

Subspace learning-based graph regularized feature selection

Shang, Ronghua; Wang, Wenbing; Stolkin, Rustam; Jiao, Licheng

DOI:

[10.1016/j.knosys.2016.09.006](https://doi.org/10.1016/j.knosys.2016.09.006)

License:

Creative Commons: Attribution-NonCommercial-NoDerivs (CC BY-NC-ND)

Document Version

Peer reviewed version

Citation for published version (Harvard):

Shang, R, Wang, W, Stolkin, R & Jiao, L 2016, 'Subspace learning-based graph regularized feature selection', *Knowledge-Based Systems*, vol. 112, pp. 152-165. <https://doi.org/10.1016/j.knosys.2016.09.006>

[Link to publication on Research at Birmingham portal](#)

Publisher Rights Statement:

Checked 11/11/2016

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

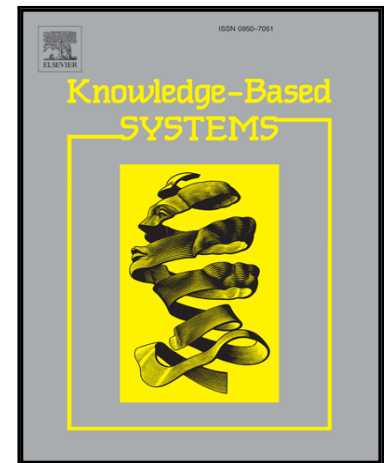
While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.

Subspace learning-based graph regularized feature selection

Ronghua Shang , Wenbing Wang , Rustam Stolkin , Licheng Jiao

PII: S0950-7051(16)30322-7
DOI: [10.1016/j.knosys.2016.09.006](https://doi.org/10.1016/j.knosys.2016.09.006)
Reference: KNOSYS 3664



To appear in: *Knowledge-Based Systems*

Received date: 27 May 2016
Revised date: 23 August 2016

Please cite this article as: Ronghua Shang , Wenbing Wang , Rustam Stolkin , Licheng Jiao , Subspace learning-based graph regularized feature selection, *Knowledge-Based Systems* (2016), doi: [10.1016/j.knosys.2016.09.006](https://doi.org/10.1016/j.knosys.2016.09.006)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Subspace learning-based graph regularized feature selection

Ronghua Shang¹, Wenbing Wang¹, Rustam Stolkin² and Licheng Jiao¹

(¹Key Laboratory of Intelligent Perception and Image Understanding of Ministry of Education, Xidian University, Xi'an, Shaanxi Province 710071, China.

²Extreme Robotics Lab, University of Birmingham UK.)

Abstract: In recent years, a variety of feature selection algorithms based on subspace learning have been proposed. However, such methods typically do not exploit information about the underlying geometry of the data. To overcome this shortcoming, we propose a novel algorithm called subspace learning-based graph regularized feature selection (SGFS). SGFS builds on the feature selection framework of subspace learning, but extends it by incorporating the idea of graph regularization, in which a feature map is constructed on the feature space in order to preserve geometric structure information on the feature manifold. Additionally, the $L_{2,1}$ -norm is used to constrain the feature selection matrix to ensure the sparsity of the feature array and avoid trivial solutions. The resulting method can provide more accurate discrimination information for feature selection. We evaluate SGFS by comparing it against five other state-of-the-art algorithms from the literature, on twelve publicly available benchmark data sets. Empirical results suggest that SGFS is more effective than the other five feature selection algorithms.

Key words: Graph regularized; subspace learning; feature manifold; sparse constraint; feature selection

1. Introduction

The increasingly rapid growth of information technology has seen a corresponding growth in the number of dimensions of gathered data. In many high dimensional data sets, only a small subset of the available features are useful, with most features being redundant, and some features even corresponding to information-less noise [1], [2], [3], [4], [5]. To facilitate the subsequent processing, it is often necessary to reduce the dimension of such high dimensional data. Processing of high dimensional data has become a challenge for researchers in many different fields [6], [7], [8], including data mining, machine learning, pattern recognition and others. Dimensionality reduction

methods can broadly be categorized into methods for feature selection and feature extraction [9], [10], [11], [12], [13]. Feature selection methods select representative features from the original set of features based on a variety of evaluation methods. In contrast, feature extraction methods map high dimensional data into a low dimensional space through a transformation matrix. Feature selection methods select a subset of the original “raw” feature data, and so retain the physical or real-world meaning of the original data. This means that the performance of the resulting classifiers can often be readily explained in terms of intuitively meaningful trends in the underlying data. In contrast, it may be difficult to explain the behaviour of feature extraction methods in terms of the relationship between the new feature and the sample class [1]. In this paper, we propose a new feature selection algorithm.

Feature selection methods can broadly be divided into: supervised [2], [14], semi supervised [15], unsupervised [16], [17], [18]. In supervised feature selection problems, the data discrimination information and also the correlation between features and the class of each data sample is available during training. However, in order to obtain large amounts of such class information, need for training such methods, a large amount of human resources are typically required, e.g. for hand annotation of data [18]. Semi-supervised feature selection requires only a smaller portion of the training data to be annotated with class label information to improve the accuracy of feature selection [15]. Unsupervised feature selection, without any class label information, only relies on the inherent information of the input data to determine the importance of features [16]. In many practical applications, the true class label information is unknown, which makes unsupervised feature selection methods more widely applicable to real problems, but also engenders greater challenges for researchers. According to various possible search strategies, unsupervised feature selection can be divided into filter, wrapper and embedded [19], [20], [21], [22], [23], [24], [25], [26] methods.

In recent years, powerful new algorithms have been proposed which exploit the advantages of matrix decomposition techniques. Well known examples of such methods include nonnegative matrix factorization (NMF) [27], [28], principal component analysis (PCA) [29], [30] and singular value decomposition (SVD) [30], [31]. However, all of these are examples of feature extraction methods. A smaller body of literature has explored how the idea of matrix decomposition can also be applied to feature selection. Wang et al. [32] proposed subspace learning for unsupervised feature

selection via matrix factorization (MFFS). In [33], Wang et al. proposed unsupervised feature selection via maximum projection and minimum redundancy (MPMR). These two algorithms find a suitable feature subspace through matrix decomposition, and the feature subspace is then used to represent the original feature space. By exploiting the advantages of the matrix factorization technique, MFFS and MPMR can both achieve good performance. However, MFFS and MPMR ignore the underlying geometry information of the data itself. In contrast, this paper shows how such geometry information can be used to further improve the quality of feature selection.

A variety of literature has shown that the distributions of high dimensional data are often sparse. Such data contain a lot of local information, which is important for mining the internal structure of such data and improving the performance of nonlinear learning [34], [35]. Some manifold learning algorithms have been proposed to discover the underlying manifold structure of data, such as Locality Preserving Projection [36], local linear embedding (LLE) [37] and Laplacian Eigenmap [38]. By analyzing the manifold structure of the data set, we can use the underlying geometric information to improve the learning efficiency of the algorithm.

Spectral graph theory [39], [40] can be used to characterize the underlying manifold structure of the data. The spectral clustering method [39] exploits spectral graph theory to obtain good clustering performance. Based on nonnegative matrix factorization (NMF) [28], Cai et al. [41] proposed graph regularized nonnegative matrix factorization (GNMF), which uses the geometry information of the data itself to greatly improve performance. Compared with concept factorization (CF) [42], locally consistent concept factorization (LCCF) [43] shows better performance, because it is able to exploit the local structure of data. In recent years, new work [44], [45], [46] has shown that the manifold information of the data is not only distributed in the data space, but also in the feature space. In [44], Shang et al. proposed a graph dual regularization non-negative matrix factorization for co-clustering algorithm (DNMF). Ye et al. [46] proposed dual-graph regularized concept factorization clustering (GCF).

Some classification algorithms also use the spectral graph theory. Belkin et al. [47] proposed manifold regularization, a geometric framework for learning from labeled and unlabeled examples, which uses graph theory to exploit the manifold structure of the data. Experimental results show that this method can use unlabeled data effectively. Based on standard SVM, Chova et al. proposed

semi-supervised image classification with Laplacian support vector machines (LapSVM) [48], which uses the geometry information of both labeled and unlabeled samples by using the graph Laplacian. Some experimental evidence suggests that LapSVM can outperform conventional SVM. In [49], Yang et al. proposed the Laplacian twin parametric-margin support vector machine for semi-supervised classification (LTPMSVM), which overcomes the shortcomings of conventional methods which are unable to effectively handle unlabeled data. LTPMSVM uses the geometric information of the unlabeled data to construct a better classifier, and experimental results have confirmed the strong performance of LTPMSVM.

Some feature selection algorithms which use local structure information have previously been proposed. Laplacian score (LapScor) [21], spectral feature selection (SPEC) [18], minimum redundancy spectral feature selection (MRSF) [50], unsupervised feature selection for multicluster data (MCFS) [51] are four well known algorithms. Extensions of MRSF and MCFS include clustering-guided sparse structural learning for unsupervised feature selection (CGSSL) [52] and joint embedding learning and sparse regression (JELSR) [53].

In this paper, we propose a new method called subspace learning-based graph regularized feature selection (SGFS). SGFS is based on the framework of subspace learning feature selection, which exploits the advantages of matrix factorization techniques. On this basis, we introduce the concept of graph regularization and preserve the local structure information of the feature space of the data. The local structure information of the feature space directly guides the learning of the coefficient matrix in the error reconstruction term, and indirectly guides the learning of the feature selection matrix. Additionally, we propose the use of the $L_{2,1}$ -norm to constrain the feature selection matrix, which guarantees its sparsity, so as to provide more accurate discrimination information for feature selection. We use an alternating iterative optimization mechanism to optimize the objective function and adjust the parameters to minimize the reconstruction error. Finally, we obtain the feature selection matrix. Through this matrix, we can calculate the scores of all the features, and select the most representative features.

The main contributions of this paper are as follows:

1. By using graph theory, the geometric structure information of the feature manifold is preserved. Through the guidance of geometry information, the learning of the feature selection matrix and

coefficient matrix are more rapid and accurate.

2. By introducing $L_{2,1}$ -norm to constrain the feature selection term, the sparsity of the feature selection matrix is guaranteed, enabling more accurate discrimination information for feature evaluation.

The structure of this paper is organized as follows. In Section 2, we introduce the framework, the iterative update rules and convergence proof of SGFS. In Section 3, we present the experimental results of comparing the performance of SGFS against five other state-of-the-art algorithms on twelve public benchmark data sets. Section 4 provides concluding remarks.

2. Subspace Learning-based Graph regularized Feature Selection

In this section, we present details of the SGFS method, which breaks down into three main parts: sparse subspace learning, local structure preserving and feature evaluation.

2.1 Related notations

First of all, we introduce some related notations. Denote $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ as the unlabeled sample data set. Where n and d respectively represent the number and dimension of the samples. We use l to indicate the number of selected features, $l \leq d$.

Given an arbitrary matrix $\mathbf{A} \in \mathbb{R}^{e \times f}$, its $L_{r,s}$ is defined as:

$$\|\mathbf{A}\|_{r,s} = (\sum_{i=1}^e (\sum_{j=1}^f |A_{ij}|^r)^{s/r})^{1/s}. \quad (1)$$

According to the definition, when $r=s=2$, it indicates Frobenius-norm or L_2 -norm. In contrast, when $r=2, s=1$, it represents sparse constraint $L_{2,1}$ -norm. We denote L_2 -norm and $L_{2,1}$ -norm respectively as $\|\cdot\|_2^2$ and $\|\cdot\|_{2,1}$ in the following.

2.2 Sparse subspace learning

2.2.1 Distance between subspaces

According to [32], we first define the distance between subspaces. Given a vector group \mathbf{X} in an m -dimensional real number space. We define $\text{span}(\mathbf{X}) = \{\mathbf{a}^T \mathbf{X} \mid \mathbf{a} \in \mathbb{R}^{|\mathbf{X}|}\}$ as the spanning subspace of \mathbf{X} , which is the set of all combinations of elements of \mathbf{X} . Where, $|\mathbf{X}|$ is the basis of \mathbf{X} . Given two

vector groups \mathbf{X}_1 and \mathbf{X}_2 in an m -dimensional real number space, the directional distance between $\text{span}(\mathbf{X}_1)$ and $\text{span}(\mathbf{X}_2)$ can be defined as:

$$\vec{d}(\text{span}(\mathbf{X}_1), \text{span}(\mathbf{X}_2)) = \sum_{\mathbf{x} \in \mathbf{X}_1} d(\mathbf{x}, \text{span}(\mathbf{X}_2)) = \min_{\mathbf{H} \in \mathbb{R}^{|\mathbf{X}_2| \times |\mathbf{X}_1|}} \|\mathbf{X}_1 - \mathbf{X}_2 \mathbf{H}\|_2^2. \quad (2)$$

where, $\mathbf{H} \in \mathbb{R}^{l \times d}$ is the coefficient matrix, which is used to calculate the directional distance between $\text{span}(\mathbf{X}_1)$ and $\text{span}(\mathbf{X}_2)$.

2.2.2 Sparse subspace learning

The main purpose of subspace learning is to find a suitable feature subspace for representing the original feature space, referred to as the process of feature selection [32]. Denote $\mathbf{X} \in \mathbb{R}^{d \times n}$, $\mathbf{X}_I \in \mathbb{R}^{l \times n}$. Where \mathbf{X}_I is the subspace of \mathbf{X} . According to the definition of equation (2), the problem of subspace learning or feature selection can be considered as solving the following problems:

$$\begin{aligned} \arg \min_I \vec{d}(\text{span}(\mathbf{X}^T), \text{span}(\mathbf{X}_I^T)) \\ \text{s.t. } |I| = l. \end{aligned} \quad (3)$$

where, I is the index set of selected features, $|I|$ represents the number of elements in set I . To minimize the distance of the two spaces, a suitable feature subset is obtained.

According to equation (2), formula (3) can be rewritten as:

$$\begin{aligned} \arg \min_I \|\mathbf{X}^T - \mathbf{X}_I^T \mathbf{H}\|_2^2 \\ \text{s.t. } |I| = l. \end{aligned} \quad (4)$$

where, $\mathbf{H} \in \mathbb{R}^{l \times d}$ is the coefficient matrix, which is used for data reconstruction. By means of the coefficient matrix \mathbf{H} , a new element of data can be approximately reconstructed as a linear combination of the features of the feature subset. The feature subset and the coefficient matrix can be obtained simultaneously by minimizing the reconstruction error. According to the idea of matrix factorization [32], equation (4) can be rewritten as:

$$\begin{aligned} \arg \min_{\mathbf{W}, \mathbf{H}} \|\mathbf{X}^T - \mathbf{X}^T \mathbf{W} \mathbf{H}\|_2^2 \\ \text{s.t. } \mathbf{W}, \mathbf{H} \geq 0, \mathbf{W}^T \mathbf{W} = \mathbf{I}_I. \end{aligned} \quad (5)$$

where, $\mathbf{W} \in \mathbb{R}^{d \times l}$ is the feature selection matrix, and $\mathbf{I}_I \in \mathbb{R}^{l \times l}$ is an identity matrix. Each row or column of \mathbf{W} contains no more than one non-zero element. It is therefore an indication of the feature

selection, i.e. $\mathbf{X}_I^T = \mathbf{X}^T \mathbf{W}$. The definition of \mathbf{W} is as follows:

$$W_{i,j} = \begin{cases} 1, & \text{the } j\text{th element of } I \text{ is } i, \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

In the objective function, the non-negative constraint and the orthogonal constraint make \mathbf{W} more close to the matrix containing 0-1 elements in the learning process. However, in the actual learning process, this process is difficult. Therefore we introduce the sparse constraint to increase the sparsity of matrix \mathbf{W} .

We use $L_{2,1}$ -norm to constrain \mathbf{W} . It can ensure the sparsity of \mathbf{W} and help avoid trivial solutions, which makes the feature selection more discriminative. Thus, the objective function (5) can be rewritten as:

$$\begin{aligned} \arg \min_{\mathbf{W}, \mathbf{H}} & \|\mathbf{X}^T - \mathbf{X}^T \mathbf{W} \mathbf{H}\|_2^2 + \beta \|\mathbf{W}\|_{2,1} \\ \text{s.t. } & \mathbf{W}, \mathbf{H} \geq 0, \mathbf{W}^T \mathbf{W} = \mathbf{I}_I. \end{aligned} \quad (7)$$

where $\|\mathbf{W}\|_{2,1} = \sum_{i=1}^d \|\mathbf{w}_i\|_2$, \mathbf{w}_i is the i -th row of the matrix \mathbf{W} , β is balance parameter.

2.3 Local structure preserving

Some studies [44], [45], [46] show that the feature manifold contains underlying geometric structure information, which is very useful for improving the performance of the algorithm. We use spectral theory [46], [44] to preserve the local structure information on the feature manifold, which can improve the efficiency and accuracy of feature selection.

According to the method of constructing graphs in [44], we construct the nearest neighbor graph in feature space. We use a row vector \mathbf{f}_i to represent the i -th feature of the sample matrix, so that the sample matrix \mathbf{X} can be rewritten as $\mathbf{X} = [\mathbf{f}_1; \mathbf{f}_2; \dots; \mathbf{f}_d] \in \mathbb{R}^{d \times n}$. We construct a k -nearest neighborhood graph $\mathbf{G}=(V, E)$, where V represents a set of feature points $\{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_d\}$. E represents the weights of the edges between the vertices. It represents the similarity between the two features: the higher the weight, the more similar the features. In this paper, we adopt the Gaussian function as the weight measurement, and its definition is as follows:

$$[S]_{ij} = \begin{cases} \exp(-\|f_i - f_j\|_2^2 / \sigma^2), & \text{if } f_i \in N(f_j) \\ & \text{or } f_j \in N(f_i) \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

where, $i, j=1,2,\dots,d$, $N(f_i)$ represents the k -nearest neighbor set of feature f_i , $[S]_{ij}$ represents the similarity between f_i and f_j . σ is the Gaussian scale parameter. L is the graph Laplacian matrix of the feature space, $L=D-S$, D is a diagonal matrix, and $D_{ii} = \sum_j [S]_{ij}$.

Denote $H = [h_1, h_2, \dots, h_d] \in \mathbb{R}^{l \times d}$ as the coefficient matrix, the objection of local structure preserving is as follows:

$$\arg \min_H \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^d \|h_i - h_j\|_2^2 [S]_{ij} = \text{Tr}(HLH^T). \quad (9)$$

If f_i and f_j are very similar, h_i and h_j are also very close to each other. Therefore, we combine equations (7) and (9), and obtain a new objective function as follows:

$$\arg \min_{H, W} \text{Tr}(HLH^T) + \alpha \|X^T - X^T WH\|_2^2 + \beta \|W\|_{2,1} + \frac{\lambda}{2} \|W^T W - I_l\|_2^2 \quad (10)$$

$s.t. W \geq 0, H \geq 0.$

where, α, β, λ are balance parameters, $\alpha, \beta, \lambda \geq 0$.

2.4 Feature evaluation

By optimizing the objective function of SGFS, we can obtain a feature selection matrix W . Where $W=[w_1; w_2; \dots; w_d]$, Using $\|w_i\|_2$ to calculate the value of each row of matrix W , these values represent the importance of the corresponding features. The higher the value, the more important the feature.

We arrange all the evaluation values of features in descending order, obtaining an index I . According to the index I , we select the first l features to represent the original data set, generating a new data matrix $X_{new} \in \mathbb{R}^{l \times n}$. In this way, the feature selection process is completed.

2.5 Connection with MFFS

By observing the objective function (10), we can see that when $\alpha=1, \beta=0$, and removing the local structure preserving regularization term, SGFS degenerates into MFFS. MFFS is mainly used to solve the following problem:

$$\begin{aligned} \arg \min_{\mathbf{H}, \mathbf{W}} & \|\mathbf{X}^T - \mathbf{X}^T \mathbf{W} \mathbf{H}\|_2^2 + \frac{\lambda}{2} \|\mathbf{W}^T \mathbf{W} - \mathbf{I}_l\|_2^2 \\ \text{s.t. } & \mathbf{W} \geq 0, \mathbf{H} \geq 0. \end{aligned} \quad (11)$$

2.6 Update rules for SGFS

We use an alternating iterative method [44] to optimize the objective function (10). We first introduce two Lagrange multipliers ψ_{ij} and ϕ_{ij} to constraints $\mathbf{W}_{ij} \geq 0$ and $\mathbf{H}_{ij} \geq 0$ respectively. Therefore, formula (10) can be rewritten as Lagrange's function:

$$\begin{aligned} \mathbf{L}(\mathbf{W}, \mathbf{H}) = & \text{Tr}(\mathbf{H} \mathbf{L} \mathbf{H}^T) + \alpha \|\mathbf{X}^T - \mathbf{X}^T \mathbf{W} \mathbf{H}\|_2^2 + \beta \|\mathbf{W}\|_{2,1} \\ & + \frac{\lambda}{2} \|\mathbf{W}^T \mathbf{W} - \mathbf{I}_l\|_2^2 + \text{Tr}(\psi \mathbf{W}^T) + \text{Tr}(\phi \mathbf{H}^T). \end{aligned} \quad (12)$$

First, we need to define a diagonal matrix $\mathbf{U} \in \mathbb{R}^{d \times d}$. The i -th diagonal element of \mathbf{U} is defined as follows:

$$\mathbf{U}_{ii} = \frac{1}{2\|\mathbf{w}_i\|_2}. \quad (13)$$

On this basis, we introduce a small enough constant ε to avoid overflow, and get the following formula:

$$\mathbf{U}_{ii} = \frac{1}{2\max(\|\mathbf{w}_i\|_2, \varepsilon)}. \quad (14)$$

Due to the definition of \mathbf{U} , we rewrite $\|\mathbf{W}\|_{2,1}$ as $\text{Tr}(\mathbf{W}^T \mathbf{U} \mathbf{W})$, so formula (12) can be rewritten as:

$$\begin{aligned} \mathbf{L}(\mathbf{W}, \mathbf{H}) = & \text{Tr}(\mathbf{H} \mathbf{L} \mathbf{H}^T) + \alpha \text{Tr}((\mathbf{X}^T - \mathbf{X}^T \mathbf{W} \mathbf{H})(\mathbf{X}^T - \mathbf{X}^T \mathbf{W} \mathbf{H})^T) + \beta \text{Tr}(\mathbf{W}^T \mathbf{U} \mathbf{W}) \\ & + \frac{\lambda}{2} \text{Tr}((\mathbf{W}^T \mathbf{W} - \mathbf{I}_l)(\mathbf{W}^T \mathbf{W} - \mathbf{I}_l)^T) + \text{Tr}(\psi \mathbf{W}^T) + \text{Tr}(\phi \mathbf{H}^T). \end{aligned} \quad (15)$$

Next, we fix \mathbf{H} and \mathbf{U} , and update \mathbf{W} . By taking the partial derivative of formula (15) with respect to \mathbf{W} , we arrive at:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{W}} = 2\alpha(\mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} \mathbf{H}^T - \mathbf{X} \mathbf{X}^T \mathbf{H}^T) + 2\beta \mathbf{U} \mathbf{W} + 2\lambda \mathbf{W} \mathbf{W}^T \mathbf{W} - 2\lambda \mathbf{W} + \psi. \quad (16)$$

Using the KKT conditions [42] $\psi_{ij} \mathbf{W}_{ij} = 0$, we get:

$$[\alpha(\mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} \mathbf{H}^T - \mathbf{X} \mathbf{X}^T \mathbf{H}^T) + \beta \mathbf{U} \mathbf{W} + \lambda \mathbf{W} \mathbf{W}^T \mathbf{W} - \lambda \mathbf{W}]_{ij} \mathbf{W}_{ij} = 0. \quad (17)$$

Hence, we get the iterative update rule for \mathbf{W} as follows:

$$\mathbf{W}_{ij} \leftarrow \mathbf{W}_{ij} \frac{[\alpha \mathbf{X} \mathbf{X}^T \mathbf{H}^T + \lambda \mathbf{W}]_{ij}}{[\alpha \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} \mathbf{H}^T + \beta \mathbf{U} \mathbf{W} + \lambda \mathbf{W} \mathbf{W}^T \mathbf{W}]_{ij}}. \quad (18)$$

Finally, we fix \mathbf{W} and \mathbf{U} , and update \mathbf{H} . By taking the partial derivatives of formula (15) with respect to \mathbf{H} , we arrive at:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{H}} = 2\alpha(\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} - \mathbf{W}^T \mathbf{X} \mathbf{X}^T) + 2\mathbf{H}(\mathbf{D} - \mathbf{S}) + \phi. \quad (19)$$

Using the KKT conditions $\phi_{ij} \mathbf{H}_{ij} = 0$, we get:

$$[\alpha(\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} - \mathbf{W}^T \mathbf{X} \mathbf{X}^T) + \mathbf{H}(\mathbf{D} - \mathbf{S})]_{ij} \mathbf{H}_{ij} = 0. \quad (20)$$

We get the iterative update rule for \mathbf{H} as follows:

$$\mathbf{H}_{ij} \leftarrow \mathbf{H}_{ij} \frac{[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T + \mathbf{H} \mathbf{S}]_{ij}}{[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} + \mathbf{H} \mathbf{D}]_{ij}}. \quad (21)$$

Table 1 shows the procedure of SGFS.

Table 1 The procedure of SGFS

Input: Data matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$; neighborhood size k ; balance parameter α, β, λ ; Gaussian scale parameter σ ; maximum number of iterations $NIter$; number of selected features l .

Output: Index of selected features I ; new data matrix $\mathbf{X}_{new} \in \mathbb{R}^{l \times n}$.

1. Construct the k -nearest neighborhood graph $G=(V, E)$ in feature space.
 2. Compute the similarity matrix \mathbf{S} , graph Laplacian matrix \mathbf{L} .
 3. Initialize $\mathbf{W}, \mathbf{H}, \mathbf{U}$.
 4. Update $\mathbf{W}, \mathbf{H}, \mathbf{U}$ according to the iterative update rules (14), (18) and (21), until the number of iterations is equal to $NIter$.
 5. Compute the evaluation values for all the features according to $\|\mathbf{w}_i\|_2$, select the features corresponding to the largest l values and get a new data matrix $\mathbf{X}_{new} \in \mathbb{R}^{l \times n}$.
-

2.7 Convergence analysis

Next, we analyze the convergence of SGFS. Based on the proof method in [54][32], we give the proof of the monotone property of formula (10) under the update rules (18) and (21).

First, we prove that formula (10) is nonincreasing under the update rule (21).

Definition 1. From Lee and Seung [54], if there is a function $G(k, k')$, which makes $F(k)$ satisfy the following conditions:

$$G(k, k') \geq F(k), G(k, k) = F(k) \quad (22)$$

then $F(k)$ is monotonically decreasing under the following updating formula:

$$k^{(t+1)} = \arg \min_h G(k, k^{(t)}) \quad (23)$$

Here $G(k, k')$ is an auxiliary function for $F(k)$.

Proof. $F(k^{(t+1)}) \leq G(k^{(t+1)}, k^{(t)}) \leq G(k^{(t)}, k^{(t)}) = F(k^{(t)})$.

We only retain the items with \mathbf{H} from the formula (10) and obtain the following function:

$$F(\mathbf{H}) = \text{Tr}(\mathbf{H}\mathbf{L}\mathbf{H}^T) + \alpha \text{Tr}(\mathbf{X}^T \mathbf{W} \mathbf{H} \mathbf{H}^T \mathbf{W}^T \mathbf{X} - 2\mathbf{X}^T \mathbf{H}^T \mathbf{W}^T \mathbf{X}). \quad (24)$$

Through taking the first-order and second-order partial derivatives of $F(\mathbf{H})$ with respect to \mathbf{H} , we arrive at:

$$F'_{ij} = \left[\frac{\partial F}{\partial \mathbf{H}} \right]_{ij} = [2\alpha(\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} - \mathbf{W}^T \mathbf{X} \mathbf{X}^T) + 2\mathbf{H}\mathbf{L}]_{ij} \quad (25)$$

$$F''_{ij} = 2\alpha[\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}]_{ii} + 2[\mathbf{L}]_{jj} \quad (26)$$

Lemma 1. From Lee and Seung [54], the auxiliary function of F_{ij} is given as follows:

$$G(\mathbf{H}_{ij}, \mathbf{H}_{ij}^{(t)}) = F_{ij}(\mathbf{H}_{ij}^{(t)}) + F'_{ij}(\mathbf{H}_{ij}^{(t)})(\mathbf{H}_{ij} - \mathbf{H}_{ij}^{(t)}) + \frac{[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} + \mathbf{H}\mathbf{D}]_{ij}}{\mathbf{H}_{ij}^{(t)}} (\mathbf{H}_{ij} - \mathbf{H}_{ij}^{(t)})^2 \quad (27)$$

Given the Taylor expansion of $F_{ij}(\mathbf{H}_{ij})$:

$$F_{ij}(\mathbf{H}_{ij}) = F_{ij}(\mathbf{H}_{ij}^{(t)}) + F'_{ij}(\mathbf{H}_{ij}^{(t)})(\mathbf{H}_{ij} - \mathbf{H}_{ij}^{(t)}) + \frac{1}{2} \{ \alpha[\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}]_{ii} + [\mathbf{L}]_{jj} \} (\mathbf{H}_{ij} - \mathbf{H}_{ij}^{(t)})^2 \quad (28)$$

According to formula (27), we know that $G(\mathbf{H}_{ij}, \mathbf{H}_{ij}^{(t)}) \geq F_{ij}(\mathbf{H}_{ij})$ is equivalent to

$$\frac{[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} + \mathbf{H}\mathbf{D}]_{ij}}{\mathbf{H}_{ij}^{(t)}} \geq \alpha[\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}]_{ii} + [\mathbf{L}]_{jj} \quad (29)$$

Clearly, it can be seen that $[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H}]_{ij} = \sum_{l=1}^d [\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}]_{il} \mathbf{H}_{lj}^{(t)} \geq \alpha[\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}]_{ii} \mathbf{H}_{ij}^{(t)}$

and $[\mathbf{H}\mathbf{D}]_{ij} = \sum_{l=1}^d \mathbf{H}_{il}^{(t)} [\mathbf{D}]_{lj} \geq \mathbf{H}_{ij}^{(t)} \mathbf{D}_{jj} \geq \mathbf{H}_{ij}^{(t)} [\mathbf{D} - \mathbf{W}]_{jj} = \mathbf{H}_{ij}^{(t)} [\mathbf{L}]_{jj}$. Thus, inequality (29) holds, i.e.

$G(\mathbf{H}_{ij}, \mathbf{H}_{ij}^{(t)}) \geq F_{ij}(\mathbf{H}_{ij})$ holds, $G(\mathbf{H}_{ij}, \mathbf{H}_{ij}) = F_{ij}(\mathbf{H}_{ij})$ also holds.

Proof. The variable \mathbf{H} satisfies the updating formula (23) that makes the F_{ij} monotonically decreasing.

By substituting $G(\mathbf{H}_{ij}, \mathbf{H}_{ij}^{(t)})$ in (27) into (23), we get the following formula:

$$\mathbf{H}_{ij}^{(t+1)} = \mathbf{H}_{ij}^{(t)} - \mathbf{H}_{ij}^{(t)} \frac{F'_{ij}(\mathbf{H}_{ij}^{(t)})}{2[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} + \mathbf{H}\mathbf{D}]_{ij}} = \mathbf{H}_{ij}^{(t)} \frac{[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T + \mathbf{H}\mathbf{S}]_{ij}}{[\alpha \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} \mathbf{H} + \mathbf{H}\mathbf{D}]_{ij}}$$

We can see that the above formula is the update rule (21). Therefore F_{ij} is non-increasing under the update rule (21).

The proof of the monotone property of formula (10) under the update rule (18) is similar to the above proof, and we can prove that F_{ij} is non-increasing under the update rule (18). Therefore, we conclude that formula (10) is also non-increasing under the update rules (18) and (21).

3. Experiments

In this section, we present the results and analysis of empirical experiments. We apply the proposed algorithm, and five comparison algorithms, to twelve public benchmark data sets, and the experimental results are analyzed. We use a clustering algorithm to cluster the feature selection results, and the clustering results are used as criteria for evaluating the performance of the feature selection algorithms. In this paper, we use *k-means* [55] as the clustering algorithm. In addition, we also analyze the parameter sensitivity of SGFS.

3.1 Data sets

In this experiment, we used twelve datasets, including digital image¹, text image¹, face image¹ and biological data¹ [32], [56]. The detailed information of the datasets is shown in table 2.

Table 2 The information of twelve datasets

| Data set | Size | Dim | Classes | Type |
|----------|------|-------|---------|---------------|
| Usps | 9298 | 256 | 10 | Digital image |
| Lung_dis | 73 | 325 | 7 | Biological |
| Isolet | 1560 | 617 | 26 | Letter image |
| Umist | 575 | 644 | 20 | Face image |
| COIL20 | 1440 | 1024 | 20 | Digital image |
| AR10P | 130 | 2400 | 10 | Face image |
| Lung | 203 | 3312 | 5 | Biological |
| Yale64 | 165 | 4096 | 15 | Face image |
| Orl64 | 400 | 4096 | 40 | Face image |
| TOX_171 | 171 | 5748 | 4 | Biological |
| Orlraws | 100 | 10304 | 10 | Face image |
| AT&T | 400 | 10304 | 40 | Face image |

¹ <http://featureselection.asu.edu/datasets.php>

3.2 Comparison algorithms

To verify the effectiveness of SGFS, we compare its performance against five feature selection algorithms, as follows:

1. Baseline: all the features of the data set are selected.
2. LapScor: Laplacian Score [21] is a classic unsupervised feature selection algorithm, it is characterized by being simple and fast, but it lacks an effective learning mechanism.
3. MCFS: unsupervised feature selection for multicluster data [51] uses a spectral embedding learning and sparse regression feature selection framework. It uses a two step strategy, and mainly solves the following problems:

$$\begin{aligned} \arg \min_{SS^T=I_m} Tr(SLS^T) \\ \arg \min_P \|P^T X - S\|_2^2 + \alpha \|P\|_1. \end{aligned} \quad (30)$$

4. UDFS: unsupervised discriminant feature selection algorithm [57] uses the discriminant information and feature correlation, which aims to find the most discriminative features. Its objective function is as follows:

$$\arg \min_{P^T P=I_m} Tr(P^T X L X^T P) + \alpha \|P\|_{2,1}. \quad (31)$$

5. MFFS: subspace learning for unsupervised feature selection via matrix factorization [32] uses a matrix decomposition technique to obtain a feature selection matrix. According to this matrix, a representative feature subspace can be found in the original feature space.

3.3 Evaluation metrics

By evaluating the clustering results of different feature selection algorithms, their performance can be compared. In this paper, we use Clustering Accuracy (ACC) [46][58][59] and Normalized Mutual Information (NMI) [46][60] to evaluate the clustering results. The higher the values of ACC and NMI, the better the performance of the algorithm.

NMI is defined as:

$$NMI(R, T) = \frac{I(R, T)}{\sqrt{H(R)H(T)}} \quad (32)$$

where R and T are two arbitrary variables, and $I(R, T)$ is the mutual information between R and T .

$H(R)$ and $H(T)$ are the information entropies of R and T respectively. When NMI is used to evaluate the clustering results, R and T represent the clustering label and the ground truth label respectively.

ACC is defined as:

$$ACC = \frac{1}{n} \sum_{i=1}^n \delta(a_i, \text{map}(b_i)) \quad (33)$$

where n represents the total number of samples, a_i and b_i are the clustering label and the ground truth label for the sample x_i respectively. $\delta(x, y)=1$, if $x=y$; $\delta(x, y)=0$, otherwise. $\text{map}(\cdot)$ is an optimal mapping function which uses Hungarian [61] to match the clustering label and the ground truth label.

3.4 Experimental results and analysis

3.4.1 Experimental setting

Before the experiment, we set up the parameters of the algorithm. For all datasets, the number of selected features l are tuned from $\{20, 30, 40, 50, 60, 70, 80, 90, 100\}$. For LapScor, MCFS, UDFS and SGFS, the nearest neighborhood parameter k is set to 5, the Gaussian scale parameter σ is set to 10. For MCFS, UDFS, MFFS and SGFS, we set the maximum number of iterations to 30. For SGFS, the balance parameter α and β are searched in the range of $\{10^{-7}, 10^{-6}, \dots, 10^{+7}\}$, the range of parameter λ is set to $\{10^{-6}, 10^{-5}, 10^{-3}, 10^{-2}, 10^{+4}, 10^{+7}, 10^{+8}\}$. On each dataset, the selection of parameters is relatively stable. By adjusting the balance parameters α, β, λ , we obtain the maximum value of ACC and NMI. We independently run 20 times, and take the average of these results as the final result. In the experiment, we used a computer with 4G memory and 2.3GHz frequency, and used Matlab as a simulation software.

3.4.2 Convergence test

First, we verify the convergence of the objective function of SGFS. We let SGFS iterate 30 times on twelve datasets, and recorded the value of the objective function after each iteration. The results are shown in **Fig. 1**.

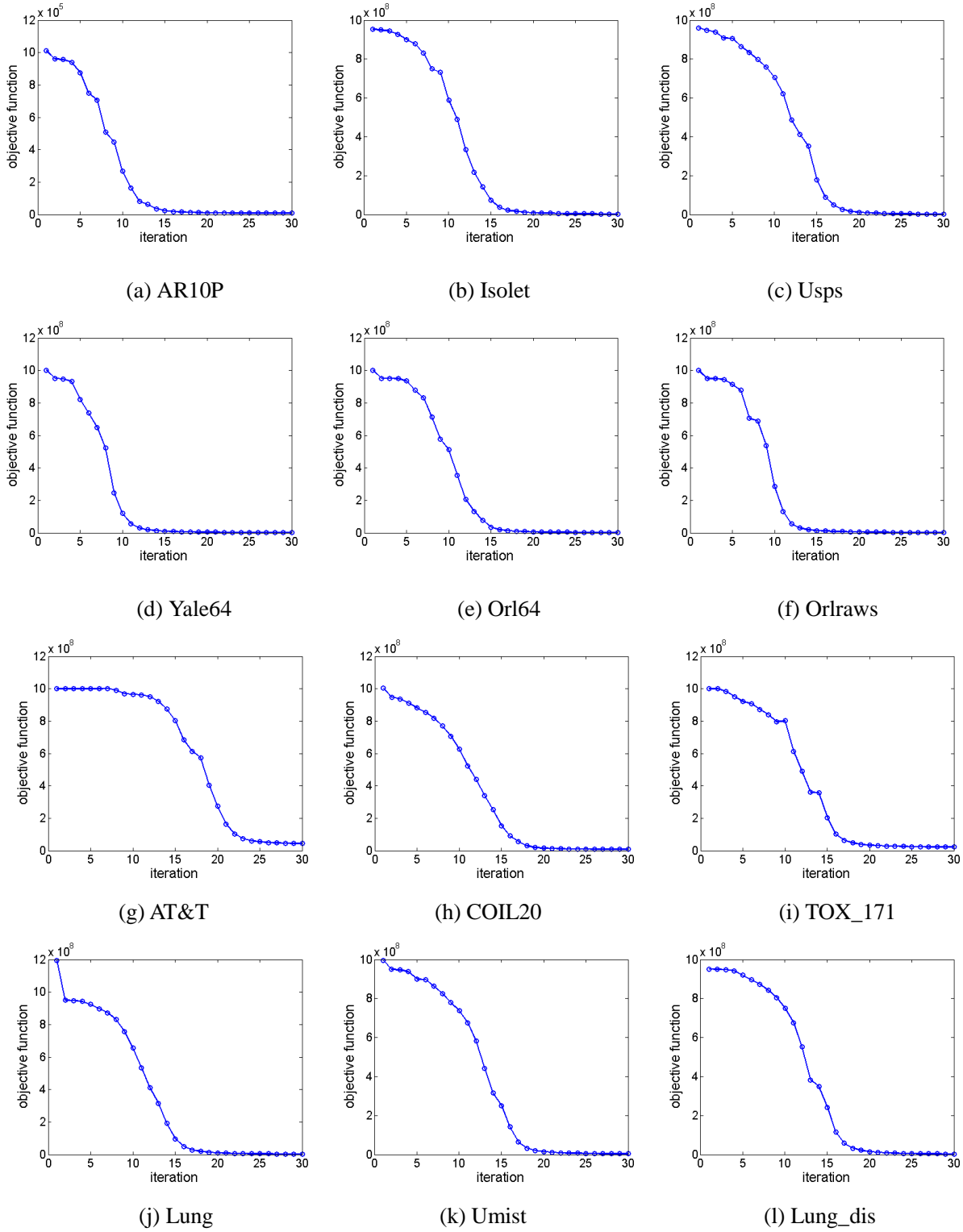


Fig. 1. Convergence of the objective function on twelve data sets.

In **Fig. 1**, the horizontal and vertical axes represent the number of iterations and the value of the objective function respectively. We can see that the value of the objective function decreases with increasing number of iterations. This is consistent with our previous proof of convergence. Also note that the value of the objective function tends to become stable within 20 iterations on most of the

datasets, which suggests that the objective function rapidly converges to a stable value.

In order to prove that the proposed algorithm is faster and more accurate than MFFS in the learning of the feature selection matrix and coefficient matrix, we analyze the number of iterations and the time required for the convergence of each objective function. Considering that the time spent in each iteration of SGFS and MFFS is almost the same, we only analyze the number of iterations required to achieve convergence of the objective function. For MFFS and SGFS, we set the maximum number of iterations to 30, the number of selected features l is fixed to 20, the range of parameter λ is set to $\{10^{-6}, 10^{-5}, 10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4, 10^5, 10^6, 10^7, 10^8\}$. For SGFS, the balance parameters α and β are searched in the range of $\{10^{-7}, 10^{-6}, \dots, 10^7\}$. The comparison results are given in Table 3.

Table 3 The number of iterations required for the convergence of the objective function

| Dataset | AR10P | Isolet | Usps | Yale64 | Orl64 | Orlraws |
|---------|-------|--------|---------|--------|-------|----------|
| MFFS | 25 | 25 | 24 | 28 | 25 | 29 |
| SGFS | 16 | 19 | 20 | 14 | 17 | 15 |
| Dataset | AT&T | COIL20 | TOX_171 | Umist | Lung | Lung_dis |
| MFFS | 27 | 27 | 22 | 24 | 25 | 23 |
| SGFS | 25 | 19 | 20 | 17 | 20 | 20 |

From Table 3, we can see that, for the objective function to converge, SGFS always requires significantly less iterations than MFFS. This suggests that SGFS can learn appropriate feature selection and coefficient matrices more quickly.

3.4.3 AT&T face dataset example

We randomly selected two different samples from the AT&T face dataset as the experimental samples, and apply the proposed algorithm to these two samples. In this experiment, the two samples were selected from the third class and the sixth class, they all contain 10304 features. The number of selected features l are tuned from $\{1280, 2560, 3840, 5120, 6400, 7680, 8960, 10240\}$. With the increase in the number of selected features, the extracted face information is also increased. We use white to represent the feature which is not selected, and the selected feature retains its original gray. The experimental results are shown in **Fig. 2**.



Fig. 2. Results of two AT&T samples under different number of selected features.

From **Fig. 2**, we can see how SGFS selects different numbers of features from the two samples. From left to right, as the number of selected features gradually increases, the reconstructed image becomes increasingly clear. We can see that a convincing outline of the human face can be seen by selecting only a very small number of features. The reason is that the most representative face features are preferentially selected, such as the eyes, nose, mouth and chin. This preferential selection of the most salient features illustrates the strong performance of the proposed algorithm for the task of feature selection.

3.4.4 Experimental results and analysis

In Table 4 and Table 5, we show the ACC and NMI values of the six algorithms on all data sets. In these results, the best ACC and NMI values are marked in bold black, and the second best ACC and NMI values are marked in underlined.

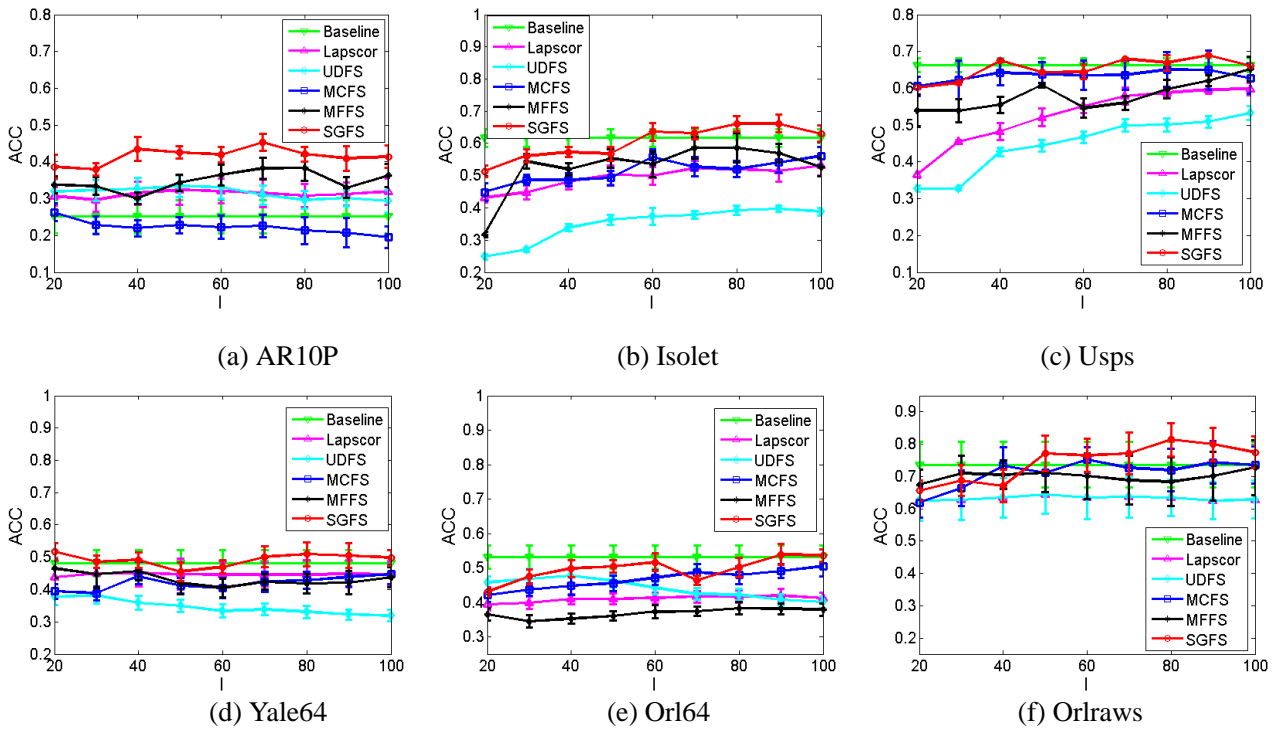
Table 4 Clustering accuracy of six algorithms on twelve datasets (ACC \pm STD%)

| Dataset | Baseline | LapScor | UDFS | MCFS | MFFS | SGFS |
|----------|----------------------------------|------------------|------------------|----------------------------------|----------------------------------|----------------------------------|
| AR10P | 25.12 \pm 4.47 | 32.42 \pm 4.00 | 33.38 \pm 2.96 | 26.31 \pm 2.45 | <u>38.31\pm3.40</u> | 45.38\pm2.34 |
| Isolet | <u>61.73\pm2.77</u> | 53.17 \pm 3.11 | 39.66 \pm 0.94 | 56.11 \pm 2.70 | 58.68 \pm 4.60 | 66.20\pm2.67 |
| Usps | <u>66.26\pm1.93</u> | 59.79 \pm 1.62 | 53.30 \pm 1.91 | 65.06 \pm 4.75 | 65.25 \pm 3.30 | 68.97\pm0.17 |
| Yale64 | 48.15 \pm 4.10 | 45.09 \pm 4.15 | 38.18 \pm 2.59 | 44.70 \pm 3.01 | <u>50.52\pm3.87</u> | 51.70\pm2.56 |
| Orl64 | <u>53.14\pm3.33</u> | 42.04 \pm 1.84 | 47.71 \pm 2.38 | 50.59 \pm 2.96 | 38.36 \pm 1.85 | 53.95\pm3.08 |
| Orlraws | 73.65 \pm 7.06 | 64.45 \pm 5.97 | 64.45 \pm 5.97 | <u>75.15\pm3.95</u> | 72.80 \pm 8.46 | 81.40\pm5.15 |
| AT&T | 60.96\pm3.30 | 47.31 \pm 1.83 | 54.46 \pm 2.05 | 55.88 \pm 1.86 | 53.46 \pm 2.73 | <u>59.65\pm2.89</u> |
| COIL20 | 65.75\pm4.16 | 60.94 \pm 2.17 | 63.38 \pm 3.04 | 62.96 \pm 2.82 | 62.00 \pm 2.40 | <u>65.61\pm2.67</u> |
| TOX_171 | <u>43.13\pm2.35</u> | 41.84 \pm 1.98 | 41.35 \pm 0.57 | 40.18 \pm 4.35 | 38.45 \pm 1.31 | 44.82\pm1.36 |
| Umist | 42.97 \pm 2.21 | 42.09 \pm 1.94 | 44.83 \pm 1.52 | <u>47.24\pm2.25</u> | 44.07 \pm 2.38 | 52.36\pm2.61 |
| Lung | 70.10 \pm 8.22 | 66.67 \pm 2.12 | 55.44 \pm 4.75 | <u>71.75\pm2.04</u> | 62.82 \pm 4.74 | 84.70\pm0.87 |
| Lung_dis | 73.63 \pm 5.26 | 62.81 \pm 5.59 | 70.07 \pm 5.83 | 74.25 \pm 4.84 | <u>74.52\pm4.86</u> | 81.03\pm4.34 |

Table 5 Normalized Mutual Information of six algorithms on twelve datasets (NMI \pm STD%)

| Dataset | Baseline | LapScor | UDFS | MCFS | MFFS | SGFS |
|----------|----------------------------------|------------------|------------------|----------------------------------|----------------------------------|----------------------------------|
| AR10P | 21.42 \pm 5.62 | 32.81 \pm 2.42 | 32.00 \pm 3.22 | 22.95 \pm 2.93 | <u>39.63\pm2.89</u> | 47.98\pm2.60 |
| Isolet | <u>76.06\pm1.26</u> | 68.40 \pm 1.34 | 52.90 \pm 0.79 | 69.92 \pm 0.96 | 72.42 \pm 2.12 | 76.78\pm0.98 |
| Usps | <u>61.13\pm0.86</u> | 56.53 \pm 0.70 | 48.01 \pm 1.28 | 58.89 \pm 2.03 | 60.02 \pm 1.72 | 62.23\pm0.14 |
| Yale64 | 54.79 \pm 3.39 | 51.46 \pm 2.58 | 44.95 \pm 1.79 | 50.11 \pm 2.23 | 56.17\pm4.77 | <u>55.20\pm1.88</u> |
| Orl64 | 73.56\pm1.50 | 63.71 \pm 0.96 | 68.68 \pm 1.64 | 70.94 \pm 1.58 | 61.55 \pm 1.03 | <u>73.41\pm0.97</u> |
| Orlraws | 79.87 \pm 5.31 | 72.64 \pm 3.80 | 72.64 \pm 3.80 | <u>82.52\pm3.07</u> | 81.44 \pm 6.57 | 83.54\pm3.56 |
| AT&T | 79.96\pm1.37 | 71.04 \pm 0.93 | 73.90 \pm 1.32 | 74.52 \pm 1.02 | 73.73 \pm 1.21 | <u>76.83\pm1.27</u> |
| COIL20 | 76.69\pm1.99 | 69.94 \pm 2.05 | 72.43 \pm 1.31 | 73.94 \pm 1.44 | 71.93 \pm 2.19 | <u>74.76\pm1.49</u> |
| TOX_171 | 14.55\pm2.64 | 11.03 \pm 1.14 | 10.12 \pm 2.03 | 10.20 \pm 1.98 | 11.31 \pm 0.13 | <u>12.45\pm1.91</u> |
| Umist | 64.83 \pm 2.04 | 62.33 \pm 1.91 | 58.31 \pm 1.47 | <u>67.06\pm1.38</u> | 62.59 \pm 1.87 | 68.91\pm1.40 |
| Lung | <u>54.47\pm2.84</u> | 47.17 \pm 1.94 | 38.01 \pm 1.87 | 52.52 \pm 0.63 | 44.53 \pm 1.14 | 58.23\pm1.01 |
| Lung_dis | 69.27 \pm 4.21 | 59.76 \pm 3.77 | 65.02 \pm 4.12 | <u>70.30\pm3.72</u> | 65.88 \pm 3.77 | 72.96\pm3.64 |

From Table 4 and Table 5, we can see that most of the results of SGFS on 12 datasets are better than the results of five comparison algorithms. In addition, the results of SGFS are even better than the results of Baseline on some datasets. **Fig. 3** shows the clustering accuracy of SGFS and five comparison algorithms on twelve datasets with different number of selected features. The horizontal axis represents the number of selected features l , the vertical axis represents clustering accuracy (ACC) and standard deviation (STD).



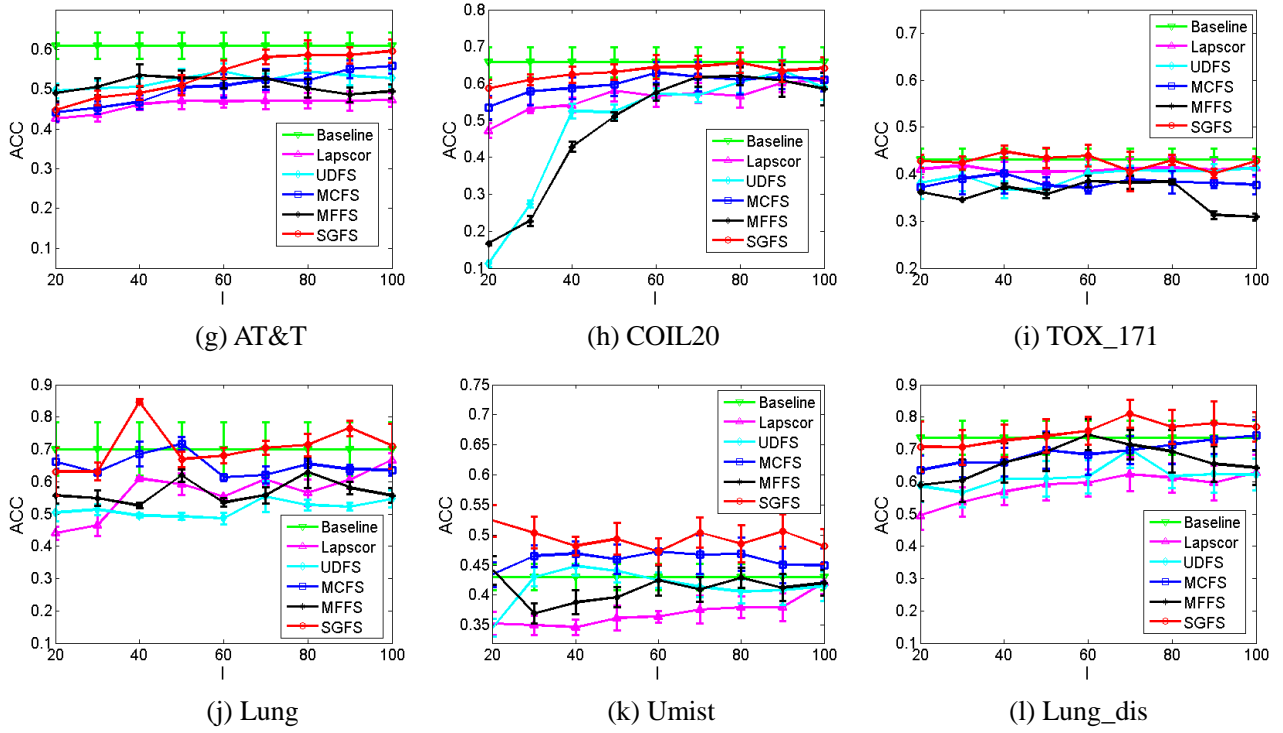
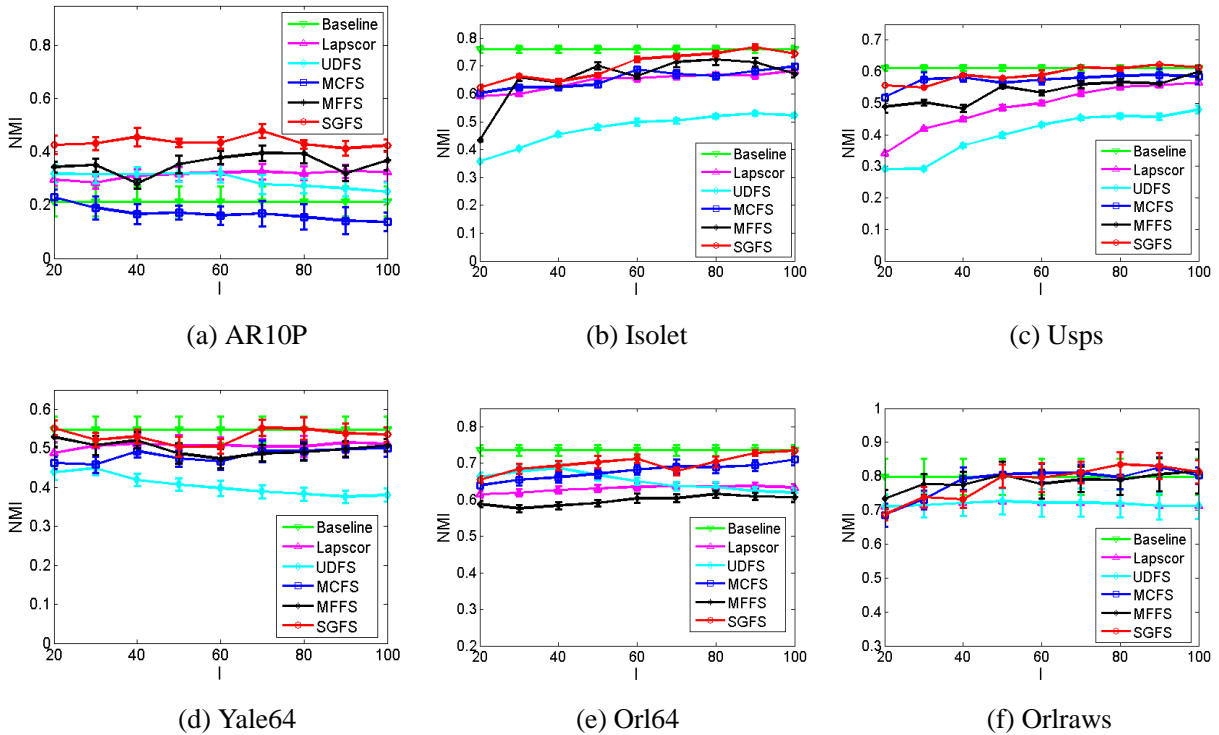


Fig. 3. Clustering accuracy of six algorithms on twelve datasets with different number of selected features.

From **Fig. 3**, we can see that the ACC of SGFS are higher than that of four comparison algorithms on all datasets, and in some cases even higher than the results of Baseline. It is worth mentioning that the results of SGFS are better than all the comparison algorithms on the AR10P dataset, which suggests that SGFS improves the clustering accuracy.



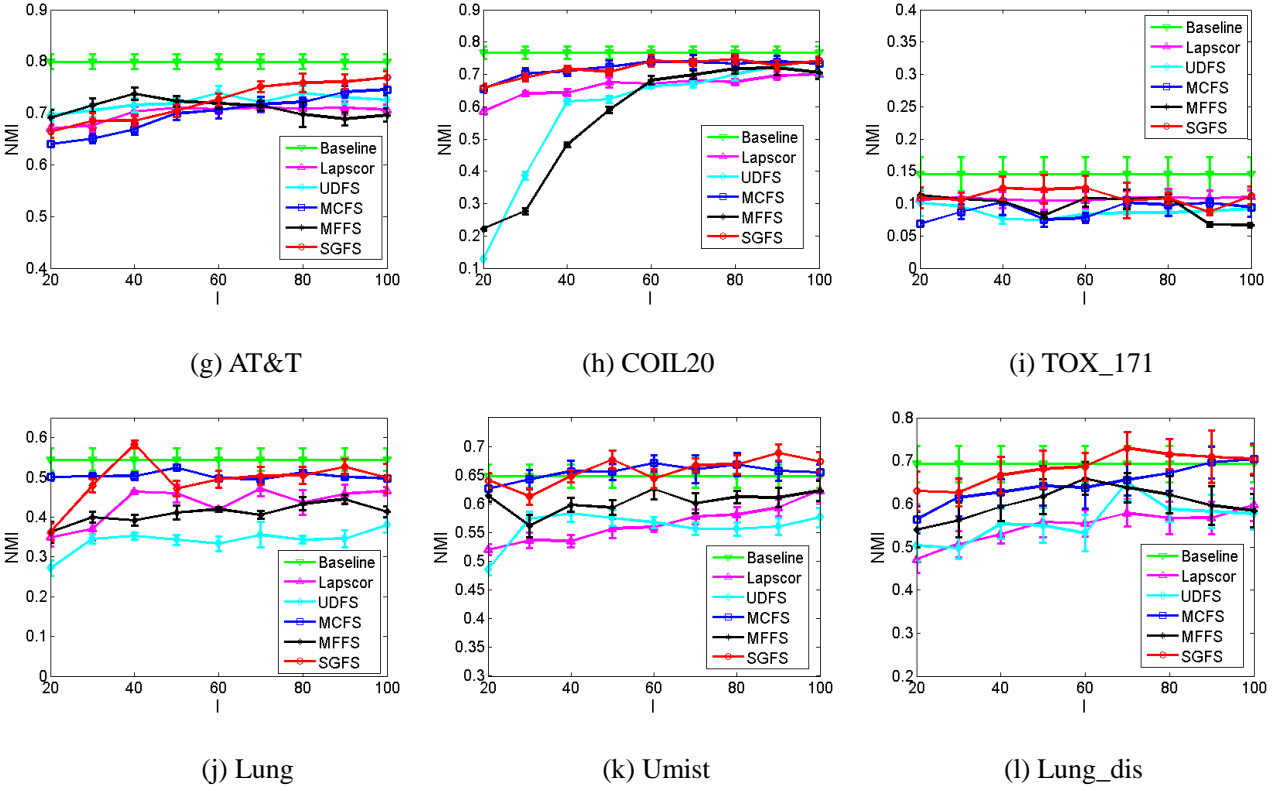


Fig. 4. Normalized Mutual Information of six algorithms on twelve datasets with different number of selected features.

Fig.4 shows the normalized mutual information of SGFS and five comparison algorithms on twelve datasets with different numbers of selected features. The horizontal coordinate also represents the number of selected features l , while the vertical coordinate represents normalized mutual information (NMI) and standard deviation (STD). We also can see that the NMI of SGFS are higher than that of four comparison algorithms on all datasets. On datasets Isolet, Usps, Yale64, Orl64 and Lung, only SGFS has results better than Baseline. It also shows that SGFS is more effective than the comparison algorithms.

3.4.5 Robustness test of algorithms

In order to verify the robustness of the proposed algorithm against noise, we artificially corrupted six datasets by adding Gaussian noise of various different magnitudes. In this experiment, we set the variance of the added Gaussian noise to 10, 20, 30. In Table 6 and Table 7, we show the ACC and NMI values of MFFS and SGFS on six test datasets.

Table 6 Clustering accuracy of SGFS and MFFS on six datasets with different variance of the Gaussian noise (ACC \pm STD%)

| Variance | 10 | | 20 | | 30 | |
|----------|------------------|------------------|------------------|------------------|------------------|------------------|
| Dataset | MFFS | SGFS | MFFS | SGFS | MFFS | SGFS |
| AR10P | 35.38 \pm 2.79 | 45.62 \pm 4.07 | 35.04 \pm 2.82 | 43.35 \pm 3.39 | 35.46 \pm 3.61 | 41.04 \pm 3.08 |
| Orl64 | 36.47 \pm 1.31 | 52.14 \pm 2.96 | 36.01 \pm 1.32 | 52.41 \pm 2.72 | 36.43 \pm 2.04 | 53.12 \pm 1.80 |
| AT&T | 53.81 \pm 2.48 | 58.80 \pm 2.81 | 53.31 \pm 1.77 | 57.98 \pm 2.26 | 53.04 \pm 2.46 | 58.27 \pm 2.72 |
| TOX_171 | 34.50 \pm 0.01 | 42.98 \pm 1.55 | 34.50 \pm 0.01 | 43.04 \pm 1.51 | 34.50 \pm 0.01 | 43.01 \pm 1.52 |
| Umist | 42.62 \pm 2.72 | 48.49 \pm 2.41 | 42.59 \pm 2.20 | 47.99 \pm 3.05 | 41.62 \pm 2.13 | 46.15 \pm 2.97 |
| Lung_dis | 69.59 \pm 7.16 | 77.19 \pm 3.78 | 68.08 \pm 6.92 | 76.51 \pm 6.79 | 68.22 \pm 6.87 | 76.16 \pm 6.61 |

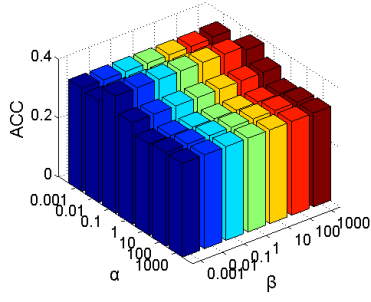
Table 7 Normalized Mutual Information of SGFS and MFFS on six datasets with different variance of the Gaussian noise (NMI \pm STD%)

| Variance | 10 | | 20 | | 30 | |
|----------|------------------|------------------|------------------|------------------|------------------|------------------|
| Dataset | MFFS | SGFS | MFFS | SGFS | MFFS | SGFS |
| AR10P | 37.75 \pm 2.97 | 46.75 \pm 4.18 | 37.07 \pm 2.94 | 43.83 \pm 3.65 | 37.73 \pm 3.36 | 43.51 \pm 3.82 |
| Orl64 | 59.56 \pm 1.31 | 72.35 \pm 1.40 | 59.01 \pm 1.11 | 71.68 \pm 1.50 | 59.11 \pm 1.22 | 72.25 \pm 0.75 |
| AT&T | 72.48 \pm 1.48 | 76.16 \pm 1.26 | 72.08 \pm 1.19 | 75.60 \pm 1.39 | 71.95 \pm 1.29 | 75.79 \pm 1.58 |
| TOX_171 | 7.96 \pm 0.01 | 10.83 \pm 1.68 | 7.96 \pm 0.01 | 10.84 \pm 1.67 | 7.96 \pm 0.01 | 10.84 \pm 1.67 |
| Umist | 62.31 \pm 1.69 | 67.55 \pm 1.65 | 62.11 \pm 1.18 | 66.46 \pm 1.71 | 61.11 \pm 1.67 | 64.88 \pm 1.34 |
| Lung_dis | 63.59 \pm 5.56 | 70.22 \pm 3.38 | 62.38 \pm 4.46 | 71.32 \pm 4.05 | 62.16 \pm 4.64 | 69.89 \pm 5.18 |

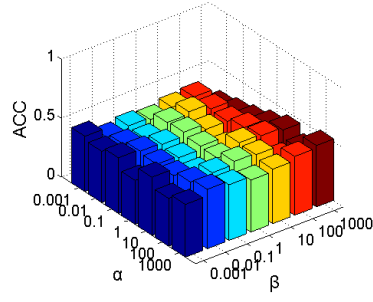
From Table 6 and Table 7, we can see that the performance of the proposed algorithm is affected very little by the added Gaussian noise. This suggests that the proposed algorithm is highly robust to noise. Furthermore, on each test data set, the results of SGFS are significantly better than those of MFFS.

3.4.6 Parameter sensitivity analysis

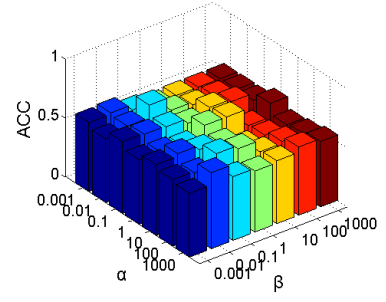
The parameters of SGFS mainly include: the neighborhood parameter k , the Gaussian scale parameter σ , and the balance parameters α , β , λ . In this paper, we only discuss the sensitivity of parameters α and β , since the other parameters are relatively stable. We explore the parameters α and β in the range of $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10^{+1}, 10^{+2}, 10^{+3}\}$, and record the clustering accuracy (ACC) and normalized mutual information (NMI) under different parameter combinations. We draw the results into 3-D diagrams, as shown in **Fig. 5** and **Fig 6**.



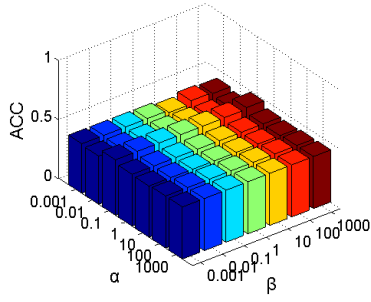
(a) AR10P



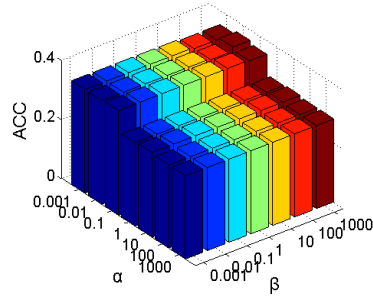
(b) Isolet



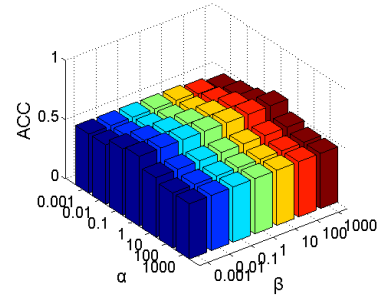
(c) Usps



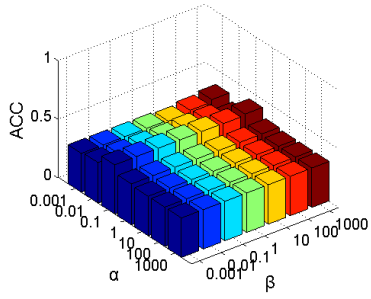
(d) Yale64



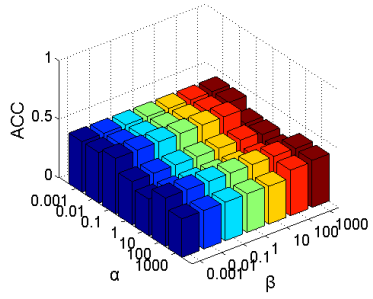
(e) Orl64



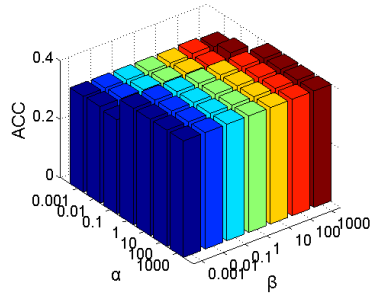
(f) Orlraws



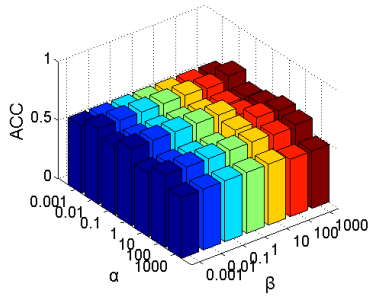
(g) AT&T



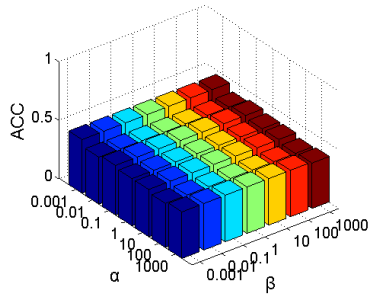
(h) COIL20



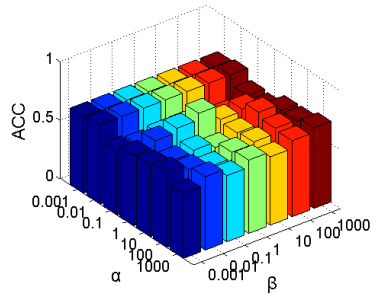
(i) TOX_171



(j) Lung



(k) Umist



(l) Lung_dis

Fig. 5. Clustering accuracy of SGFS on twelve datasets under different parameter combinations.

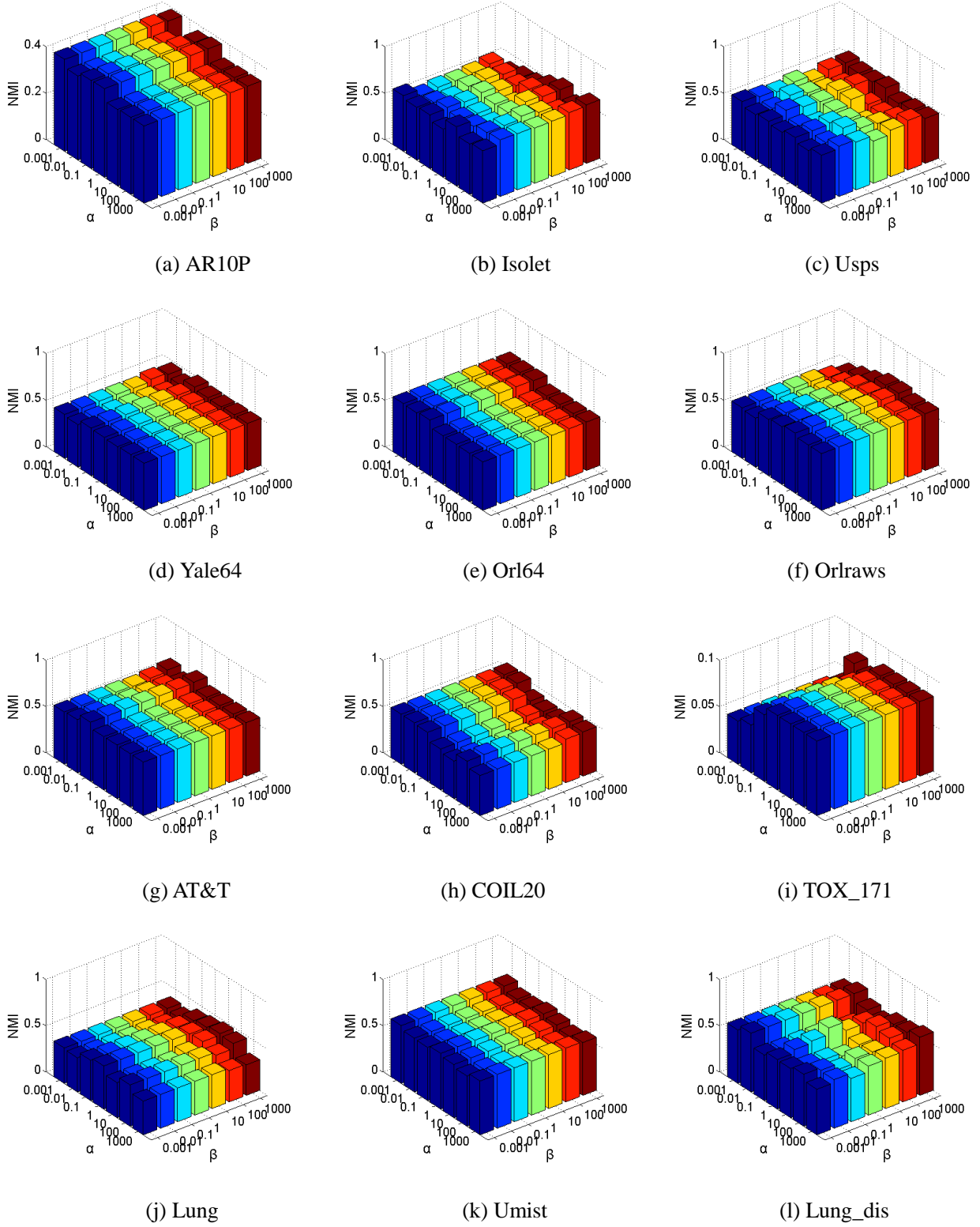


Fig. 6. Normalized Mutual Information of SGFS on twelve datasets under different parameter combinations.

From **Fig.5** and **Fig.6**, we can see that the parameters α and β on most datasets are relatively stable, especially on datasets Yale64, AT&T and Umist. This suggests that the proposed algorithm is not sensitive to the values of parameters α and β .

To further analyze the parameter sensitivity of the proposed algorithm, we compare the clustering results of SGFS under the default parameters and the best parameters. We set the default parameters α and β to a fixed value of 0.1. The best parameters are searched in the range of α and β . The comparison results are given in **Fig. 7**.

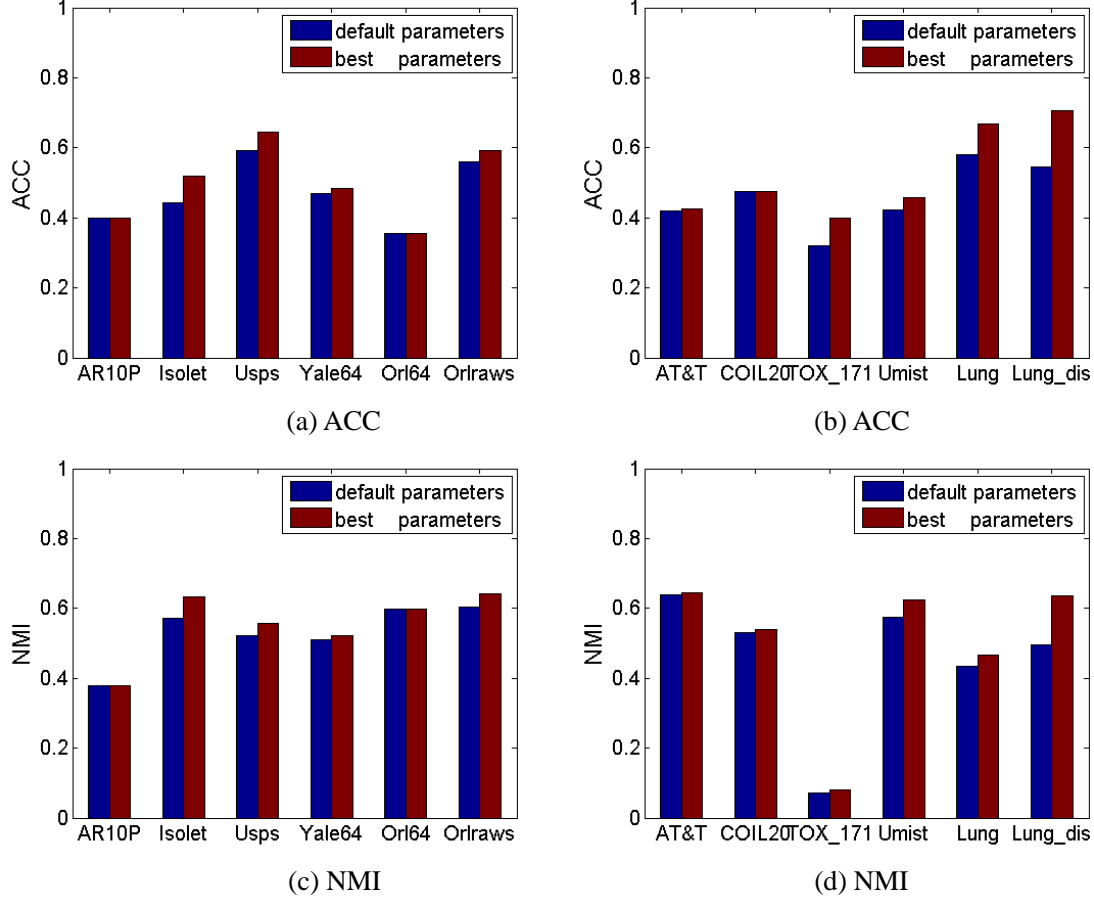


Fig. 7. Clustering results under the default parameters and the best parameters (default $\alpha=0.1$, $\beta=0.1$).

Fig. 7 shows that the clustering results on all datasets under the default parameters and the best parameters lead to only small differences in performance. On datasets AR10P and Orl64, the performance that results from default parameters and the best parameters is the same. These results suggest that overall performance is not strongly dependent on the choice of parameter values. Good performance can be achieved with arbitrary parameter values.

4. Conclusions

In this paper, we have proposed a novel algorithm, called subspace learning-based graph regularized feature selection (SGFS). The proposed algorithm is based on the framework of subspace

learning feature selection via matrix factorization (MFFS). We have extended this approach, by first introducing an $L_{2,1}$ -norm sparse constraint, which ensures the sparsity of the feature selection matrix \mathbf{W} and increases the discriminating ability of the selected features. Furthermore, we have shown how the idea of graph regularization can be incorporated into feature selection, to preserve the local structure information of the feature space. This structure is then used to guide the learning of the coefficient matrix \mathbf{H} and the feature selection matrix \mathbf{W} . We have presented a variety of experimental results, which suggest that SGFS outperforms several well known comparison algorithms from the literature.

A deficiency of the proposed algorithm is that the alternative iterative optimization method can sometimes be prone to converging on local optima. In future work, we hope to incorporate a mechanism for optimizing \mathbf{W} and \mathbf{H} simultaneously to obtain the globally optimal solution. We also intend to apply the idea of dual-graph regularization to the feature selection framework, and use the geometric structure information of the data space and the feature space simultaneously.

Acknowledgement

We would like to express our sincere appreciation to the editors and the anonymous reviewers for their insightful comments, which have greatly helped us in improving the quality of the paper. This work was partially supported by the National Natural Science Foundation of China, under Grants 61371201, the National Basic Research Program (973 Program) of China under Grant 2013CB329402, the Program for Cheung Kong Scholars and Innovative Research Team in University under Grant IRT_15R53.

References

- [1] H. Yan, J. Yang, Sparse discriminative feature selection, *Pattern Recogn.* 48 (5) (2015) 1827-1835.
- [2] F. Nie, H. Huang, X. Cai, C. Ding, Efficient and robust feature selection via joint $L_{2,1}$ -norms minimization, in: *Advances in Neural Information Processing Systems*, 2010.
- [3] I. Koprinska, M. Rana, V. G. Agelidis, Correlation and instance based feature selection for electricity load forecasting, *Knowl.-Based Syst.* 82 (2015) 29-40.
- [4] A. G. Wang, N. An, G. L. Chen, L. Li, G. Alterovitz, Accelerating wrapper-based feature

selection with k-nearest-neighbor, *Knowl.-Based Syst.* 83 (2015) 81-91.

- [5] Y. J. Li, H. X. Guo, X. Liu, Y. N. Li, J. L. Li, Adapted ensemble classification algorithm based on multiple classifier system and feature selection for classifying multi-class imbalanced data, *Knowl.-Based Syst.* 94 (2016) 88-104.
- [6] B. Gu, V. S. Sheng, K. Y. Tay, W. Romano, and S. Li, Incremental Support Vector Learning for Ordinal Regression, *IEEE Transactions on Neural Networks and Learning Systems*, 26(7) (2015): 1403-1416.
- [7] B. Gu and V. S. Sheng, A Robust Regularization Path Algorithm for v-Support Vector Classification, *IEEE Transactions on Neural Networks and Learning Systems*, DOI : 10.1109/TNNLS.2016.2527796, 2016
- [8] Z. J. Chen, C. Z. Wu, Y. S. Zhang, Z. Huang, B. Ren, Feature selection with redundancy complementariness dispersion, *Knowl.-Based Syst.* 89 (2015) 203-217.
- [9] P. Moradi, M. Rostami, Integration of graph clustering with ant colony optimization for feature selection, *Knowl.-Based Syst.* 84 (2015) 144-161.
- [10] F. Nie, S. Xiang, Y. Liu, C. Hou, C. Zhang, Orthogonal vs. uncorrelated least squares discriminant analysis for feature extraction, *Pattern Recogn. Lett.* 33 (5) (2012) 485-491.
- [11] B. Gu, X. M. Sun, and V. S. Sheng, Structural Minimax Probability Machine, *IEEE Transactions on Neural Networks and Learning Systems*, DOI : 10.1109/TNNLS.2016.2544779, 2016
- [12] F. Nie, D. Xu, I. W. Tsang, C. Zhang, Flexible manifold embedding: a framework for semi-supervised and unsupervised dimension reduction, *IEEE Trans. Image Process.* 19 (7) (2010) 1921-1932.
- [13] B. Gu, V. S. Sheng, Z. Wang, D. Ho, S. Osman, and S. Li, Incremental learning for v-Support Vector Regression, *Neural Networks*, 67(2015):140-150.
- [14] M. Sikonja, I. Kononenko, Theoretical and empirical analysis of relief and relieff, *Mach Learn.* 53 (1-2) (2003) 23-69.
- [15] Z. Xu, I. King, M. R.-T. Lyu R. Jin, Discriminative semi-supervised feature selection via manifold regularization, *IEEE Trans. Neural Netw.* 21 (7) (2010) 1033-1047.
- [16] J. G. Dy, C. E. Brodley, Feature selection for unsupervised learning, *J. Mach Learn. Res.* 5 (2004) 845-889.
- [17] M. Banerjee, N. R. Pal, Unsupervised Feature Selection with Controlled Redundancy (UFESCoR), *IEEE Trans. Knowl. Data Eng.* 27 (12) (2015) 3390-3403.

- [18]Z. Zhao, H. Liu, Spectral feature selection for supervised and unsupervised learning, in: Proceedings International Conference on Machine Learning, 2007, pp. 1151-1157.
- [19]I. Guyon, A. Elisseeffi, An introduction to variable and feature selection, J. Mach Learn. Res. 3 (2003) 1157-1182.
- [20]R. Kohavi, G.H. John, Wrappers for feature subset selection, Artif. Intell. 97 (1-2) (1997) 273-324.
- [21]X. He, D. Cai, P. Niyogi, Laplacian score for feature selection, in: Advances in Neural Information Processing Systems, 2005, pp. 507-514.
- [22]C. Hou, F. Nie, D. Yi, Y. Wu, Feature selection via joint embedding learning and sparse regression, in IJCAI Proceedings-International Joint Conference on Artificial Intelligence, 2011, pp. 1324-1329.
- [23]V. Vapnik, Statistical Learning Theory, Wiley, New York, NY, USA, 1998.
- [24]Y. W. Wang, Y. N. Liu, L. Z. Feng, X. D. Zhu, Novel feature selection method based on harmony search for email classification, Knowl.-Based Syst. 73 (2015) 311-323.
- [25]L. G. Zhou, D. Lu, H. Fujita, The performance of corporate financial distress prediction models with feature selection guided by domain knowledge and data mining approaches, Knowl.-Based Syst. 85 (2015) 52-61.
- [26]Q. F. Zhou, H. Zhou, T. Li, Cost-sensitive feature selection using random forest: Selecting low-cost subsets of informative features, Knowl.-Based Syst. 95 (2016) 1-11.
- [27]P. Paatero, U. Tapper, Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values, Environmetrics. 5 (2) (1994) 111-126.
- [28]D. D. Lee, H. S. Seung, Learning the parts of objects by non-negative matrix factorization, Nature. 401 (1999) 788-791.
- [29]I. Jolliffe, Principle Component Analysis, Springer, 1986.
- [30]S. Lipovetsky, PCA and SVD with nonnegative loadings, Pattern Recogn. 42 (1) (2009) 68-76.
- [31]G. Golub, C. Reinsch, Singular value decomposition and least squares solutions, Numer. Math. 14 (5) (1970) 403-420.
- [32]S. Wang, W. Pedrycz, Q. Zhu, W. Zhu, Subspace learning for unsupervised feature selection via matrix factorization, Pattern Recogn. 48 (1) (2015) 10-19.
- [33]S. Wang, W. Pedrycz, Q. Zhu, W. Zhu, Unsupervised feature selection via maximum projection and minimum redundancy, Knowl.-Based Syst. 75 (2015) 19-29.

- [34]X. Fang, Y. Xu, X. Li, Z. Fan, H. Liu, Y. Chen, Locality and similarity preserving embedding for feature selection, *Neurocomputing*. 128 (2014) 304-315.
- [35]K. Yu, T. Zhang, Y. H. Gong, Nonlinear learning using local coordinate coding, in: *Advances in Neural Information Processing Systems*, 2009, pp. 2223-2231.
- [36]X. He, P. Niyogi, Locality Preserving Projections, in: *Advances in Neural Information Processing Systems*, 2003, pp. 153-160.
- [37]S. Roweis, L. Saul, Nonlinear Dimensionality Reduction by Locally Linear Embedding, *Science*. 290 (5500) (2000) 2323-2326.
- [38]M. Belkin, P. Niyogi, Laplacian Eigenmaps and Spectral Techniques for Embedding and Clustering, in: *Advances in Neural Information Processing Systems*, 2001, pp. 585-591.
- [39]F. Nie, Z. Zeng, I. W. Tsang, D. Xu, C. Zhang, Spectral embedded clustering: a framework for in-sample and out-of-sample spectral clustering, *IEEE Trans. Neural Netw.* 22 (11) (2011) 1796-1808.
- [40]F. Shang, Y. Liu, F. Wang, Learning spectral embedding for semi-supervised clustering, 2011 *IEEE 11th International Conference on Data Mining*, 2011, pp. 597-606.
- [41]D. Cai, X. He, J. Han, T. Huang, Graph regularized nonnegative matrix factorization for data representation, *IEEE Trans. Pattern Anal. Mach. Intell.* 33 (8) (2011) 1548-1560.
- [42]W. Xu, Y. Gong, Document clustering by concept factorization, in: *Proceedings of International Conference on Research and Development in Information. Retrieval (SIGIR'04)*, Sheffield, UK, 2004, pp. 202-209.
- [43]D. Cai, X. He, J. Han, Locally consistent concept factorization for document clustering, *IEEE Trans. Knowl. Data Eng.* 23 (6) (2011) 902-913.
- [44]F. Shang, L.C. Jiao, F. Wang, Graph dual regularization non-negative matrix factorization for co-clustering, *Pattern Recogn.* 45 (6) (2012) 2237-2250.
- [45]J. Bu, P. Li, C. Chen, Z. He, D. Cai, Relational multi-manifold co-clustering, *IEEE Trans. Cybern.* 43 (6) (2013) 1871-1881.
- [46]J. Ye, Z. Jin, Dual-graph regularized concept factorization for clustering, *Neurocomputing*. 138 (2014) 120-130.
- [47]M. Belkin, P. Niyogi, V. Sindhwani, Manifold regularization: A geometric framework for learning from labeled and unlabeled examples, *J. Mach. Learn. Res.*, 7 (2006) 2399-2434.
- [48]L. Chova, G. Valls, J. Mari, J. Calpe, Semi-supervised image classification with Laplacian

support vector machines, *IEEE Geosci. Remote Sens. Lett.*, 5 (2008) 336-340.

- [49]Z. Yang and Y. Xu, Laplacian twin parametric-margin support vector machine for semi-supervised classification, *Neurocomputing*, 171 (2016) 325-334.
- [50]Z. Zhao, L. Wang, H. Liu, Efficient spectral feature selection with minimum redundancy, in: *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010, pp. 673-678.
- [51]D. Cai, C. Zhang, X. He, Unsupervised feature selection for multicluster data, in: *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2010, pp. 333-342.
- [52]Z. Li, J. Liu, Y. Yang, X. Zhou, H. Lu, Clustering-Guided sparse structural learning for unsupervised feature selection, *IEEE Trans. Knowl. Data Eng.* 26 (9) (2014) 2138-2150.
- [53]C. Hou, F. Nie, X. Li, D. Yi, Y. Wu, Joint embedding learning and sparse regression: A framework for unsupervised feature selection, *IEEE Trans. Cybern.* 44 (6) (2014) 2168-2267.
- [54]D. D. Lee, H. Seung, Algorithms for non-negative matrix factorization, in: *Advances in Neural Information Processing Systems*, 2001, pp. 556-562.
- [55]A. Rakhlin, A. Caponnetto, Stability of K-Means clustering, in: *Advances in Neural Information Processing Systems*, 2007, pp. 216-222.
- [56]H. Lu, Z. Fu, X. Shu, Non-negative and sparse spectral clustering, *Pattern Recogn.* 47 (1) (2014) 418-426.
- [57]Y. Yang, H. T. Shen, Z. Ma, Z. Huang, X. Zhou, $\ell_{2,1}$ -Norm regularized discriminative feature selection for unsupervised learning, in: *Proceedings of the 22nd International Joint Conference on Artificial Intelligence*, 2011, pp. 1589-1594.
- [58]F. Bach, Consistency of the group lasso and multiple kernel learning, *J. Mach Learn. Res.* 9 (2008) 1179-1225.
- [59]M. Wu, B. Schölkopf, A local learning approach for clustering. In: *Advances in Neural Information Processing Systems*, 2007, pp. 1529-1536.
- [60]A. Strehl, J. Ghosh, Cluster ensembles-a knowledge reuse framework for combining multiple partitions. *J. Mach Learn. Res.* 3 (2002) 583-617.
- [61]C. Papadimitriou, K. Steiglitz, *Combinatorial optimization: Algorithms and complexity*, Dover, New York, NY, USA, 1998.