

Flexible Subspace Clustering: A Joint Feature Selection and K-Means Clustering Framework[☆]

Zhong-Zhen Long^a, Guoxia Xu^b, Jiao Du^c, Hu Zhu^d, Taiyu Yan^e, Yu-Feng Yu^{f,*}

^a*Shenzhen Securities Communication Co., Ltd., Shenzhen 518041, China.*

^b*Department of Computer Science, Norwegian University of Science and Technology, 2815 Gjøvik, Norway.*

^c*School of Computer Science and Educational Software, Guangzhou University, Guangzhou 510006, China.*

^d*College of Telecommunication and Information Engineering, Nanjing University of Posts and Telecommunications, Nanjing 210003, China.*

^e*Department of Imaging and Interventional Radiology, Chinese University of Hong Kong, Hong Kong, China.*

^f*Department of Statistics, Guangzhou University, Guangzhou 510006, China.*

Abstract

Regarding as an important computing paradigm, cloud computing is to address big and distributed databases and rather simple computation. In this paradigm, data mining is one of the most important and fundamental problems. A large amount of data is generated by sensors and other intelligent devices. Data mining for these big data is crucial in various applications. K-means clustering is a typical technique to group the similar data into the same clustering, and has been commonly used in data mining. However, it is still a challenge to the data con-

[☆]This work was supported in part by National Natural Science Foundation of China under Grant 62006056 and Grant 61802148, in part by National Statistical Science Research Project of China under Grant 2020LY090, and in part by the Natural Science Foundation of Guangdong Province under Grant 2019A1515011266.

*Corresponding author

Email addresses: longzhzh@mail2.sysu.edu.cn (Zhong-Zhen Long), gxxu.re@gmail.com (Guoxia Xu), djiao@gzhu.edu.cn (Jiao Du), peter.hu.zhu@gmail.com (Hu Zhu), taiyuyan@163.com (Taiyu Yan), yuyufeng220@163.com (Yu-Feng Yu)

taining a large amount of noise, outliers and redundant features. In this paper, we propose a robust K-means clustering algorithm, namely, flexible subspace clustering. The proposed method incorporates feature selection and K-means clustering into a unified framework, which can select the refined features and improve the clustering performance. Moreover, for the purpose of enhancing the robustness, the $l_{2,p}$ -norm is embedded into the objective function. We can flexibly choose appropriate p according to the different data and thus obtain more robust performance. Experimental results verify the presented method has more robust and better performance on benchmark databases compared to the existing approaches.

Keywords: Big data, K-means clustering, cloud computing, subspace learning.

1. Introduction

The cloud computing with big data is regarded as an important paradigm, which handles big and distributed databases and rather simple computation. Many interesting studies concentrate on cloud security [3], smart service [2] and mobile cloud computing [1]. However, most of these papers focus on hardware, data storage and management in clouds. Recently, the internet of things (IoT) is gaining increasing attention and many related studies are proposed in various applications such as quality prediction [4][5], dynamic resource discovery [6], bearing test [7] and feature recognition [8][9]. IoT and big data have an increasing impact on the future development of cloud computing. In IoT, enormous amount of data generated by sensors and other intelligent devices contain valuable information, but also encompass a large amount of noise, outliers and redundant features. Thus, data mining for these big data is crucial to be suitable for various applications. As one of the most important and fundamental technique in data mining, clustering

15 has been studied a lot and applied in many fields, such as resource scheduling in
16 cloud computing [10], abnormal behavior detection in cloud [11], clinical obser-
17 vation [12], heterogeneous data analysis [13] and so on. Clustering is a kind of
18 unsupervised learning, which groups the similar data points into the same cluster.
19 As for the similarity, the most common used criterion is the distance, and K-means
20 (KM) is a typical algorithm of this criterion.

21 The classical K-means distributes data points to k different clusters using l_2 -
22 norm distance. It's simple and easy to be solved, but easily affected by outliers
23 and noises [14]. To overcome this problem, one direction is to use a distance
24 measure that can be more robust. The use of l_p -norm is a successful extension.
25 Hathaway et al. [15] conclude that $p = 1$ shows its property of robustness and
26 choosing the value of p can provide better clustering results than fixing p as 1
27 or 2 but the model could be difficult to be solved. Salem et al. [16] adopt l_1 -
28 norm to evaluate the similarity between the observation and the centroid, which is
29 shown efficiency and suitable to noisy data and outliers. Cai et al. [17] propose
30 a multi-view K-means clustering based on $l_{2,1}$ -norm. Liang et al. [18] propose
31 a robust K-means using $l_{2,1}$ -norm in the feature space and then extend it to the
32 kernel space. The reform of the distance metric can improve the performance of
33 K-means algorithm, which has been demonstrated in the above literatures.

34 However, with the development of science and technology, the data in re-
35 al life is explosive. Big data sets generated from many fields contains a large
36 amount of attributes, and some of which are noise and redundant attributes. It
37 poses a remarkable challenge on the traditional clustering methods. For example,
38 in face recognition applications, given a face image data of 128×128 resolution
39 which is relatively small, it will generate a 16384-dimensional feature vector. This

40 kind of high-dimensional data always contains a large amount of noises, outliers
41 and redundant features. It is difficult to cluster directly, and sometimes leads to
42 high computational complexity and performance degradation [19], especially in
43 K-means and its extensions. To deal with the curse of dimensionality and reduce
44 the noise, outliers and redundant features, an intuitive approach is to conduct di-
45 mensionality reduction processing on the data before clustering. Many dimension
46 reduction methods have been studied in the past decades, such as Principal Com-
47 ponent Analysis (PCA) [20], Linear Discriminant Analysis (LDA) [21], sparse ap-
48 proximation to discriminant projection learning (SADPL) [22] and Locally Linear
49 Embedding (LLE) [23]. PCAKM is a typical method that sequentially conducts
50 PCA for dimension reduction and K-means for clustering [24]. Yin et al. [25]
51 apply LLE to preprocess the data before performing K-means to make better use
52 of the manifold information. These sequential methods can improve the com-
53 putational efficiency, but the subspace got from the dimension reduction process
54 may not be the optimal one for the clustering process, so that some researchers
55 believe that the separation of dimension reduction and clustering may result in
56 worse clustering performance [26].

57 Intuitively, if clustering is embedded into the process of dimension reduction,
58 the performance of clustering may be improved. This kind of methods try to find
59 the optimal structure of data in the low-dimensional feature space for clustering.
60 They perform K-means and the subspace learning process simultaneously. For
61 example, Ding et al. [28] construct an adaptive framework LDAKM, in which
62 LDA and K-means are jointly implemented, that is, labels are generated by K-
63 means algorithm, and the obtained labels are used by LDA to learn the subspace.
64 Since LDA may fail when the number of samples is very small, several LDA's

65 extensions have been used to replace LDA, and get better results than LDAK-
66 M [29], such as Maximum Margin Criterion (MMC) [30], Orthogonal Centroid
67 Method (OCM) [31] and Orthogonal Least Squares Discriminant Analysis (OLS-
68 DA) [32]. Hou et al. [29] consider the relation between PCA and K-means, and
69 propose a general subspace clustering framework. This kind of algorithms have
70 been proved to get better results than the sequential algorithms, but they also have
71 some drawbacks. These algorithms all need to compute an approximate solution
72 by eigenvalue decomposition, which will increase the computational burden so
73 that when facing the high-dimensional data, these algorithms may fail. And s-
74 ince the optimal subspace is found by orthogonal linear transformation, it may
75 have difficulty to understand the meaning of the obtained low-dimensional fea-
76 tures. Wang et al. [27] construct a special feature selection matrix and propose
77 a fast adaptive subspace clustering algorithm FAKM based on DEC, which can
78 effectively select the most representative subspace without requiring eigenvalue
79 decomposition. FAKM also performs adaptive learning to the K-means part.

80 Most methods mentioned above are based on the l_2 -norm distance metric,
81 which is known to be very sensitive to data outliers and noise. Therefore, it is
82 meaningful to build a model with robust distance metric. Recently, $l_{2,p}$ -norm is
83 successfully used to replace l_2 -norm as distance metric for improving the robust-
84 ness, such as DCM [33] and $l_{2,p}$ -PCA [34]. In $l_{2,p}$ -PCA, $l_{2,p}$ -norm is incorporated
85 into PCA, and it is robust to outliers and can retain the desirable properties from
86 big data. Inspired by FAKM and $l_{2,p}$ -PCA, we propose a flexible subspace cluster-
87 ing method. Our method flexibly chooses appropriate p according to the data and
88 thus obtains more robust clustering performance. Several experimental results on
89 various datasets prove the effectiveness of the proposed algorithm.

90 The main contributions of our paper are listed as follows.

- 91 • The proposed algorithm combines the feature selection and clustering into
92 a single framework jointly.
- 93 • The use of $l_{2,p}$ -norm on K-means makes our algorithm robust to noise and
94 redundant features of big data.
- 95 • The proposed approach is neither convex nor Lipschitz continuous, thus it is
96 difficult to be solved directly. We propose an iterative algorithm to optimize
97 it.

98 The rest of the paper is organized as follows. We propose our model and
99 derive an efficient algorithm to optimize the model in Section 2. In Section 3, the
100 proposed model is evaluated on the benchmark databases. Finally, we draw the
101 conclusion in Section 4.

102 **2. The Proposed Method**

103 In this section, we introduce the details about the proposed method for cluster-
104 ing. The main content will be separated into the following several parts including
105 the formulation of the proposed approach, an efficient algorithm, convergence and
106 computational complexity analysis.

107 *2.1. Formulation*

Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in R^{D \times n}$ be a high-dimensional data matrix, and $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_c] \in R^{D \times c}$ be c centroid vectors. $\mathbf{F} \in \{0, 1\}^{n \times c}$ denotes the indicator matrix, here $F_{ik} = 1$ if \mathbf{x}_i belongs to the k -th cluster, otherwise $F_{ik} = 0$.

Following [37][38], we can obtain the K-means formulation as

$$\begin{aligned} \min_{\mathbf{F}, \mathbf{Z}} \sum_{i=1}^n \sum_{k=1}^c F_{ik} \|\mathbf{x}_i - \mathbf{z}_k\|_2^2, \\ \text{s.t. } \mathbf{F} \in \{0, 1\}^{n \times c}, \mathbf{F}\mathbf{1} = \mathbf{1}. \end{aligned} \quad (1)$$

Taking a simple algebra, the objective in (1) becomes

$$\begin{aligned} \min_{\mathbf{F}, \mathbf{Z}} \|\mathbf{X} - \mathbf{Z}\mathbf{F}^T\|_F^2, \\ \text{s.t. } \mathbf{F} \in \{0, 1\}^{n \times c}, \mathbf{F}\mathbf{1} = \mathbf{1}. \end{aligned} \quad (2)$$

Considering that the high-dimensional data could contain a large amount of noises, outliers and redundant features. It leads to high computational complexity and performance degradation. The direct idea is to find a transformation matrix $\mathbf{W} \in R^{D \times d}$ which transforms the high-dimensional features to a low-dimensional feature space $\mathbf{Y} = \mathbf{W}^T \mathbf{X}$, where $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in R^{d \times n}$. Following the feature selection [39][27], we use the column vectors \mathbf{w}_i as follow

$$\mathbf{w}_i = \underbrace{[0, \dots, 0]_{i-1}}_{i-1}, 1, \underbrace{[0, \dots, 0]_{D-i}}_{D-i}. \quad (3)$$

Then the feature selection matrix \mathbf{W} can be represented as

$$\mathbf{W} = [\mathbf{w}_{I(1)}, \mathbf{w}_{I(2)}, \dots, \mathbf{w}_{I(d)}], \quad (4)$$

108 where I is a permutation of $\{1, 2, \dots, D\}$. It can be seen that the transformation
109 matrix \mathbf{W} is sparse and column-full-rank.

To achieve the goal of feature selection and K-means clustering simultaneously, we incorporate the subspace learning and K-means clustering into a unified

framework as

$$\begin{aligned}
& \max_{\mathbf{W}, \mathbf{G}, \mathbf{F}} Tr(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) - \lambda \|\mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T\|_2^p \\
& \text{s.t. } \mathbf{W} \in \{0, 1\}^{D \times d}, \text{rank}(\mathbf{W}) = d, \mathbf{W}^T \mathbf{1} = \mathbf{1}, \\
& \mathbf{F} \in \{0, 1\}^{n \times c}, \mathbf{F} \mathbf{1} = \mathbf{1},
\end{aligned} \tag{5}$$

110 where $\mathbf{S}_t = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$ is the total scatter matrix. $\mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_c] \in R^{d \times c}$ is c
111 centroid vectors in the low-dimensional space. It should be noted that FAKM [27]
112 is also a joint model of subspace learning and clustering. It uses the pattern of $\|\mathbf{M}\|_\sigma = \sum_i \frac{(1+\sigma)\|\mathbf{m}_i\|_2^2}{\|\mathbf{m}_i\|_2 + \sigma}$
113 to construct the K-means clustering, here \mathbf{M} is an arbitrary
114 matrix, \mathbf{m}_i is the i -th column and σ is a parameter. Different from FAKM, the
115 model in (5) has more robust performance since it adopts $l_{2,p}$ -norm to construct
116 the K-means clustering and can flexibly choose appropriate p according to the
117 different data.

118 2.2. Optimization

119 Since our objective function (5) involves $l_{2,p}$ -norm, it is difficult to get its
120 closed-form solution directly. In [34], an iterative algorithm is proposed to solve
121 the objective function in the form of $l_{2,p}$ -norm. Similar techniques are used in
122 [35] to solve the problem of the minimization of LDA with regular term based on
123 $l_{2,p}$ -norm ($0 < p \leq 2$). Inspired by these papers, we propose an effective iterative
124 algorithm to solve our objective function.

Let $d_i = \frac{p}{2} \|\mathbf{W}^T \mathbf{x}_i - \mathbf{G} \mathbf{f}_i\|_2^{p-2}$, then (5) can be transformed to

$$\max_{\mathbf{W}, \mathbf{G}, \mathbf{F}} Tr(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) - 2\lambda/p \sum_{i=1}^n d_i \|\mathbf{W}^T \mathbf{x}_i - \mathbf{G} \mathbf{f}_i\|_2^2. \tag{6}$$

125 Since λ is an arbitrary constant, for convenience, we will still mark $2\lambda/p$ as λ .

Denote Δ as a diagonal matrix with its i -th diagonal element as d_i , and $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n] = \mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T$, where $\mathbf{u}_i \in \mathbb{R}^d$ is the i -th column of \mathbf{U} . We have

$$\max_{\mathbf{W}, \mathbf{G}, \mathbf{F}, \Delta} Tr(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) - \lambda Tr(\mathbf{U}^T \Delta \mathbf{U}). \quad (7)$$

126 Since the objective function in (7) is not jointly convex with all the variables,
 127 and Δ is dependent on \mathbf{W} , \mathbf{F} and \mathbf{G} , we propose the following iterative algorithm
 128 to alternatively update \mathbf{W} , \mathbf{G} , \mathbf{F} and Δ .

129 **Step 1: Fixing \mathbf{W} , \mathbf{G} and Δ and Optimizing \mathbf{F} .**

When \mathbf{W} , \mathbf{G} and Δ are fixed, the first term in (7) is constant, and we only need to minimize the second term. The optimization problem becomes

$$\min_{\mathbf{F}} \sum_{i=1}^n d_i \|\mathbf{W}^T \mathbf{x}_i - \mathbf{G} \mathbf{f}_i\|_2^2 = \min_{\mathbf{F}} \sum_{i=1}^n d_i \sum_{k=1}^c \|\mathbf{W}^T \mathbf{x}_i - \mathbf{g}_k\|_2^2 F_{ik}, \quad (8)$$

Since \mathbf{G} is fixed and \mathbf{F} is the cluster indicator matrix, according to the algorithm in [29] and [27], the optimized \mathbf{F} can be derived from

$$\mathbf{F}_{ij} = \begin{cases} 1, & j = \arg \min_k \|\mathbf{W}^T \mathbf{x}_i - \mathbf{g}_k\|_2^2, \\ 0, & \text{Otherwise.} \end{cases} \quad (9)$$

130 **Step 2: Fixing Δ and \mathbf{F} and Optimizing \mathbf{W} and \mathbf{G} .**

When Δ and \mathbf{F} are fixed, the closed-form solution of \mathbf{W} and \mathbf{G} can be derived as follows. Denote

$$\mathcal{L}(\mathbf{W}, \mathbf{G}) = Tr(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) - \lambda Tr(\mathbf{U}^T \Delta \mathbf{U}), \quad (10)$$

131 where $\mathbf{U} = \mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T$.

We take a derivative of $\mathcal{L}(\mathbf{W}, \mathbf{G})$ over \mathbf{G}

$$\begin{aligned}
\frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{G})}{\partial \mathbf{G}} &= -\lambda \frac{\partial Tr((\mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T)^T \Delta (\mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T))}{\partial \mathbf{G}}, \\
&= -\lambda \frac{\partial Tr((\mathbf{G} \mathbf{F}^T - \mathbf{W}^T \mathbf{X})^T \Delta (\mathbf{G} \mathbf{F}^T - \mathbf{W}^T \mathbf{X}))}{\partial \mathbf{G}}, \\
&= -\lambda \frac{\partial Tr(\mathbf{G} \mathbf{F}^T \Delta \mathbf{F} \mathbf{G}^T) - 2Tr(\mathbf{G} \mathbf{F}^T \Delta \mathbf{X}^T \mathbf{W})}{\partial \mathbf{G}}, \\
&= -2\lambda (\mathbf{G} \mathbf{F}^T \Delta \mathbf{F} - \mathbf{W}^T \mathbf{X} \Delta \mathbf{F}).
\end{aligned} \tag{11}$$

Let the above equation equals to zero, and we have

$$\mathbf{G} = \mathbf{W}^T \mathbf{X} \Delta \mathbf{F} (\mathbf{F}^T \Delta \mathbf{F})^{-1}. \tag{12}$$

Substituting \mathbf{G} into $\mathcal{L}(\mathbf{W}, \mathbf{G})$, we have

$$\begin{aligned}
\mathcal{L}(\mathbf{W}) &= Tr(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) - \lambda Tr((\mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T)^T \Delta (\mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T)), \\
&= Tr(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) - \lambda Tr(\mathbf{W}^T \mathbf{X} \Delta \mathbf{X}^T \mathbf{W} - \mathbf{W}^T \mathbf{X} \mathbf{F} (\mathbf{F}^T \Delta \mathbf{F})^{-1} \mathbf{F}^T \Delta^T \mathbf{X}^T \mathbf{W}), \\
&= Tr(\mathbf{W}^T (\mathbf{S}_t - \lambda \mathbf{X} \Delta \mathbf{X}^T + \lambda \mathbf{X} \Delta \mathbf{F} (\mathbf{F}^T \Delta \mathbf{F})^{-1} \mathbf{F}^T \Delta \mathbf{X}^T) \mathbf{W}), \\
&= Tr(\mathbf{W}^T \mathbf{M} \mathbf{W}),
\end{aligned} \tag{13}$$

132 where $\mathbf{M} = \mathbf{S}_t - \lambda \mathbf{X} \Delta \mathbf{X}^T + \lambda \mathbf{X} \Delta \mathbf{F} (\mathbf{F}^T \Delta \mathbf{F})^{-1} \mathbf{F}^T \Delta \mathbf{X}^T$.

Therefore the problem to optimize \mathbf{W} becomes

$$\max_{\mathbf{W}} Tr(\mathbf{W}^T \mathbf{M} \mathbf{W}) = \max_{\mathbf{W}} \sum_{i=1}^d Tr(\mathbf{w}_i^T \mathbf{M} \mathbf{w}_i), \tag{14}$$

133 According to the definition of \mathbf{W} in (4), we can optimize \mathbf{W} by locating the
134 first d largest diagonal elements of matrix \mathbf{M} .

Step 3: Updating Δ by calculating its i -th diagonal element as

$$d_i = \frac{p}{2} \|\mathbf{W}^T \mathbf{x}_i - \mathbf{G} \mathbf{f}_i\|_2^{p-2}. \tag{15}$$

135 It is important to note that there is a problem when using the above alternative
136 algorithm. Although the above solving strategy can guarantee convergence, its
137 result is not satisfactory. Like the traditional K-means method, there are a lot of
138 local optimizations which depend on initialization. Considering the above update
139 rules, when \mathbf{F} is fixed, the algorithm can quickly adjust \mathbf{W} and \mathbf{G} to adapt to the
140 \mathbf{F} . In other words, when we need to update the \mathbf{F} in the next step, the optimal \mathbf{F}
141 is the same as before. That is to say, the algorithm has fast convergence speed
142 and the optimal solution depends on the initial value. In order to avoid the local
143 optimal problem, the update rule proposed in [29] and [27] is employed. In each
144 step of updating \mathbf{F} , we will randomly initialize \mathbf{F} several times (20 times in our
145 experiment). If the value of the objective function $\|\mathbf{W}^T \mathbf{X} - \mathbf{G} \mathbf{F}^T\|_F^2$ is smaller
146 than that of the previous \mathbf{F} , then updating \mathbf{F} according to the random initialization.
147 Otherwise, updating \mathbf{F} by (9). That is, assume that in the i -th iteration, we have
148 gotten \mathbf{F}_i^* , \mathbf{W}_i^* and \mathbf{G}_i^* . In the $(i+1)$ -th iteration, we will get $\mathbf{F}_{i+1}^1, \mathbf{F}_{i+1}^2, \dots, \mathbf{F}_{i+1}^t$ by
149 random initialization, where t is the number of random initialization. We update
150 \mathbf{F} according to the following rules

$$\mathbf{F}_{i+1}^* = \begin{cases} \mathbf{F}_{i+1}^j, & \|\mathbf{W}_i^*\mathbf{X} - \mathbf{G}_i^*(\mathbf{F}_{i+1}^j)^T\|_F^2 < \|\mathbf{W}_i^*\mathbf{X} - \mathbf{G}_i^*(\mathbf{F}_i^*)^T\|_F^2, \\ \mathbf{F}^*, & \text{Otherwise,} \end{cases} \quad (16)$$

where \mathbf{F}^* is defined as

$$\mathbf{F}_{ij}^* = \begin{cases} 1, & j = \arg \min_k \|\mathbf{W}_i^* x_i - (\mathbf{g}_i^*)_k\|_2^2, \\ 0, & \text{Otherwise.} \end{cases} \quad (17)$$

151 The pseudo code of optimizing the proposed algorithm is listed in Algorithm
152 1.

Algorithm 1 The Algorithm to Solve Problem (5)

Input: The input data $\mathbf{X} \in \mathbb{R}^{D \times n}$, the reduced dimension number d , the number of clusters c , regularization parameter λ , and the distance metric parameter p .

Output: Transformation matrix \mathbf{W} , cluster indicator matrix \mathbf{F} , and cluster centroid matrix \mathbf{G} .

1: Initialize Δ as identity matrix, and randomly initialize \mathbf{W} and \mathbf{G} .

2: **while** Not convergent **do**

3: Update \mathbf{F} by (16);

4: Update \mathbf{G} by (12);

5: Update \mathbf{W} by locating the d largest diagonal elements of the matrix \mathbf{M} in (14);

6: Update Δ by calculating its diagonal elements by $d_i = \frac{p}{2} \|\mathbf{W}^T \mathbf{x}_i - \mathbf{G} \mathbf{f}_i\|_2^{p-2}$;

7: **end while**

153 *2.3. Convergence Analysis*

154 In this section, we prove the convergence of the proposed algorithm. First, we
155 give the following Lemma:

Lemma 1 [34]: For any nonzero vectors $\mathbf{e}^{t+1}, \mathbf{e}^t \in \mathbb{R}^m$, when $0 < p \leq 2$, we have:

$$\frac{\|\mathbf{e}^{t+1}\|_2^p}{\|\mathbf{e}^t\|_2^p} - \frac{p}{2} \frac{\|\mathbf{e}^{t+1}\|_2^2}{\|\mathbf{e}^t\|_2^2} - 1 + \frac{p}{2} \leq 0. \quad (18)$$

156 *Theorem 1:* When \mathbf{W}, \mathbf{G} and Δ are fixed, the derived \mathbf{F} in (9) is the global
157 solution to the problem (7). Similarly, when \mathbf{F} and Δ are fixed, the derived \mathbf{G}
158 in (12) and the derived \mathbf{W} by locating the d largest diagonal elements of $\mathbf{S}_t -$
159 $\lambda \mathbf{X} \Delta \mathbf{X}^T + \lambda \mathbf{X} \Delta \mathbf{F} (\mathbf{F}^T \Delta \mathbf{F})^{-1} \mathbf{F}^T \Delta \mathbf{X}^T$ are also the global solutions to the problem
160 in (7).

161 *Proof:* When \mathbf{W} , \mathbf{G} and Δ are fixed, optimizing the problem in (7) is equal to
 162 solving the traditional K-means on $\mathbf{W}^T \mathbf{X}$ with fixed centroid. Thus the optimized
 163 solution is unique.

164 According to (3), \mathbf{w}_i is a vector with only one element being 1 and the rest
 165 being 0. Obviously, the derived \mathbf{W} by locating the d largest diagonal elements
 166 of $\mathbf{S}_t - \lambda \mathbf{X} \Delta \mathbf{X}^T + \lambda \mathbf{X} \Delta \mathbf{F} (\mathbf{F}^T \Delta \mathbf{F})^{-1} \mathbf{F}^T \Delta \mathbf{X}^T$ maximizes the objective function in
 167 (14). When \mathbf{F} and Δ are fixed, \mathbf{G} is dependent on \mathbf{W} , and the global solution of
 168 \mathbf{W} can be derived from the process above.

169 To sum up, the theorem is proved.

170 *Theorem 2:* The procedure in Algorithm 1 monotonically increases the objec-
 171 tive function of the problem in (5) in each iteration.

Proof: Assume that we have derived the updated \mathbf{W}_t , \mathbf{G}_t in the t -th iteration.
 In the $(t + 1)$ -th iteration, we fix \mathbf{W}_t , \mathbf{G}_t and Δ_t , and get the optimized \mathbf{F}_{t+1} by
 (16). According to Theorem 1 and the updating rule in (9), we have

$$\begin{aligned} & Tr(\mathbf{W}_t^T \mathbf{S}_t \mathbf{W}_t) - \lambda \|\mathbf{W}_t^T \mathbf{X} - \mathbf{G}_t \mathbf{F}_t^T\|_{2,p}^p \\ & \leq Tr(\mathbf{W}_t^T \mathbf{S}_t \mathbf{W}_t) - \lambda \|\mathbf{W}_t^T \mathbf{X} - \mathbf{G}_t \mathbf{F}_{t+1}^T\|_{2,p}^p. \end{aligned} \quad (19)$$

Then we fix Δ_t and \mathbf{F}_{t+1} , and update \mathbf{G} and \mathbf{W} by maximizing (10). Let
 $f(\mathbf{W}) = Tr(\mathbf{W}_t^T \mathbf{S}_t \mathbf{W}_t)$, $\mathbf{u}_i^t = \mathbf{W}_t^T \mathbf{x}_i - \mathbf{G}_t(\mathbf{f}_i)_{t+1}$, and $\mathbf{u}_i^{t+1} = \mathbf{W}_{t+1}^T \mathbf{x}_i - \mathbf{G}_{t+1}(\mathbf{f}_i)_{t+1}$,
 we have

$$f(\mathbf{W}_t) - \lambda \sum_i d_i^t \|\mathbf{u}_i^t\|_2^2 \leq f(\mathbf{W}_{t+1}) - \lambda \sum_i d_i^t \|\mathbf{u}_i^{t+1}\|_2^2. \quad (20)$$

Since $d_i^t = \frac{p}{2} \|\mathbf{W}_t^T \mathbf{x}_i - \mathbf{G}_t(\mathbf{f}_i)_t\|_2^{p-2}$, thus we have

$$f(\mathbf{W}_t) - \lambda \sum_i \frac{p}{2} \|\mathbf{u}_i^t\|_2^p \leq f(\mathbf{W}_{t+1}) - \lambda \sum_i \frac{p}{2} \|\mathbf{u}_i^t\|_2^{p-2} \|\mathbf{u}_i^{t+1}\|_2^2, \quad (21)$$

which can be wrote as

$$f(\mathbf{W}_t) - \lambda \sum_i \frac{p}{2} \frac{\|\mathbf{u}_i^t\|_2^2}{\|\mathbf{u}_i^t\|_2^{2-p}} \leq f(\mathbf{W}_{t+1}) - \lambda \sum_i \frac{p}{2} \frac{\|\mathbf{u}_i^{t+1}\|_2^2}{\|\mathbf{u}_i^t\|_2^{2-p}}, \quad (22)$$

According to Lemma 1, we have

$$\frac{p}{2} \frac{\|\mathbf{u}_i^{t+1}\|_2^2}{\|\mathbf{u}_i^t\|_2^2} \|\mathbf{u}_i^t\|_2^p \geq \|\mathbf{u}_i^{t+1}\|_2^p - (1 - \frac{p}{2}) \|\mathbf{u}_i^t\|_2^p, \quad (23)$$

which holds for each index i , thus we have

$$\frac{p}{2} \sum_i \frac{\|\mathbf{u}_i^{t+1}\|_2^2}{\|\mathbf{u}_i^t\|_2^2} \|\mathbf{u}_i^t\|_2^p \geq \sum_i \|\mathbf{u}_i^{t+1}\|_2^p - (1 - \frac{p}{2}) \sum_i \|\mathbf{u}_i^t\|_2^p, \quad (24)$$

that is

$$- \sum_i \|\mathbf{u}_i^t\|_2^p + \frac{p}{2} \sum_i \frac{\|\mathbf{u}_i^t\|_2^2}{\|\mathbf{u}_i^t\|_2^{2-p}} \leq - \sum_i \|\mathbf{u}_i^{t+1}\|_2^p + \frac{p}{2} \sum_i \frac{\|\mathbf{u}_i^{t+1}\|_2^2}{\|\mathbf{u}_i^t\|_2^{2-p}}, \quad (25)$$

Combining (22) and (25), we have

$$f(\mathbf{W}_t) - \lambda \sum_i \|\mathbf{u}_i^t\|_2^p \leq f(\mathbf{W}_{t+1}) - \lambda \sum_i \|\mathbf{u}_i^{t+1}\|_2^p. \quad (26)$$

172 To sum up, Algorithm 1 monotonically increases the objective function of the
 173 problem in (5) in each iteration. Since (5) has an obvious upper bound $Tr(\mathbf{X}\mathbf{X}^T)$,
 174 Algorithm 1 will monotonically increase the objective function until it converges.

175 2.4. Complexity Analysis

176 First we consider the computation complexity of Algorithm 1. It contains
 177 three main components, i.e., K-means in the subspace with computation com-
 178 plexity $O(dcn)$, the process of computing matrix \mathbf{G} with computation complexity
 179 $O(dcn + c^2n)$ and computing matrix \mathbf{M} 's diagonal elements to optimize \mathbf{W} with
 180 computation complexity $O(Dn + D + d \log d)$. Denote the repeated initialization

181 times of \mathbf{F} in (16) as T_k , and the number of iterations in the whole algorithm as T_t ,
 182 then the computational complexity of our algorithm is $O(T_t(T_k(DCN) + DCN +$
 183 $c^2n + Dn + d \log d) \sim O(Dn)$. Next we consider the memory cost of Algorithm
 184 1. Algorithm 1 mainly involves matrices such as \mathbf{X} , \mathbf{F} , \mathbf{G} , etc. $O(Dn + cn + dc)$
 185 is needed for storage. Thus, the calculation cost of our algorithm has a linear re-
 186 lationship with the dimension of the data. According to the above analysis, our
 187 algorithm can deal with high-dimensional data well.

188 2.5. Parameter Determination

189 Our method mainly involves three important parameters: the reduced dimen-
 190 sion d , the balance parameter λ , and the p value of the $l_{2,p}$ -norm used in the dis-
 191 tance metric. Since the determination of the parameters is still an open problem
 192 in the related fields, we use heuristic and empirical methods to determine the pa-
 193 rameters.

194 The first parameter d represents the number of features that can best repre-
 195 sent the original data. When d is too large, the representation of the original data
 196 is still redundant and the curse of dimension still exists. When d is too small,
 197 there may be loss of information so that different clusters cannot be separated.
 198 In this paper, by changing the value of d , the parameters with the best accuracy
 199 are selected through grid search. The second parameter is the balance parame-
 200 ter λ . Obviously, this parameter balances the effect of dimensionality reduction
 201 and clustering on the value of objective function. The larger λ , the greater the
 202 impact of clustering is. Following the setting in [27], we search λ in the range
 203 of $[10^{-6}, 10^{-4}, 10^{-2}, \dots, 10^2, 10^4, 10^6]$. The third parameter p affects the distance
 204 between data points in KM, and then influences the clustering results. We adjust
 205 the value between 0 and 2. The influence of different parameter values will be

206 discussed in the experimental section.

207 **3. Experiments**

208 *3.1. Data Description and Evaluation Metric*

209 *3.1.1. Data Description*

210 We conduct analytical experiments on seven datasets to evaluate the perfor-
211 mance. For each dataset, we preprocess all the values by centralization. These
212 datasets include:

213 UCI datasets ¹: We evaluate our algorithm on four datasets: Cars, Wine, Iono-
214 sphere, and Ecoli.

215 USPS Digit Dataset ²: The dataset includes 9298 handwritten digital images,
216 all of which are grayscale images of 16 pixels. We select 20% of the dataset for
217 the experiment.

218 Umist Face Dataset ³: 575 images in total, corresponding to 20 different peo-
219 ple. Each category consists of 19 to 48 images.

220 COIL-20 Object Dataset [40]: It contains 20 objects and each object has 72
221 samples taken at pose intervals of five degrees. We first extract LBP features with
222 3076 dimensions and reduce the dimension to 300 for evaluating the performance
223 of our method.

224 The detailed description of the aforementioned datasets is displayed in Table
225 1.

¹<http://archive.ics.uci.edu/ml/datasets.html>

²<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multiclass.html>

³<http://images.ee.umist.ac.uk/danny/database.html>

Table 1: Summary of the different datasets.

Datasets	Classes (c)	Samples (n)	Total features (D)
Cars	3	392	8
Wine	3	178	13
Ionosphere	2	351	34
Ecoli	8	366	343
USPS	10	1854	256
Umist	20	575	644
COIL-20	20	1440	3076

226 3.1.2. Evaluation Metric

227 In order to evaluate the effectiveness of the proposed method, we will compare
 228 it with some relevant subspace clustering methods. Meanwhile, in order to express
 229 the effect of dimensionality reduction, we will also provide the results of K-means
 230 clustering for comparison. The detailed introduction is as follows:

- 231 - **KM** represents the traditional K-means algorithm, and its results will be
 232 used as the benchmark in the experiment.
- 233 - **PCA KM** means that PCA is first used to reduce the dimension of data, and
 234 then KM clustering is used for clustering.
- 235 - **DEC** [29] is a general discriminant subspace learning framework, which
 236 optimizes both PCA and KM simultaneously.
- 237 - **TRACK** [36] adopts LDA and KM clustering methods, and uses regulariza-
 238 tion technique of structured sparse induction criterion to select discriminant
 239 features.

240 - **FAKM** [27] combines feature selection with KM clustering, and uses an
 241 adaptive loss function in the objective function.

242 All the compared methods are implemented in MATLAB (R2016a). The com-
 243 puter processor is Intel(R)Core(TM) i7-7500T CPU @ 2.70GHz, and the memory
 244 is 8-GB. We used three indicators of accuracy (ACC), normalized mutual infor-
 245 mation (NMI) and purity to evaluate the clustering performance of all methods.

Denote g_i as the real label of x_i , q_i as the result of cluster process. Accuracy (ACC) is defined as follow

$$ACC = \frac{\sum_{i=1}^n \sigma(g_i, \text{map}(q_i))}{n}, \quad (27)$$

where $\text{map}(\cdot)$ is a mapping function to obtain the matching between real tags and clustering tags by Kuhn-Munkres algorithm. $\delta(x, y)$ is the Kronecker function

$$\delta(x, y) = \begin{cases} 1, & x = y, \\ 0, & \text{Otherwise.} \end{cases} \quad (28)$$

246 A larger value of accuracy (ACC) indicates a better clustering result.

Denote C as the real classes tag set of the sample, C' as the classes tag set obtained by clustering algorithm. Normalized mutual information (NMI) can be defined by the following formula

$$NMI(C, C') = \frac{MI(C, C')}{\max(H(L), H(C))}, \quad (29)$$

where $H(\cdot)$ represents the entropy. $MI(C, C')$ is the mutual information between C and C' , as defined below

$$MI(C, C') = \sum_{c_i \in C, c'_j \in C'} p(c_i, c'_j) \cdot \log_2 \frac{p(c_i, c'_j)}{p(c_i) \cdot p(c'_j)}, \quad (30)$$

247 here $p(c_i, c'_j)$ is the probability of a randomly selected sample belongs to both
248 cluster c_i and c'_j . It is easy to observe that the value of normalized mutual infor-
249 mation (NMI) is between 0 and 1. Similar to the accuracy rate (ACC), the larger
250 the NMI, the better the clustering result.

Purity is a very simple clustering evaluation method, which is calculated by assigning the labels of a cluster to the most frequent classes. The mathematical definition is as follows

$$purity(C, C') = \frac{1}{N} \sum_j \max_i |c'_j \cap c_j|, \quad (31)$$

251 where N represents the total number of samples. Similarly, $purity \in [0, 1]$, the
252 closer the value is to 1, the better the result.

253 3.2. Toy Example on Iris

254 To show the visual effectiveness, we first conduct a small experiment on Iris
255 dataset⁴. The dataset consists of three categories (setosa, versicolor and Virginia).
256 The petal length and petal width are chosen for experiment to show a visualization
257 example. DEC is used to compare with our method, and d is set as 2. We first
258 cluster the iris data, and then use the obtained optimal transformation matrix to
259 project the original data into a two-dimensional space. The clustering results are
260 shown in Fig 1, where the samples of the wrong cluster are marked with red 'x'.

261 As we can see, our method has fewer error markers than DEC. In addition,
262 From Fig. 1.(c), it can be seen that the features selected by our method are con-
263 sistent with the two features that can distinguish the various types of samples
264 visually, namely, the length and width of petals. From the Fig. 1.(b), we can see

⁴<http://archive.ics.uci.edu/ml/datasets/Iris>

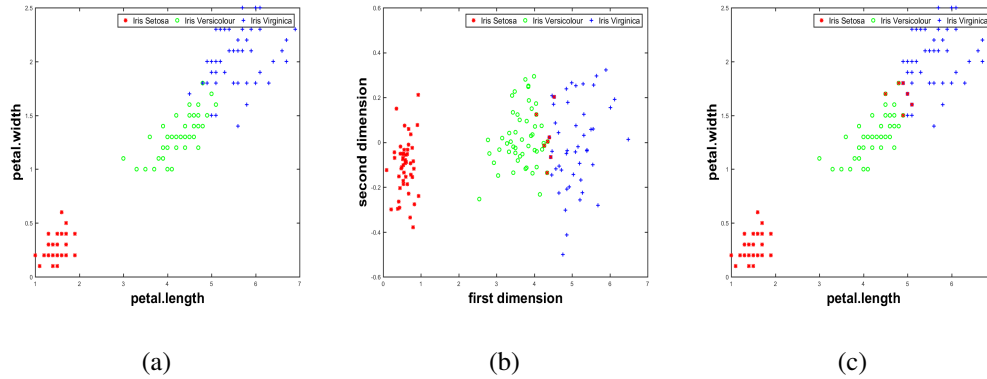


Figure 1: Clustering results on the Iris dataset, the dimension is reduced to 2. (a) Original data. (b) Clustering results of DEC. (c) Clustering results of our methods.

Table 2: Comparison of clustering results (ACC%)

Methods	Cars	Wine	Ionosphere	Ecoli	USPS	Umist	COIL-20
KM	44.79 \pm 0.13	64.80 \pm 6.44	70.75 \pm 1.60	55.67 \pm 7.69	62.03 \pm 3.89	41.67 \pm 2.23	62.78 \pm 0.04
PCAKM	44.82 \pm 0.12	67.64 \pm 5.44	71.11 \pm 0.14	68.93 \pm 6.41	63.91 \pm 1.64	42.10 \pm 2.32	59.38 \pm 3.27
TRACK	45.66 \pm 0.00	70.22 \pm 0.00	71.88 \pm 0.14	63.01 \pm 5.42	65.70 \pm 0.27	47.97 \pm 4.02	54.04 \pm 3.08
DEC	47.68 \pm 0.08	70.22 \pm 0.00	71.23 \pm 0.00	62.08 \pm 3.85	64.96 \pm 0.09	44.54 \pm 1.73	67.74 \pm 2.81
FAKM	59.18 \pm 0.17	88.20 \pm 0.00	72.31 \pm 3.63	69.73 \pm 6.11	66.98 \pm 3.37	48.43 \pm 1.98	67.02 \pm 3.33
OURS	62.50 \pm 1.31	88.20 \pm 0.00	74.93 \pm 0.00	72.05 \pm 2.25	67.49 \pm 2.79	48.54 \pm 3.10	67.23 \pm 1.98

265 that DEC has a completely different structure. Therefore, our approach better p-
 266 reserves the structure of the original data than that of DEC by selecting the most
 267 representative features.

268 3.3. Comparison of Clustering Results

269 In this section, we show the clustering results of different methods on differ-
 270 ent datasets. Grid search is conducted for different parameters according to the
 271 above mentioned, and the best combination of parameters is selected to repeat the
 272 experiment for 10 times and the average value is taken. The results are shown in

Table 3: Comparison of clustering results (NMI%)

Methods	Cars	Wine	Ionosphere	Ecoli	USPS	Umist	COIL-20
KM	19.35 ± 0.33	41.61 ± 1.49	12.30 ± 2.99	49.09 ± 4.00	61.73 ± 2.67	62.93 ± 2.27	73.28 ± 1.97
PCAKM	19.45 ± 0.32	42.27 ± 1.57	13.01 ± 0.00	57.38 ± 3.51	62.34 ± 0.68	64.07 ± 2.27	71.82 ± 2.09
TRACK	30.39 ± 3.78	43.56 ± 2.68	13.49 ± 0.48	55.29 ± 6.66	63.60 ± 0.79	64.43 ± 2.27	66.30 ± 1.92
DEC	19.10 ± 0.00	42.87 ± 0.00	13.12 ± 0.00	56.54 ± 2.58	62.90 ± 0.66	65.77 ± 2.27	75.94 ± 1.29
FAKM	19.10 ± 7.17	65.69 ± 0.00	12.85 ± 9.72	57.59 ± 1.58	63.60 ± 0.88	66.84 ± 2.27	75.60 ± 1.63
OURS	30.39 ± 0.00	65.69 ± 0.00	18.86 ± 0.00	58.52 ± 1.46	64.08 ± 1.12	66.74 ± 2.27	75.65 ± 1.13

Table 4: Comparison of clustering results (purity%)

Methods	Cars	Wine	Ionosphere	Ecoli	USPS	Umist	COIL-20
KM	65.05 ± 0.00	69.52 ± 0.84	70.75 ± 1.60	76.60 ± 3.17	70.76 ± 3.22	49.90 ± 2.36	66.06 ± 2.90
PCAKM	65.05 ± 0.00	69.89 ± 1.12	71.11 ± 0.00	80.59 ± 2.89	71.51 ± 1.49	50.63 ± 2.36	63.28 ± 2.82
TRACK	65.03 ± 0.00	70.22 ± 0.00	71.88 ± 0.27	80.86 ± 7.43	73.19 ± 1.54	52.89 ± 2.36	57.78 ± 2.28
DEC	65.05 ± 0.00	70.22 ± 0.00	71.23 ± 0.00	82.17 ± 2.39	72.40 ± 1.58	52.94 ± 2.36	70.39 ± 2.57
FAKM	67.85 ± 7.17	88.20 ± 0.00	72.30 ± 6.38	81.69 ± 7.87	73.40 ± 1.80	56.94 ± 2.36	70.04 ± 2.32
OURS	69.03 ± 0.12	88.20 ± 0.00	75.73 ± 0.00	82.83 ± 0.14	73.59 ± 2.21	56.50 ± 2.36	70.24 ± 1.52

273 Tables 2-4.

274 From the tables, we can get the following observations:

- 275 - Most KM-based subspace clustering algorithms have better performance
276 than KM on each dataset, which shows the effectiveness of this kind of
277 algorithm. Although the NMIs of DEC and FAKM on the Cars dataset are
278 lower than KM, these methods still achieve a smaller gap with KM when the
279 dimension is reduced and the calculation cost of subsequent learning tasks
280 is greatly reduced.
- 281 - DEC achieves better results than PCAKM on all datasets except Ecoli be-
282 cause it builds a more general discriminant clustering framework.

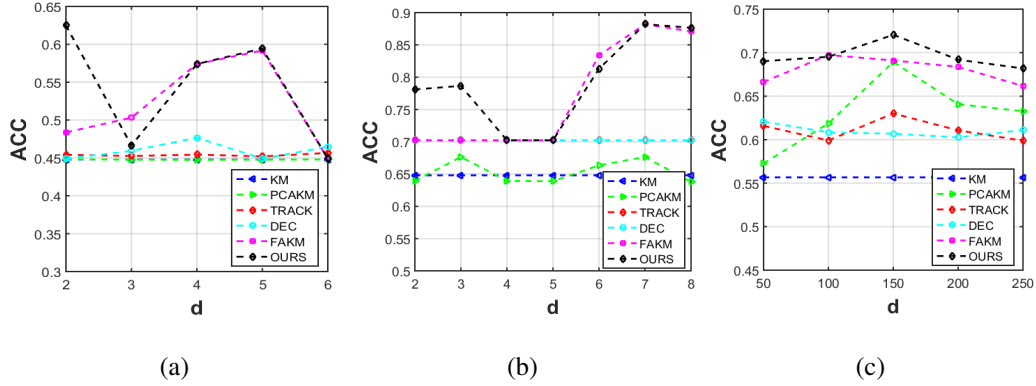


Figure 2: Clustering results (ACC) of the compared methods on the different d . (a) Cars. (b) Wine. (c) Ecoli.

- 283 - Compared to DEC, we can see that our method achieves better results on
- 284 the most of datasets due to the robustness of $l_{2,p}$ -norm as a distance metric.
- 285 - Compared with TRACK, which also combines feature selection and clus-
- 286 tering, our method also has better performance. The reason may be that our
- 287 method is more flexible in balancing the scatter matrix.
- 288 - FAKM defines an adaptive objective function to improve the robustness of
- 289 the method. In comparison, our method has similar or better results, which
- 290 indicates that the objective function based on $l_{2,p}$ -norm is more robust.

291 3.4. Impact of Dimension Reduction

292 In addition, we also study the effect of the reduced dimension d on different
 293 datasets by different methods, and the parameter setting is the same as above, each
 294 experiment is repeated ten times, and the mean value is recorded. The results are
 295 shown in Fig. 2 and Fig. 3.

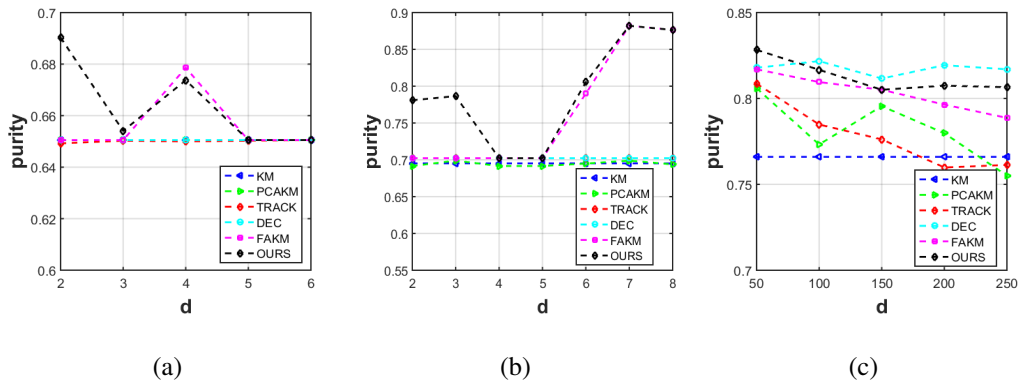
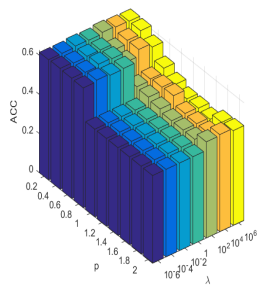


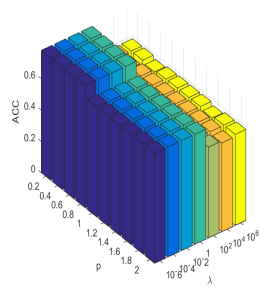
Figure 3: Clustering results (purity) of the compared methods on the different d . (a) Cars. (b) Wine. (c) Ecoli.

296 Through observation, the following conclusions can be drawn:

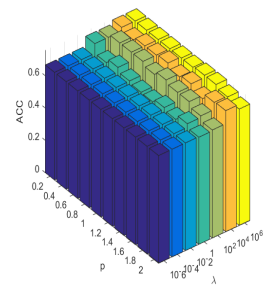
- 297 - Not all the methods can achieve better results when d is increased, which
 298 indicates that dimension reduction can effectively improve the performance
 299 of clustering.
- 300 - When only a small dimension is reserved, the performance of some sub-
 301 space clustering methods will decline because of the excessive information
 302 loss.
- 303 - Our method tends to perform better on smaller dimensions than other meth-
 304 ods. In addition, the optimal results are usually obtained on the smaller
 305 dimensions, which indicates that our method can effectively select the most
 306 important features in the data.
- 307 - Our method can get the better results in most cases.



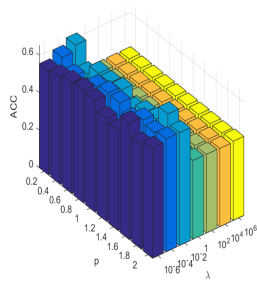
(a)



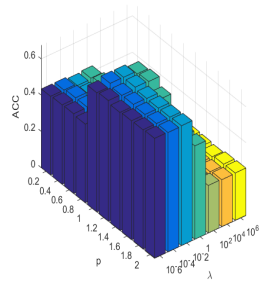
(b)



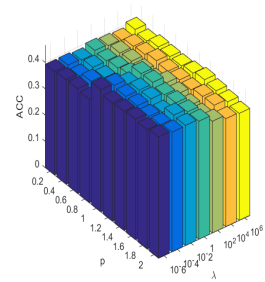
(c)



(d)



(e)



(f)

Figure 4: Parameters sensitivity analysis. (a) Cars; (b) Wine; (c) Ionosphere; (d) Ecoli; (e) Usps; (f) Umist

308 3.5. Parameter Analysis

309 In order to understand how the parameters λ and p affect the results of the
310 clustering experiment, we fix the value of d in the experiment and carry out the
311 parameters sensitivity experiment. The results are shown in Fig. 4.

312 As can be seen from Fig. 4, λ and p both have great influence on the final
313 clustering accuracy. Let's first discuss the impact of λ . From the experimental
314 results, we notice that the clustering performance is sensitive to λ . For example,
315 in Fig. 4.(b), i.e., the Wine dataset, the result of $\lambda < 1$ is significantly better
316 than that of $\lambda > 1$. At the same time, it can be found that if we can choose a
317 value close to the λ that get the optimal result, we can get a good result, but it is
318 also affected by the value of parameter p . We can see that the p also affects the
319 result by the different value range. Take Fig. 4.(a) and (b) as examples, when p
320 approximately belongs to $(0, 1)$, the clustering results are better. In Fig. 4.(c) and
321 (e), it is approximately within the range of $(1, 2)$ for higher accuracies. The above
322 observation is very helpful for parameter selection, that is, the parameter value
323 can be approximately determined by finding which range of results are better.

324 4. Conclusion

325 In this paper, we propose a flexible subspace clustering model. Specifically,
326 we first incorporate feature selection and K-means clustering into a single frame-
327 work, which can select the refined features and improve the clustering perfor-
328 mance. Second, we embed the $l_{2,p}$ -norm into the framework to enhance the ro-
329 bustness and retain the desirable properties from big data. Finally, considering the
330 proposed model is neither convex nor Lipschitz continuous, we develop an effec-
331 tive algorithm to solve it. In addition, we also theoretically prove the convergence

332 of the proposed algorithm. Experimental results verify the presented method has
333 more robust and better performance on benchmark databases compared to the ex-
334 isting approaches.

335 It should be noted that the proposed method could obtain more robust results
336 than the existing methods due to the flexibility of selecting p value. However, the
337 proposed model can only choose the parameter p manually for different dataset-
338 s. Recently, many adaptive learning approaches [41] [42] are successfully used
339 in data mining and pattern recognition. Can the idea be used for our clustering
340 model to adjust the parameter p automatically according to characters of different
341 datasets? If the answer is yes, how to design the adaptive scheme? Our future
342 work will focus on the topic.

343 **References**

- 344 [1] Yao D, Yu C, Yang L T, Jin H. Using Crowdsourcing to Provide QoS for
345 Mobile Cloud Computing. IEEE Transactions on Cloud Computing, 2019,
346 7(2): 344-356.
- 347 [2] Sakr S, Elgammal A. Towards a Comprehensive Data Analytics Framework
348 for Smart Healthcare Services. Big Data Research, 2016, 4: 44-58.
- 349 [3] Hendre A, Joshi K P. A Semantic Approach to Cloud Security and Compli-
350 ance. International Conference on Cloud Computing, 2015: 1081-1084.
- 351 [4] Ren L, Meng Z, Wang X, Zhang L, Yang L T , A Data-driven Approach of
352 Product Quality Prediction for Complex Production Systems, IEEE Trans-
353 actions on Industrial Informatics, doi: 10.1109/TII.2020.3001054, 2020.

- 354 [5] Ren L, Meng Z, Wang X, Luan R, Yang L T, A Wide-Deep-Sequence
355 Model based Quality Prediction Method in Industrial Process Analy-
356 sis, IEEE Transactions on Neural Networks and Learning Systems, doi:
357 10.1109/TNNLS.2020.3001602, 2020.
- 358 [6] Li Z, Chen R, Liu L, Min G. Dynamic Resource Discovery Based on Pref-
359 erence and Movement Pattern Similarity for Large-Scale Social Internet of
360 Things. IEEE Internet of Things Journal, 2016, 3(4): 581-589.
- 361 [7] Wang X, Yang L T, Wang Y, Ren L, Deen M J. ADTT: A Highly-Efficient
362 Distributed Tensor-Train Decomposition Method for IIoT Big Data. IEEE
363 Transactions on Industrial Informatics, doi: 10.1109/tii.2020.2967768,
364 2020.
- 365 [8] Li G, Wu H, Jiang G, Xu S, Liu H. Dynamic Gesture Recognition in the
366 Internet of Things. IEEE Access, 2019, 7: 23713-23724.
- 367 [9] Wang X, Yang L T, Song L, Wang H, Ren L, Deen M J. A Tensor-based
368 Multi-Attributes Visual Feature Recognition Method for Industrial Intel-
369 ligence. IEEE Transactions on Industrial Informatics, doi: 10.1109/TI-
370 I.2020.2999901, 2020.
- 371 [10] Wang X, Wang Y, Zhe H, Juan D. The Research on Resource Scheduling
372 Based on Fuzzy Clustering in Cloud Computing. International Conference
373 on Intelligent Computation Technology and Automation, 2015: 1025-1028.
- 374 [11] Zhang X, Meng F, Xu J. PerfInsight: A Robust Clustering-Based Abnor-
375 mal Behavior Detection System for Large-Scale Cloud. International Con-
376 ference on Cloud Computing, 2018:896-899.

- 377 [12] Estiri H , Omran B A , Murphy S N . kluster : An Efficient Scalable Proce-
378 dure for Approximating the Number of Clusters in Unsupervised Learning.
379 Big Data Research, 2018, 13: 38-51.
- 380 [13] Zhang Q, Yang L T, Chen Z, Li P. High-order possibilistic c-means al-
381 gorithms based on tensor decompositions for big data in IoT. Information
382 Fusion, 2018, 39: 72-80.
- 383 [14] Hodge V J, Austin J. A Survey of Outlier Detection Methodologies. Artifi-
384 cial Intelligence Review, 2004, 22(2): 85-126.
- 385 [15] Hathaway R J, Bezdek J C, Hu Y. Generalized fuzzy c-means clustering
386 strategies using $L_{p/p}$ norm distances. IEEE Transactions on Fuzzy Sys-
387 tems, 2000, 8(5): 576-582.
- 388 [16] Salem S B, Naouali S, Chtourou Z, et al. A fast and effective partitional
389 clustering algorithm for large categorical datasets using a k-means based
390 approach. Computers & Electrical Engineering, 2018, 68: 463-483.
- 391 [17] Cai X, Nie F, Huang H. Multi-view K-means clustering on big data. Inter-
392 national Joint Conference on Artificial Intelligence, 2013: 2598-2604.
- 393 [18] Liang D, Peng Z, Lei S, Wang H, Fan M, Wang W, Shen Y. Robust mul-
394 tiple kernel K-means using $\ell_{2,p}$ norm. International Joint Conference on
395 Artificial Intelligence, 2015, 3476-3482.
- 396 [19] Chang X, Nie F, Wang S, Yang Y, Zhou X, Zhang C. Compound rank- k
397 projections for bilinear analysis. IEEE Transactions on Neural Networks
398 and Learning Systems, 2016, 27(7): 1502-1513.

- 399 [20] Jolliffe, IT (1986). Principal Component Analysis. New York: Springer-
400 Verlag.
- 401 [21] Duda, Richard O, Hart, Peter E, Stork, David G. Pattern Classification (2nd
402 Edition). John Wiley, 2000.
- 403 [22] Yu Y F, Ren C X, Jiang M, et al. Sparse approximation to discriminant
404 projection learning and application to image classification. Pattern Recog-
405 nition, 2019, 96: 1-10.
- 406 [23] Roweis S T and Saul L K, Nonlinear dimensionality reduction by locally
407 linear embedding, Science, 2000, vol.290, no.5500, pp.2323-2326.
- 408 [24] Hou C, Nie F, Jiao Y, Zhang C, Wu Y. Learning a subspace for clustering via
409 pattern shrinking. Information Processing and Management, 2013, 49(4):
410 871-883.
- 411 [25] Yin X, Chen S, Hu E. Regularized soft K-means for discriminant analysis.
412 Neurocomputing, 2013, 103: 29-42.
- 413 [26] Wang X, Chen R, Hong C, Zhang Z. Unsupervised feature analysis with
414 sparse adaptive learning. Pattern Recognition Letters, 2018, 102: 89-94.
- 415 [27] Wang X, Chen R, Yan F, et al. Fast adaptive K-means subspace clustering
416 for high-dimensional data. IEEE Access, 2019, 7: 42639-42651.
- 417 [28] Ding C, Li T. Adaptive dimension reduction using discriminant analysis
418 and K -means clustering. International Conference on Machine Learning,
419 2007: 521-528.

- 420 [29] Hou C, Nie F, Yi D, Tao D. Discriminative embedded clustering: A frame-
421 work for grouping high-dimensional data. *IEEE Transactions on Neural*
422 *Networks and Learning Systems*, 2015, 26(6): 1287-1299.
- 423 [30] Li H, Jiang T, Zhang K. Efficient and robust feature extraction by maxi-
424 mum margin criterion. *IEEE Transactions on Neural Networks and Learn-*
425 *ing Systems*, 2006, 17(1): 157-165.
- 426 [31] Park H, Jeon M, Rosen J B, et al. Lower dimensional representation of
427 text data based on centroids and least squares. *Bit Numerical Mathematics*,
428 2003, 43(2): 427-448.
- 429 [32] Nie F, Xiang S, Liu Y, Hou C, Zhang C. Orthogonal vs. uncorrelated least
430 squares discriminant analysis for feature extraction. *Pattern Recognition*
431 *Letters*, 2012, 33(5): 485-491.
- 432 [33] Yu Y F, Xu G, Huang K K, Zhu H, Chen L, Wang H. Dual
433 Calibration Mechanism based $L_{2,p}$ -Norm for Graph Matching. *IEEE*
434 *Transactions on Circuits and Systems for Video Technology*, doi:
435 10.1109/TCSVT.2020.3023781.
- 436 [34] Wang Q, Gao Q, Gao X, Nie F. $\ell_{2,p}$ -Norm Based PCA for Image Recogni-
437 tion. *IEEE Transactions on Image Processing*, 2018, 27(3): 1336-1346.
- 438 [35] Tao H, Hou C, Nie F, Jiao Y, Yi D. Effective discriminative feature selec-
439 tion with nontrivial solution. *IEEE Transactions on Neural Networks and*
440 *Learning Systems*, 2016, 27(4): 796-808.
- 441 [36] Wang D, Nie F, Huang H. Unsupervised feature selection via unified trace

- 442 ratio formulation and K -means clustering (TRACK). European Conference
443 on Machine Learning, 2014, 306-321.
- 444 [37] Nguyen B, De Baets B. Kernel-based distance metric learning for su-
445 pervised k -means clustering. IEEE Transactions on Neural Networks and
446 Learning Systems, 2019, 30(10): 3084-3095.
- 447 [38] Peng J, Wei Y. Approximating K-means-type clustering via semidefinite
448 programming. Siam Journal on Optimization, 2007, 18(1): 186-205.
- 449 [39] Luo M, Nie F, Chang X, Yang Y, Hauptmann A. G. and Zheng Q. Adaptive
450 unsupervised feature selection with structure regularization. IEEE Transac-
451 tions on Neural Networks and Learning Systems, 2018, 29(4): 944-956.
- 452 [40] Yu K, Zhang T, Gong Y. Nonlinear Learning using Local Coordinate Cod-
453 ing. Advances in Neural Information Processing Systems, 2009, 2223-
454 2231.
- 455 [41] Ren C X, Liang B H, Ge P, Zhai Y M, Lei Z. Domain Adaptive Person Re-
456 Identification via Camera Style Generation and Label Propagation. IEEE
457 Transactions on Information Forensics and Security, 2020, 15: 1290-1302.
- 458 [42] Yan F , Wang X D , Zeng Z Q , et al. Adaptive Multi-view Subspace Clus-
459 tering for High-dimensional Data. Pattern Recognition Letters, 2019, 130:
460 299-305.