT [36] consists of 3,755 faces of 276 individuals, each image has 76 manual landmarks. This dataset provides a diversity of lighting, age, and ethnicity.

ORL [37] dataset contains 10 images of each of 40 distinct subjects. The images were taken at different times, varying the lighting, facial expressions and facial details.

YaleFace [38] contains 165 grayscale images of 15 individuals. There are 11 images per subject under different facial expressions or configurations.

Face96 [39] are the open source datasets of face recognition provided by the University of Essex. Face96 has 3,016 images of 151 individuals. Each sample is cropped into 16×16 scale.

Isolet [40], **Mpeg7** [41], and **Palm**³ are widely used datasets for clustering and classification. Isolet consists of 7797 samples of 26 letters represented by 617 features spoken by 150 speakers. Mpeg7 consists of 1400 shapes grouped into 70 classes. Palm contains 2000 fingerprint images generated by 100 different palms.

2. <http://biometrics.idealtest.org/dbDetailForUser.do?id=9>
3. <http://www.escience.cn/people/chenxiaojun>

TABLE 3: Benchmark datasets

	FaceV5	CFPW	FEI	Isolet	Grimace
# Samples	2500	7000	700	7797	360
# Features	256	256	147	617	256
# Subjects	500	500	50	26	18
	Mpeg7	Face96	IMM	YaleFace	MSRA25
# Samples	1400	3016	240	165	1799
# Features	728	256	256	256	256
# Subjects	70	151	40	15	12
	Palm	MUCT	ORL	PEAL	CACD
# Samples	2000	3755	400	30863	163446
# Features	256	256	256	256	256
# Subjects	100	276	40	1040	2000
	WebFace	CelebA			
# Samples	494414	202599			
# Features	256	256			
# Subjects	10575	10177			

To verify the efficiency of k -sums, four large-scale datasets are used. They are **PEAL** [42]⁴, **CACD** [43], **WebFace** [44], and **CelebA** [45]. The statistics of all real-world datasets are summarized in Table 3.

5.2.2 Face representation

Because of the difference in lighting, posture, and background, it is very unwise to directly measure the Euclidean distance between two face images. Therefore, we use a pre-trained deep neural network model [46] to extract the features of the face images and then performing clustering algorithms in the feature space.

As shown in Figure 5, before feeding the pictures to the model, we use Dlib to find the location of 68 facial landmarks. Then the face image is rotated so that 20- and 25-th landmarks are located on the same horizontal line. Next, with 31-th landmark as the center, we cut out a square with a length of $3d$ (d is the vertical distance from the 31-th landmark to the 25-th landmark).

5.2.3 Baseline

To verify the effectiveness of the proposed k -sums, we compare it with five graph-based clustering methods including: Anchor-based Graph Clustering-improved algorithm

4. The subset named CAS-PEAL-R1 is used, which is publicly available.

TABLE 4: Performance of algorithms taking graph as input (mean (\pm std))

		AGCI	FCDMF	FINCH	RCAE	RCC	k -sums
FaceV5	ACC	0.734 (\pm 9.0e-03)	0.546 (\pm 1.2e-02)	0.866 (\pm 0.0e+00)	0.834 (\pm 2.1e-02)	0.389 (\pm 0.0e+00)	0.956 (\pm5.0e-03)
	NMI	0.932 (\pm 2.9e-03)	0.847 (\pm 4.1e-03)	0.935 (\pm 0.0e+00)	0.882 (\pm 9.6e-03)	0.632 (\pm 0.0e+00)	0.986 (\pm1.5e-03)
	ARI	0.626 (\pm 2.0e-02)	0.342 (\pm 1.4e-02)	0.236 (\pm 0.0e+00)	0.115 (\pm 1.6e-02)	0.029 (\pm 0.0e+00)	0.931 (\pm6.8e-03)
	Time	16.823 (\pm 1.2e-02)	5.530 (\pm 7.4e-03)	4.943 (\pm 0.0e+00)	1.7e+2 (\pm 8.1e+1)	8.926 (\pm 1.8e-02)	0.295 (\pm4.4e-04)
CFPW	ACC	0.541 (\pm 9.5e-03)	0.357 (\pm 9.1e-03)	0.665 (\pm 0.0e+00)	0.726 (\pm 3.3e-02)	0.749 (\pm 0.0e+00)	0.761 (\pm5.4e-03)
	NMI	0.770 (\pm 6.4e-03)	0.686 (\pm 5.3e-03)	0.679 (\pm 0.0e+00)	0.736 (\pm 2.8e-02)	0.883 (\pm 3.8e-04)	0.895 (\pm1.8e-03)
	ARI	0.206 (\pm 1.7e-02)	0.158 (\pm 9.1e-03)	0.016 (\pm 0.0e+00)	0.029 (\pm 8.3e-03)	0.323 (\pm 5.8e-03)	0.647 (\pm6.4e-03)
	Time	50.718 (\pm 1.1e-01)	17.581 (\pm 3.1e-02)	40.038 (\pm 0.0e+00)	1.4e+3 (\pm 1.1e+3)	21.132 (\pm 2.5e-01)	2.346 (\pm4.0e-03)
Face96	ACC	0.655 (\pm 2.7e-02)	0.758 (\pm 1.5e-02)	0.966 (\pm 0.0e+00)	0.970 (\pm 1.3e-02)	0.975 (\pm 0.0e+00)	0.992 (\pm3.3e-16)
	NMI	0.875 (\pm 1.5e-02)	0.921 (\pm 4.6e-03)	0.979 (\pm 0.0e+00)	0.980 (\pm 5.6e-03)	0.989 (\pm 0.0e+00)	0.996 (\pm2.0e-16)
	ARI	0.446 (\pm 6.6e-02)	0.725 (\pm 1.5e-02)	0.888 (\pm 0.0e+00)	0.898 (\pm 4.4e-02)	0.961 (\pm 0.0e+00)	0.986 (\pm2.2e-16)
	Time	7.988 (\pm 1.2e-02)	4.050 (\pm 1.4e-02)	0.648 (\pm 0.0e+00)	1.1e+2 (\pm 5.2e+1)	9.095 (\pm 6.9e-02)	0.391 (\pm7.8e-04)
FEI	ACC	0.459 (\pm 2.1e-02)	0.420 (\pm 1.7e-02)	0.516 (\pm 0.0e+00)	0.518 (\pm 6.0e-02)	0.043 (\pm 0.0e+00)	0.540 (\pm1.6e-02)
	NMI	0.689 (\pm 1.1e-02)	0.653 (\pm 7.6e-03)	0.734 (\pm0.0e+00)	0.687 (\pm 3.6e-02)	0.047 (\pm 0.0e+00)	0.728 (\pm 9.0e-03)
	ARI	0.337 (\pm 1.8e-02)	0.289 (\pm 1.4e-02)	0.363 (\pm 0.0e+00)	0.316 (\pm 6.4e-02)	0.003 (\pm 0.0e+00)	0.407 (\pm1.8e-02)
	Time	0.330 (\pm 7.9e-04)	0.215 (\pm 6.8e-05)	0.034 (\pm 0.0e+00)	3.452 (\pm 1.4e+00)	1.239 (\pm 1.6e-02)	0.019 (\pm4.3e-04)
Grimace	ACC	0.791 (\pm 7.2e-02)	0.754 (\pm 4.9e-02)	1.000 (\pm0.0e+00)	0.984 (\pm 2.1e-02)	1.000 (\pm0.0e+00)	1.000 (\pm0.0e+00)
	NMI	0.908 (\pm 4.7e-02)	0.868 (\pm 2.2e-02)	1.000 (\pm0.0e+00)	0.988 (\pm 1.1e-02)	1.000 (\pm0.0e+00)	1.000 (\pm0.0e+00)
	ARI	0.768 (\pm 1.3e-01)	0.721 (\pm 4.1e-02)	1.000 (\pm0.0e+00)	0.975 (\pm 2.6e-02)	1.000 (\pm0.0e+00)	1.000 (\pm0.0e+00)
	Time	0.093 (\pm 2.8e-04)	0.081 (\pm 1.5e-04)	0.026 (\pm 0.0e+00)	0.629 (\pm 4.1e-01)	0.872 (\pm 2.7e-02)	0.004 (\pm8.9e-07)
IMM	ACC	0.761 (\pm 5.2e-02)	0.734 (\pm 2.6e-02)	0.996 (\pm 0.0e+00)	0.873 (\pm 3.0e-02)	0.871 (\pm 0.0e+00)	1.000 (\pm0.0e+00)
	NMI	0.919 (\pm 2.4e-02)	0.887 (\pm 1.2e-02)	0.997 (\pm 0.0e+00)	0.915 (\pm 2.3e-02)	0.953 (\pm 0.0e+00)	1.000 (\pm0.0e+00)
	ARI	0.721 (\pm 7.5e-02)	0.668 (\pm 3.1e-02)	0.991 (\pm 0.0e+00)	0.701 (\pm 8.8e-02)	0.862 (\pm 1.1e-16)	1.000 (\pm0.0e+00)
	Time	0.057 (\pm 2.4e-04)	0.052 (\pm 4.9e-05)	0.130 (\pm 0.0e+00)	0.622 (\pm 3.1e-01)	0.353 (\pm 1.3e-02)	0.002 (\pm1.4e-06)
Isolet	ACC	0.535 (\pm 2.5e-02)	0.260 (\pm 1.5e-02)	0.469 (\pm 0.0e+00)	0.555 (\pm 1.6e-02)	0.038 (\pm 0.0e+00)	0.619 (\pm1.1e-02)
	NMI	0.710 (\pm 1.3e-02)	0.425 (\pm 1.3e-02)	0.700 (\pm 0.0e+00)	0.703 (\pm 2.0e-02)	0.000 (\pm 0.0e+00)	0.733 (\pm5.2e-03)
	ARI	0.476 (\pm 2.9e-02)	0.183 (\pm 1.4e-02)	0.415 (\pm 0.0e+00)	0.456 (\pm 2.7e-02)	0.000 (\pm 0.0e+00)	0.533 (\pm9.5e-03)
	Time	23.445 (\pm 4.8e-02)	20.949 (\pm 2.7e-03)	1.617 (\pm0.0e+00)	1.5e+3 (\pm 3.6e+2)	8.2e+2 (\pm 3.6e+0)	3.370 (\pm 9.5e-02)
Mpeg7	ACC	0.466 (\pm 2.4e-02)	0.439 (\pm 1.3e-02)	0.161 (\pm 0.0e+00)	0.466 (\pm 5.2e-02)	0.328 (\pm 1.2e-01)	0.541 (\pm7.1e-03)
	NMI	0.662 (\pm 1.6e-02)	0.653 (\pm 8.1e-03)	0.218 (\pm 0.0e+00)	0.629 (\pm 4.3e-02)	0.667 (\pm 3.5e-02)	0.725 (\pm3.4e-03)
	ARI	0.288 (\pm 3.0e-02)	0.303 (\pm 1.4e-02)	0.007 (\pm 0.0e+00)	0.155 (\pm 5.5e-02)	0.336 (\pm 1.4e-01)	0.415 (\pm6.1e-03)
	Time	3.999 (\pm 5.5e-03)	2.609 (\pm 2.5e-04)	2.921 (\pm 0.0e+00)	1.6e+1 (\pm 1.0e+1)	10.938 (\pm 1.3e-02)	0.088 (\pm7.3e-04)
MSRA25	ACC	0.520 (\pm 4.4e-02)	0.353 (\pm 2.2e-02)	0.444 (\pm 0.0e+00)	0.514 (\pm 2.2e-02)	0.139 (\pm 0.0e+00)	0.585 (\pm3.0e-02)
	NMI	0.591 (\pm 3.1e-02)	0.367 (\pm 2.7e-02)	0.495 (\pm 0.0e+00)	0.564 (\pm 1.8e-02)	0.065 (\pm 0.0e+00)	0.592 (\pm1.8e-02)
	ARI	0.368 (\pm 3.8e-02)	0.200 (\pm 3.3e-02)	0.220 (\pm 0.0e+00)	0.328 (\pm 3.6e-02)	0.034 (\pm 0.0e+00)	0.413 (\pm2.2e-02)
	Time	1.978 (\pm 3.4e-03)	1.564 (\pm 1.5e-02)	0.112 (\pm0.0e+00)	26.848 (\pm 8.9e+00)	34.898 (\pm 3.5e-02)	0.167 (\pm 8.3e-03)
MUCT	ACC	0.745 (\pm 2.2e-02)	0.734 (\pm 1.6e-02)	0.986 (\pm0.0e+00)	0.976 (\pm 1.3e-02)	0.885 (\pm 1.1e-16)	0.984 (\pm 2.3e-03)
	NMI	0.935 (\pm 1.1e-02)	0.918 (\pm 4.2e-03)	0.993 (\pm0.0e+00)	0.988 (\pm 6.3e-03)	0.963 (\pm 0.0e+00)	0.993 (\pm9.5e-04)
	ARI	0.661 (\pm 1.0e-01)	0.691 (\pm 1.5e-02)	0.968 (\pm 0.0e+00)	0.919 (\pm 7.4e-02)	0.848 (\pm 1.1e-16)	0.979 (\pm2.5e-03)
	Time	18.924 (\pm 1.6e-02)	5.605 (\pm 2.0e-02)	0.473 (\pm0.0e+00)	1.9e+2 (\pm 1.1e+2)	10.260 (\pm 7.5e-03)	0.632 (\pm 8.0e-04)
ORL	ACC	0.687 (\pm 7.8e-02)	0.758 (\pm 3.4e-02)	1.000 (\pm0.0e+00)	0.988 (\pm 9.9e-03)	0.998 (\pm 1.1e-16)	1.000 (\pm0.0e+00)
	NMI	0.871 (\pm 5.8e-02)	0.897 (\pm 1.2e-02)	1.000 (\pm0.0e+00)	0.992 (\pm 4.8e-03)	0.998 (\pm 0.0e+00)	1.000 (\pm0.0e+00)
	ARI	0.578 (\pm 1.9e-01)	0.714 (\pm 3.6e-02)	1.000 (\pm0.0e+00)	0.977 (\pm 1.7e-02)	0.997 (\pm 0.0e+00)	1.000 (\pm0.0e+00)
	Time	0.134 (\pm 3.0e-04)	0.112 (\pm 1.4e-04)	0.031 (\pm 0.0e+00)	0.924 (\pm 6.1e-01)	0.852 (\pm 1.8e-02)	0.006 (\pm9.5e-07)
Palm	ACC	0.695 (\pm 2.1e-02)	0.593 (\pm 1.3e-02)	0.591 (\pm 0.0e+00)	0.824 (\pm4.3e-02)	0.020 (\pm 0.0e+00)	0.807 (\pm 1.9e-02)
	NMI	0.892 (\pm 6.4e-03)	0.824 (\pm 5.8e-03)	0.778 (\pm 0.0e+00)	0.927 (\pm 1.0e-02)	0.012 (\pm 0.0e+00)	0.931 (\pm6.6e-03)
	ARI	0.649 (\pm 2.1e-02)	0.516 (\pm 1.3e-02)	0.148 (\pm 0.0e+00)	0.736 (\pm 3.8e-02)	0.000 (\pm 0.0e+00)	0.778 (\pm2.1e-02)
	Time	3.473 (\pm 3.3e-03)	2.092 (\pm 1.8e-03)	0.648 (\pm 0.0e+00)	3.1e+1 (\pm 1.6e+1)	8.815 (\pm 1.4e-02)	0.167 (\pm8.3e-04)
YaleFace	ACC	0.695 (\pm 9.1e-02)	0.736 (\pm 5.8e-02)	0.976 (\pm 0.0e+00)	0.976 (\pm 1.1e-16)	0.976 (\pm 1.1e-16)	0.994 (\pm3.3e-16)
	NMI	0.811 (\pm 6.4e-02)	0.834 (\pm 2.4e-02)	0.968 (\pm 0.0e+00)	0.968 (\pm 0.0e+00)	0.968 (\pm 0.0e+00)	0.992 (\pm2.2e-16)
	ARI	0.562 (\pm 1.4e-01)	0.663 (\pm 4.9e-02)	0.942 (\pm 0.0e+00)	0.942 (\pm 0.0e+00)	0.942 (\pm 0.0e+00)	0.986 (\pm2.2e-16)
	Time	0.033 (\pm 1.8e-04)	0.031 (\pm 1.5e-04)	0.017 (\pm 0.0e+00)	0.173 (\pm 1.7e-01)	0.213 (\pm 2.6e-02)	0.001 (\pm1.0e-06)

TABLE 5: Performance of algorithms taking features as input (mean (\pm std))

		<i>k</i> -means	CDKM	LKM	RKM	SPKM	<i>k</i> -sums-x
FaceV5	ACC	0.730 ($\pm 1.2\text{e-}02$)	0.931 ($\pm 5.4\text{e-}03$)	0.664 ($\pm 2.1\text{e-}01$)	0.814 ($\pm 4.2\text{e-}02$)	0.757 ($\pm 1.2\text{e-}02$)	0.963 ($\pm 3.0\text{e-}03$)
	NMI	0.930 ($\pm 2.6\text{e-}03$)	0.978 ($\pm 1.2\text{e-}03$)	0.897 ($\pm 7.5\text{e-}02$)	0.947 ($\pm 1.1\text{e-}02$)	0.941 ($\pm 3.8\text{e-}03$)	0.986 ($\pm 8.3\text{e-}04$)
	ARI	0.608 ($\pm 1.5\text{e-}02$)	0.814 ($\pm 1.3\text{e-}02$)	0.521 ($\pm 2.7\text{e-}01$)	0.709 ($\pm 8.1\text{e-}02$)	0.698 ($\pm 1.7\text{e-}02$)	0.915 ($\pm 6.3\text{e-}03$)
	Time	0.084 ($\pm 2.5\text{e-}02$)	17.311 ($\pm 3.3\text{e+}00$)	4.729 ($\pm 2.1\text{e+}00$)	3.121 ($\pm 8.8\text{e-}01$)	0.405 ($\pm 6.7\text{e-}02$)	0.254 ($\pm 3.8\text{e-}02$)
CFPW	ACC	0.547 ($\pm 1.0\text{e-}02$)	0.717 ($\pm 2.7\text{e-}03$)	0.480 ($\pm 2.4\text{e-}01$)	0.659 ($\pm 1.6\text{e-}02$)	0.623 ($\pm 1.0\text{e-}02$)	0.673 ($\pm 4.8\text{e-}03$)
	NMI	0.771 ($\pm 6.0\text{e-}03$)	0.830 ($\pm 2.2\text{e-}03$)	0.763 ($\pm 1.1\text{e-}01$)	0.844 ($\pm 1.8\text{e-}02$)	0.865 ($\pm 2.8\text{e-}03$)	0.856 ($\pm 1.5\text{e-}03$)
	ARI	0.208 ($\pm 1.6\text{e-}02$)	0.318 ($\pm 1.2\text{e-}02$)	0.289 ($\pm 1.9\text{e-}01$)	0.450 ($\pm 9.4\text{e-}02$)	0.528 ($\pm 8.0\text{e-}03$)	0.518 ($\pm 4.8\text{e-}03$)
	Time	0.561 ($\pm 2.3\text{e-}01$)	1.8e+2 ($\pm 5.6\text{e+}01$)	14.208 ($\pm 6.9\text{e+}00$)	1.8e+1 ($\pm 3.7\text{e+}00$)	1.355 ($\pm 2.0\text{e-}01$)	0.944 ($\pm 8.9\text{e-}02$)
Face96	ACC	0.666 ($\pm 3.0\text{e-}02$)	0.949 ($\pm 6.3\text{e-}03$)	0.620 ($\pm 2.1\text{e-}01$)	0.898 ($\pm 7.7\text{e-}02$)	0.789 ($\pm 1.5\text{e-}02$)	0.989 ($\pm 4.0\text{e-}03$)
	NMI	0.884 ($\pm 1.7\text{e-}02$)	0.977 ($\pm 1.8\text{e-}03$)	0.784 ($\pm 9.8\text{e-}02$)	0.968 ($\pm 2.1\text{e-}02$)	0.948 ($\pm 4.9\text{e-}03$)	0.993 ($\pm 1.1\text{e-}03$)
	ARI	0.477 ($\pm 7.8\text{e-}02$)	0.915 ($\pm 9.6\text{e-}03$)	0.581 ($\pm 2.5\text{e-}01$)	0.890 ($\pm 7.6\text{e-}02$)	0.796 ($\pm 1.8\text{e-}02$)	0.980 ($\pm 3.9\text{e-}03$)
	Time	0.046 ($\pm 1.1\text{e-}02$)	6.413 ($\pm 2.8\text{e+}00$)	10.960 ($\pm 1.1\text{e+}01$)	1.465 ($\pm 4.6\text{e-}01$)	0.218 ($\pm 4.9\text{e-}02$)	0.091 ($\pm 1.8\text{e-}02$)
FEI	ACC	0.465 ($\pm 2.6\text{e-}02$)	0.490 ($\pm 1.8\text{e-}02$)	0.413 ($\pm 2.8\text{e-}02$)	0.465 ($\pm 1.6\text{e-}02$)	0.460 ($\pm 1.6\text{e-}02$)	0.515 ($\pm 1.8\text{e-}02$)
	NMI	0.692 ($\pm 1.6\text{e-}02$)	0.713 ($\pm 1.0\text{e-}02$)	0.631 ($\pm 2.6\text{e-}02$)	0.677 ($\pm 1.1\text{e-}02$)	0.684 ($\pm 9.4\text{e-}03$)	0.721 ($\pm 7.5\text{e-}03$)
	ARI	0.342 ($\pm 2.8\text{e-}02$)	0.372 ($\pm 1.9\text{e-}02$)	0.298 ($\pm 1.9\text{e-}02$)	0.329 ($\pm 1.3\text{e-}02$)	0.322 ($\pm 1.4\text{e-}02$)	0.396 ($\pm 1.6\text{e-}02$)
	Time	0.007 ($\pm 6.2\text{e-}04$)	0.567 ($\pm 1.5\text{e-}01$)	1.081 ($\pm 1.2\text{e-}01$)	0.100 ($\pm 1.5\text{e-}02$)	0.023 ($\pm 3.4\text{e-}03$)	0.048 ($\pm 5.6\text{e-}03$)
Grimace	ACC	0.727 ($\pm 9.9\text{e-}02$)	0.992 ($\pm 2.3\text{e-}02$)	0.653 ($\pm 3.1\text{e-}01$)	1.000 ($\pm 0.0\text{e+}00$)	0.790 ($\pm 5.0\text{e-}02$)	1.000 ($\pm 0.0\text{e+}00$)
	NMI	0.871 ($\pm 7.5\text{e-}02$)	0.997 ($\pm 8.0\text{e-}03$)	0.710 ($\pm 3.2\text{e-}01$)	1.000 ($\pm 0.0\text{e+}00$)	0.919 ($\pm 1.9\text{e-}02$)	1.000 ($\pm 0.0\text{e+}00$)
	ARI	0.680 ($\pm 1.7\text{e-}01$)	0.993 ($\pm 2.2\text{e-}02$)	0.607 ($\pm 3.8\text{e-}01$)	1.000 ($\pm 0.0\text{e+}00$)	0.803 ($\pm 4.5\text{e-}02$)	1.000 ($\pm 0.0\text{e+}00$)
	Time	0.006 ($\pm 2.4\text{e-}04$)	0.084 ($\pm 2.8\text{e-}02$)	0.240 ($\pm 1.0\text{e-}01$)	0.016 ($\pm 2.3\text{e-}03$)	0.008 ($\pm 1.5\text{e-}03$)	0.022 ($\pm 3.0\text{e-}03$)
IMM	ACC	0.748 ($\pm 5.7\text{e-}02$)	0.970 ($\pm 2.5\text{e-}02$)	0.648 ($\pm 2.3\text{e-}01$)	0.932 ($\pm 4.4\text{e-}02$)	0.800 ($\pm 5.0\text{e-}02$)	1.000 ($\pm 0.0\text{e+}00$)
	NMI	0.914 ($\pm 2.0\text{e-}02$)	0.991 ($\pm 7.8\text{e-}03$)	0.790 ($\pm 1.4\text{e-}01$)	0.979 ($\pm 1.4\text{e-}02$)	0.937 ($\pm 1.6\text{e-}02$)	1.000 ($\pm 0.0\text{e+}00$)
	ARI	0.703 ($\pm 7.0\text{e-}02$)	0.966 ($\pm 2.8\text{e-}02$)	0.572 ($\pm 3.3\text{e-}01$)	0.927 ($\pm 4.7\text{e-}02$)	0.792 ($\pm 4.8\text{e-}02$)	1.000 ($\pm 0.0\text{e+}00$)
	Time	0.006 ($\pm 2.9\text{e-}04$)	0.100 ($\pm 2.2\text{e-}02$)	0.169 ($\pm 3.4\text{e-}02$)	0.020 ($\pm 3.2\text{e-}03$)	0.006 ($\pm 1.3\text{e-}03$)	0.026 ($\pm 3.2\text{e-}03$)
Isolet	ACC	0.529 ($\pm 3.8\text{e-}02$)	0.536 ($\pm 3.0\text{e-}02$)	0.450 ($\pm 1.2\text{e-}01$)	0.592 ($\pm 3.5\text{e-}02$)	0.546 ($\pm 4.8\text{e-}02$)	0.610 ($\pm 1.3\text{e-}02$)
	NMI	0.709 ($\pm 1.5\text{e-}02$)	0.711 ($\pm 1.0\text{e-}02$)	0.594 ($\pm 1.3\text{e-}01$)	0.720 ($\pm 1.1\text{e-}02$)	0.711 ($\pm 1.8\text{e-}02$)	0.731 ($\pm 7.4\text{e-}03$)
	ARI	0.468 ($\pm 3.7\text{e-}02$)	0.471 ($\pm 2.5\text{e-}02$)	0.355 ($\pm 1.4\text{e-}01$)	0.506 ($\pm 3.0\text{e-}02$)	0.475 ($\pm 4.0\text{e-}02$)	0.524 ($\pm 1.3\text{e-}02$)
	Time	0.196 ($\pm 5.1\text{e-}02$)	31.072 ($\pm 8.1\text{e+}00$)	1.851 ($\pm 7.4\text{e-}01$)	1.603 ($\pm 7.4\text{e-}01$)	0.746 ($\pm 2.1\text{e-}01$)	0.349 ($\pm 6.5\text{e-}02$)
Mpeg7	ACC	0.460 ($\pm 1.8\text{e-}02$)	0.517 ($\pm 1.6\text{e-}02$)	0.276 ($\pm 1.9\text{e-}01$)	0.516 ($\pm 2.4\text{e-}02$)	0.480 ($\pm 1.5\text{e-}02$)	0.551 ($\pm 9.9\text{e-}03$)
	NMI	0.661 ($\pm 1.6\text{e-}02$)	0.721 ($\pm 7.9\text{e-}03$)	0.500 ($\pm 1.7\text{e-}01$)	0.709 ($\pm 1.5\text{e-}02$)	0.689 ($\pm 7.6\text{e-}03$)	0.738 ($\pm 5.5\text{e-}03$)
	ARI	0.278 ($\pm 3.1\text{e-}02$)	0.364 ($\pm 2.1\text{e-}02$)	0.156 ($\pm 1.7\text{e-}01$)	0.383 ($\pm 2.8\text{e-}02$)	0.342 ($\pm 1.7\text{e-}02$)	0.434 ($\pm 9.4\text{e-}03$)
	Time	0.033 ($\pm 5.1\text{e-}03$)	7.875 ($\pm 2.0\text{e+}00$)	1.338 ($\pm 9.7\text{e-}01$)	0.267 ($\pm 1.4\text{e-}02$)	0.214 ($\pm 3.8\text{e-}02$)	0.120 ($\pm 1.3\text{e-}02$)
MSRA25	ACC	0.522 ($\pm 5.9\text{e-}02$)	0.499 ($\pm 3.2\text{e-}02$)	0.404 ($\pm 1.1\text{e-}01$)	0.517 ($\pm 4.1\text{e-}02$)	0.474 ($\pm 3.1\text{e-}02$)	0.530 ($\pm 3.4\text{e-}02$)
	NMI	0.584 ($\pm 4.3\text{e-}02$)	0.577 ($\pm 2.4\text{e-}02$)	0.430 ($\pm 1.6\text{e-}01$)	0.574 ($\pm 3.1\text{e-}02$)	0.553 ($\pm 3.9\text{e-}02$)	0.600 ($\pm 2.7\text{e-}02$)
	ARI	0.345 ($\pm 5.5\text{e-}02$)	0.344 ($\pm 2.9\text{e-}02$)	0.250 ($\pm 1.1\text{e-}01$)	0.360 ($\pm 3.2\text{e-}02$)	0.323 ($\pm 5.5\text{e-}02$)	0.392 ($\pm 3.2\text{e-}02$)
	Time	0.019 ($\pm 3.5\text{e-}03$)	0.694 ($\pm 1.3\text{e-}01$)	0.135 ($\pm 9.3\text{e-}02$)	0.179 ($\pm 4.8\text{e-}02$)	0.037 ($\pm 1.1\text{e-}02$)	0.034 ($\pm 6.3\text{e-}03$)
MUCT	ACC	0.741 ($\pm 2.5\text{e-}02$)	0.970 ($\pm 5.3\text{e-}03$)	0.746 ($\pm 9.3\text{e-}02$)	0.875 ($\pm 5.3\text{e-}02$)	0.791 ($\pm 2.2\text{e-}02$)	0.975 ($\pm 4.1\text{e-}03$)
	NMI	0.935 ($\pm 1.4\text{e-}02$)	0.990 ($\pm 1.1\text{e-}03$)	0.868 ($\pm 6.5\text{e-}02$)	0.963 ($\pm 1.1\text{e-}02$)	0.954 ($\pm 5.1\text{e-}03$)	0.991 ($\pm 1.4\text{e-}03$)
	ARI	0.672 ($\pm 1.2\text{e-}01$)	0.956 ($\pm 5.6\text{e-}03$)	0.717 ($\pm 9.9\text{e-}02$)	0.860 ($\pm 4.5\text{e-}02$)	0.794 ($\pm 2.3\text{e-}02$)	0.969 ($\pm 4.6\text{e-}03$)
	Time	0.074 ($\pm 1.4\text{e-}02$)	11.546 ($\pm 2.3\text{e+}00$)	6.550 ($\pm 4.2\text{e+}00$)	2.204 ($\pm 1.9\text{e-}01$)	0.323 ($\pm 3.7\text{e-}02$)	0.177 ($\pm 1.6\text{e-}02$)
ORL	ACC	0.708 ($\pm 6.3\text{e-}02$)	0.969 ($\pm 2.4\text{e-}02$)	0.564 ($\pm 2.6\text{e-}01$)	0.948 ($\pm 4.5\text{e-}02$)	0.792 ($\pm 3.3\text{e-}02$)	0.998 ($\pm 7.5\text{e-}03$)
	NMI	0.887 ($\pm 4.3\text{e-}02$)	0.992 ($\pm 6.1\text{e-}03$)	0.696 ($\pm 1.8\text{e-}01$)	0.984 ($\pm 1.4\text{e-}02$)	0.935 ($\pm 1.1\text{e-}02$)	0.999 ($\pm 2.8\text{e-}03$)
	ARI	0.624 ($\pm 1.4\text{e-}01$)	0.971 ($\pm 2.2\text{e-}02$)	0.517 ($\pm 3.3\text{e-}01$)	0.947 ($\pm 4.6\text{e-}02$)	0.795 ($\pm 3.3\text{e-}02$)	0.997 ($\pm 8.5\text{e-}03$)
	Time	0.006 ($\pm 3.1\text{e-}04$)	0.151 ($\pm 3.2\text{e-}02$)	0.267 ($\pm 1.4\text{e-}01$)	0.029 ($\pm 2.4\text{e-}03$)	0.011 ($\pm 1.9\text{e-}03$)	0.050 ($\pm 5.6\text{e-}03$)
Palm	ACC	0.685 ($\pm 2.4\text{e-}02$)	0.750 ($\pm 1.9\text{e-}02$)	0.569 ($\pm 1.9\text{e-}01$)	0.703 ($\pm 4.5\text{e-}02$)	0.691 ($\pm 1.9\text{e-}02$)	0.799 ($\pm 1.8\text{e-}02$)
	NMI	0.890 ($\pm 8.4\text{e-}03$)	0.916 ($\pm 5.2\text{e-}03$)	0.770 ($\pm 1.6\text{e-}01$)	0.889 ($\pm 1.9\text{e-}02$)	0.891 ($\pm 7.8\text{e-}03$)	0.925 ($\pm 5.7\text{e-}03$)
	ARI	0.641 ($\pm 2.6\text{e-}02$)	0.711 ($\pm 1.9\text{e-}02$)	0.506 ($\pm 2.1\text{e-}01$)	0.658 ($\pm 5.5\text{e-}02$)	0.643 ($\pm 2.2\text{e-}02$)	0.766 ($\pm 1.8\text{e-}02$)
	Time	0.023 ($\pm 3.3\text{e-}03$)	3.332 ($\pm 7.4\text{e-}01$)	2.856 ($\pm 1.7\text{e+}00$)	0.532 ($\pm 7.8\text{e-}02$)	0.177 ($\pm 3.8\text{e-}02$)	0.070 ($\pm 1.6\text{e-}02$)
YaleFace	ACC	0.767 ($\pm 9.4\text{e-}02$)	0.968 ($\pm 2.4\text{e-}02$)	0.649 ($\pm 2.7\text{e-}01$)	0.944 ($\pm 5.2\text{e-}02$)	0.833 ($\pm 4.6\text{e-}02$)	0.988 ($\pm 0.0\text{e+}00$)
	NMI	0.864 ($\pm 5.6\text{e-}02$)	0.964 ($\pm 1.2\text{e-}02$)	0.700 ($\pm 2.6\text{e-}01$)	0.963 ($\pm 2.8\text{e-}02$)	0.925 ($\pm 2.1\text{e-}02$)	0.984 ($\pm 2.2\text{e-}16$)
	ARI	0.700 ($\pm 1.3\text{e-}01$)	0.932 ($\pm 3.1\text{e-}02$)	0.567 ($\pm 3.5\text{e-}01$)	0.922 ($\pm 6.1\text{e-}02$)	0.824 ($\pm 4.5\text{e-}02$)	0.972 ($\pm 1.1\text{e-}16$)
	Time	0.005 ($\pm 1.8\text{e-}04$)	0.032 ($\pm 5.6\text{e-}03$)	0.130 ($\pm 9.3\text{e-}02$)	0.006 ($\pm 2.7\text{e-}03$)	0.005 ($\pm 2.3\text{e-}03$)	0.017 ($\pm 2.6\text{e-}03$)

TABLE 6: Performance on large-scale datasets

Data	Met.	FINCH	k -sums	k -means	k -sums-x
CACD	ACC	0.745	0.832	0.68	0.819
	NMI	0.79	0.903	0.858	0.897
	ARI	0.021	0.755	0.461	0.736
	Time	2217.65	20.119	114.032	114.449
WebFace	ACC	-	0.565	0.529	0.489
	NMI	-	0.847	0.836	0.83
	ARI	-	0.391	0.329	0.334
	Time	-	75.114	1841.865	4232.793
CelebA	ACC	-	0.467	0.365	0.44
	NMI	-	0.802	0.52	0.555
	ARI	-	0.298	0	0
	Time	-	7.476	301.742	828.339
PEAL	ACC	0.884	0.877	0.696	0.849
	NMI	0.906	0.942	0.894	0.929
	ARI	0.125	0.806	0.587	0.747
	Time	75.844	1.401	3.112	9.706

- Out of memory

(AGCI) [47], Fast clustering with Co-clustering via Discrete non-negative Matrix Factorization (FCDMF) [48], First Integer Neighbor Clustering Hierarchy (FINCH) [49], Rank-Constrained clustering with Adaptive Embedding (RCAE), and Robust continuous clustering (RCC) [50].

We verify the superiority of k -sums-x by comparing it with five k -means based algorithms, which are: k -means (KM) [51], Coordinate Descent Method for k -means(CDKM) [52], lasso k -means (LKM) [53], regularized k -means (RKM) [54], and spherical k -means (SPKM) [55].

5.2.4 Settings

The parameters involved in these algorithms are described here. In all methods that required an anchor graph as input, i.e., AGCI and FCDMF, the anchors are always generated by k -means (the cluster center is regarded as an anchor point, in the final clustering result). The number of anchors, m , is set as follows:

$$m = \min(n/2, 1024) \quad (63)$$

For these methods that require a k -NN graph or anchor graph, heat-kernel is used to construct the graph, and the number of neighbours, k , is determined by $k = \lfloor n/c * 1.2 \rfloor$. In RCAE, the regularization parameter γ is tuned from $\{0.01, 0.1, 1, 10, 100\}$. the clustering threshold in RCC is tuned from $\{0.1, 0.3, 0.5, 0.7, 0.9\}$. For a fair comparison, the number of neighbours, k , is also determined by $k = \lfloor n/c * 1.2 \rfloor$, in RCC and k -sums.

The methods shown in Table 5 are all initialized in a random way. The regularization parameter λ is tuned from $\{10e-6, 10e-4, 10e-2, 10e0, 10e2, 10e4, 10e6\}$. RKM used here is the version with the strict balance constraint. Algorithms whose performance is affected by initialization were run 50 times, and average performance was reported.

5.2.5 Experimental results

From the results shown in Table 4, 5 and Figure 6, it is not difficult to find that:

1. First of all, it can be clearly seen that our algorithms have achieved the highest performance on most datasets showing the effectiveness and robustness of them. Take k -sums as an example, it achieves 32.4%, 30.5%, 7.9%, 5.7%, 4.5%, and 4.4% gain on CFPW, FaceV5, Mpeg7, Isolet, MSRA25, and YaleFace over the second-best method respectively, in terms of ARI score. Similar conclusions can be obtained with the metrics of recall and precision.

2. In Table 4, the time consumed by constructing the k -NN graph is also included. Here, we construct it by brute force since the size of these datasets is small. As can be seen, on medium-scale data sets, k -sums shows a speed advantage compared to other graph-based algorithms, but it is not as efficient as the k -means. This is why we do experiments on large-scale datasets.

3. Algorithms that require a hyperparameter have poor performance on some data, which may be caused by inappropriate parameters. k -sums algorithm requires only one parameter, k , which makes the algorithm easier to use. Moreover, it has satisfactory performance, even if the parameter is not tuned, as shown in Table 4.

4. For graph-based algorithms, their performance usually is effected by the value of k . However, from Table 4, we can see that k -sums has achieved satisfactory performance on Isolet ($n/c \approx 300$), MSRA25($n/c \approx 150$), CFPW($n/c \approx 14$) and YaleFace ($n/c \approx 11$), which shows that the formula used to determine the value of k is reasonable. The effect of parameter k on performance has also been studied, and the results are shown in Figure 6.

5. Again, the time of construction the k -NN graph is included in Table 6. As we have seen, in terms of efficiency, k -sums outperforms k -means on large-scale datasets. On large-scale datasets, most graph-based algorithms (AGCI, FCDMF, RCAE and RCC) cannot produce results because of their high computational complexity or high memory overhead. EFANNA [56] is the algorithm that is used to generate the k -nearest neighbors graph.

From the results shown in Table 4, 5, and 6, in most benchmarks, k -sums achieves higher performance over k -sums-x. It seems that k -sums-x is not very necessary. Therefore, we think it is necessary to list some scenarios where k -sums-x is applicable.

- In many practical scenarios, the cluster centers are also very important information. In other words, we sometimes need meaningful cluster centers instead of a partition [57], [58]. In k -sums-x, the means of each cluster can be used as cluster centers. However, in k -sums, the means of each cluster are meaningless.
- If some new training samples are obtained, it is difficult for k -sums to update the model, which is also a disadvantage of most graph-based algorithms. Usually, it will be retrained on whole datasets. But it is easy for k -sums-x to update the model.

In summary, k -sums should be used if a partition is all you need, k -sums-x should be used if clustering centers are needed or you need to update the model often.

6 CONCLUSION

Based on the unified framework of k -means and graph cut models, we proposed a novel clustering algorithm,

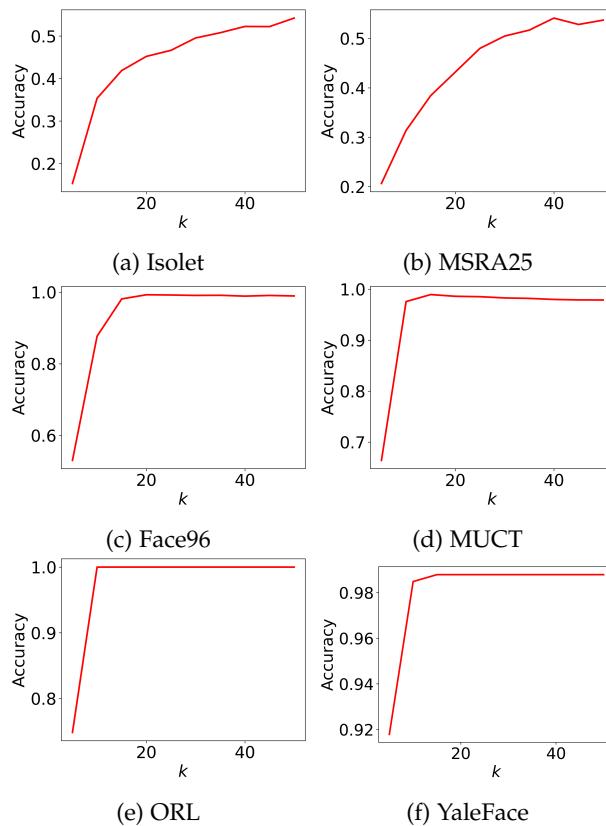


Fig. 6: Performance with different k

called k -sums. The exciting thing is that our model can be efficiently optimized if the standard coordinate descent algorithm is adopted, and its performance is comparable to most graph-based clustering algorithms. Specifically, both computational and memory overhead are $O(nk)$ that is irrelevant to the product of the number of samples and subjects.

To deal with the situation in which the k -NN graph is not available, we proposed a variant of k -sums, called k -sums-x. It is a k -means-like algorithm, and only takes features and the number of clusters as input. By the optimization method based on coordinate descent proposed in this paper, k -sums-x can also be solved very efficiently. The computational complexity of it is the same as that of k -means.

Extensive experiments have been conducted on 10 synthetic datasets, 13 middle-scale benchmark datasets, and 4 large-scale datasets. to verify the effectiveness and efficiency of the proposed models, i.e., k -sums and k -sums-x, and the properties of them are also analysed.

ACKNOWLEDGMENTS

This work was supported in part by the National Key Research and Development Program of China under Grant 2018AAA0101902, in part by the Natural Science Basic Research Program of Shaanxi (Program No. 2021JM-071), in part by the National Natural Science Foundation of China under Grant 62176212, Grant 61936014 and Grant 61772427, and in part by the Fundamental Research Funds for the Central Universities under Grant G2019KY0501.

REFERENCES

- [1] M. Tapaswi, M. T. Law, and S. Fidler, "Video face clustering with unknown number of clusters," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5027–5036.
- [2] Y. Shi, C. Otto, and A. K. Jain, "Face clustering: representation and pairwise constraints," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 7, pp. 1626–1640, 2018.
- [3] R. Janani and S. Vijayarani, "Text document clustering using spectral clustering algorithm with particle swarm optimization," *Expert Systems with Applications*, vol. 134, pp. 192–200, 2019.
- [4] L. Krawczyk, "Comparing applicability of prevalent clustering algorithms for document clustering," B.S. thesis, Humboldt-Universität zu Berlin, 2019.
- [5] H.-W. Lee, N. Malik, F. Shi, and P. Mucha, "Social clustering in epidemic spread on coevolving networks," *Physical Review E*, vol. 99, no. 6, p. 062301, 2019.
- [6] M. Xu, Y. Li, R. Li, F. Zou, and X. Gu, "Eadp: An extended adaptive density peaks clustering for overlapping community detection in social networks," *Neurocomputing*, vol. 337, pp. 287–302, 2019.
- [7] R. Jain and R. Sharma, "Image segmentation through fuzzy clustering: A survey," in *Harmony Search and Nature Inspired Optimization Algorithms*, 2019, pp. 497–508.
- [8] T. Lei, X. Jia, Y. Zhang, S. Liu, H. Meng, and A. K. Nandi, "Superpixel-based fast fuzzy c-means clustering for color image segmentation," *IEEE Transactions on Fuzzy Systems*, vol. 27, no. 9, pp. 1753–1766, 2018.
- [9] A. Ng, M. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," in *Advances in neural information processing systems*, 2002, pp. 849–856.
- [10] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 22, no. 8, pp. 888–905, 2000.
- [11] C. Fowlkes, S. Belongie, F. Chung, and J. Malik, "Spectral grouping using the nyström method," *IEEE transactions on pattern analysis and machine intelligence*, vol. 26, no. 2, pp. 214–225, 2004.
- [12] A. Choromanska, T. Jebara, H. Kim, M. Mohan, and C. Monteleoni, "Fast spectral clustering via the nyström method," in *Algorithmic Learning Theory*, 2013, pp. 367–381.
- [13] G. Zhong and C.-M. Pun, "Revisiting nyström extension for hypergraph clustering," *Neurocomputing*, vol. 403, pp. 247–256, 2020.
- [14] M. Li, J. Kwok, and B. Lü, "Making large-scale nyström approximation possible," in *ICML 2010-Proceedings, 27th International Conference on Machine Learning*, 2010, p. 631.
- [15] T. Sakai and A. Imaia, "Fast spectral clustering with random projection and sampling," in *Machine Learning and Data Mining in Pattern Recognition*, 2009, pp. 372–384.
- [16] X. Yang, W. Yu, R. Wang, G. Zhang, and F. Nie, "Fast spectral clustering learning with hierarchical bipartite graph for large-scale data," *Pattern Recognition Letters*, vol. 130, pp. 345 – 352, 2020.
- [17] Y. Li, J. Huang, and W. Liu, "Scalable sequential spectral clustering," in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, 2016, p. 1809–1815.
- [18] A. Hassanzadeh, A. Kaarna, and T. Kauranne, "Sequential spectral clustering of hyperspectral remote sensing image over bipartite graph," *Applied Soft Computing*, vol. 73, pp. 727 – 734, 2018.
- [19] B. Schölkopf, A. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural computation*, vol. 10, pp. 1299–1319, 1998.
- [20] D. Kim, K. Lee, D. Lee, and K. Lee, "Evaluation of the performance of clustering algorithms in kernel-induced feature space," *Pattern Recognition*, vol. 38, pp. 607–611, 2005.
- [21] R. Chitta, R. Jin, T. Havens, and A. Jain, "Approximate kernel k-means: Solution to large scale kernel clustering," in *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2011, pp. 895–903.
- [22] D. Arthur and S. Vassilvitskii, "k-means++: The advantages of careful seeding," Stanford, Tech. Rep., 2006.
- [23] O. Bachem, M. Lucic, S. Hassani, and A. Krause, "Approximate k-means++ in sublinear time," in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, 2016, pp. 1459–1467.
- [24] O. Bachem, M. Lucic, H. Hassani, and A. Krause, "Fast and provably good seedings for k-means," in *Advances in Neural Information Processing Systems 29*, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, Eds., 2016, pp. 55–63.
- [25] I. Dhillon, Y. Guan, and B. Kulis, "Kernel k-means: spectral clustering and normalized cuts," in *Proceedings of the tenth ACM SIGKDD*

- international conference on Knowledge discovery and data mining*, 2004, pp. 551–556.
- [26] I. S. Dhillon, Y. Guan, and B. Kulis, *A unified view of kernel k-means, spectral clustering and graph cuts*. Citeseer, 2004.
- [27] C. Ding, X. He, and H. Simon, “On the equivalence of nonnegative matrix factorization and spectral clustering,” in *Proceedings of the 2005 SIAM international conference on data mining*, 2005, pp. 606–610.
- [28] R. Zass and A. Shashua, “Doubly stochastic normalization for spectral clustering,” in *Advances in Neural Information Processing Systems*, B. Schölkopf, J. Platt, and T. Hoffman, Eds., vol. 19, 2007.
- [29] X. Wang, F. Nie, and H. Huang, “Structured doubly stochastic matrix for graph based clustering: Structured doubly stochastic matrix,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, p. 1245–1254.
- [30] P. E. Hart, D. G. Stork, and R. O. Duda, *Pattern classification*. Wiley Hoboken, 2000.
- [31] C. C. V. P. R. C. D. J. Sengupta, J.C. Cheng, “Frontal to profile face verification in the wild,” in *Proc. WACV*, 2016, pp. 1–9.
- [32] C. E. Thomaz and G. A. Giraldi, “A new ranking method for principal components analysis and its application to face image analysis,” *Image Vision Comput.*, vol. 28, no. 6, pp. 902–913, 2010.
- [33] U. A. Kamerikar and M. Chavan, “Experimental assessment of lda and kllda for face recognition,” *International Journal*, vol. 2, no. 2, pp. 137–146, 2014.
- [34] M. Nordstrom, M. Larsen, J. Sierakowski, and M. Stegmann, *The IMM Face Database - An Annotated Dataset of 240 Face Images*. Technical University of Denmark, DTU Informatics, 2004.
- [35] L. Zhang, Q. Zhang, B. Du, J. You, and D. Tao, “Adaptive manifold regularized matrix factorization for data clustering,” in *Proc. IJCAI*, 2017, pp. 3399–3405.
- [36] S. Milborrow, J. Morkel, and F. Nicolls, “The muct landmarked face database,” *Pattern Recognition Association of South Africa*, pp. 32–34, 2010.
- [37] D. Cai, X. He, Y. Hu, J. Han, and T. Huang, “Learning a spatially smooth subspace for face recognition,” in *Proc. CVPR*, 2007, pp. 1–7.
- [38] A. S. Georghiades, P. N. Belhumeur, and D. J. Kriegman, “From few to many: illumination cone models for face recognition under variable lighting and pose,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 6, pp. 643–660, 2001.
- [39] D. Virmani, P. Girdhar, P. Jain, and P. Bamdev, “Fdnet: Face detection and recognition pipeline,” *Engineering, Technology and Applied Science Research*, vol. 9, no. 2, pp. 3933–3938, 2019.
- [40] M. Fanty and R. Cole, “Spoken letter recognition,” *Proc. NIPS*, vol. 3, pp. 220–226, 1990.
- [41] L. J. Latecki and R. Lakamper, “Shape similarity measure based on correspondence of visual parts,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 10, pp. 1185–1190, 2000.
- [42] W. Gao, B. Cao, S. Shan, X. Chen, D. Zhou, X. Zhang, and D. Zhao, “The cas-peal large-scale chinese face database and baseline evaluations,” *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, vol. 38, no. 1, pp. 149–161, 2008.
- [43] B.-C. Chen, C.-S. Chen, and W. H. Hsu, “Face recognition and retrieval using cross-age reference coding with cross-age celebrity dataset,” *IEEE Transactions on Multimedia*, vol. 17, no. 6, pp. 804–815, 2015.
- [44] D. Yi, Z. Lei, S. Liao, and S. Z. Li, “Learning face representation from scratch,” *arXiv preprint arXiv:1411.7923*, 2014.
- [45] Z. Liu, P. Luo, X. Wang, and X. Tang, “Deep learning face attributes in the wild,” in *Proceedings of International Conference on Computer Vision (ICCV)*, 2015, pp. 3730–3738.
- [46] X. Wu, R. He, Z. Sun, and T. Tan, “A light cnn for deep face representation with noisy labels,” *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 11, pp. 2884–2896, 2018.
- [47] Y. Zhao, Y. Yuan, and Q. Wang, “Fast spectral clustering for unsupervised hyperspectral image classification,” *Remote Sensing*, vol. 11, no. 4, p. 399, 2019.
- [48] F. Nie, S. Pei, R. Wang, and X. Li, “Fast clustering with co-clustering via discrete non-negative matrix factorization for image identification,” in *Proc. ICASSP*, 2020, pp. 2073–2077.
- [49] V. S. M. Saquib Sarfraz and R. Stiefelhagen, “Efficient parameter-free clustering using first neighbor relations,” in *Proc. CVPR*, 2019, pp. 8934–8943.
- [50] S. A. Shah and V. Koltun, “Robust continuous clustering,” *Proceedings of the National Academy of Sciences*, vol. 114, no. 37, pp. 9814–9819, 2017.
- [51] J. A. Hartigan and M. A. Wong, “Algorithm as 136: A k -means clustering algorithm,” *Journal of the Royal Statistical Society*, vol. 28, no. 1, pp. 100–108, 1979.
- [52] F. Nie, J. Xue, D. Wu, R. Wang, H. Li, and X. Li, “Coordinate descent method for k -means,” *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021.
- [53] Z. Li, F. Nie, X. Chang, Z. Ma, and Y. Yang, “Balanced clustering via exclusive lasso: A pragmatic approach,” in *Proc. AAAI*, 2018, pp. 3596–3603.
- [54] W. Lin, Z. He, and M. Xiao, “Balanced clustering: A uniform model and fast algorithm,” in *Proc. IJCAI*, 7 2019, pp. 2987–2993.
- [55] C. Buchta, M. Kober, I. Feinerer, and K. Hornik, “Spherical k -means clustering,” *Journal of statistical software*, vol. 50, no. 10, pp. 1–22, 2012.
- [56] C. Fu and D. Cai, “Efanna: An extremely fast approximate nearest neighbor search algorithm based on knn graph,” *arXiv preprint arXiv:1609.07228*, 2016.
- [57] Y.-N. Zhu and Y.-F. Li, “Semi-supervised streaming learning with emerging new labels,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 7015–7022.
- [58] H. Jegou, M. Douze, and C. Schmid, “Product quantization for nearest neighbor search,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 1, pp. 117–128, 2010.

Shenfei Pei is now working toward his PhD degree in Northwestern Polytechnical University, Xi'an, China. His research interests include machine learning and its application fields, such as data mining, computer vision and Clustering.

Huimin Chen Huimin Chen received the bachelor's degree in software engineering from Northwestern Polytechnical University, Xi'an, China, in 2019. Her current research interests in machine learning and its applications, including clustering and dimensionality reduction.

Feiping Nie received the Ph.D. degree in Computer Science from Tsinghua University, China in 2009, and currently is full professor in Northwestern Polytechnical University, China. His research interests are machine learning and its applications, such as pattern recognition, data mining, computer vision, image processing and information retrieval. He has published more than 100 papers in the following journals and conferences: TPAMI, IJCV, TIP, TNNLS, TKDE, ICML, NIPS, KDD, IJCAI, AAAI, ICCV, CVPR, ACM MM. His papers have been cited more than 20000 times and the H-index is 78. He is now serving as Associate Editor or PC member for several prestigious journals and conferences in the related fields.

Rong Wang received the B.S. degree in information engineering, the M.S. degree in signal and information processing, and the Ph.D. degree in computer science from Xi'an Research Institute of Hi-Tech, Xi'an, China, in 2004, 2007 and 2013, respectively. During 2007 and 2013, he also studied in the Department of Automation, Tsinghua University, Beijing, China for his Ph.D. degree. He is currently an associate professor at the School of Cybersecurity and School of Artificial Intelligence, Optics and Electronics (iOPEN), Northwestern Polytechnical University, Xi'an, China. His research interests focus on machine learning and its applications.

Xuelong Li (Fellow, IEEE) is a Full Professor with the school of Artificial Intelligence, Optics and Electronics (iOPEN), Northwestern Polytechnical University, Xi'an, China.