

Low-rank Dictionary Learning for Unsupervised Feature Selection

Mohsen Ghassemi Parsa^a, Hadi Zare^{a,*}, Mehdi Ghatee^b

^a*Faculty of New Sciences and Technologies, University of Tehran, Tehran, Iran*

^b*Department of Mathematics and Computer Science, Amirkabir University of Technology, Tehran, Iran*

Abstract

There exist many high-dimensional data in real-world applications such as biology, computer vision, and social networks. Feature selection approaches are devised to confront with high-dimensional data challenges with the aim of efficient learning technologies as well as reduction of models complexity. Due to the hardship of labeling on these datasets, there are a variety of approaches on feature selection process in an unsupervised setting by considering some important characteristics of data. In this paper, we introduce a novel unsupervised feature selection approach by applying dictionary learning ideas in a low-rank representation. Dictionary learning in a low-rank representation not only enables us to provide a new representation, but it also maintains feature correlation. Then, spectral analysis is employed to preserve sample similarities. Finally, a unified objective function for unsupervised feature selection is proposed in a sparse way by an $\ell_{2,1}$ -norm regularization. Furthermore, an efficient numerical algorithm is designed to solve the corresponding optimization problem. We demonstrate the performance of the proposed method based on a variety of standard datasets from different applied domains. Our experimental findings reveal that the proposed method outperforms the state-of-the-art algorithm.

Keywords: Unsupervised feature selection, Dictionary Learning, Sparse Learning, Spectral analysis, Low-rank representation

*Corresponding author

Email addresses: mgparsa@ut.ac.ir (Mohsen Ghassemi Parsa), h.zare@ut.ac.ir (Hadi Zare), ghatee@aut.ac.ir (Mehdi Ghatee)

1. Introduction

Technological advancement and popularity of social networks provide many huge and high-dimensional data and information sources. High-dimensional data are available in many applications, including machine vision (Shi et al., 2015), text mining (Rogati & Yang, 2002), and biology (Hoseini & Mansoori, 2019). High-dimensionality not only increases the complexity of the training process and the learned model but also degrades performance, that is called as the curse of dimensionality (Murphy, 2012). To address the issue, dimensionality reduction can be considered in two main approaches, feature extraction (FE) and feature selection (FS) (Pandit et al., 2020). The new features are constituted by a linear or non-linear transformation of the original features in FE approaches, while FS methods aim to select appropriate features by considering some evaluation criteria. FS attains more attraction than FE in some situations, specifically when the primary aim is to take advantage of more interpretable and understandable features (Li et al., 2017a).

FS methods can be classified to filter, wrapper, and embedded approaches according to feature evaluation. Filters (Krzanowski, 1987; He et al., 2005; Zare & Niazi, 2016) exploit data properties to find out the importance of the features, while wrappers (Dy & Brodley, 2004) evaluate the feature subsets by a learning algorithm. In embedded methods (Li et al., 2012), the feature selection process is embedded in a learning algorithm. Recently, unsupervised feature selection (UFS) has attracted many efforts among researchers due to the unavailability of the right answers on practical domains and real-world applications (Parsa et al., 2020; Zare et al., 2020).

UFS algorithms are mainly categorized into similarity preserving, sparse learning, reconstruction, and dictionary learning methods. Similarity preserving methods (He et al., 2005; Zhao & Liu, 2007) are tried to maintain the local geometric structures among the selected features. Sparse learning methods (Parsa et al., 2020; Zare et al., 2020) considered selecting more relevant features in a regularized way. Data reconstruction methods (Masaeli et al., 2010; Farahat et al., 2013) re-express the features to eliminate uninformative ones. One of the most important approaches in UFS is based on dictionary learning (Zhu et al., 2016, 2017; Ding et al., 2020), in which a new sparse representation of the data matrix is obtained on a dictionary basis space.

Most of the earlier dictionary learning approaches proposed to build the dictionary matrix without any restriction on the rank of the basis matrix. A

more natural assumption is to employ the low-rank representation on high-dimensional data to alleviate the noisy and redundant features (Chen & Huang, 2012). In addition, the basis matrix can be learned in a parsimonious way by imposing the rank constraint, which can be improved the learning process. In this paper, we propose a dictionary learning-based unsupervised feature selection method, named DLUFS, to provide a sparse representation of the original data. Furthermore, we employ a low-rank constraint to eliminate the noisy and redundant features. The local sample structure is also considered by exploiting a spectral analysis.

We summarize the main contributions of this paper as,

- A dictionary learning method is proposed to select features to obtain a sparse representation of data.
- Low-rank constraint on the basis matrix is imposed to eliminate the noisy and redundant features.
- Spectral analysis is employed to preserve the local similarities among samples.

This paper is organized as follows. The existing UFS methods are reviewed in Section 2. The proposed method, an illustrative example, and the corresponding optimization algorithm are presented in Section 3. We analyze the convergence behavior of the proposed algorithm in Section 4. Section 5 presents the experimental results on benchmark datasets based on state-of-the-art methods. The conclusions are given in Section 6.

2. Related Works

In this section, we review unsupervised feature selection methods in four categories, similarity preserving, sparse machine learning, data reconstruction, and dictionary learning methods.

Similarity preserving methods consider the sample structure in the selected features, such as Laplacian score, LS (He et al., 2005), spectral feature selection, SPEC (Zhao & Liu, 2007), and trace ratio criterion for feature selection, TrRatio (Nie et al., 2008). The learning models are not employed in this category of methods which results in the selection of less relevant features.

In sparse machine learning methods, feature selection is performed based on learning a regularized model, such as low-dimensional embedding and

sparse regression, JELSR (Hou et al., 2014), local discriminative sparse subspace learning, LDSSL (Shang et al., 2019), multi-cluster feature selection, MCFS (Cai et al., 2010), non-negative discriminative feature selection, NDFS (Li et al., 2012), similarity preserving feature selection, SPFS (Zhao et al., 2013), structure preservation robust spectral feature selection, SRFS (Zhu et al., 2018), unsupervised discriminative feature selection, UDFS (Yang et al., 2011). These methods commonly select features based on learning a regularized regression matrix without involving data reconstruction.

Data reconstruction approaches were proposed to select features based on their explanation on linear and non-linear transformation of the data, including sparse principal component analysis, CPFS (Masaeli et al., 2010), greedy unsupervised feature selection, GreedyFS (Farahat et al., 2013), graph regularized feature selection, GRFS (Zhao et al., 2016), embedded reconstruction based unsupervised feature selection, REFS (Li et al., 2017b), structure preserving unsupervised feature selection, SPUFS (Lu et al., 2018), reconstruction error minimization, REMFS (Yang et al., 2019).

While dictionary learning and reconstruction based methods can be regarded as similar techniques to learn the basis matrix, the dictionary learning methods enable us to provide a new data representation along with the elimination of redundant features. Feature selection process in dictionary learning methods is conducted in two main phases, learning the basis matrix and sparse new data representation. Most of the dictionary learning methods were proposed by considering a two-step procedure such as (Zheng et al., 2011; Zhu et al., 2017). Graph sparse coding, GSC (Zheng et al., 2011) performed dictionary learning and spectral analysis to yield a sparse representation by an ℓ_1 -norm regularization. Robust joint graph sparse coding, RJGSC (Zhu et al., 2017), extended GSC by using an $\ell_{2,1}$ -norm regularizer. On the other hand, DGL (Ding et al., 2020) jointly learns the basis and sparse data matrix in a unified framework. Earlier dictionary-based methods have not constructed the basis matrix by considering the natural low-rank assumption of it, which is justified on many real-world high-dimensional data (Chen & Huang, 2012).

Table 1 presents a summary of the related methods by considering important characteristics in a UFS process. Sparse learning employs a regularization approach to the learning model. Subspace learning indicates the low-dimensional representation of the original data. In spectral analysis, the local structure of the samples is taken into account. Joint learning refers to a unified objective function. In data reconstruction, features are expressed by

Table 1: Summary of the state-of-the-art unsupervised feature selection methods.

Algorithm	Sparse learning	Subspace learning	Spectral analysis	Joint learning	Data reconstruction	Dictionary learning	Low-rank representation
LS (He et al., 2005)	×	×	✓	✓	×	×	×
MCFS (Cai et al., 2010)	✓	✓	✓	×	×	×	×
UDFS (Yang et al., 2011)	✓	✓	×	✓	×	×	×
NDFS (Li et al., 2012)	✓	✓	✓	✓	×	×	×
SPFS (Zhao et al., 2013)	✓	✓	×	✓	×	×	×
JELSR (Hou et al., 2014)	✓	✓	✓	✓	×	×	×
LDSSL (Shang et al., 2019)	✓	✓	✓	✓	✓	×	×
SRFS (Zhu et al., 2018)	✓	✓	✓	✓	✓	×	✓
RJGSC (Zhu et al., 2017)	✓	×	✓	×	✓	✓	×
DGL (Ding et al., 2020)	✓	×	✓	✓	✓	✓	×
DLUFS	✓	✓	✓	✓	✓	✓	✓

a linear or non-linear combination of all features to discard redundant ones. In dictionary learning, a new representation of the original data is learned in a basis space. By low-rank representation, the reconstruction matrix is decomposed to low-rank matrices to consider the correlation among features. In this paper, we propose a unified UFS method based on all of these main characteristics to yield an efficient and robust procedure.

3. The Proposed Method

In this section, at first notations are presented. Then, the proposed method and its algorithm details are introduced. Finally, the proposed algorithm is illustrated through an example.

3.1. Notations

In this paper, the vectors and the matrices are denoted by bold lowercase and bold uppercase characters. For a given vector \mathbf{v} , its ℓ_2 -norm is denoted by $\|\mathbf{v}\|_2$. Suppose \mathbf{M} is an arbitrary matrix, M_{ij} represents its (i, j) -th element, \mathbf{m}_i is the i -th row, \mathbf{m}^j is the j -th column, $\text{tr}(\mathbf{M})$ is the trace, and \mathbf{M}^\top is the transpose of the matrix. The Frobenius norm is denoted by $\|\mathbf{M}\|_F$, and the $\ell_{2,1}$ -norm is defined as,

$$\|\mathbf{M}\|_{2,1} = \sum_i \sqrt{\sum_j M_{ij}^2}.$$

Let $\mathbf{X} \in \mathbb{R}^{p \times n}$ represents the data matrix, where p is the number of features and n is the number of samples.

3.2. The Proposed Method

At first, a new data representation matrix can be learned based on dictionary learning approach as,

$$\min_{\mathbf{Q}, \mathbf{Z}} \|\mathbf{X} - \mathbf{QZ}\|_F^2, \quad (1)$$

where $\mathbf{Z} \in \mathbb{R}^{p \times n}$ is a representation of the data matrix \mathbf{X} in the space of the dictionary matrix $\mathbf{Q} \in \mathbb{R}^{p \times p}$. The rank of high-dimensional data can be increased by the noisy and outlier features (Chen & Huang, 2012). Based on this fact, the data can be represented in a low-rank space. In this regard, a low-rank constraint on the basis matrix is imposed as,

$$\begin{aligned} \min_{\mathbf{Q}, \mathbf{Z}} \quad & \|\mathbf{X} - \mathbf{QZ}\|_F^2 \\ \text{s.t.} \quad & \text{rank}(\mathbf{Q}) = r, \end{aligned} \quad (2)$$

where $r \ll \{n, p\}$ is the induced rank to \mathbf{Q} . The low-rank constraint on Eq. (2) is equivalent to multiply two rank r matrices as,

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{Z}} \|\mathbf{X} - \mathbf{ABZ}\|_F^2, \quad (3)$$

where $\mathbf{A} \in \mathbb{R}^{p \times r}$, $\mathbf{B} \in \mathbb{R}^{r \times p}$. By $\mathbf{Q} = \mathbf{AB}$ decomposition, the feature correlation is considered in a low-rank space. The matrix \mathbf{Z} is transformed to a low-dimensional matrix $\mathbf{BZ} \in \mathbb{R}^{r \times n}$ to perform subspace learning, while \mathbf{A} re-transforms \mathbf{BZ} to the original space. More specifically, as further expressed in Eq. (11), the mentioned subspace learning is calculated based on LDA (Fukunaga, 1990).

The global feature correlations are maintained by low-rank constraint. Furthermore, the spectral analysis is applied to take the local sample structure into account as,

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{Z}} \|\mathbf{X} - \mathbf{ABZ}\|_F^2 + \alpha \text{tr}(\mathbf{ZLZ}^\top), \quad (4)$$

where α is a tuning parameter. The Laplacian matrix \mathbf{L} is calculated as $\mathbf{L} = \mathbf{D} - \mathbf{S}$, where the diagonal matrix \mathbf{D} is defined as $D_{ii} = \sum_j S_{ij}$ and the similarity matrix \mathbf{S} is calculated as follows,

$$S_{ij} = \begin{cases} \exp\left(-\frac{\|\mathbf{x}^i - \mathbf{x}^j\|_2^2}{\sigma^2}\right), & \text{if } \mathbf{x}^i \in N_k(\mathbf{x}^j) \text{ or } \mathbf{x}^j \in N_k(\mathbf{x}^i) \\ 0, & \text{otherwise,} \end{cases} \quad (5)$$

where $N_k(\mathbf{x}^i)$ represents the set of k -nearest neighbors of \mathbf{x}^i , and σ is the width parameter for the Gaussian kernel.

A feature selection framework can be provided by inducing a sparse learning on the new representation matrix \mathbf{Z} . Hence, the final objective function is proposed as,

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{Z}} \|\mathbf{X} - \mathbf{ABZ}\|_F^2 + \alpha \text{tr}(\mathbf{ZLZ}^\top) + \lambda \|\mathbf{Z}\|_{2,1}, \quad (6)$$

where λ is the regularization parameter. The $\ell_{2,1}$ -norm regularizer provides sparsity on the rows of \mathbf{Z} , inspired by the group lasso penalty (Yuan & Lin, 2006). The rows are closer to zero, then the corresponding features are more likely regarded as uninformative features.

3.3. Optimization

We consider the main objective function as the following optimization problem,

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{Z}} f(\mathbf{A}, \mathbf{B}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{ABZ}\|_F^2 + \alpha \text{tr}(\mathbf{ZLZ}^\top) + \lambda \|\mathbf{Z}\|_{2,1}, \quad (7)$$

First, by fixing \mathbf{Z} in the main optimization problem in Eq. (7) to get,

$$\min_{\mathbf{A}, \mathbf{B}} f(\mathbf{A}, \mathbf{B}) = \|\mathbf{X} - \mathbf{ABZ}\|_F^2. \quad (8)$$

By setting the derivative of the Eq. (8) with respect to \mathbf{A} to zero,

$$\mathbf{A} = \mathbf{XZ}^\top \mathbf{B}^\top (\mathbf{BS}_w \mathbf{B}^\top)^{-1}, \quad (9)$$

where $\mathbf{S}_w = \mathbf{ZZ}^\top$ is the correlation among features in the new representation matrix \mathbf{Z} .

We rewrite Eq. (8) as,

$$\min_{\mathbf{A}, \mathbf{B}} \text{tr}(\mathbf{XX}^\top) - 2\text{tr}(\mathbf{ABZ}\mathbf{X}^\top) + \text{tr}(\mathbf{ABZZ}^\top \mathbf{B}^\top \mathbf{A}^\top). \quad (10)$$

Using obtained \mathbf{A} from Eq. (9) in Eq. (10) to derive the objective function of \mathbf{B} as,

$$\begin{aligned} & \min_{\mathbf{B}} -\text{tr}(\mathbf{XZ}^\top \mathbf{B}^\top (\mathbf{BS}_w \mathbf{B}^\top)^{-1} \mathbf{BZ}\mathbf{X}^\top), \\ \Leftrightarrow & \max_{\mathbf{B}} \text{tr}((\mathbf{BS}_w \mathbf{B}^\top)^{-1} \mathbf{BS}_b \mathbf{B}^\top), \end{aligned} \quad (11)$$

where $\mathbf{S}_b = \mathbf{Z}\mathbf{X}^\top\mathbf{X}\mathbf{Z}^\top$. Similar to discriminant analysis (Fukunaga, 1990), \mathbf{S}_w and \mathbf{S}_b can be interpreted as within-class and between-class scatter matrices. Therefore, \mathbf{B}^\top can be learned by r eigenvectors of $\mathbf{S}_w^{-1}\mathbf{S}_b$ corresponding to top r eigenvalues.

By rewriting the objective function in Eq. (6),

$$f(\mathbf{Z}) = \|\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{Z}\|_F^2 + \alpha \text{tr}(\mathbf{Z}\mathbf{L}\mathbf{Z}^\top) + \lambda \|\mathbf{Z}\|_{2,1}. \quad (12)$$

Let \mathbf{A} and \mathbf{B} are fixed in Eq. (12). Then, by setting the derivative of $f(\mathbf{Z})$ to zero, the following Sylvester equation (Bartels & Stewart, 1972) can be obtained,

$$((\mathbf{A}\mathbf{B})^\top\mathbf{A}\mathbf{B} + \lambda\mathbf{D})\mathbf{Z} + \mathbf{Z}(\alpha\mathbf{L}) = (\mathbf{A}\mathbf{B})^\top\mathbf{X}, \quad (13)$$

where \mathbf{D} is a diagonal matrix as,

$$D_{ii} = \frac{1}{2\|\mathbf{z}_i\|_2 + \epsilon}. \quad (14)$$

Here, Algorithm 1 summarizes the iterative procedure of obtaining the main optimization variables in Eq. (7). By descending order of $\|\mathbf{z}_i\|_2$'s, the importance of the corresponding features are determined.

Algorithm 1 DLUFS algorithm.

Input: Data matrix $\mathbf{X} \in \mathbb{R}^{p \times n}$ and parameters α and λ .

1: $t = 0$.

2: Initialize $\mathbf{Z}^t = \mathbf{X}$.

3: **repeat**

4: Update \mathbf{B}^{t+1} by solving Eq. (11).

5: Update \mathbf{A}^{t+1} by Eq. (9).

6: Update the diagonal matrix \mathbf{D}^{t+1} by Eq. (14).

7: Update \mathbf{Z}^{t+1} by solving the Sylvester equation in (13).

8: **until** the convergence of the objective function in Eq. (7).

Output: Sorting features in descending order of $\|\mathbf{z}_i\|_2$'s.

3.4. Computational Complexity

The computational complexity of Algorithm 1 consists of computing \mathbf{B} , \mathbf{A} and \mathbf{Z} in each iteration. By considering (11), cost of updating \mathbf{B} equals

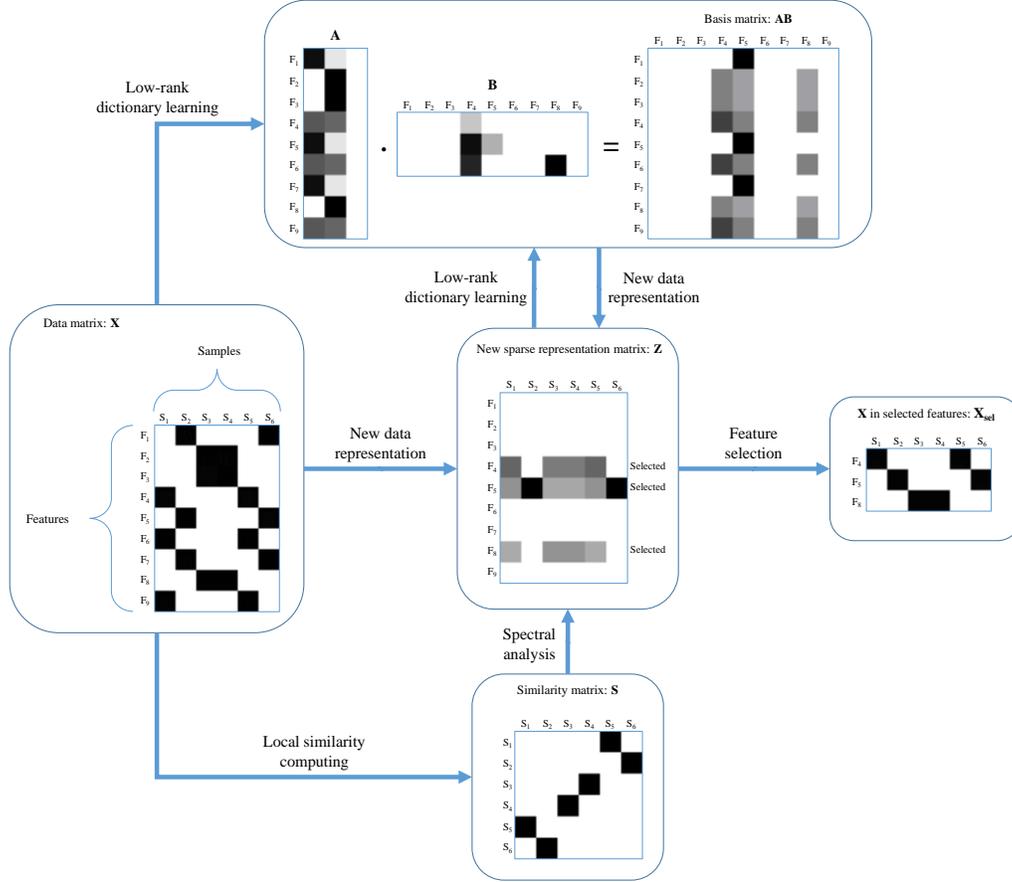


Figure 1: An illustrative example for describing the proposed method.

to $\max\{O(p^3), O(p^2n)\}$. Next by (9), the time complexity of computing \mathbf{A} is $\max\{O(p^2n), O(p^2r), O(r^3)\}$. Furthermore, the time complexity of updating \mathbf{Z} contains two elements, computing the input of Sylvester equation in (13), and solving the Sylvester equation, which derives as $\max\{O(p^3), O(p^2n), O(p^2r)\}$. Since the assumption of $r \ll \{p, n\}$ holds on high-dimensional settings, computational complexity of Algorithm 1 reduces to $\max\{O(p^3), O(p^2n)\}$.

3.5. An illustrative example

Fig. 1 describes the proposed algorithm through an illustrative example. All matrices are assumed to be non-negative, where the brighter elements indicate more closeness to zero and the darker ones are far from zero. \mathbf{X} is

an artificial data matrix with nine features and six samples. First, a local similarity matrix \mathbf{S} is calculated based on the sample similarities in \mathbf{X} . A low-rank basis matrix \mathbf{AB} is obtained by dictionary learning which can be interpreted as the correlation among features based on the low-dimensional matrices \mathbf{A} and \mathbf{B} . Then, a new data representation matrix \mathbf{Z} is computed based on spectral analysis in the basis space. The \mathbf{B} , \mathbf{A} , and \mathbf{Z} are iteratively updated until convergence. Features are ranked by calculating the ℓ_2 -norm on their corresponding rows of \mathbf{Z} as our new data matrix. Finally, selected features \mathbf{X}_{sel} are given in the output.

In this example, samples are classified to three sets as $\{\mathbf{S}_1, \mathbf{S}_5\}$, $\{\mathbf{S}_2, \mathbf{S}_6\}$, and $\{\mathbf{S}_3, \mathbf{S}_4\}$ in terms of similarities. In addition, there are three categories of features, $\{\mathbf{F}_1, \mathbf{F}_5, \mathbf{F}_7\}$, $\{\mathbf{F}_2, \mathbf{F}_3, \mathbf{F}_8\}$, and $\{\mathbf{F}_4, \mathbf{F}_6, \mathbf{F}_9\}$ by their similarities. The basis matrix \mathbf{AB} is learned based on low-rank dictionary learning with rank = 3. Fig. 1 indicates that the top three remaining rows of \mathbf{Z} are the \mathbf{F}_4 , \mathbf{F}_5 , and \mathbf{F}_8 . Therefore, the output \mathbf{X}_{sel} is formed by selected features.

4. Convergence Analysis

Our aim is to show the non-increasing behavior of Algorithm 1 based on the objective function in Eq. (7). Initially, a lemma is given, then the main theorem is presented.

Lemma 1. *Let \mathbf{u} and \mathbf{v} are two non-zero vectors, then this inequality holds,*

$$\|\mathbf{u}\|_2 - \frac{\|\mathbf{u}\|_2^2}{2\|\mathbf{v}\|_2} \leq \|\mathbf{v}\|_2 - \frac{\|\mathbf{v}\|_2^2}{2\|\mathbf{v}\|_2}. \quad (15)$$

The proof of Lemma 1 derived in (Nie et al., 2010).

Theorem 1. *Algorithm 1 behaves non-increasingly in each update through the primary objective function in (7).*

Proof. In the following, the t -th iteration of a vector \mathbf{v} and a matrix \mathbf{M} is denoted by \mathbf{v}^t and \mathbf{M}^t .

The non-increasing behavior of the objective function of \mathbf{Z} in (12) is derived by assuming \mathbf{A}^t and \mathbf{B}^t to be fixed. Since the non-smooth $\|\mathbf{Z}\|_{2,1}$ is iteratively optimized by updating \mathbf{D} and \mathbf{Z} , the following inequality can be

presented,

$$\begin{aligned} & \|\mathbf{X} - \mathbf{Q}^t \mathbf{Z}^{t+1}\|_F^2 + \alpha \operatorname{tr}(\mathbf{Z}^{t+1} \mathbf{L}(\mathbf{Z}^{t+1})^\top) + \lambda \sum_{i=1}^p \frac{\|\mathbf{z}_i^{t+1}\|_2^2}{2\|\mathbf{z}_i^t\|_2} \\ & \leq \|\mathbf{X} - \mathbf{Q}^t \mathbf{Z}^t\|_F^2 + \alpha \operatorname{tr}(\mathbf{Z}^t \mathbf{L}(\mathbf{Z}^t)^\top) + \lambda \sum_{i=1}^p \frac{\|\mathbf{z}_i^t\|_2^2}{2\|\mathbf{z}_i^t\|_2}. \end{aligned} \quad (16)$$

Where $\mathbf{Q} = \mathbf{A}\mathbf{B}$. Then, the inequality (16) can be rewritten as,

$$\begin{aligned} & \|\mathbf{X} - \mathbf{Q}^t \mathbf{Z}^{t+1}\|_F^2 + \alpha \operatorname{tr}(\mathbf{Z}^{t+1} \mathbf{L}(\mathbf{Z}^{t+1})^\top) + \lambda \|\mathbf{Z}^{t+1}\|_{2,1} - \lambda \sum_{i=1}^p \left(\|\mathbf{z}_i^{t+1}\|_2 - \frac{\|\mathbf{z}_i^{t+1}\|_2^2}{2\|\mathbf{z}_i^t\|_2} \right) \\ & \leq \|\mathbf{X} - \mathbf{Q}^t \mathbf{Z}^t\|_F^2 + \alpha \operatorname{tr}(\mathbf{Z}^t \mathbf{L}(\mathbf{Z}^t)^\top) + \lambda \|\mathbf{Z}^t\|_{2,1} - \lambda \sum_{i=1}^p \left(\|\mathbf{z}_i^t\|_2 - \frac{\|\mathbf{z}_i^t\|_2^2}{2\|\mathbf{z}_i^t\|_2} \right). \end{aligned} \quad (17)$$

According to Lemma 1,

$$\|\mathbf{z}_i^{t+1}\|_2 - \frac{\|\mathbf{z}_i^{t+1}\|_2^2}{2\|\mathbf{z}_i^t\|_2} \leq \|\mathbf{z}_i^t\|_2 - \frac{\|\mathbf{z}_i^t\|_2^2}{2\|\mathbf{z}_i^t\|_2}, \quad (18)$$

we obtain,

$$\begin{aligned} & \|\mathbf{X} - \mathbf{Q}^t \mathbf{Z}^{t+1}\|_F^2 + \alpha \operatorname{tr}(\mathbf{Z}^{t+1} \mathbf{L}(\mathbf{Z}^{t+1})^\top) + \lambda \|\mathbf{Z}^{t+1}\|_{2,1} \\ & \leq \|\mathbf{X} - \mathbf{Q}^t \mathbf{Z}^t\|_F^2 + \alpha \operatorname{tr}(\mathbf{Z}^t \mathbf{L}(\mathbf{Z}^t)^\top) + \lambda \|\mathbf{Z}^t\|_{2,1}. \end{aligned} \quad (19)$$

Therefore,

$$f(\mathbf{A}^t, \mathbf{B}^t, \mathbf{Z}^{t+1}) \leq f(\mathbf{A}^t, \mathbf{B}^t, \mathbf{Z}^t). \quad (20)$$

In the same way, by fixing \mathbf{Z}^{t+1} , it can be shown that,

$$f(\mathbf{A}^{t+1}, \mathbf{B}^{t+1}, \mathbf{Z}^{t+1}) \leq f(\mathbf{A}^t, \mathbf{B}^t, \mathbf{Z}^{t+1}). \quad (21)$$

By considering inequalities (20) and (21),

$$f(\mathbf{A}^{t+1}, \mathbf{B}^{t+1}, \mathbf{Z}^{t+1}) \leq f(\mathbf{A}^{t+1}, \mathbf{B}^{t+1}, \mathbf{Z}^t) \leq f(\mathbf{A}^t, \mathbf{B}^t, \mathbf{Z}^t). \quad (22)$$

Therefore, the non-increasing behavior of Algorithm 1 based on the primary objective function in Eq. (7) is given. \square

Table 2: The statistics of datasets.

Dataset	samples	features	classes	Type	Category
BA	1404	320	36	Binary	Image
Colon	62	2000	2	Discrete	Biology
GLIOMA	50	4434	4	Continuous	Biology
Madelon	2600	500	2	Continuous	Artificial
ORL	400	1024	40	Discrete	Image
PCMAC	1943	3289	2	Discrete	Text
WarpAR10P	130	2400	10	Discrete	Image
Yale	165	1024	15	Discrete	Image

5. Experiments

This section is divided into five subsections to describe our experimental setup. A brief summary of the applied datasets in our study, evaluation measures, the details of parameters setting, the obtained results, and the sensitivity analysis are presented in the following.

5.1. Datasets

A variety of domains of applications are considered in employed datasets such as images of digits (BA (Belhumeur et al., 1997)), colon cancer (Colon (Alon et al., 1999)), malignant brain tumor (GLIOMA (et al., 2003)), non-sparse artificial dataset (Madelon (Li et al., 2017a)), image of faces (ORL, Yale (Cai et al., 2006), and WarpAR10P (Li et al., 2017a)), and PC versus MAC from 20-newsgroups dataset (PCMAC (Lang, 1995)). The BA is available on <https://cs.nyu.edu/~roweis/data.html>, while all others are accessed from (Li et al., 2017a). Table 2 reports the main characteristics of datasets.

5.2. Evaluation measures

The clustering techniques are usually applied to evaluate the UFS methods (Li et al., 2017a). Based on the attained clustering results and the ground truth information, two common evaluation measures are used frequently, Accuracy and Normalized Mutual Information.

Let the ground truth and the predicted one based on clustering approach are shown by \mathbf{y} and $\hat{\mathbf{y}}$.

The Accuracy (called as ACC) is defined as,

$$\text{ACC}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_{i=1}^n \delta(y_i, \text{map}(z_i)),$$

where the function $\delta(a, b)$ equals to 1, when $a = b$, and 0 elsewhere. For the $\text{map}(\cdot)$ function, the Kuhn-Munkres approach (Lovasz, 1986) is employed to find the best permutation for matching the categories in vectors \mathbf{y} and $\hat{\mathbf{y}}$. Based on definitions of entropy measure $H(\cdot)$ and the mutual information of \mathbf{y} and $\hat{\mathbf{y}}$ denoted by $I(\mathbf{y}, \hat{\mathbf{y}})$ (Bishop, 2006), Normalized Mutual Information (called as NMI) is given as,

$$\text{NMI}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{I(\mathbf{y}, \hat{\mathbf{y}})}{\max(H(\mathbf{y}), H(\hat{\mathbf{y}}))},$$

5.3. The experimental setting

We compare the proposed method, DLUFS, with the state-of-the-art unsupervised feature selection algorithms, including JELSR (Hou et al., 2014), LDSSL (Shang et al., 2019), LS (He et al., 2005), MCFS (Cai et al., 2010), NDFS (Li et al., 2012), SPFS (Zhao et al., 2013), SRFS (Zhu et al., 2018), UDFS (Yang et al., 2011), and selecting all features namely Baseline.

The number of neighborhoods in k-nearest neighbor algorithm is set to $k = 5$. We set $\sigma = 1$ in our method and others requiring a similarity matrix based on a Gaussian kernel. For NDFS method, the default $\gamma = 10^8$ is considered. Finding the suitable choices of the tuning parameters α and λ in our method is performed by a grid search approach from the set of $\{10^{-4}, 10^{-2}, 1, 10^2, 10^4\}$ candidates, where a similar approach is used to set the tuning parameters in the other methods. The stopping convergence condition for all of the iterative algorithms, including our method, is given as $\frac{|f(t) - f(t-1)|}{f(t)} < 10^{-3}$, where $f(t)$ equals to the objective function in the t -th iteration. The NMI and ACC measures are reported on 20 times repetitions. In each repetition, the k-means algorithm is applied by setting the numbers of features from $\{50, 100, 150, 200, 250, 300\}$. We report two main descriptive statistical quantities of NMI and ACC in these repeated experiments, mean and standard deviation (STD).

5.4. Experimental results

We demonstrate the performance of the proposed method, DLUFS, by comparing to benchmark UFS methods defined in the earlier subsection. Table 3 and 4 represent the obtained results based on ACC and NMI (mean \pm STD). The bold numbers represent the best attained results. We used the underlined numbers to indicate the second-best results. We summarize the main findings from the obtained results of Table 3 and Table 4 in the following,

Table 3: The results of clustering (ACC% \pm std) of UFS methods on standard datasets. The best and the second-best are presented in bold and underlined numbers.

Dataset	BA	Colon	GLIOMA	Madelon	ORL	PCMAC	WarpAR10P	Yale
Baseline	43.04 \pm 1.18	54.84 \pm 0.00	<u>61.30</u> \pm 4.11	50.30 \pm 0.07	59.24 \pm 2.17	50.54 \pm 0.04	21.04 \pm 2.93	41.91 \pm 2.36
JELSR	39.97 \pm 2.14	56.76 \pm 1.34	<u>52.53</u> \pm 1.00	<u>57.53</u> \pm 1.21	57.22 \pm 3.15	50.55 \pm 0.03	<u>33.96</u> \pm 1.50	37.56 \pm 1.05
LDSSL	40.73 \pm 3.28	<u>58.01</u> \pm 0.47	56.98 \pm 0.70	<u>52.23</u> \pm 1.57	43.09 \pm 5.51	<u>50.69</u> \pm 0.07	33.60 \pm 0.51	40.23 \pm 0.85
LS	41.59 \pm 3.47	57.80 \pm 0.60	55.75 \pm 2.28	50.35 \pm 0.04	48.13 \pm 2.92	50.45 \pm 0.03	32.60 \pm 3.43	<u>44.07</u> \pm 2.42
MCFS	41.75 \pm 2.16	54.66 \pm 1.44	60.55 \pm 4.73	52.30 \pm 0.99	57.72 \pm 1.27	50.46 \pm 0.14	23.38 \pm 3.39	39.76 \pm 0.49
NDFS	39.42 \pm 4.06	57.39 \pm 1.45	57.37 \pm 2.14	57.06 \pm 1.64	53.43 \pm 2.73	50.59 \pm 0.03	31.38 \pm 1.14	37.15 \pm 1.47
SPFS	41.39 \pm 2.07	58.33 \pm 1.61	55.68 \pm 4.48	51.58 \pm 0.12	58.08 \pm 1.28	50.52 \pm 0.04	30.96 \pm 5.55	41.33 \pm 0.71
SRFS	37.78 \pm 4.73	<u>60.32</u> \pm 0.92	61.20 \pm 1.32	52.00 \pm 1.65	58.17 \pm 1.24	50.54 \pm 0.02	31.19 \pm 2.61	42.06 \pm 1.88
UDFS	40.69 \pm 3.05	55.67 \pm 0.86	54.17 \pm 3.05	57.47 \pm 1.66	54.74 \pm 2.10	50.57 \pm 0.05	28.60 \pm 3.78	33.53 \pm 1.33
DLUFS	<u>42.38</u> \pm 1.91	64.83 \pm 7.83	64.22 \pm 2.18	59.13 \pm 0.43	<u>59.22</u> \pm 1.25	50.80 \pm 0.12	35.65 \pm 2.23	47.31 \pm 0.98

Table 4: The results of clustering (NMI% \pm std) of UFS methods on standard datasets. The best and the second-best are presented in bold and underlined numbers.

Dataset	BA	Colon	GLIOMA	Madelon	ORL	PCMAC	WarpAR10P	Yale
Baseline	58.21 \pm 0.69	00.62 \pm 0.00	<u>50.93</u> \pm 2.47	00.00 \pm 0.00	<u>77.81</u> \pm 0.83	00.04 \pm 0.04	17.41 \pm 3.60	48.93 \pm 1.85
JELSR	55.62 \pm 1.60	02.76 \pm 0.86	31.11 \pm 2.10	01.60 \pm 0.43	75.70 \pm 1.94	00.90 \pm 0.05	33.85 \pm 2.17	44.90 \pm 1.13
LDSSL	56.04 \pm 3.26	01.52 \pm 0.30	49.74 \pm 0.37	00.22 \pm 0.20	66.24 \pm 3.99	01.08 \pm 0.12	<u>35.00</u> \pm 1.17	47.08 \pm 0.80
LS	57.14 \pm 3.02	01.80 \pm 0.20	49.77 \pm 2.14	00.00 \pm 0.00	71.32 \pm 2.13	<u>01.23</u> \pm 0.30	33.10 \pm 4.42	<u>49.83</u> \pm 2.15
MCFS	56.64 \pm 1.58	00.23 \pm 0.13	37.64 \pm 7.80	00.18 \pm 0.13	76.78 \pm 0.61	00.89 \pm 0.37	20.30 \pm 4.59	47.10 \pm 0.68
NDFS	54.65 \pm 4.23	01.71 \pm 1.40	49.87 \pm 0.46	01.51 \pm 0.54	73.43 \pm 2.00	00.67 \pm 0.07	29.73 \pm 1.19	44.38 \pm 1.31
SPFS	56.91 \pm 1.62	01.21 \pm 0.48	35.41 \pm 6.24	00.07 \pm 0.01	77.22 \pm 0.77	00.90 \pm 0.06	27.09 \pm 6.04	48.24 \pm 0.76
SRFS	53.44 \pm 4.68	<u>02.96</u> \pm 0.89	49.62 \pm 0.59	00.21 \pm 0.32	77.27 \pm 0.84	00.09 \pm 0.09	29.16 \pm 1.37	49.60 \pm 1.10
UDFS	56.25 \pm 2.72	01.13 \pm 0.26	29.13 \pm 4.92	<u>02.13</u> \pm 1.14	74.87 \pm 1.52	00.83 \pm 0.26	25.43 \pm 3.63	42.00 \pm 1.00
DLUFS	<u>57.63</u> \pm 1.84	07.59 \pm 9.66	51.75 \pm 0.34	02.43 \pm 0.23	78.02 \pm 0.71	02.24 \pm 0.48	36.65 \pm 1.07	56.20 \pm 0.62

- In most cases, DLUFS outperforms the Baseline, which would lead to the superiority of selecting features over learning on all features. Therefore, feature selection process improves the efficiency (i.e. ease of computation) and performance (i.e. better results).
- DLUFS obtain good results as the best or the second-best (after Baseline) aligned with the good performance of LDSSL, SRFS which could be due to data reconstruction.
- Our method achieves good performance in most datasets. It can be attributed to low-rank representation as SRFS.

Furthermore, the performance of the proposed method is demonstrated by selecting various number of features from 50 to 300. Fig. 2 and Fig. 3 show the results for different numbers of selected features, ignoring LS, MCFS, SPFS, and UDFS due to weak results. Obviously, DLUFS achieves the best performance compared to the competitors even when the number of selected features is varied.

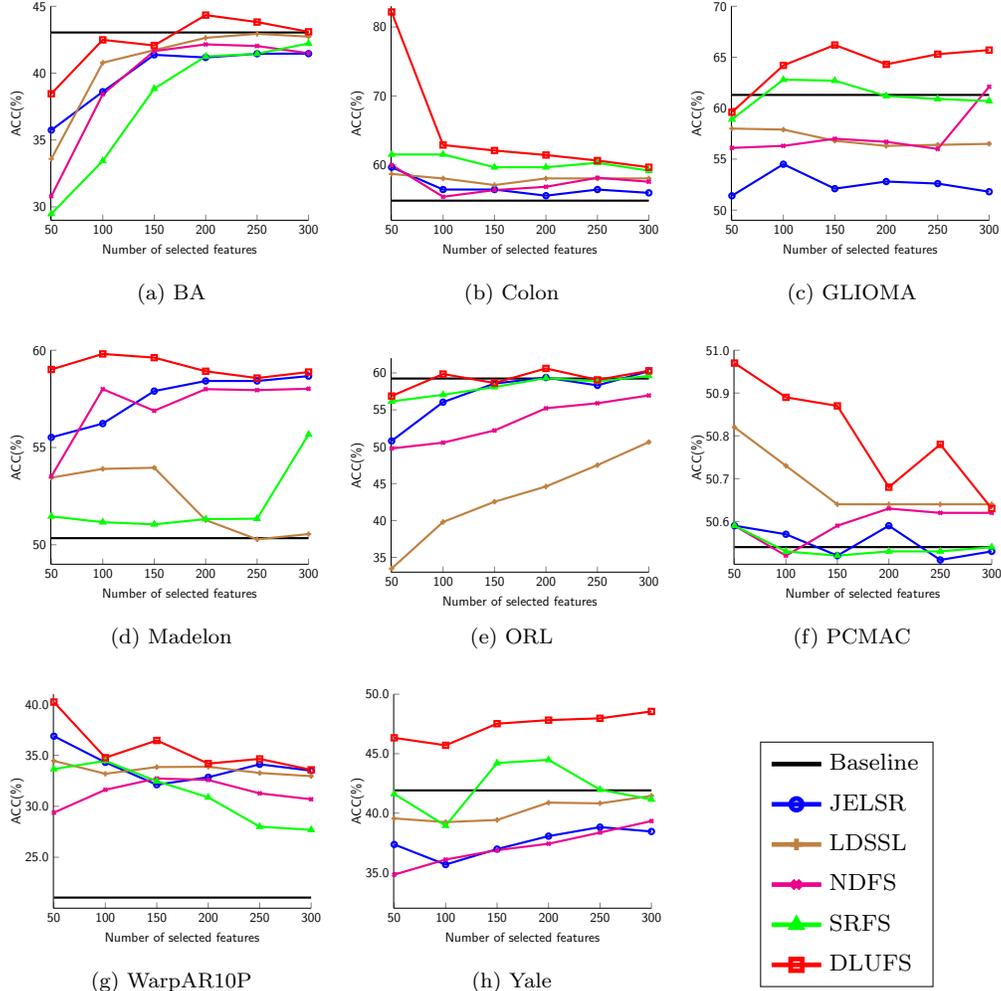


Figure 2: The achieved results in ACC measure by selecting different numbers of features.

5.5. Parameter sensitivity and convergence study

First, the sensitivity of the parameters is investigated in DLUFS. The main parameters of the method are α and λ , by considering the objective function of DLUFS (Eq. (7)). We explore the set of candidate values for tuning parameters α and λ from $\{10^{-4}, 10^{-2}, 1, 10^2, 10^4\}$. The datasets of Table 2 are used to investigate the affect of variations in the main tuning parameters. We report the obtained results of DLUFS by considering ACC and NMI measures. Fig. 4 shows ACC values based on different settings

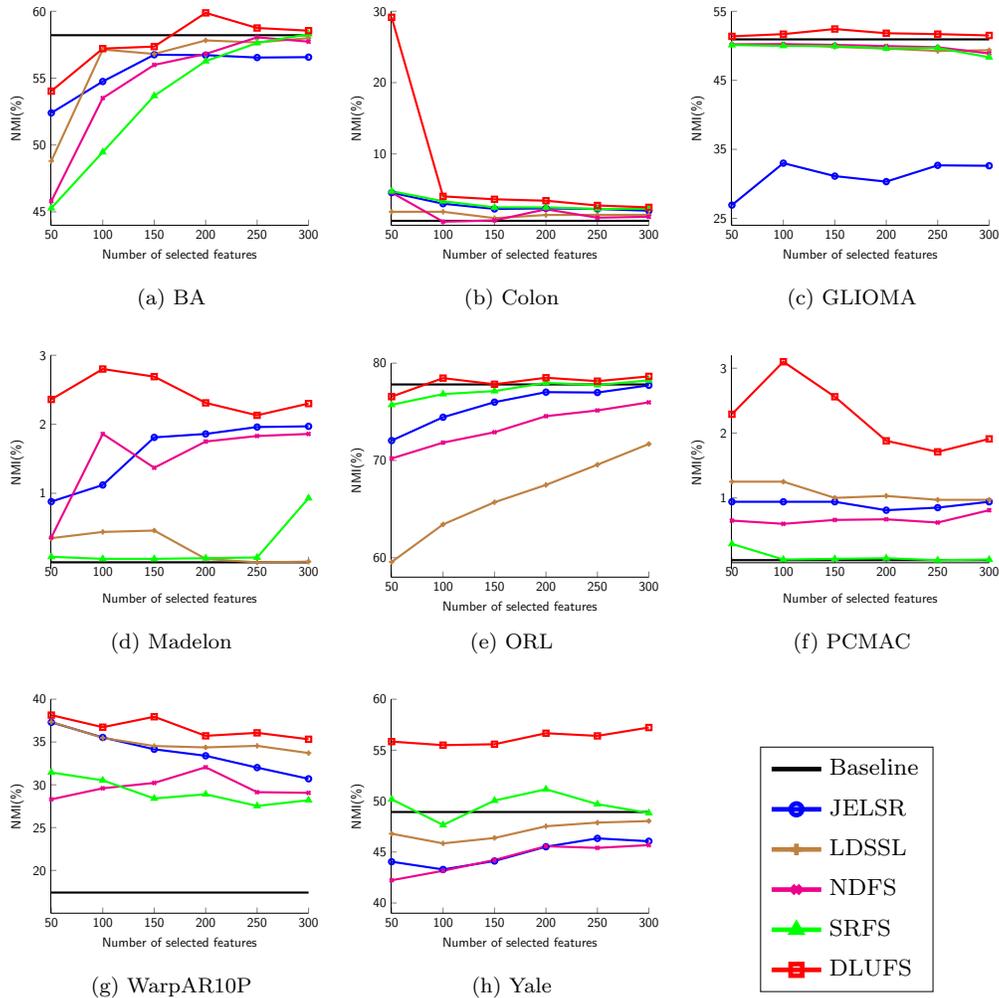


Figure 3: The achieved results in NMI measure by selecting different numbers of features.

of α and λ parameters, and Fig. 5 represents NMI values. The results indicate that there exists a slight sensitivity to these main parameters in the performance of DLUFS.

Next, the convergence behavior of the DLUFS algorithm is experimentally demonstrated on the datasets of Table 2. The convergence speed of the objective function is depicted versus the number of iterations in Fig. 6. These findings from Fig. 6 reveal the efficiency of the proposed algorithm based on the convergence speed.

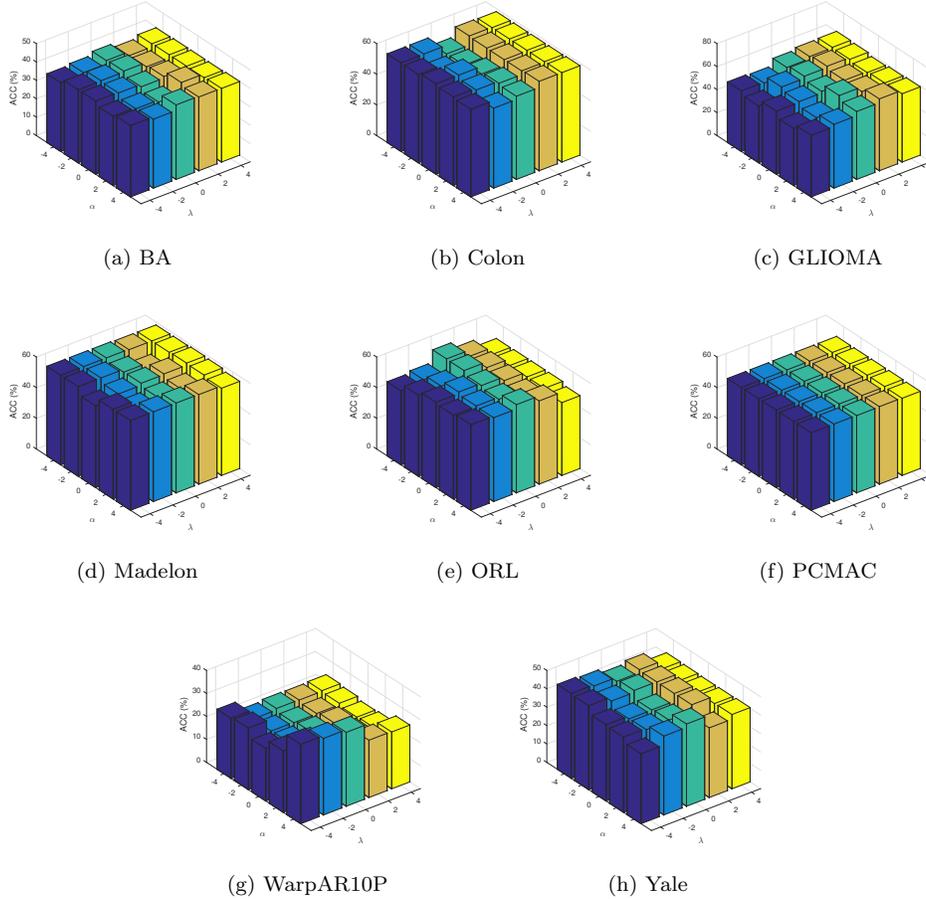


Figure 4: Performance of DLUFS in ACC measure using different values of the tuning parameters α and β in \log_{10} .

6. Conclusion

In this work, we proposed a new unsupervised feature selection method based on dictionary learning to learn a new data representation in a basis space. We devise a low-rank constraint on the basis matrix to preserve the feature correlation along with subspace learning. Spectral analysis was employed to consider the sample similarities in the learned data representation matrix. Moreover, an $\ell_{2,1}$ -norm regularization was applied in our primary objective function to discard uninformative features. We presented a unified framework based on the important characteristics to select features. An

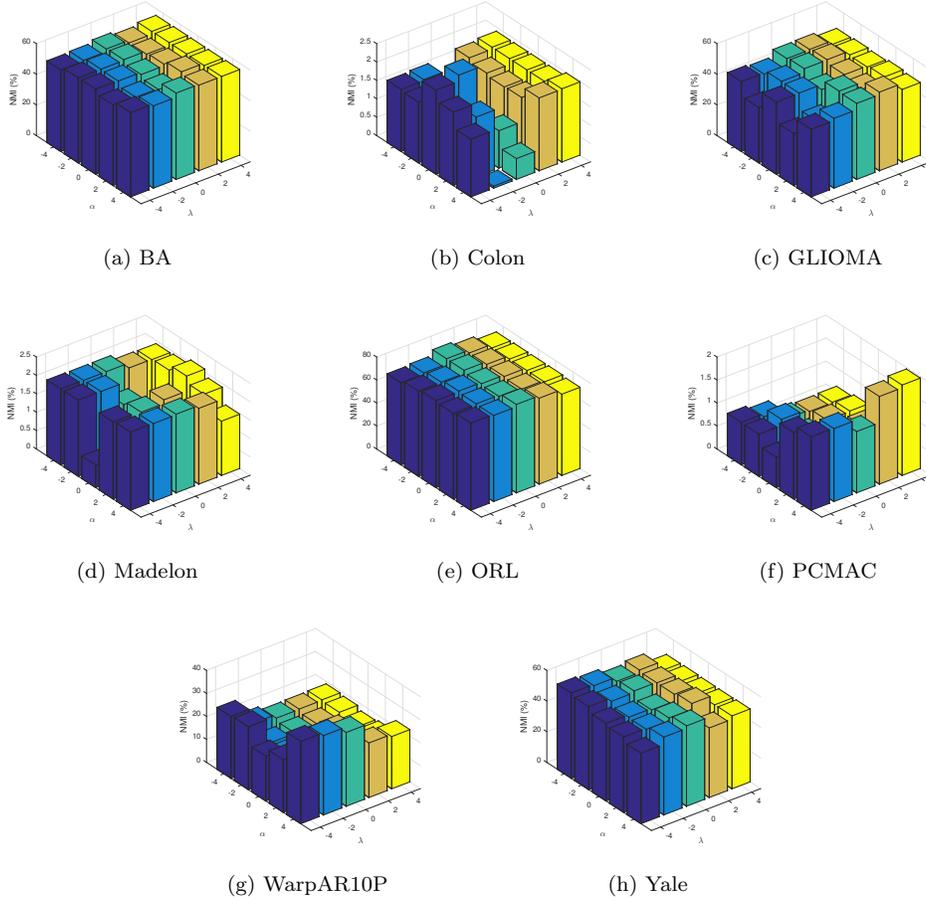


Figure 5: Performance of DLUFS in NMI measure using different values of the tuning parameters α and β in \log_{10} .

efficient numerical algorithm is introduced for the proposed method. The performance of the proposed method was investigated through state-of-the-art methods by the aid of a variety of standard datasets. The attained results verified the strength of the proposed approach in terms of accuracy and speed of convergence.

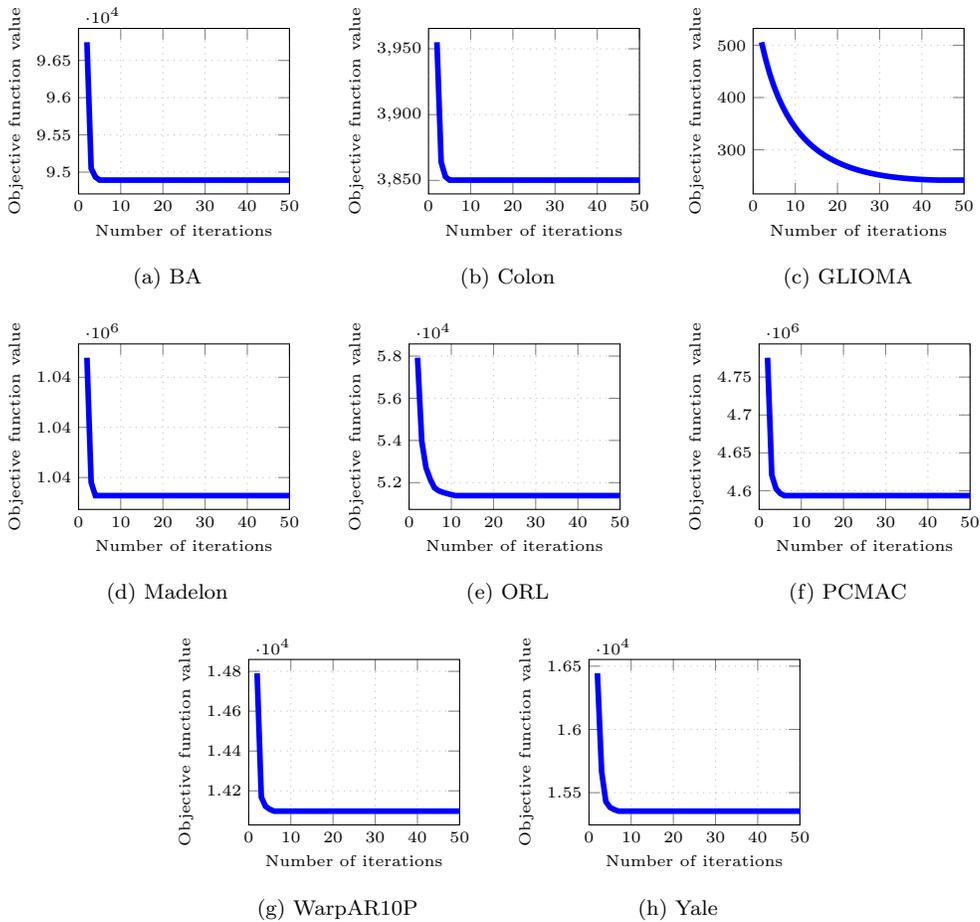


Figure 6: Convergence curve of DLUFS on different datasets.

References

- et al., C. L. N. (2003). Gene Expression-based Classification of Malignant Gliomas Correlates Better with Survival than Histological Classification. *Cancer Research*, 63, 1602–1607.
- Alon, U., Barkai, N., Notterman, D. A., Gish, K., Ybarra, S., Mack, D., & Levine, A. J. (1999). Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proceedings of the National Academy of Sciences*, 96, 6745–6750.

- Bartels, R. H., & Stewart, G. W. (1972). Solution of the matrix equation $AX + XB = C$ [F4]. *Communications of the ACM*, 15, 820–826.
- Belhumeur, P., Hespanha, J., & Kriegman, D. (1997). Eigenfaces vs. Fisherfaces: recognition using class specific linear projection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19, 711–720.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag.
- Cai, D., He, X., Han, J., & Zhang, H.-J. (2006). Orthogonal Laplacianfaces for Face Recognition. *IEEE Transactions on Image Processing*, 15, 3608–3614.
- Cai, D., Zhang, C., & He, X. (2010). Unsupervised Feature Selection for Multi-cluster Data. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining KDD '10* (pp. 333–342). New York, NY, USA: ACM.
- Chen, L., & Huang, J. Z. (2012). Sparse Reduced-Rank Regression for Simultaneous Dimension Reduction and Variable Selection. *Journal of the American Statistical Association*, 107, 1533–1545.
- Ding, D., Xia, F., Yang, X., & Tang, C. (2020). Joint dictionary and graph learning for unsupervised feature selection. *Applied Intelligence*, 50, 1379–1397.
- Dy, J. G., & Brodley, C. E. (2004). Feature Selection for Unsupervised Learning. *Journal of Machine Learning Research*, 5, 845–889.
- Farahat, A. K., Ghodsi, A., & Kamel, M. S. (2013). Efficient greedy feature selection for unsupervised learning. *Knowledge and Information Systems*, 35, 285–310.
- Fukunaga, K. (1990). *Introduction to statistical pattern recognition*. Computer science and scientific computing (2nd ed.). Boston: Academic Press.
- He, X., Cai, D., & Niyogi, P. (2005). Laplacian Score for Feature Selection. In *Proceedings of the 18th International Conference on Neural Information Processing Systems NIPS'05* (pp. 507–514). Cambridge, MA, USA: MIT Press.

- Hoseini, E., & Mansoori, E. G. (2019). Unsupervised feature selection in linked biological data. *Pattern Analysis & Applications*, 22, 999–1013.
- Hou, C., Nie, F., Li, X., Yi, D., & Wu, Y. (2014). Joint Embedding Learning and Sparse Regression: A Framework for Unsupervised Feature Selection. *IEEE Transactions on Cybernetics*, 44, 793–804.
- Krzanowski, W. J. (1987). Selection of Variables to Preserve Multivariate Data Structure, Using Principal Components. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 36, 22–33.
- Lang, K. (1995). Newsweeder: Learning to filter netnews. In *Proceedings of the Twelfth International Conference on Machine Learning* (pp. 331–339).
- Li, J., Cheng, K., Wang, S., Morstatter, F., Trevino, R. P., Tang, J., & Liu, H. (2017a). Feature Selection: A Data Perspective, <http://featureselection.asu.edu/>. *ACM Computing Surveys*, 50, 94:1–94:45.
- Li, J., Tang, J., & Liu, H. (2017b). Reconstruction-based Unsupervised Feature Selection: An Embedded Approach. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence IJCAI'17* (pp. 2159–2165). AAAI Press.
- Li, Z., Yang, Y., Liu, J., Zhou, X., & Lu, H. (2012). Unsupervised Feature Selection Using Nonnegative Spectral Analysis. In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence AAAI'12* (pp. 1026–1032). AAAI Press.
- Lovasz, L. (1986). *Matching Theory (North-Holland Mathematics Studies)*. Oxford, UK, UK: Elsevier Science Ltd.
- Lu, Q., Li, X., & Dong, Y. (2018). Structure Preserving Unsupervised Feature Selection. *Neurocomputing*, 301, 36–45.
- Masaeli, M., Yan, Y., Cui, Y., Fung, G., & Dy, J. (2010). Convex Principal Feature Selection. In *Proceedings of the 2010 SIAM International Conference on Data Mining Proceedings* (pp. 619–628). Society for Industrial and Applied Mathematics.

- Murphy, K. P. (2012). *Machine learning: a probabilistic perspective*. MIT press.
- Nie, F., Huang, H., Cai, X., & Ding, C. (2010). Efficient and Robust Feature Selection via Joint L_{2,1}-norms Minimization. In *Proceedings of the 23rd International Conference on Neural Information Processing Systems - Volume 2 NIPS'10* (pp. 1813–1821). USA: Curran Associates Inc.
- Nie, F., Xiang, S., Jia, Y., Zhang, C., & Yan, S. (2008). Trace ratio criterion for feature selection. In *Proceedings of the 23rd national conference on Artificial intelligence - Volume 2 AAAI'08* (pp. 671–676). Chicago, Illinois: AAAI Press.
- Pandit, A. A., Pimpale, B., & Dubey, S. (2020). A Comprehensive Review on Unsupervised Feature Selection Algorithms. In G. Singh Tomar, N. S. Chaudhari, J. L. V. Barbosa, & M. K. Aghwariya (Eds.), *International Conference on Intelligent Computing and Smart Communication 2019 Algorithms for Intelligent Systems* (pp. 255–266). Singapore: Springer.
- Parsa, M. G., Zare, H., & Ghatee, M. (2020). Unsupervised feature selection based on adaptive similarity learning and subspace clustering. *Engineering Applications of Artificial Intelligence*, *95*, 103855.
- Rogati, M., & Yang, Y. (2002). High-performing feature selection for text classification. In *Proceedings of the eleventh international conference on Information and knowledge management CIKM '02* (pp. 659–661). New York, NY, USA: Association for Computing Machinery.
- Shang, R., Meng, Y., Wang, W., Shang, F., & Jiao, L. (2019). Local discriminative based sparse subspace learning for feature selection. *Pattern Recognition*, *92*, 219–230.
- Shi, C., Ruan, Q., Guo, S., & Tian, Y. (2015). Sparse feature selection based on L_{2,1/2}-matrix norm for web image annotation. *Neurocomputing*, *151*, 424–433.
- Yang, S., Zhang, R., Nie, F., & Li, X. (2019). Unsupervised Feature Selection Based on Reconstruction Error Minimization. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 2107–2111).

- Yang, Y., Shen, H. T., Ma, Z., Huang, Z., & Zhou, X. (2011). L2,1-norm Regularized Discriminative Feature Selection for Unsupervised Learning. In *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence - Volume Volume Two IJCAI'11* (pp. 1589–1594). AAAI Press.
- Yuan, M., & Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, *68*, 49–67.
- Zare, H., & Niazi, M. (2016). Relevant based structure learning for feature selection. *Engineering Applications of Artificial Intelligence*, *55*, 93–102.
- Zare, H., Parsa, M. G., Ghatee, M., & Alizadeh, S. H. (2020). Similarity Preserving Unsupervised Feature Selection based on Sparse Learning. In *2020 10th International Symposium on Telecommunications (IST)* (pp. 50–55).
- Zhao, Z., He, X., Cai, D., Zhang, L., Ng, W., & Zhuang, Y. (2016). Graph Regularized Feature Selection with Data Reconstruction. *IEEE Transactions on Knowledge and Data Engineering*, *28*, 689–700.
- Zhao, Z., & Liu, H. (2007). Spectral Feature Selection for Supervised and Unsupervised Learning. In *Proceedings of the 24th International Conference on Machine Learning ICML '07* (pp. 1151–1157). New York, NY, USA: ACM.
- Zhao, Z., Wang, L., Liu, H., & Ye, J. (2013). On Similarity Preserving Feature Selection. *IEEE Transactions on Knowledge and Data Engineering*, *25*, 619–632.
- Zheng, M., Bu, J., Chen, C., Wang, C., Zhang, L., Qiu, G., & Cai, D. (2011). Graph Regularized Sparse Coding for Image Representation. *IEEE Transactions on Image Processing*, *20*, 1327–1336.
- Zhu, P., Hu, Q., Zhang, C., & Zuo, W. (2016). Coupled dictionary learning for unsupervised feature selection. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence AAAI'16* (pp. 2422–2428). Phoenix, Arizona: AAAI Press.

- Zhu, X., Li, X., Zhang, S., Ju, C., & Wu, X. (2017). Robust Joint Graph Sparse Coding for Unsupervised Spectral Feature Selection. *IEEE Transactions on Neural Networks and Learning Systems*, *28*, 1263–1275.
- Zhu, X., Zhang, S., Hu, R., Zhu, Y., & song, j. (2018). Local and Global Structure Preservation for Robust Unsupervised Spectral Feature Selection. *IEEE Transactions on Knowledge and Data Engineering*, *30*, 517–529.