

# Half-Quadratic Minimization for Unsupervised Feature Selection on Incomplete Data

Heng Tao Shen<sup>1</sup>, Senior Member, IEEE, Yonghua Zhu, Wei Zheng, and Xiaofeng Zhu<sup>2</sup>, Senior Member, IEEE

**Abstract**—Unsupervised feature selection (UFS) is a popular technique of reducing the dimensions of high-dimensional data. Previous UFS methods were often designed with the assumption that the whole information in the data set is observed. However, incomplete data sets that contain unobserved information can be often found in real applications, especially in industry. Thus, these existing UFS methods have a limitation on conducting feature selection on incomplete data. On the other hand, most existing UFS methods did not consider the sample importance for feature selection, i.e., different samples have various importance. As a result, the constructed UFS models easily suffer from the influence of outliers. This article investigates a new UFS method for conducting UFS on incomplete data sets to investigate the abovementioned issues. Specifically, the proposed method deals with unobserved information by using an indicator matrix to filter it out the process of feature selection and reduces the influence of outliers by employing the half-quadratic minimization technique to automatically assigning outliers with small or even zero weights and important samples with large weights. This article further designs an alternative optimization strategy to optimize the proposed objective function as well as theoretically and experimentally prove the convergence of the proposed optimization strategy. Experimental results on both real and synthetic incomplete data sets verified the effectiveness of the proposed method compared with previous methods, in terms of clustering performance on the low-dimensional space of the high-dimensional data.

**Index Terms**—Feature selection, half-quadratic minimization, incomplete data, robust statistics, sparse learning.

## I. INTRODUCTION

WHILE high-dimensional representation is widely applied in real applications to build high-quality machine learning models for data analysis, it still suffers some issues, such as high computation time, storage cost, and the

issue of curse of dimensionality [1], [2]. Moreover, the data set with high-dimensional representation often contains redundant features and outliers to possibly degrade the performance of data analysis. To address these issues, dimensionality reduction was designed to reduce the dimensions of high-dimensional data by removing unimportant features [3], [4].

Unsupervised feature selection (UFS) is one of the dimensionality reduction techniques to handle high-dimensional data since UFS has interpretability and obtaining labels of the samples is usually difficult in real applications. By considering the way of the model construction [5], [6], previous UFS methods can be partitioned into three groups, i.e., filter model, wrapper model, and embedded model. The filter model is the most simple and efficient, whereas the embedded model usually achieves the best effectiveness [1]. Hence, this article is focusing on investigating a new embedded model of UFS methods to deal with the issue of high-dimensional data under the setting of unsupervised learning.

Feature selection first assumes that different features have different importance for data analysis, i.e., feature importance for short, and then keeps important features and removes unimportant features to reduce the dimensions of high-dimensional data. Actually, samples also have such characteristics. Specifically, different samples produce different influences for data analysis, i.e., sample importance for short. For example, outliers may have smaller importance compared with other samples because outliers usually result in larger residual than other samples [3], [7]. Moreover, both feature importance and sample importance are necessary for data analysis [8]. Otherwise, the performance of data analysis will be influenced if only considering each of them [9], [10]. However, most previous UFS methods do not take sample importance into account for feature selection [1], [11].

The incomplete data set with unobserved information in some samples is very normal in practical applications. In particular, an industrial data set may miss 90% of the whole information [12], [13]. On the other hand, most techniques of machine learning are designed to conduct data analysis on complete data sets where the whole information is observed. Hence, these techniques cannot be directly applied in incomplete data sets. Recently, a number of solutions have been proposed to deal with unobserved data by outputting the complete data set with different strategies so that existing machine learning techniques can be employed [12], [14], [15].

The classic solutions to deal with unobserved information include case-deletion methods and imputation methods [12], [15]. The case-deletion method removes incomplete samples

Manuscript received September 17, 2019; revised May 11, 2020; accepted July 11, 2020. Date of publication July 30, 2020; date of current version July 7, 2021. This work was supported in part by the National Key Research and Development Program of China under Grant 2018AAA0102200, in part by the National Natural Science Foundation of China under Grant 61632007 and Grant 61876046, and in part by the Sichuan Science and Technology Program under Grant 2018GZDZX0032. (Corresponding author: Xiaofeng Zhu.)

Heng Tao Shen and Xiaofeng Zhu are with the Center for Future Media, University of Electronic Science and Technology of China, Chengdu 611731, China, and also with the School of Computer Science and Technology, University of Electronic Science and Technology of China, Chengdu 611731, China (e-mail: shenhengtao@hotmail.com; seanzhuxf@gmail.com).

Yonghua Zhu and Wei Zheng are with the Guangxi Key Laboratory of Multi-Source Information Mining and Security, Guangxi Normal University, Guilin 541004, China.

This article has supplementary downloadable material available at <https://ieeexplore.ieee.org>, provided by the authors.

Color versions of one or more of the figures in this article are available online at <https://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2020.3009632

2162-237X © 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.

See <https://www.ieee.org/publications/rights/index.html> for more information.

(i.e., the samples with unobserved information) from the data set and only uses the complete samples (i.e., the whole information of the samples is observed) for data analysis. Obviously, the case-deletion method overlooks the observed information in incomplete samples as the incomplete samples include observed information as well as unobserved information. Moreover, the case-deletion method uses limited information by discarding incomplete samples so that it cannot guarantee the effectiveness of data analysis [14], [15]. To address these issues, the imputation method first guesses a value for unobserved information (i.e., imputation) and then utilizes all complete samples (including the incomplete samples whose unobserved information has been imputed) for data analysis. Although the imputation method utilizes more information (i.e., the incomplete samples) compared with the case-deletion method, the imputed values of unobserved information may be unuseful or even noisy. The reason is that there is no ground truth for unobserved information so that the correctness of the imputed information cannot be evaluated [14].

This article extends our conference version [16] to propose a new UFS method by involving two following components: 1) using an indicator matrix to utilize all observed information in both complete samples and incomplete samples for conducting feature selection on incomplete data sets and 2) employing the half-quadratic minimization technique [9], [11] to flexibly reduce the influence of outliers. Specifically, the use of the indicator matrix prohibits unobserved information involving the process of feature selection and makes the use of observed information in incomplete samples, whereas the use of the half-quadratic minimization technique pushes outliers to reduce the influence for feature selection. As a result, the proposed method could simultaneously remove sample-level outliers and feature-level redundancy for conducting UFS on incomplete data sets.

Compared with previous UFS methods, we summarize the contribution of our proposed method as follows.

- 1) Previous UFS methods (see [17], [18]) use an  $\ell_{2,1}$ -norm loss function to reduce the influence of outliers. The use of the  $\ell_{2,1}$ -norm loss function is fixed and has no theoretical guarantee. On the contrary, the proposed method employs half-quadratic minimization [3], [19] to achieve the same goal so that it is flexible and has theoretical guarantee by robust statistics [9]. Specifically, the proposed method first automatically selects easy samples (i.e., the samples outputting top smallest residuals between the predictions and the ground truths<sup>1</sup>) to construct an initial feature selection model and then automatically detect other easy samples from the left samples to improve the robustness and generalization ability of this initial model. This process of easy sample selection will be repeated until this model achieves convergence or all samples have been selected. In this way, outliers will either be selected later than easy

samples or never be selected. As a result, the influence of outliers is reduced. Moreover, the number of easy samples in each iteration is data-driven to make our method flexible.

- 2) To deal with incomplete data sets, the case-deletion method uses a part of the whole information, whereas the imputation method uses the whole information but it possibly introduces noise. By contrast, our proposed method introduces an indicator matrix to remove unobserved information from the process of feature selection so that all observed information can be used and less noise is involved.
- 3) To the best of our knowledge, the proposed method is the first work to integrate feature selection, dealing with unobserved data, and half-quadratic minimization, in a unified framework since previous works were only designed to focus on a part of these three tasks [20], [21]. Moreover, experimental results on real and synthetic incomplete data sets verified the effectiveness of the proposed method compared with previous state-of-the-art methods. Furthermore, it is easy to extend the proposed method for dealing with unobserved data to other tasks, such as other UFS methods, supervised feature selection, and semisupervised feature selection.

## II. RELATED WORK

In this section, we review previous methods of feature selection as well as analyze the difference between our proposed method and them.

### A. Feature Selection on Complete Data

The goal of feature selection is to find a useful subset from all original features. This makes feature selection interpretable [22], [23]. Specifically, the filter models conduct a two-step strategy for feature selection by first ranking the importance of every feature and then selecting the most important features. In the two-step strategy, the optimal results of the first step cannot guarantee to output the optimal of the second step [3], [10]. The wrapper model first employs an exhaustive search method to generate the subsets of all features, i.e.,  $2^d - 1$  subsets for  $d$  features, and then applies the model of data analysis to evaluate the importance of different subsets. Obviously, exhaustive search makes the wrapper model possibly obtain an optimized subset of features, but the time complexity of the wrapper model is very high. Thus, the wrapper model prohibits applying in data sets with a large number of features. Furthermore, the wrapper model is a two-step strategy for feature selection. The embedded model integrates feature ranking with the process of feature selection in a unified framework, i.e., a one-step strategy, to automatically select important features.

Based on the label information, the embedded model can be partitioned into three subgroups, i.e., unsupervised embedded model, supervised embedded model, and semisupervised embedded model. Specifically, unsupervised embedded model employs either dictionary learning methods or self-representation methods to conduct feature selection without

<sup>1</sup>Easy samples are defined as the samples outputting top smallest residuals in each iteration. Moreover, easy samples may change in different iterations and the selection of easy samples is dependent on the residual between the ground truth and the prediction for each sample. Furthermore, easy samples in the last iteration are called important samples.

the supervision of labels. For example, the method of regularized self-representation (RSR) utilizes the  $\ell_{2,1}$ -norm loss function based on the self-representation property of features as well as uses the  $\ell_{2,1}$ -norm regularization to select features [24]. Supervised embedded model applies the  $\ell_{2,r}$ -norm loss function based on the difference between the labels and the predictions to reduce the influence of outliers as well as applies the  $\ell_{2,p}$ -norm regularization to conduct feature selection where both  $r$  and  $p$  are nonnegative tuning parameters. Moreover, their ranges are set as  $0 < r \leq 2$  and  $0 < p \leq 1$ . For example, the method of general sparsity feature selection employs the  $\ell_{2,r}$ -norm loss function and the  $\ell_{2,1}$ -norm regularization [25], whereas the method of simultaneous capped  $\ell_2$ -norm loss function and  $\ell_{2,p}$ -norm regularization minimization uses the capped  $\ell_2$ -norm loss function and the  $\ell_{2,p}$ -norm regularization [26]. In semisupervised embedded model, the method of robust feature selection (RFS) with linear discriminant analysis [8] simultaneously finds sample- outliers and feature noises by using both labeled samples and unlabeled samples to conduct robust and discriminative semisupervised feature selection.

The aforementioned feature selection methods propose different strategies of robustness to reduce the influence of outliers, i.e., the  $\ell_{2,r}$ -norm loss function and the capped  $\ell_2$ -norm loss function. In this article, we focus on employing half-quadratic minimization to conduct UFS.

### B. Feature Selection on Incomplete Data

The literature (see [15], [20]) usually designs a two-stage strategy to conduct feature selection on incomplete data sets, i.e., using imputation methods to guess a value for unobserved information and then conducting feature selection with previous techniques on imputed data sets. For example, the method of feature selection based on conditional entropy [27] first uses entropy-based uncertainty measures to deal with unobserved information and then conducts feature selection. Recently, Zhang *et al.* [15] proposed first conducting feature selection and then imputing unobserved information by decision tree techniques. The method of multicriteria-based feature selection on cost-sensitive data with missing values [28] proposed first selecting features and then using a rough set theory to deal with unobserved information.

Different from previous literature which separately imputes unobserved information and conducts feature selection, our proposed method prohibits unobserved information to involve the process of feature selection by introducing an indicator matrix so that it simultaneously handles unobserved information, reducing the influence of outliers, and conducts feature selection by utilizing all observed information.

## III. APPROACH

In this article, we use boldface uppercase letters, boldface lowercase letters, and normal italic letters, respectively, to denote matrices, vectors, and scalars.

### A. Robust Feature Selection

Sparse theory has widely been used in embedded models since it selects features by pushing the weight coefficients

of unimportant features to have small values or even zeros and the weight coefficients of important features to have large values [22]. Specifically, given a feature matrix  $\mathbf{X} = [\mathbf{x}^1; \dots; \mathbf{x}^n] = [\mathbf{x}_1, \dots, \mathbf{x}_d] \in \mathbb{R}^{n \times d}$ , traditional unsupervised sparse feature selection methods [24], [29], [30] can be formulated as

$$\min_{\mathbf{W}} \|\mathbf{X} - \mathbf{X}\mathbf{W}\|_F^2 + \lambda \|\mathbf{W}\|_{2,1} \quad (1)$$

where  $\lambda$  is a nonnegative tuning parameter,  $\mathbf{W} \in \mathbb{R}^{d \times d}$  is the weight coefficient matrix,  $\|\cdot\|_F$  is a Frobenius norm, and  $\|\mathbf{W}\|_{2,1}$  (i.e.,  $\|\mathbf{W}\|_{2,1} = \sum_i (\sum_j w_{ij}^2)^{1/2}$ ) is an  $\ell_{2,1}$ -norm sparsity regularizer, which selects important features by automatically assigning their corresponding rows of the weight coefficients  $\mathbf{W}$  with nonzero values.

Equation (1) equivalently considers every sample, so the resulting model is easily influenced by outliers [9], [31]. To address this issue, Nie *et al.* [31] proposed to replace the Frobenius norm in (1) by an  $\ell_{2,1}$ -norm loss function to have

$$\min_{\mathbf{W}} \|\mathbf{X} - \mathbf{X}\mathbf{W}\|_{2,1} + \lambda \|\mathbf{W}\|_{2,1}. \quad (2)$$

In (2), the residual of  $\mathbf{X} - \mathbf{X}\mathbf{W}$  (i.e.,  $\|\mathbf{X} - \mathbf{X}\mathbf{W}\|_{2,1}$ ) does not have a squared operator, and thus, outliers have less influence compared with Frobenius norm [31]. Compared with (1), (2) considers sample importance for feature selection. However, (2) is less flexible as it cannot guarantee the effectiveness without considering to reduce the influence of sample-level outliers. Moreover, both (1) and (2) cannot be used for dealing with incomplete data sets.

In this article, we employ the theory of half-quadratic minimization [9], [32], such as self-paced learning [33], [34] and M-estimation (i.e., maximum-likelihood-type estimation) [35] to replace the  $\ell_{2,1}$ -norm loss function in (2) with predefined robust loss functions, aimed at achieving flexibly RFS. In half-quadratic minimization, functions can be constructed so that they can first select easy samples to construct an initial model, and the left easy samples are automatically detected from the left samples joining with the former easy samples to revise the robustness and generalization of this initial model gradually until either this model is not further improved or all the samples are used up [33], [34]. By replacing the  $\ell_{2,1}$ -norm loss function in (2) with a general robust loss function  $\phi(\cdot)$ , the RFS framework can be formulated as

$$\min_{\mathbf{W}} \phi(\|\mathbf{X} - \mathbf{X}\mathbf{W}\|_F) + \lambda \|\mathbf{W}\|_{2,1}. \quad (3)$$

A number of robust loss functions have been designed in the literature [19], [36], such as Cauchy function,  $\ell_1$ - $\ell_2$  function, Welsch M-estimator, and GemanMcClure estimator. Every robust loss function  $\phi(\cdot)$  has individual characteristics different from the others and can be flexibly selected in real applications, e.g., changing  $\phi(\cdot)$  to achieve flexibility. In this way, (2) can be regarded as a special issue of (3).

### B. Robust Feature Selection on Incomplete Data

A few feature selection methods have been designed to deal with incomplete data sets. A possible method of conducting feature selection on incomplete data sets is to construct feature



selection models only using complete samples by discarding incomplete samples. Obviously, the larger the number of incomplete samples, the lower the robustness and generalization ability of the feature selection model. To address this issue, an alternative method is to first impute unobserved values with imputation methods (such as mean-value method [37] and  $k$  nearest neighbors (kNNs) imputation method [15]) and then conduct feature selection using all the samples, including the imputed incomplete samples. However, the imputation strategy for feature selection is usually unpractical. Specifically, if the imputation models are estimated perfectly, then the constructed feature selection models may be overfitting. Otherwise, the models will be underfitting with a bad imputation. Moreover, we have no idea on the ground truth of unobserved information so that it is difficult to evaluate imputation models. Furthermore, the imputed information is obtained by the complete samples, so the imputation method does not utilize the observed information in incomplete samples, so does the case-deletion method.

Based on the above observations, it may be unnecessary to impute unobserved information for feature selection. To make the use of the observed information in incomplete samples, we propose to use an indicator matrix  $\mathbf{D} \in \mathbb{R}^{n \times d}$  (which has the same size as the feature matrix  $\mathbf{X}$ ) to avoid unobserved values involving the process of feature selection. We thus propose the following objective function to conduct RFS on incomplete data sets:

$$\min_{\mathbf{W}} \phi(\|\mathbf{D} \circ (\mathbf{X} - \mathbf{XW})\|_F) + \lambda \|\mathbf{W}\|_{2,1} \quad (4)$$

where  $\mathbf{D} = [\mathbf{d}^1; \dots; \mathbf{d}^n]$  is an indicator matrix. Specifically, if the  $i$ th row and the  $j$ th column element  $x_{ij}$  are unobserved, the value of  $d_{ij}$  is 0, otherwise 1. The symbol  $\circ$  is a Hadamard product operator that conducts the elementwise multiplication between two same size matrices.

According to (4), our proposed method conducts feature selection by using all available information (without imputing unobserved values) as well as reducing the influence of outliers. Moreover, the definition of  $\mathbf{D}$  can be applied in any sparse feature selection models for conducting robust UFS using other embedded models, supervised feature selection, and semisupervised feature selection, on incomplete data sets.

### C. Proposed Objective Function

Although a number of robust loss functions have been reported in the literature, the resulting objective functions in these predefined robust functions may be optimized difficultly or inefficiently or even nonconvex. To address this, the half-quadratic minimization technique introduces an auxiliary variable  $\mathbf{v}$ .

*Lemma 1:* Given a fixed scalar  $z$ , if a differentiable function  $\phi(z)$  satisfies four conditions listed in [19], then the following holds:

$$\phi(z) = \inf_{z \in \mathbb{R}} \{ \mathbf{v}z^2 + \psi(\mathbf{v}) \} \quad (5)$$

where  $\mathbf{v}$  is a variable determined by the minimization function of the dual potential function  $\psi(\mathbf{v})$  of the differentiable loss function  $\phi(z)$ .

Based on Lemma 1, each robust loss function  $\phi(z)$  theoretically has a corresponding potential function  $\psi(\mathbf{v})$ , which is optimized by its minimization function or directly optimized by other optimization methods. In this article, for a fixed scalar  $z$ , we select a specific robust loss function, i.e., GemanMcClure loss function [38]

$$\phi(z) = \frac{\mu z^2}{\mu + z^2} \quad (6)$$

to replace  $\phi(\cdot)$  in (4) to have

$$\min_{\mathbf{W}} \frac{\mu \|\mathbf{D} \circ (\mathbf{X} - \mathbf{XW})\|_F^2}{\mu + \|\mathbf{D} \circ (\mathbf{X} - \mathbf{XW})\|_F^2} + \lambda \|\mathbf{W}\|_{2,1} \quad (7)$$

where  $\mu$  is used to tune the number of easy samples in each iteration [39]. Equation (7) can be optimized by any gradient methods. By the consideration of efficient and scalable optimization, we apply Lemma 1 to (7) and then obtain our final objective function

$$\min_{\mathbf{W}, \mathbf{v}} \sum_{i=1}^n (v_i \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W})\|_2^2 + \mu(\sqrt{v_i} - 1)^2) + \lambda \|\mathbf{W}\|_{2,1} \quad (8)$$

where  $\psi(\mathbf{v}) = \mu(\sqrt{\mathbf{v}} - 1)^2$ ,  $\mathbf{v} = [v_1, \dots, v_n] \in \mathbb{R}^n$  is an auxiliary variable and  $\mu$  and  $\lambda$  are two nonnegative tuning parameters. According to Lemma 1, (7) and (8) are equivalent with respect to the optimization of  $\mathbf{W}$  and are more flexible than either (2) or (1) since  $\mathbf{v}$  controls the number of easy samples as well as enables to define additional constraints to make the model [see (8)] more flexible. Specifically, by considering the optimization result of  $v_i$  in (10), if the residual (i.e.,  $\|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W})\|_2^2$ ) is small, then the  $i$ th sample  $\mathbf{x}^i$  can be regarded an easy sample and, thus, its weight  $v_i$  will be large. By contrast, the weight of an outlier is small due to its large residual. Moreover, in alternative iterations of our proposed Algorithm 1 to solve (8), the value of every  $v_i$  will vary with the iteratively updated  $\mathbf{W}$ . In this way, the initial model constructed in the first iteration will be updated gradually with different weights for each sample until the model achieves stability. Furthermore, the number of easy samples in every iteration can be flexibly controlled by the tuning parameter  $\mu$  according to the practical demand [39], i.e.,  $\mu \gg \bar{f}$  where  $\bar{f}$  is the average of  $\{f^i = \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W})\|_2^2, i = 1, \dots, n\}$ .

### D. Optimization

The optimization of (8) is challenging due to the introduction of an auxiliary variable  $\mathbf{v}$ , the convex but no-smooth constraint on  $\mathbf{W}$  (i.e.,  $\|\mathbf{W}\|_{2,1}$ ), and the indicator matrix  $\mathbf{D}$ . In this article, we propose an alternative optimization strategy based on the framework of iteratively reweighted least squares (IRLS) [40] to optimize (8): 1) update  $\mathbf{v}$  by fixing  $\mathbf{W}$  and 2) update  $\mathbf{W}$  by fixing  $\mathbf{v}$ . We list the pseudo in Algorithm 1.

1) *Update  $\mathbf{v}$  by Fixing  $\mathbf{W}$ :* While  $\mathbf{W}$  is fixed, (8) can be changed to

$$\min_{\mathbf{v}} \sum_{i=1}^n (v_i \|\mathbf{d}^i \circ \mathbf{x}^i - \mathbf{d}^i \circ (\mathbf{x}^i \mathbf{W})\|_2^2 + \mu(\sqrt{v_i} - 1)^2). \quad (9)$$

**Algorithm 1** Proposed Algorithm for Optimizing (8)

**Input:**  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $\lambda$ ;  
**Output:**  $\mathbf{W} \in \mathbb{R}^{d \times d}$ ,  $\mathbf{v} \in \mathbb{R}^{n \times 1}$ ;  
1. Initialize the binary matrix  $\mathbf{D}$ ;  
2. Initialize  $\mathbf{W}$  by conducting ridge regression on observed samples;  
3. Initialize  $\mathbf{Q} = \mathbf{I} \in \mathbb{R}^{d \times d}$ ;  
4. Initialize  $\tilde{\mathbf{f}}$  as  $\tilde{\mathbf{f}} = \sum_{i=1}^n \mathbf{f}^i$ ;  
5. Initialize  $\mu$  as  $\mu = 1$ ;  
6. **repeat:**  
6.1. Update  $\mathbf{v}$  via Eq. (10);  
6.2. Update  $\mathbf{Q}$  via Eq. (13);  
6.3. Update  $\mathbf{W}$  via Eq. (16);  
6.4. Update  $\mu$  via  $\mu = \max(\frac{\mu}{2}, \frac{\tilde{\mathbf{f}}}{2})$ ;  
**until** Eq. (8) converges

According to either [39] or directly conducting the derivative with respect to  $v_i$ , the closed-form solution of  $v_i$  ( $i = 1, \dots, n$ ) is

$$v_i = \left( \frac{\mu}{\mu + \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W})\|_2^2} \right)^2. \quad (10)$$

2) Update  $\mathbf{W}$  by Fixing  $\mathbf{v}$ : While  $\mathbf{v}$  is fixed, (8) becomes

$$\min_{\mathbf{W}} \sum_{i=1}^n v_i \|\mathbf{d}^i \circ \mathbf{x}^i - \mathbf{d}^i \circ (\mathbf{x}^i \mathbf{W})\|_2^2 + \lambda \|\mathbf{W}\|_{2,1}. \quad (11)$$

By defining  $\tilde{\mathbf{X}} = \text{diag}(\sqrt{\mathbf{v}})\mathbf{X}$  and  $\mathbf{B} = \mathbf{D} \circ \tilde{\mathbf{X}}$  where  $\text{diag}(\mathbf{v})$  transfers the vector  $\mathbf{v}$  to a diagonal matrix whose  $i$ th diagonal element is the value of  $v_i$ , (11) is thus changed to its matrix form as follows:

$$\min_{\mathbf{W}} \|\mathbf{B} - \mathbf{D} \circ (\tilde{\mathbf{X}}\mathbf{W})\|_F^2 + \lambda \text{tr}(\mathbf{W}^T \mathbf{Q} \mathbf{W}) \quad (12)$$

where  $\mathbf{Q}$  is a diagonal matrix whose diagonal element is defined as

$$q_{i,i} = \frac{1}{2\|\mathbf{w}^i\|_2}, \quad i = 1, \dots, d. \quad (13)$$

According to the IRLS framework, (12) can be rewritten as

$$\begin{aligned} \min_{\mathbf{W}} & \text{tr}(\mathbf{B}^T \mathbf{B} - 2\mathbf{D}^T \circ (\tilde{\mathbf{X}}\mathbf{W})^T \mathbf{B}) \\ & + \text{tr}(\mathbf{D}^T \circ (\tilde{\mathbf{X}}\mathbf{W})^T ((\tilde{\mathbf{X}}\mathbf{W}) \circ \mathbf{D})) + \lambda \text{tr}(\mathbf{W}^T \mathbf{Q} \mathbf{W}). \end{aligned} \quad (14)$$

The derivative of (14) with respect to  $\mathbf{W}$  is

$$-2\mathbf{X}^T (\mathbf{D} \circ \mathbf{B}) + 2(\tilde{\mathbf{X}} \circ \mathbf{D})^T \tilde{\mathbf{X}} \mathbf{W} + 2\lambda \mathbf{Q} \mathbf{W} = 0. \quad (15)$$

Hence, the solution of  $\mathbf{W}$  is

$$\mathbf{W} = (\mathbf{B}^T \tilde{\mathbf{X}} + \lambda \mathbf{Q})^{-1} \tilde{\mathbf{X}}^T \mathbf{B}. \quad (16)$$

### E. Convergence Analysis

By denoting  $\mathbf{v}^{(t)}$  and  $\mathbf{W}^{(t)}$ , respectively, as the  $t$ th iteration result of  $\mathbf{v}$  and  $\mathbf{W}$ , according to Algorithm 1, (8) can be

changed to

$$\begin{aligned} J(\mathbf{W}^{(t)}, \mathbf{v}^{(t)}) &= \sum_{i=1}^n (v_i^{(t)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t)})\|_2^2 \\ &+ \mu (\sqrt{v_i^{(t)}} - 1)^2) + \lambda \|\mathbf{W}^{(t)}\|_{2,1}. \end{aligned} \quad (17)$$

To prove the convergence of the proposed optimization Algorithm 1, we first list the following lemma and theorem based on [10] and [22].

*Lemma 2:* For any positive real numbers  $p$  and  $q$ , the following inequality always holds:

$$\sqrt{p} - \frac{p}{2\sqrt{q}} \leq \sqrt{q} - \frac{q}{2\sqrt{q}}. \quad (18)$$

*Theorem 1:* The objective function value of (8) monotonically decreases until Algorithm 1 converges.

*Proof:* With the fixed  $\mathbf{W}^{(t)}$  and according to the half-quadratic minimization theory [32], we have

$$\begin{aligned} & \sum_{i=1}^n (v_i^{(t+1)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t)})\|_2^2 \\ &+ \mu (\sqrt{v_i^{(t+1)}} - 1)^2) + \lambda \|\mathbf{W}^{(t)}\|_{2,1} \\ &\leq \sum_{i=1}^n (v_i^{(t)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t)})\|_2^2 \\ &+ \mu (\sqrt{v_i^{(t)}} - 1)^2) + \lambda \|\mathbf{W}^{(t)}\|_{2,1}. \end{aligned} \quad (19)$$

While fixing  $\mathbf{v}^{(t+1)}$ , we have the following inequality according to the IRLS framework and (13):

$$\begin{aligned} & \sum_{i=1}^n (v_i^{(t+1)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t+1)})\|_2^2 \\ &+ \mu (\sqrt{v_i^{(t+1)}} - 1)^2) + \lambda \sum_{j=1}^d \frac{\|\mathbf{w}^{j(t+1)}\|_2^2}{2\|\mathbf{w}^{j(t)}\|_2} \\ &\leq \sum_{i=1}^n (v_i^{(t+1)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t)})\|_2^2 \\ &+ \mu (\sqrt{v_i^{(t+1)}} - 1)^2) + \lambda \sum_{j=1}^d \frac{\|\mathbf{w}^{j(t)}\|_2^2}{2\|\mathbf{w}^{j(t)}\|_2} \end{aligned} \quad (20)$$

where  $\mathbf{w}^{j(t)}$  and  $\mathbf{w}^{j(t+1)}$ , respectively, are the  $j$ th row of  $\mathbf{W}^{(t)}$  and  $\mathbf{W}^{(t+1)}$ . According to Lemma 2, we have

$$\|\mathbf{w}^{j(t+1)}\|_2 - \frac{\|\mathbf{w}^{j(t+1)}\|_2^2}{2\|\mathbf{w}^{j(t)}\|_2} \leq \|\mathbf{w}^{j(t)}\|_2 - \frac{\|\mathbf{w}^{j(t)}\|_2^2}{2\|\mathbf{w}^{j(t)}\|_2}. \quad (21)$$

TABLE I

SUMMARIZATION OF USED DATA SETS, WHERE 0.26% (FEATURES) INDICATES THAT THE RATIO OF INCOMPLETE FEATURES IS 0.26% OF ALL FEATURES AND 28.60% (SAMPLES) INDICATES THAT THE RATIO OF INCOMPLETE SAMPLES IS 28.60% OF ALL SAMPLES

Name	# (Samples)	# (Features)	# (Classes)	# (Incomplete feature/sample ratio)	# (Data types)
Advertisement	3279	1558	2	0.26% (features), 28.60% (samples)	Image data
Arrhythmia	452	279	13	1.79% (features), 84.96% (samples)	Medical data
Cvpu	2215	145	9	28.28% (features), 94.99% (samples)	Image data
Mice	1567	590	2	91.19% (features), 48.89% (samples)	Gene data
Cane	1080	856	9	Complete data set	Text data
Usps	2007	256	10	Complete data set	Digit data
Yale	165	1024	15	Complete data set	Face image data
Yeast	1484	1470	10	Complete data set	Gene data
Palm	2000	256	100	Complete data set	Text data
Cifar	60000	3072	10	Complete data set	Images data
Connect	67557	126	3	Complete data set	Game data
Mnist	70000	784	10	Complete data set	Images data
Vehicle	78823	100	3	Complete data set	Sensor data
Aloi	108000	128	1000	Complete data set	Image data

By combining (20) with (21), we have

$$\begin{aligned}
& \sum_{i=1}^n (v_i^{(t+1)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t+1)})\|_2^2 \\
& + \mu(\sqrt{v_i^{(t+1)}} - 1)^2) + \lambda \sum_{j=1}^d \|\mathbf{w}^{j(t+1)}\|_2 \\
& \leq \sum_{i=1}^n (v_i^{(t+1)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t)})\|_2^2 \\
& + \mu(\sqrt{v_i^{(t+1)}} - 1)^2) + \lambda \sum_{j=1}^d \|\mathbf{w}^{j(t)}\|_2. \quad (22)
\end{aligned}$$

Thus, we obtain final inequality by integrating (19) with (22) to have

$$\begin{aligned}
& \sum_{i=1}^n (v_i^{(t+1)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t+1)})\|_2^2 \\
& + \mu(\sqrt{v_i^{(t+1)}} - 1)^2) + \lambda \|\mathbf{W}^{(t+1)}\|_2 \\
& \leq \sum_{i=1}^n (v_i^{(t)} \|\mathbf{d}^i \circ (\mathbf{x}^i - \mathbf{x}^i \mathbf{W}^{(t)})\|_2^2 \\
& + \mu(\sqrt{v_i^{(t)}} - 1)^2) + \lambda \|\mathbf{W}^{(t)}\|_2. \quad (23)
\end{aligned}$$

According to (23), (17) is nonincreasing at each iteration in Algorithm 1. Hence, Algorithm 1 converges.  $\square$

#### F. Complexity Analysis

The complexity of our proposed Algorithm 1 is focused on the variables of  $\mathbf{v}$  and  $\mathbf{W}$ , which have closed-form solutions. Specifically, the optimization of the variables  $\mathbf{v}$  is with linear complexity [i.e.,  $O(n)$ ] to the sample size  $n$  according to the literature [32], and the complexity of  $\mathbf{W}$  is  $O(r^3)$ , where  $r$  denotes the number of the rank of the feature matrix, usually  $r < d$ , where  $d$  is the number of features. Moreover, our experimental results showed that our Algorithm 1 usually achieved convergence within 20 iterations, so the complexity of our proposed optimization Algorithm 1 is  $\max\{O(tr^3), O(tn)\}$ , where  $t$  is the iteration time. The space complexity of our method is linear to the sample size.

In our experiments, we found that our proposed Algorithm 1 is insensitive to the initialization of variables and our method achieved fast convergence with the initialization settings in Algorithm 1.

#### IV. EXPERIMENTS

In this section, we evaluated our proposed method by comparing with four previous feature selection methods on four real incomplete data sets and ten synthetic data sets in terms of their clustering performance on the low-dimensional data.

##### A. Data Sets

In our experiments, we used two kinds of data sets to evaluate the effectiveness of our proposed method, i.e., incomplete data sets and complete data sets. Specifically, we randomly missed different ratios of information to generate different data sets to demonstrate the robustness of our method, while we used real incomplete data sets to analyze the generalization ability of our proposed method.

The complete data sets (i.e., Cane, Usps, Yale, Yeast, Palm, Cifar, Connect, Mnist, Vehicle, and Aloi) and incomplete data sets (i.e., Advertisement, Arrhythmia, Cvpu, and Mice) were downloaded from the LIBSVM data website<sup>2</sup> and the public UCI website<sup>3</sup>. The used data sets had different scales and have been applied in different domains. The diversity of used data sets makes our experimental results convincing. We listed the details of used data sets in Table I.

In this article, we denoted the incomplete feature ratio as the percentage of incomplete features out of all features and the incomplete sample ratio as the percentage of incomplete samples out of all samples. The incomplete sample ratios on incomplete data sets, i.e., Advertisement, Arrhythmia, Cvpu, and Mice, are 28.60%, 84.96%, 94.99%, and 48.89%, respectively, whereas the incomplete feature ratios on these data sets are 0.26%, 1.79%, 28.28%, and 91.19%, respectively. For complete data sets, we randomly marked observed information as unobserved information by setting the incomplete sample

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

<sup>3</sup><https://archive.ics.uci.edu/ml/datasets.php>

ratio as 0% (i.e., complete data set), 10%, 30%, 50%, 70%, and 90% and then analyzed the results of feature selection by removing 50% of all features as such a range of features makes all methods achieve good results in our experiments. It is noteworthy that we did not report the results of all methods with different incomplete feature ratios as such results are similar to the ones with different incomplete sample ratios.

### B. Comparison Methods

We listed the details of the comparison methods as follows.

- 1) Laplacian score (LaPscore) [41] is a filter model to evaluate the importance of every feature based on the Laplacian score of every feature.
- 2) RSR [24] first uses the self-representation of the feature level to reconstruct each feature by all features and then uses an  $\ell_{2,1}$ -norm regularization to conduct feature selection.
- 3) General framework for sparsity regularized (GSR) [25] is a general sparse embedded model to simultaneously conduct feature selection and outliers reduce by tuning parameters.
- 4) RFS, i.e., (2), is a classic embedded model of feature selection. RFS is used to verify the effectiveness of reducing outlier influence compared with our proposed method.

In our experiments, we separated every incomplete data set into two subsets, i.e., incomplete set (IS) including all incomplete samples and observed set (OS) including all observed samples. First, we denoted the method of conducting  $k$ -means clustering with original features of OS as Baseline. We also employed a filter model (i.e., LaPscore) and three embedded feature selection methods, i.e., RSR, GSR, and RFS, to conduct feature selection on OS, and then conducted  $k$ -means clustering on OS with selected features. Second, we used the information in OS to impute unobserved values in IS by the imputation method, i.e., mean-value imputation method (Mean) and kNN imputation method (kNN) [12], [15], and then conducted feature selection on the combination of OS and IS (i.e.,  $OS \cup IS$ ) by GSR, i.e., GSR\_mean and GSR\_knn, followed by conducting  $k$ -means clustering on OS with the selected features. Third, we used our proposed (8) to conduct feature selection on  $OS \cup IS$  and then conducted  $k$ -means clustering on OS with the selected features.

### C. Experimental Setting

The codes of all methods in our experiments were implemented with MATLAB, and all experiments were run under the environment of a Windows server with 128-GB RAM and Intel Core i9-7900x CPU.

We employed the tenfold cross-validation scheme to repeat every method on every data set ten times. We reported the average results of these ten times, each of which is the average of ten results of  $k$ -means clustering.<sup>4</sup> We set the ranges of the

parameters of the comparison methods according to the corresponding literature to guarantee that all comparison methods achieved their best results and set the ranges of the parameter (i.e.,  $\lambda$ ) of our method in (8) as  $\{10^{-3}, 10^{-2}, \dots, 10^3\}$ . We further set the number of clusters in  $k$ -means clustering as the number of real classes of the data sets.

### D. Evaluation Metrics

We used two evaluation metrics [i.e., accuracy (ACC) and normalized mutual information (NMI)] to evaluate the clustering performance of all methods. ACC indicates the percentage of correctly classified samples, that is

$$ACC = \frac{N_c}{N} \quad (24)$$

where  $N$  denotes the number of the samples and  $N_c$  is the number of the correctly classified samples.

NMI uncovers a correlation between the predictions and the ground truths, that is

$$NMI = \frac{2I(\mathbf{X}, \mathbf{Y})}{H(\mathbf{X}) + H(\mathbf{Y})} \quad (25)$$

where  $I(\mathbf{X}, \mathbf{Y})$  denotes mutual information between the predicted labels and the real labels, i.e.,  $I(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^{|\mathbf{X}|} \sum_{j=1}^{|\mathbf{Y}|} P(i, j) \log((P(i, j)/P(i)P'(j)))$ , where  $P(i, j) = (|\mathbf{X}_i \cap \mathbf{Y}_j|/N)$ , and  $|\mathbf{X}|$  and  $|\mathbf{Y}|$  imply the cardinality of  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.  $H(\cdot)$  is the entropy, i.e.,  $H(\cdot) = -\sum_{i=1}^{|\cdot|} P(i) \log(P(i))$ .

Besides, we conducted the paired-sample t-tests (at 95% significance level) between our method and every comparison method, in terms of ACC and NMI. Moreover, the symbols “\*” and “\*\*\*,” respectively, indicate that our method has statistically significant difference with  $p < 0.05$  and  $p < 0.001$  on the paired-sample t-tests at 95% significance level compared with the comparison method.

### E. Real Incomplete Data Sets

We reported clustering performance of all methods with different feature ratios (i.e.,  $\{10\%, 20\%, \dots, 90\%\}$  of all the features) on four real incomplete data sets, i.e., Advertisement, Arrhythmia, Cypu, and Mice, in Figs. 1 and 2.

- 1) Our method achieved the best clustering performance, followed by GSR\_knn, GSR\_mean, RFS, GSR, RSR, LapScore, and Baseline. Moreover, the results of our method are statistically significant based on the statistical analysis. For example, our proposed method on average improved by 2.15% and 0.48%, compared with the best comparison method (i.e., GSR\_knn), as well as on average improved by 18.70% and 14.42%, compared with the worst comparison method (i.e., Baseline), in terms of ACC and NMI results. The reason may be that our method uses all the available information as well as considers to flexibly reduce the influence of outliers. It is noteworthy that imputation methods (i.e., GSR\_mean and GSR\_knn) used all the available information in complete samples, but they do not take the influence of outliers into account, whereas feature

<sup>4</sup>Since  $k$ -means clustering is sensitive to the initialization, we need to repeat the clustering process many times to report the average result for avoiding the randomness.



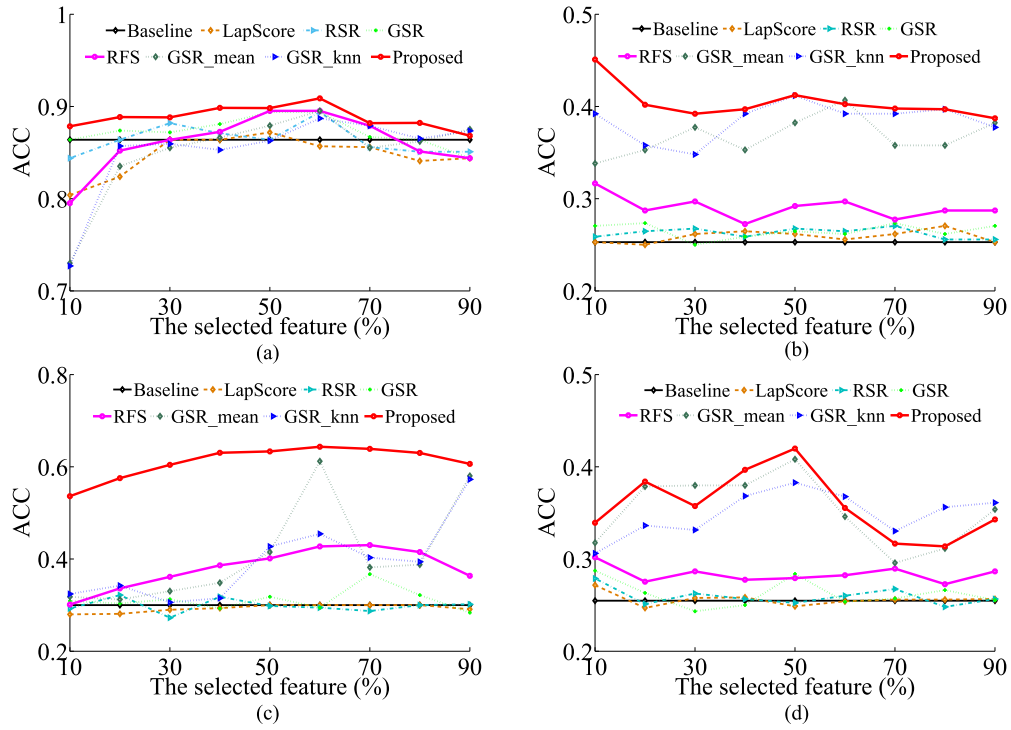


Fig. 1. ACC results of all methods with different features on different real incomplete data sets. (a) Advertisement. (b) Arrhythmia. (c) Cvpu. (d) Mice.

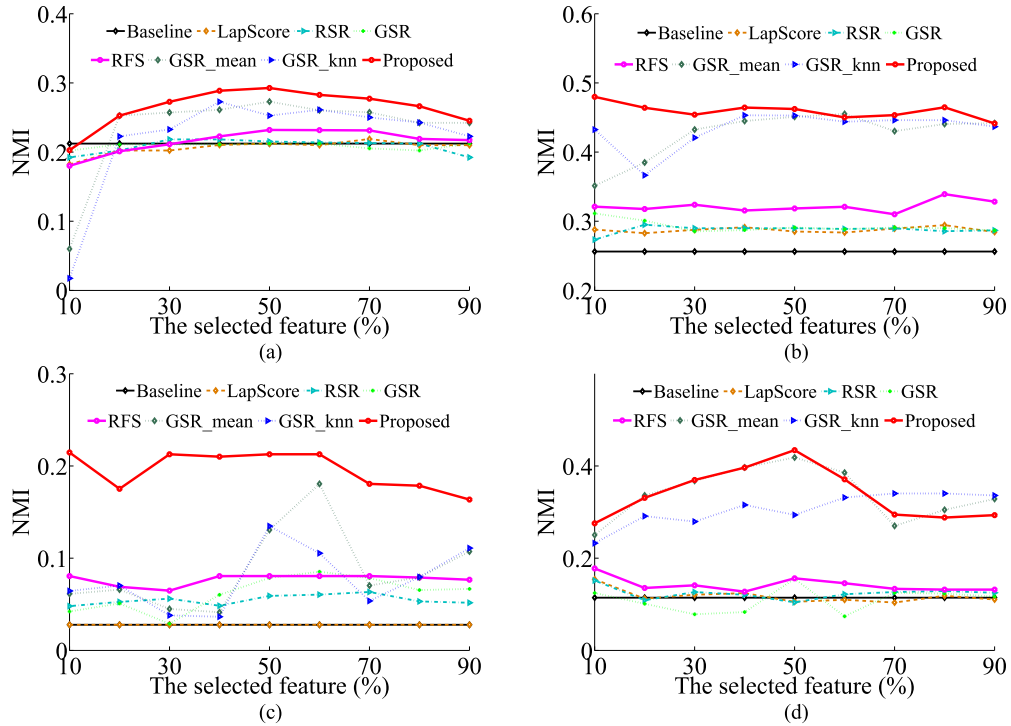


Fig. 2. NMI results of all methods with different features on different real incomplete data sets. (a) Advertisement. (b) Arrhythmia. (c) Cvpu. (d) Mice.

selection methods (i.e., LapScore, RSR, GSR, and RFS) only used the information in observed samples and do not use the observed information in incomplete samples.

- 2) Our method achieved the maximal improvement compared with four comparison methods on data set

Cvpu because this data set only has 5.01% complete samples, i.e., the minimal ratio of incomplete samples on all data sets. This indicates that reliable feature selection models need enough training samples.

- 3) In all, our proposed method achieved the best performance for two evaluation metrics on the data sets



TABLE II

CLUSTERING RESULTS OF SYNTHETIC INCOMPLETE DATA SETS WITH DIFFERENT INCOMPLETE SAMPLE RATIOS WHILE SELECTING 50% OF ALL FEATURES FOR FEATURE SELECTION METHODS. THE BOLD NUMBERS INDICATE THE BEST RESULTS THROUGH THE WHOLE ROW FOR THE SAME EVALUATION METRIC AND THE UNDERLINED BOLD NUMBERS IMPLY THE BEST RESULT IN ONE BLOCK FOR THE SAME EVALUATION METRIC. THE SYMBOLS “\*” AND “\*\*,” RESPECTIVELY, INDICATE THE STATISTICALLY SIGNIFICANT DIFFERENCE WITH  $p < 0.05$  AND  $p < 0.001$ , ON THE PAIRED-SAMPLE T-TESTS AT 95% SIGNIFICANCE LEVEL OF THE RESULTS BETWEEN OUR METHOD AND EVERY COMPARISON METHOD

Data set	Ratio	Baseline		LapScore		RSR		GSR		RFS		GSR_mean		GSR_knn		Proposed	
		ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
Cane	0	48.2**	41.4**	49.9**	42.3**	48.5**	41.3**	48.8**	42.0**	54.5*	46.9**	46.2**	40.2**	49.3**	41.5**	<b>55.8</b>	<b>48.9</b>
		±2.6	±2.6	±2.6	±2.2	±3.0	±2.5	±3.1	±2.4	±2.1	±1.7	±2.6	±2.1	±2.5	±2.2	±3.1	±2.9
	0.1	46.4**	37.6**	49.0**	40.2**	46.2**	41.3**	50.6**	43.6**	55.2**	47.1*	49.1**	41.3**	50.7**	40.8**	<b>58.1</b>	<b>48.2</b>
		±3.0	±2.3	±3.0	±2.9	±2.2	±1.9	±2.2	±2.1	±2.0	±2.9	±1.4	±2.2	±2.1	±2.1	±1.9	±2.6
	0.3	45.3**	37.4**	45.7**	39.8**	53.2*	47.3**	52.5**	45.6**	60.2*	52.8	48.8**	42.9**	52.9**	44.6**	<u><b>61.1</b></u>	<u><b>53.3</b></u>
		±3.0	±2.7	±1.9	±1.9	±2.2	±2.4	±1.9	±2.1	±3.2	±2.9	±3.7	±2.8	±2.9	±1.5	±3.5	±2.6
	0.5	37.1**	30.2**	51.7**	44.6**	48.9**	41.3**	49.0**	42.1**	55.0**	47.0	42.6**	35.2**	45.5**	37.0**	<b>56.9</b>	<b>47.9</b>
		±1.7	±1.7	±3.7	±3.7	±2.6	±2.1	±3.8	±3.5	±3.3	±2.5	±3.8	±0.3	±3.7	±3.2	±3.5	±3.3
Uspst	0	42.9**	36.6**	49.6**	42.7**	41.7**	36.0**	44.1**	39.7**	<b>54.5</b>	<b>47.8*</b>	43.3**	33.1**	46.2**	37.8**	54.3	46.8
		±3.6	±3.3	±2.8	±2.4	±3.1	±2.9	±3.3	±3.6	±3.1	±3.0	±0.8	±0.2	±1.8	±0.8	±4.2	±3.9
	0.7	42.1**	39.1**	45.3**	44.6**	48.7**	45.2**	44.9**	41.9**	53.4	<b>51.8**</b>	45.1**	37.0**	45.8**	37.4**	<b>54.6</b>	48.9
		±2.1	±2.1	±2.0	±2.0	±2.0	±2.4	±2.9	±2.6	±3.2	±3.0	±1.1	±2.8	±3.0	±2.9	±2.2	±2.1
	0.9	64.4**	59.5**	59.7**	57.2**	62.4**	58.2**	64.3**	58.6**	64.2**	58.2**	62.4**	58.2**	64.3**	58.6**	<b>67.7</b>	<b>60.9</b>
		±1.2	±0.7	±1.7	±1.0	±2.3	±0.9	±1.9	±0.9	±1.1	±0.7	±2.3	±0.9	±1.9	±0.9	±2.2	±0.8
	0.1	60.9**	58.1**	59.2**	56.9**	63.5**	59.2*	64.6**	59.1	64.6**	58.7	65.1**	58.3**	<b>67.5*</b>	59.3	67.0	<b>59.4</b>
		±1.9	±1.0	±1.9	±1.0	±1.4	±0.8	±1.8	±1.3	±1.6	±0.7	±2.2	±0.5	±3.4	±2.1	±0.8	±0.7
Yale	0	62.0**	57.9**	58.6**	56.7	63.7**	57.6**	61.8**	55.5**	65.5**	57.6**	63.9**	<b>61.3**</b>	63.0**	58.2**	<b>67.9</b>	56.7
		±1.9	±0.9	±2.1	±0.9	±1.4	±0.9	±1.9	±1.2	±1.1	±0.8	±3.1	±0.7	±2.2	±1.6	±0.7	±0.2
	0.3	61.0**	59.9**	59.0**	58.7**	63.4**	<b>62.0*</b>	56.5**	53.6**	65.5*	61.2**	61.4**	57.8**	57.6**	54.7**	<b>65.9</b>	61.3
		±1.5	±0.7	±2.2	±1.1	±2.7	±1.0	±1.3	±0.8	±1.8	±0.5	±2.1	±2.0	±1.3	±1.0	±1.8	±0.6
	0.5	61.3**	58.9**	59.5**	58.7**	61.5**	56.4**	59.4**	56.6	67.8	60.0**	65.4*	58.8**	63.8**	57.6**	<u><b>68.1</b></u>	<b>61.0</b>
		±1.8	±0.7	±1.0	±0.7	±1.9	±1.0	±2.1	±1.5	±1.1	±0.6	±2.5	±1.2	±3.2	±1.7	±2.3	±1.0
	0.7	61.1**	59.4**	58.0**	59.2**	56.5**	57.5**	58.6**	59.3	62.0**	61.9	61.6**	57.6**	64.3**	59.4**	<b>66.7</b>	<u><b>61.9</b></u>
		±2.0	±0.9	±2.3	±1.2	±1.5	±1.1	±2.2	±1.4	±1.3	±0.5	±1.9	±1.4	±2.8	±2.6	±1.5	±0.7
Yeast	0	48.4**	57.8*	53.6	57.7*	54.4	<b>61.0**</b>	46.3**	55.2**	46.9**	54.4**	51.4**	59.0*	46.3**	55.2**	<b>53.6</b>	58.6
		±2.1	±1.9	±2.4	±2.0	±2.3	±2.7	±1.6	±1.4	±2.1	±1.7	±2.3	±1.3	±1.6	±1.6	±4.0	±3.0
	0.1	50.4**	61.9**	51.0**	67.0*	<b>55.7*</b>	<b>71.7**</b>	55.1**	68.8**	41.8**	48.0**	48.7**	57.9**	54.0	59.5**	53.3	65.6
		±1.7	±1.8	±1.9	±1.7	±1.2	±1.8	±1.3	±1.3	±1.6	±1.6	±0.5	±1.5	±2.0	±1.9	±1.9	±1.3
	0.3	47.0**	58.0**	60.8**	75.8**	68.7**	85.8**	72.9**	85.8**	77.0**	86.9	67.7**	72.4**	70.8**	74.4**	<u><b>79.1</b></u>	<b>88.0</b>
		±1.8	±1.7	±2.3	±2.0	±1.5	±2.6	±1.4	±1.6	±1.3	±1.4	±4.0	±3.1	±1.1	±1.3	±1.7	±1.5
	0.5	46.3**	57.5**	48.3**	60.5**	48.3**	60.8**	47.9**	61.1**	50.8*	60.2**	48.7**	55.4**	47.7**	54.6**	<b>51.6</b>	<b>62.2</b>
		±1.9	±2.1	±1.5	±1.4	±2.7	±3.3	±2.2	±1.5	±1.4	±1.5	±1.5	±1.9	±2.5	±2.0	±1.6	±1.7
Palm	0	40.1**	50.6*	44.1*	51.0**	48.1**	58.3**	49.7**	58.8**	<b>56.4**</b>	<b>69.0**</b>	40.8*	47.6**	44.2**	48.9*	47.4	49.6
		±2.3	±2.1	±1.7	±1.7	±3.0	±2.7	±1.5	±1.5	±1.8	±1.8	±4.6	±4.6	±2.8	±3.7	±1.6	±1.3
	0.7	41.5**	51.9**	44.3*	<b>53.7**</b>	43.7*	53.0**	43.7*	49.5**	44.6	51.2**	38.3**	44.8**	<b>47.2**</b>	49.5	45.1	50.2
		±2.0	±2.2	±1.6	±1.8	±2.5	±3.0	±1.7	±1.0	±1.5	±1.1	±2.1	±0.6	±1.5	±1.9	±1.4	±1.0
	0.9	40.8**	<b>30.3**</b>	37.0**	18.7*	41.7**	23.2**	43.4*	11.9**	<b>44.8**</b>	30.3**	40.0**	23.4**	41.3**	14.6**	42.8	17.4
		±1.3	±3.6	±0.3	±0.6	±2.8	±2.7	±2.4	±2.9	±4.7	±3.6	±3.2	±0.4	±1.6	±2.3	±0.7	±1.7
	0.1	38.6**	27.5**	37.8**	29.9**	42.3	<b>32.7**</b>	43.2**	31.8**	<b>45.0**</b>	32.4**	38.1**	24.6**	40.9**	16.1**	42.0	19.7
		±1.6	±4.0	±0.6	±1.1	±2.1	±2.1	±2.3	±1.3	±3.1	±1.6	±2.4	±0.5	±0.1	±3.1	±1.0	±1.6
Fruit	0	40.5**	26.5**	37.6**	23.3**	40.3**	28.0**	42.9**	27.1**	43.1**	34.0**	34.0**	15.7**	36.7**	12.8**	<u><b>45.7</b></u>	<u><b>41.2</b></u>
		±1.6	±3.4	±0.5	±1.4	±2.1	±1.9	±2.4	±2.8	±2.5	±1.1	±1.2	±0.1	±1.2	±1.1	±0.7	±1.0
	0.3	35.9**	14.8	31.0**	14.1	37.0**	25.7**	<b>42.2**</b>	24.3**	38.2	<b>30.9**</b>	27.9**	12.9**	31.0**	5.9**	38.4	15.1
		±1.5	±3.2	±0.7	±1.0	±2.1	±1.7	±2.5	±1.3	±2.4	±1.5	±1.2	±0.1	±0.6	±0.1	±0.4	±1.5
	0.5	36.1**	13.4**	30.6**	18.0**	32.1**	16.0**	33.1**	10.8**	<b>42.8**</b>	<b>40.7**</b>	38.5*	26.9**	35.2**	11.6**	38.9	24.4
		±1.0	±4.0	±0.5	±0.7	±2.1	±1.7	±2.4	±1.5	±2.4	±1.3	±1.3	±0.2	±0.9	±0.3	±0.5	±0.6
	0.7	36.5**	15.5*	31.6**	15.6**	32.6**	13.2**	34.7**	<b>17.3**</b>	38.0**	12.3**	29.4**	8.1**	35.0**	11.0**	<b>38.6</b>	16.5
		±1.8	±2.5	±0.3	±0.2	±2.1	±1.9	±1.7	±1.9	±3.0	±1.0	±0.9	±0.2	±0.8	±0.8	±0.6	±0.6
Palm	0	70.2**	90.1	70.7**	90.0**	70.9**	89.8**	70.3**	89.6**	71.2**	89.9**	71.2**	<b>90.3*</b>	73.2	90.0**	<b>73.3</b>	90.2
		±0.5	±0.3	±1.5	±0.4	±1.0	±0.4	±1.1	±0.4	±1.5	±0.5	±0.5	±0.2	±0.7	±0.3	±0.9	±0.3
	0.1	68.8**	89.3**	68.9**	89.9*	70.7**	89.7**	71.3**	90.0*	72.1**	89.9**	73.8**	90.4	<b>74.4</b>	<b>90.6</b>	74.0	90.3
		±1.2	±0.1	±1.2	±0.5	±1.4	±0.5	±1.0	±0.4	±1.2	±0.4	±0.1	±0.1	±0.4	±0.1	±1.3	±0.5
	0.3	68.7**	89.7**	71.2**	90.4**	71.0**	89.5**	73.1**	90.6**	70.7**	89.5**	74.0**	90.4**	74.4**	90.5**	<b>75.1</b>	<b>90.8</b>
		±1.6	±0.4	±0.6	±0.2	±1.2	±0.5	±0.9	±0.3	±1.2	±0.4	±0.6	±0.2	±0.6	±0.2	±0.9	±0.3
	0.5	68.5**	89.2**	72.1**	90.3**	70.3**	89.5**	72.8**	90.4**	72.6**	90.5**	75.0**	90.9**	75.2*	91.0**	<b>76.9</b>	<b>91.3</b>
		±0.6	±0.6	±1.3	±0.4	±1.5	±0.5	±1.2	±0.4	±1.1	±0.5	±0.5	±0.1	±0.6	±0.4	±0.7	±0.3
	0.7	71.5**	90.7**	73.2**	91.2	72.3**	91.4	76.4	<b>92.6**</b>	74.1**	91.6	76.5	91.9**	<b>77.1*</b>	92.2**	76.3	91.3
		±0.3	±0.3	±0.9	±0.3	±1.5	±0.6	±1.0	±0.4	±0.9	±0.3	±0.5	±0.8	±0.6	±0.3	±0.6	±0.3
	0.9	71.0**	91.1**	74.5**	92.2**	77.0**	92.4**	75.0**	92.3**	77.5**	93.2**	78.7**	93.3**	77.3**	92.8**	<u><b>79.1</b></u>	<u><b>94.1</b></u>
		±0.4	±0.4	±1.0	±0.4	±1.7	±0.6	±0.9	±0.4	±0.7	±0.3	±0.8	±0.1	±0.6	±0.2	±0.8	±0.7

with different numbers of features. This indicated the robustness of our proposed method.

#### F. Synthetic Incomplete Data Sets

We randomly missed the observed information (i.e., randomly marked observed information as unobserved) in ten complete data sets to obtain different incomplete sample ratios,

i.e., 0% (i.e., complete data sets), 10%, 30%, 50%, 70%, 90%, and for every data set, and we kept 50% of all the features for all feature selection methods on these incomplete data sets. We listed the clustering results in Tables II and III.

Obviously, our method achieved the best clustering results, in terms of ACC and NMI, at different missing ratios. Moreover, the clustering results of our method statistically significant outperform the ones of all comparison methods.

TABLE III

CLUSTERING RESULTS OF SYNTHETIC INCOMPLETE DATA SETS WITH DIFFERENT INCOMPLETE SAMPLE RATIOS WHILE SELECTING 50% OF ALL FEATURES FOR FEATURE SELECTION METHODS. THE BOLD NUMBERS INDICATE THE BEST RESULTS THROUGH THE WHOLE ROW FOR THE SAME EVALUATION METRIC AND THE UNDERLINED BOLD NUMBERS IMPLY THE BEST RESULT IN ONE BLOCK FOR THE SAME EVALUATION METRIC. THE SYMBOLS “\*\*” AND “\*\*\*,” RESPECTIVELY, INDICATE THE STATISTICALLY SIGNIFICANT DIFFERENCE WITH  $p < 0.05$  AND  $p < 0.001$ , ON THE PAIRED-SAMPLE  $t$ -TESTS AT 95% SIGNIFICANCE LEVEL OF THE RESULTS BETWEEN OUR METHOD AND EVERY COMPARISON METHOD

Data set	Ratio	Baseline		LapScore		RSR		GSR		RFS		GSR_mean		GSR_knn		Proposed	
		ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
Cifar	0	21.0**	7.5**	18.4**	4.8**	21.2	7.6**	19.4**	7.0**	18.6**	4.9**	19.5**	6.9**	19.7**	7.4**	<b>21.3</b>	<b>8.1</b>
	0.1	$\pm 0.2$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 0.4$	$\pm 0.3$	$\pm 0.3$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.3$	$\pm 0.2$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$
		21.2**	7.5**	18.3**	4.7**	20.9**	7.2**	19.7**	6.6**	18.6**	4.9**	20.6**	7.1**	21.0**	7.9**	<b>21.7</b>	<b>8.4</b>
	0.3	$\pm 0.2$	$\pm 0.1$	$\pm 0.3$	$\pm 0.1$	$\pm 0.2$	$\pm 0.2$	$\pm 0.5$	$\pm 0.2$	$\pm 0.2$	$\pm 0.1$	$\pm 0.4$	$\pm 0.2$	$\pm 0.2$	$\pm 0.1$	$\pm 0.2$	$\pm 0.2$
		20.6**	7.3**	17.3**	1.7**	20.8**	7.2	21.3**	<b>7.7</b>	18.6**	4.9**	21.3**	7.5**	21.3**	7.5**	<b>21.6</b>	<b>7.7</b>
	0.5	$\pm 0.1$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 0.1$	$\pm 0.2$
		20.2**	7.0**	17.7**	4.5**	20.5**	6.9**	20.9**	7.0**	17.9**	6.0**	20.6**	7.0**	20.8**	7.6**	<b>21.3</b>	<b>7.8</b>
	0.7	$\pm 0.2$	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.2$	$\pm 0.2$	$\pm 0.3$	$\pm 0.3$	$\pm 0.1$	$\pm 0.3$	$\pm 0.2$	$\pm 0.3$	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$
0.9	19.8**	6.2**	17.9**	4.3**	19.6**	6.5**	20.0**	6.4**	17.6**	5.9**	20.2**	6.2**	20.1**	6.6*	<b>20.8</b>	6.7	
	$\pm 0.3$	$\pm 0.1$	$\pm 0.4$	$\pm 0.3$	$\pm 0.2$	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.2$	$\pm 0.3$	$\pm 0.3$	$\pm 0.5$	$\pm 0.2$	$\pm 0.4$	$\pm 0.4$	$\pm 0.3$	
		18.9	5.7**	17.3**	3.6**	18.8*	5.5**	19.4**	6.1	17.0**	4.6**	<b>20.0**</b>	6.0	<b>20.0**</b>	6.1	<b>20.0</b>	<b>6.2</b>
		$\pm 0.7$	$\pm 0.4$	$\pm 0.1$	$\pm 0.2$	$\pm 0.5$	$\pm 0.4$	$\pm 0.7$	$\pm 0.4$	$\pm 0.6$	$\pm 0.5$	$\pm 0.6$	$\pm 0.2$	$\pm 0.4$	$\pm 0.2$	$\pm 0.1$	$\pm 0.3$
Connect	0	37.5**	10.2**	38.5**	10.3**	38.1**	10.2**	42.6**	<b>11.7**</b>	40.6**	10.4**	42.7	11.4**	42.3**	11.4**	<b>42.8</b>	11.5
	0.1	$\pm 0.6$	$\pm 0.1$	$\pm 1.3$	$\pm 0.1$	$\pm 0.6$	$\pm 0.1$	$\pm 1.0$	$\pm 0.3$	$\pm 1.3$	$\pm 0.1$	$\pm 2.1$	$\pm 0.2$	$\pm 1.4$	$\pm 0.5$	$\pm 1.6$	$\pm 0.2$
		37.8**	10.2**	38.8	10.3**	38.9**	10.2**	41.2**	<b>11.5*</b>	40.4**	10.5**	<b>42.6*</b>	<b>11.5</b>	41.5**	11.4	41.7	<b>11.5</b>
	0.3	$\pm 0.7$	$\pm 0.1$	$\pm 0.9$	$\pm 0.1$	$\pm 0.6$	$\pm 0.1$	$\pm 2.9$	$\pm 0.2$	$\pm 1.3$	$\pm 0.1$	$\pm 3.1$	$\pm 0.1$	$\pm 1.7$	$\pm 0.5$	$\pm 0.8$	$\pm 0.2$
		37.6**	10.2**	38.8	10.7**	35.8*	10.2**	42.2**	11.2	35.6**	10.1**	<b>42.5**</b>	11.4**	41.8	11.5**	41.4	<b>11.8</b>
	0.5	$\pm 0.2$	$\pm 0.1$	$\pm 1.1$	$\pm 0.2$	$\pm 0.7$	$\pm 0.1$	$\pm 0.8$	$\pm 0.2$	$\pm 0.5$	$\pm 0.1$	$\pm 2.8$	$\pm 0.2$	$\pm 2.9$	$\pm 0.3$	$\pm 1.4$	$\pm 0.3$
		38.8**	10.3**	37.8**	10.2**	36.3**	10.2**	41.3**	<b>12.2**</b>	35.4*	10.1**	41.5**	11.5*	41.3**	11.4	<b>41.7</b>	11.3
	0.7	$\pm 1.5$	$\pm 0.1$	$\pm 0.6$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 2.6$	$\pm 0.4$	$\pm 0.4$	$\pm 0.1$	$\pm 2.6$	$\pm 0.1$	$\pm 2.5$	$\pm 0.4$	$\pm 1.1$	$\pm 0.2$
37.5**		10.1**	37.0**	10.2**	35.9*	10.1**	40.8**	11.3	35.9*	10.1**	41.7**	11.6**	<b>41.8**</b>	<b>11.7**</b>	<b>41.8</b>	11.3	
0.9	$\pm 0.6$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 2.0$	$\pm 0.2$	$\pm 0.3$	$\pm 0.1$	$\pm 2.1$	$\pm 0.4$	$\pm 2.3$	$\pm 0.1$	$\pm 1.1$	$\pm 0.2$	
	37.5**	10.2**	37.0	10.2**	35.2**	10.2**	40.5*	11.2	35.6**	10.1**	40.7**	11.4*	40.7**	<b>11.5**</b>	<b>40.9</b>	11.2	
		$\pm 0.4$	$\pm 0.1$	$\pm 0.7$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 2.0$	$\pm 0.2$	$\pm 0.6$	$\pm 0.1$	$\pm 3.0$	$\pm 0.3$	$\pm 2.8$	$\pm 0.1$	$\pm 0.8$	$\pm 0.1$
Mnist	0	55.0**	47.7**	57.2**	<b>48.6</b>	54.8**	42.5**	55.0**	44.7**	53.9**	42.4**	54.7**	43.8**	55.0**	44.9**	<b>58.8</b>	<b>48.6</b>
	0.1	$\pm 1.4$	$\pm 0.6$	$\pm 1.4$	$\pm 1.0$	$\pm 2.5$	$\pm 3.0$	$\pm 0.7$	$\pm 0.8$	$\pm 2.0$	$\pm 0.7$	$\pm 2.1$	$\pm 1.1$	$\pm 0.6$	$\pm 0.6$	$\pm 1.4$	$\pm 0.7$
		55.6**	47.9**	58.2**	48.5	51.7**	42.8**	56.0**	45.0**	52.6**	42.3**	53.9**	43.5**	54.2**	43.3**	<b>58.8</b>	<b>48.8</b>
	0.3	$\pm 1.4$	$\pm 0.6$	$\pm 2.3$	$\pm 1.1$	$\pm 6.5$	$\pm 4.8$	$\pm 15$	$\pm 1.4$	$\pm 1.5$	$\pm 0.4$	$\pm 2.0$	$\pm 1.0$	$\pm 1.2$	$\pm 0.6$	$\pm 1.3$	$\pm 0.9$
		54.5**	47.4**	56.4**	47.9*	51.1**	40.5**	54.3**	45.5**	52.0**	41.6**	53.8**	43.4**	52.1**	42.6**	<b>58.2</b>	<b>48.4</b>
	0.5	$\pm 1.6$	$\pm 0.6$	$\pm 1.6$	$\pm 0.9$	$\pm 0.2$	$\pm 1.8$	$\pm 1.0$	$\pm 0.7$	$\pm 1.7$	$\pm 0.6$	$\pm 2.0$	$\pm 1.0$	$\pm 0.5$	$\pm 0.7$	$\pm 1.2$	$\pm 0.5$
		54.8**	47.1*	56.0*	47.0**	52.7**	42.0**	55.4**	44.4**	49.6**	41.1**	51.5**	41.1**	51.4**	40.8**	<b>57.2</b>	<b>47.9</b>
	0.7	$\pm 1.6$	$\pm 0.7$	$\pm 1.7$	$\pm 0.8$	$\pm 2.8$	$\pm 3.2$	$\pm 1.4$	$\pm 0.7$	$\pm 0.9$	$\pm 0.4$	$\pm 1.4$	$\pm 1.8$	$\pm 1.0$	$\pm 1.1$	$\pm 1.5$	$\pm 0.7$
53.3**		46.9**	54.2**	45.7**	50.9**	40.5**	55.8**	45.0**	49.0**	41.6**	52.0**	40.6**	53.2**	41.3**	<b>57.2</b>	<b>48.1</b>	
0.9	$\pm 1.4$	$\pm 0.6$	$\pm 1.5$	$\pm 0.7$	$\pm 0.5$	$\pm 2.5$	$\pm 1.5$	$\pm 0.9$	$\pm 1.9$	$\pm 0.6$	$\pm 3.0$	$\pm 2.0$	$\pm 0.6$	$\pm 0.6$	$\pm 1.3$	$\pm 0.8$	
	50.0**	43.9**	53.9**	45.3	49.1*	41.2**	54.7**	44.5*	52.4*	<b>45.5*</b>	51.1**	38.6**	51.8**	38.8**	<b>55.8</b>	44.8	
		$\pm 2.3$	$\pm 0.8$	$\pm 2.5$	$\pm 1.0$	$\pm 3.1$	$\pm 3.9$	$\pm 2.3$	$\pm 1.1$	$\pm 5.3$	$\pm 3.0$	$\pm 0.5$	$\pm 0.6$	$\pm 0.7$	$\pm 0.4$	$\pm 1.7$	$\pm 0.9$
Vehicle	0	54.5**	17.1**	57.2**	15.7**	45.6**	9.5**	55.8**	15.0**	55.3**	15.1**	55.9**	14.7**	57.5**	16.1*	<b>57.6</b>	<b>16.2</b>
	0.1	$\pm 1.2$	$\pm 0.5$	$\pm 0.3$	$\pm 0.1$	$\pm 1.6$	$\pm 1.0$	$\pm 2.7$	$\pm 0.9$	$\pm 1.4$	$\pm 1.1$	$\pm 1.7$	$\pm 1.0$	$\pm 0.4$	$\pm 0.2$	$\pm 1.4$	$\pm 0.2$
		54.6**	16.7**	57.1	15.6*	46.3**	9.5**	57.0*	15.9	56.0**	15.6*	56.4**	15.8	54.3**	<b>19.7**</b>	<b>57.1</b>	15.7
	0.3	$\pm 1.2$	$\pm 0.8$	$\pm 0.8$	$\pm 0.5$	$\pm 1.4$	$\pm 0.9$	$\pm 1.0$	$\pm 0.4$	$\pm 1.4$	$\pm 0.8$	$\pm 0.8$	$\pm 0.5$	$\pm 4.1$	$\pm 2.2$	$\pm 1.2$	$\pm 0.6$
		56.3**	15.4**	56.5	15.3*	55.7**	14.9*	56.0**	15.4	53.5**	15.5**	53.8**	<b>15.5*</b>	56.0	<b>15.5</b>	<b>56.4</b>	<b>15.5</b>
	0.5	$\pm 1.3$	$\pm 0.3$	$\pm 0.1$	$\pm 0.2$	$\pm 0.8$	$\pm 0.7$	$\pm 0.5$	$\pm 0.2$	$\pm 2.5$	$\pm 0.8$	$\pm 2.5$	$\pm 0.9$	$\pm 1.5$	$\pm 1.3$	$\pm 1.9$	$\pm 0.9$
		54.7**	15.0**	54.3**	13.5**	53.3**	14.6	54.0**	14.2	54.4**	14.9**	53.0**	14.3	54.6**	<b>16.1**</b>	<b>56.3</b>	14.3
	0.7	$\pm 1.6$	$\pm 1.0$	$\pm 2.3$	$\pm 1.3$	$\pm 0.2$	$\pm 0.7$	$\pm 0.7$	$\pm 0.4$	$\pm 1.5$	$\pm 0.5$	$\pm 2.3$	$\pm 1.1$	$\pm 3.2$	$\pm 1.6$	$\pm 0.8$	$\pm 0.5$
<b>55.2**</b>		<b>14.8**</b>	50.2**	11.5**	53.3*	14.4**	53.6	13.2**	52.1*	13.8	49.5**	12.0**	51.2**	14.5**	53.5	13.8	
0.9	$\pm 1.3$	$\pm 0.8$	$\pm 2.9$	$\pm 2.3$	$\pm 1.6$	$\pm 0.7$	$\pm 1.2$	$\pm 0.6$	$\pm 1.4$	$\pm 0.9$	$\pm 2.8$	$\pm 1.3$	$\pm 2.4$	$\pm 2.3$	$\pm 1.3$	$\pm 1.0$	
	51.2**	12.7*	50.9	11.6*	50.1**	11.9	50.3*	11.7	50.6**	<b>13.0**</b>	49.6**	10.9**	48.5**	11.6*	<b>51.5</b>	11.8	
		$\pm 1.5$	$\pm 0.6$	$\pm 0.5$	$\pm 0.2$	$\pm 1.0$	$\pm 1.0$	$\pm 1.5$	$\pm 0.7$	$\pm 0.7$	$\pm 0.2$	$\pm 1.2$	$\pm 0.8$	$\pm 2.3$	$\pm 3.0$	$\pm 1.8$	$\pm 1.5$
Aloi	0	50.4**	79.5**	48.5**	73.1**	50.5*	80.3**	50.3**	80.6**	50.6*	<b>80.5*</b>	50.1**	80.3	<b>50.7</b>	80.4**	<b>50.7</b>	<b>80.5</b>
	0.1	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 0.6$	$\pm 0.2$	$\pm 0.5$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 0.5$	$\pm 0.2$
		50.5**	79.6**	48.5**	73.1**	50.0**	80.4*	50.7*	<b>80.8**</b>	50.5**	80.7**	50.3**	80.5**	50.0**	80.4	<b>50.8</b>	80.4
	0.3	$\pm 0.4$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 0.5$	$\pm 0.2$	$\pm 0.2$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 0.2$	$\pm 0.1$	$\pm 0.5$	$\pm 0.2$
		48.4**	78.0**	<b>48.5*</b>	73.1**	48.0**	79.2	47.8**	79.3	48.2**	<b>79.4*</b>	47.9**	79.2	47.6**	78.9**	<b>48.5</b>	<b>79.4</b>
	0.5	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 0.3$	$\pm 0.1$	$\pm 0.3$	$\pm 0.1$	$\pm 0.8$	$\pm 0.3$	$\pm 0.2$	$\pm 0.1$	$\pm 0.3$	$\pm 0.1$
		<b>41.9**</b>	74.0**	40.1**	<b>76.3**</b>	41.2**	75.9**	41.1**	75.8**	41.1**	76.0**	41.2**	75.9**	40.7**	75.7	41.5	75.6
	0.7	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.1$	$\pm 0.3$	$\pm 0.2$	$\pm 0.1$	$\pm 0.1$	$\pm 0.4$	$\pm 0.1$	$\pm 0.5$	$\pm 0.1$	$\pm 0.5$	$\pm 0.2$	$\pm 0.2$	$\pm 0.1$
22.5**		69.7**	32.0**	77.5**	32.8**	73.0**</											

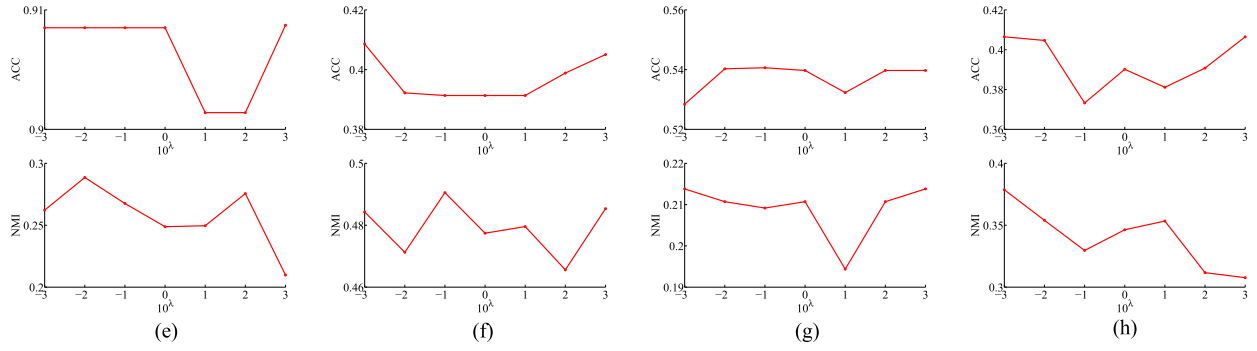


Fig. 3. ACC (top row) and NMI (bottom row) results of our method with varied parameter' setting (i.e.,  $\lambda$ ) on real incomplete data sets keeping 50% of all features. (a) Advertisement. (b) Arrhythmia. (c) Cvpu. (d) Mice.

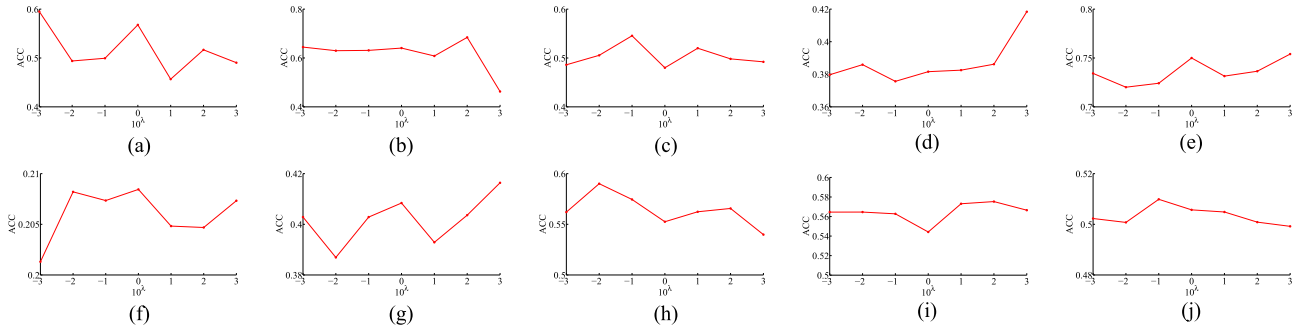


Fig. 4. ACC results of our method with varied parameter' setting (i.e.,  $\lambda$ ) on synthetic incomplete data sets keeping 50% of all features and 90% of all observed information. (a) Cane. (b) Uspst. (c) Yale. (d) Yeast. (e) Palm. (f) Cifar. (g) Connect. (h) Mnist. (i) Vehicle. (j) AloI.

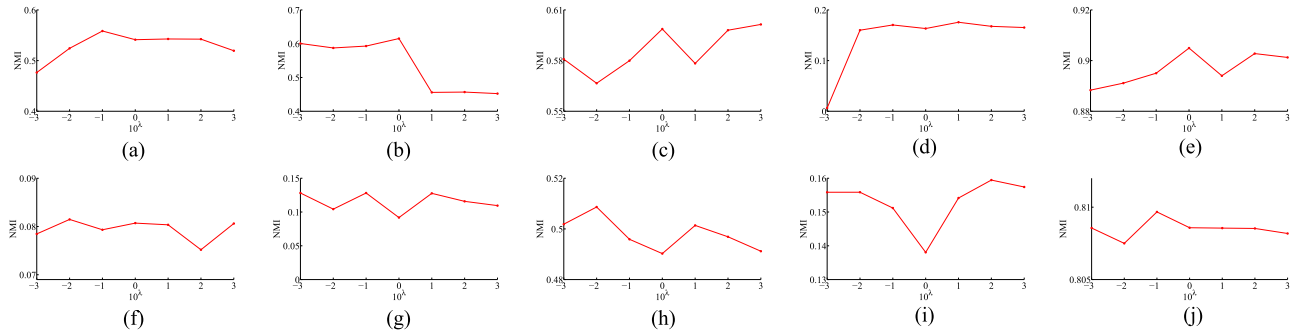


Fig. 5. NMI results of our method with varied parameter' setting (i.e.,  $\lambda$ ) on synthetic incomplete data sets keeping 50% of all features and 90% of all observed information. (a) Cane. (b) Uspst. (c) Yale. (d) Yeast. (e) Palm. (f) Cifar. (g) Connect. (h) Mnist. (i) Vehicle. (j) AloI.

selection models. Second, with the increase of the incomplete sample ratio, the difference between our method and every comparison feature selection method was larger than the difference between our proposed method and every two imputation methods. For example, by considering the ACC results on data set Uspst, GSR\_knn improved by 2.9% and 5.7% at the incomplete sample ratio 0.1 and 0.9, respectively, compared with GSR, whereas our method improved by 2.4% and 4.7% compared with RFS in the same incomplete sample ratios. The reason may be that both our method and imputation methods have more information to be used. Third, in imputation methods, GSR\_knn outperformed GSR\_mean. However, both of them were worse than our method. The reason is that imputed information in both GSR\_knn and GSR\_mean does not add the information of the whole data set (as the imputed information

was derived from observed samples) and possibly introduces noise to degrade the effectiveness of feature selection.

#### G. Parameters' Sensitivity Analysis

Our proposed objective function in (8) only needs to tune one parameter, i.e.,  $\lambda$ . We listed the variations of ACC and NMI, respectively, of our proposed method, at different values of the parameter  $\lambda$  (i.e.,  $\lambda \in \{10^{-3}, 10^{-2}, \dots, 10^3\}$ ) in Figs. 3–5.

Our method is sensitive to parameters' setting, but it may achieve stability and good clustering in some ranges. For example, the ACC results of our method varied from 68.43% to 46.27%, but it achieved the best results while  $\lambda \in [10^1, 10^2]$  on data set Uspst. This enables our method to achieve reliable clustering performance via easy parameter tuning.

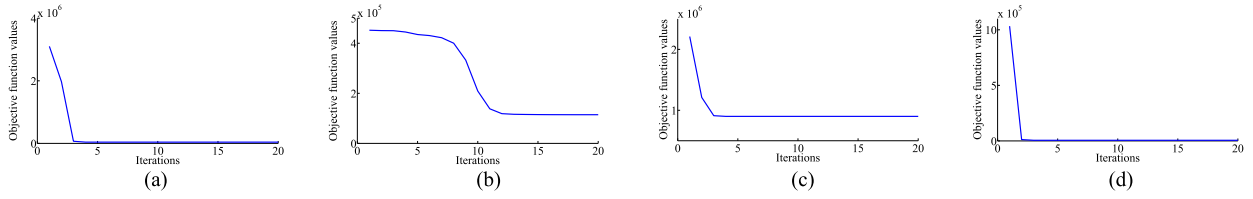


Fig. 6. Convergence analysis of our Algorithm 1 on real incomplete data sets keeping 50% of all features. (a) Advertisement. (b) Arrhythmia. (c) Cypu. (d) Mice.

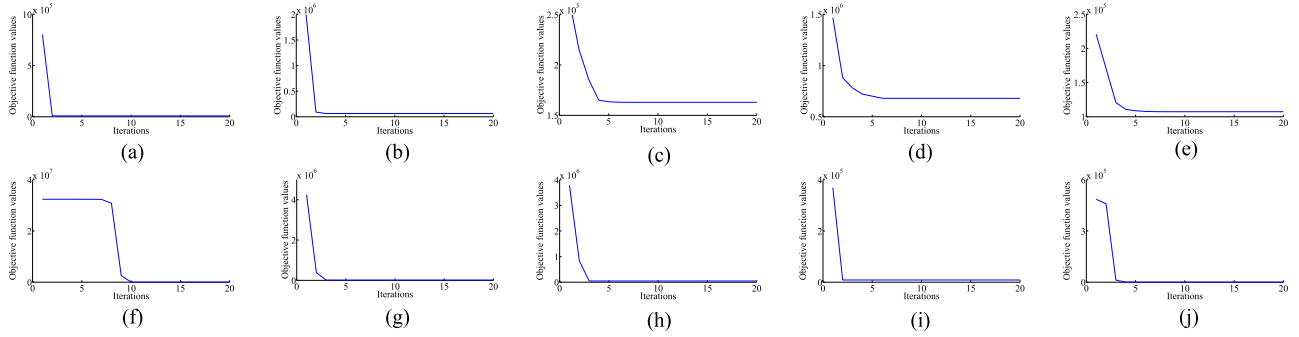


Fig. 7. Convergence analysis of our Algorithm 1 on synthetic incomplete data sets keeping 50% of all features and 90% of all observed information. (a) Cane. (b) Uspst. (c) Yale. (d) Yeast. (e) Palm. (f) Cifar. (g) Connect. (h) Mnist. (i) Vehicle. (j) Aloj.

## H. Convergence Analysis

Section III-E theoretically proved the convergence of our proposed Algorithm 1 to solve (8). Figs. 6 and 7 reported the variations of the objective function value of (8) at different iterations. In our experiments, we set the stop criterion of Algorithm 1 as  $(|\text{obj}(t+1) - \text{obj}(t)|/\text{obj}(t)) \leq 10^{-5}$ , where  $\text{obj}(t)$  represents the objective function value of the  $t$ th iteration.

Based on the results, Algorithm 1 achieved convergence as it monotonically decreased the objective function values of (8). Moreover, our experimental results indicated that our proposed Algorithm 1 can effectively optimize (8) with a fast convergence, i.e., within 20 iterations.

## I. Comparison of Robust Loss Functions

We selected GemanMcClure loss function [see (6)] [38] in our proposed method. In this section, we compared it with other two robust loss functions (such as  $\ell_1 - \ell_2$  estimator (i.e.,  $\phi(z) = (\mu + z^2)^{1/2}$ ) and the Welsch estimator (i.e.,  $\phi(z) = 1 - \exp(-(z^2/\mu^2))$ ) [42]) in this article. To do this, we denoted our (6) using the  $\ell_1 - \ell_2$  estimator and the Welsch estimator, respectively, as Proposed1 and Proposed2. It is noteworthy that the difference among different robust loss functions is not either the main focus or the contribution of this work.

Specifically, we manually generated outliers by replacing the original values with the random maximum/minimum values of the features [43] at different outlier ratios in the range of  $\{0, 5\%, 10\%, 15\%, \dots, 45\%, 50\%\}$  on four data sets, i.e., Cane, Uspst, Connect, and Vehicle. The comparison methods include Baseline, GSR\_knn, Proposed, Proposed1, and Proposed2. Moreover, “Proposed 0” and “Proposed 0.7,” respectively, indicate that our proposed method with zero and 70% missing ratios. Baseline and GSR\_knn, respectively, are

the worst and the best comparison methods. We reported the ACC results of all methods in Fig. 8.

First, all three methods with the robust loss function outperformed either GSR\_knn or Baseline. Specifically, these methods improved on average by 4.71% and 2.12%, respectively, compared with Baseline and GSR\_knn on four data sets. Second, with the increase of the outlier ratio, the clustering performance of every method with the robust loss function gradually decreases, while the clustering performance of both Baseline and GSR\_knn drastically decreases. This implies that the methods with the robust loss function can effectively reduce the effect of outliers. Third, each method with the robust loss function has a tiny difference to the other two in terms of clustering performance. This indicates that there is not a significant difference among them in our proposed method. Moreover, it is difficult to find the best one among them since they achieved the best clustering performance on different data sets, different outlier ratios, and different missing ratios. This verified the conclusion again that different robust loss functions could result in different robustness [43].

## J. Time and Space Costs

We listed the time cost (second) and space cost (MB) of the training process of all methods on the complete data sets in Table IV. It is noteworthy that three methods (i.e., GSR, GSR\_knn, and GSR\_mean) have the same time and space costs on the complete data sets. We did not list the time and space costs at a different missing ratio as the case-deletion methods (i.e., LPscore, RSR, GSR, and RFS) has less training samples compared with the proposed method for the incomplete data sets. As a result, RSR took the most time cost and GSR needed the most space cost as the time complexity of RSR is cubic to the sample size and the space complexity of GSR is cubic to the feature number.



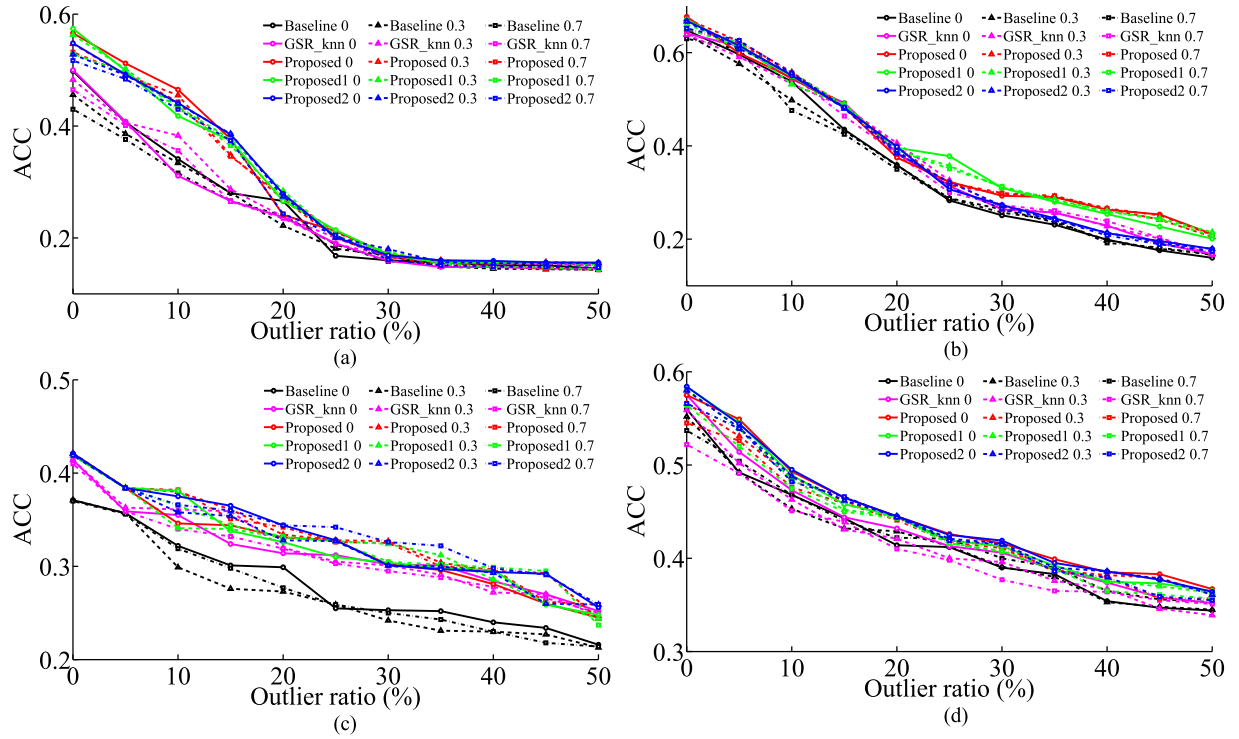


Fig. 8. Clustering results of different robust loss functions on four synthetic data sets keeping 50% of all features. (a) Cane. (b) Usps. (c) Connect. (d) Vehicle.

TABLE IV  
TIME AND SPACE COSTS OF ALL METHODS USING ALL SAMPLES FOR TRAINING

	Cane	Usps	Yale	Yeast	Palm	Cifar	Connect	Mnist	Vehicle	Aloi
Time cost (second)										
LPscore	0.86	0.63	0.38	1.91	0.72	203.80	13.62	45.381	14.33	22.34
RSR	2.81	2.56	1.23	5.49	2.77	6505.83	6674.03	7490.86	9752.71	61646.22
GSR	1.90	1.17	0.78	3.55	0.73	396.98	250.63	547.10	137.96	705.42
RFS	1.19	0.78	0.57	2.27	0.74	275.77	11.72	46.57	12.84	11.83
Proposed	3.20	1.35	0.62	3.16	0.80	593.50	21.14	157.5	23.82	30.88
Space cost (MB)										
LPscore	7.23	4.02	1.32	17.08	4.01	1442.82	66.64	387.13	81.37	108.21
RSR	95.03	7.76	75.94	108.35	4.01	1442.75	102.69	515.65	61.70	396.70
GSR	9.90	25.53	8.20	17.08	25.36	22825.63	28937.47	31068.50	39393.22	7355.70
RFS	42.67	4.02	8.21	17.08	4.01	1442.75	66.64	429.60	61.70	108.20
Proposed	10.08	4.02	8.21	17.08	6.17	1442.75	66.64	431.56	61.70	109.67

## V. CONCLUSION

This article proposed a novel method to conduct UFS on incomplete data sets. Specifically, we employed half-quadratic minimization to consider sample importance for conducting feature selection and used an indicator matrix to avoid unobserved information taking participation in the process of feature selection as well as making the use of all the observed information. Experimental results showed that our method achieved the best clustering results compared with the comparison methods.

We will extend our proposed unsupervised framework to conduct supervised/semisupervised feature selection on incomplete data sets in our future work.

## REFERENCES

- [1] X. Zhu, X. Li, S. Zhang, C. Ju, and X. Wu, "Robust joint graph sparse coding for unsupervised spectral feature selection," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 6, pp. 1263–1275, Jun. 2017.
- [2] Z. Zhang, L. Liu, F. Shen, H. T. Shen, and L. Shao, "Binary multi-view clustering," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 7, pp. 1774–1782, Jul. 2019.
- [3] X. Zhu, S. Zhang, Y. Zhu, W. Zheng, and Y. Yang, "Self-weighted multi-view fuzzy clustering," *ACM Trans. Knowl. Discovery Data*, vol. 14, no. 4, p. 48, 2020, doi: [10.1145/3396238](https://doi.org/10.1145/3396238).
- [4] R. Hu, X. Zhu, Y. Zhu, and J. Gan, "Robust SVM with adaptive graph learning," *World Wide Web*, vol. 23, pp. 1945–1968, Dec. 2019, doi: [10.1007/s11280-019-00766-x](https://doi.org/10.1007/s11280-019-00766-x).
- [5] G. Chandrashekar and F. Sahin, "A survey on feature selection methods," *Comput. Electr. Eng.*, vol. 40, no. 1, pp. 16–28, Jan. 2014.
- [6] S. Solorio-Fernández, J. A. Carrasco-Ochoa, and J. F. Martínez-Trinidad, "A review of unsupervised feature selection methods," *Artif. Intell. Rev.*, vol. 53, no. 2, pp. 907–948, Feb. 2020.
- [7] Y. Yang, Z. Ma, Y. Yang, F. Nie, and H. T. Shen, "Multitask spectral clustering by exploring intertask correlation," *IEEE Trans. Cybern.*, vol. 45, no. 5, pp. 1069–1080, May 2015.
- [8] E. Adeli *et al.*, "Semi-supervised discriminative classification robust to sample-outliers and feature-noises," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 2, pp. 515–522, Feb. 2019.
- [9] P. J. Huber, "Robust statistics," in *International Encyclopedia of Statistical Science*. 2011, pp. 1248–1251.

- [10] X. Zhu, Y. Zhu, and W. Zheng, "Spectral rotation for deep one-step clustering," *Pattern Recognit.*, vol. 105, Sep. 2020, Art. no. 107175, doi: [10.1016/j.patcog.2019.107175](https://doi.org/10.1016/j.patcog.2019.107175).
- [11] X. Zhu, J. Gan, G. Lu, J. Li, and S. Zhang, "Spectral clustering via half-quadratic optimization," *World Wide Web*, vol. 23, pp. 1969–1988, Nov. 2020, doi: [10.1007/s11280-019-00731-8](https://doi.org/10.1007/s11280-019-00731-8).
- [12] X. Zhu, J. Yang, C. Zhang, and S. Zhang, "Efficient utilization of missing data in cost-sensitive learning," *IEEE Trans. Knowl. Data Eng.*, early access, Nov. 28, 2019, doi: [10.1109/TKDE.2019.2956530](https://doi.org/10.1109/TKDE.2019.2956530).
- [13] Y. Zhou, L. Tian, C. Zhu, X. Jin, and Y. Sun, "Video coding optimization for virtual reality 360-degree source," *IEEE J. Sel. Topics Signal Process.*, vol. 14, no. 1, pp. 118–129, Jan. 2020.
- [14] J. Van Hulse and T. M. Khoshgoftaar, "Incomplete-case nearest neighbor imputation in software measurement data," *Inf. Sci.*, vol. 259, pp. 596–610, Feb. 2014.
- [15] S. Zhang, X. Li, M. Zong, X. Zhu, and R. Wang, "Efficient kNN classification with different numbers of nearest neighbors," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 5, pp. 1774–1785, May 2018.
- [16] W. Zheng, X. Zhu, Y. Zhu, and S. Zhang, "Robust feature selection on incomplete data," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, Jul. 2018, pp. 3191–3197.
- [17] C. Hou, F. Nie, X. Li, D. Yi, and Y. Wu, "Joint embedding learning and sparse regression: A framework for unsupervised feature selection," *IEEE Trans. Cybern.*, vol. 44, no. 6, pp. 793–804, Jun. 2014.
- [18] P. Zhu, Q. Hu, C. Zhang, and W. Zuo, "Coupled dictionary learning for unsupervised feature selection," in *Proc. 13th AAAI Conf. Artif. Intell.*, 2016, pp. 2422–2428.
- [19] M. Nikolova and R. H. Chan, "The equivalence of half-quadratic minimization and the gradient linearization iteration," *IEEE Trans. Image Process.*, vol. 16, no. 6, pp. 1623–1627, Jun. 2007.
- [20] G. Doquire and M. Verleysen, "Feature selection with missing data using mutual information estimators," *Neurocomputing*, vol. 90, pp. 3–11, Aug. 2012.
- [21] D. Huang, R. Cabral, and F. De la Torre, "Robust regression," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 363–375, Feb. 2016.
- [22] X. Zhu, X. Li, S. Zhang, Z. Xu, L. Yu, and C. Wang, "Graph PCA hashing for similarity search," *IEEE Trans. Multimedia*, vol. 19, no. 9, pp. 2033–2044, Sep. 2017.
- [23] H. T. Shen *et al.*, "Exploiting subspace relation in semantic labels for cross-modal hashing," *IEEE Trans. Knowl. Data Eng.*, early access, Jan. 29, 2020, doi: [10.1109/TKDE.2020.2970050](https://doi.org/10.1109/TKDE.2020.2970050).
- [24] P. Zhu, W. Zuo, L. Zhang, Q. Hu, and S. C. K. Shiu, "Unsupervised feature selection by regularized self-representation," *Pattern Recognit.*, vol. 48, no. 2, pp. 438–446, Feb. 2015.
- [25] H. Peng and Y. Fan, "A general framework for sparsity regularized feature selection via iteratively reweighted least square minimization," in *Proc. 31st AAAI Conf. Artif. Intell.*, 2017, pp. 2471–2477.
- [26] G. Lan, C. Hou, F. Nie, T. Luo, and D. Yi, "Robust feature selection via simultaneous sapped norm and sparse regularizer minimization," *Neurocomputing*, vol. 283, pp. 228–240, Mar. 2018.
- [27] L. Sun, J. Xu, and Y. Tian, "Feature selection using rough entropy-based uncertainty measures in incomplete decision systems," *Knowl.-Based Syst.*, vol. 36, pp. 206–216, Dec. 2012.
- [28] W. Shu and H. Shen, "Multi-criteria feature selection on cost-sensitive data with missing values," *Pattern Recognit.*, vol. 51, pp. 268–280, Mar. 2016.
- [29] J. Huang, T. Zhang, and D. Metaxas, "Learning with structured sparsity," *J. Mach. Learn. Res.*, vol. 12, pp. 3371–3412, Jan. 2011.
- [30] S. Bengio, F. Pereira, Y. Singer, and D. Strelow, "Group sparse coding," in *Proc. NIPS*, 2009, pp. 82–89.
- [31] F. Nie, H. Huang, X. Cai, and C. H. Ding, "Efficient and robust feature selection via joint  $\ell_{2,1}$ -norms minimization," in *Proc. NIPS*, 2010, pp. 1813–1821.
- [32] M. Nikolova and M. K. Ng, "Analysis of half-quadratic minimization methods for signal and image recovery," *SIAM J. Sci. Comput.*, vol. 27, no. 3, pp. 937–966, Jan. 2005.
- [33] M. P. Kumar, B. Packer, and D. Koller, "Self-paced learning for latent variable models," in *Proc. NIPS*, 2010, pp. 1189–1197.
- [34] Y. J. Lee and K. Grauman, "Learning the easy things first: Self-paced visual category discovery," in *Proc. CVPR*, Jun. 2011, pp. 1721–1728.
- [35] S. Negahban, B. Yu, M. J. Wainwright, and P. K. Ravikumar, "A unified framework for high-dimensional analysis of  $M$ -estimators with decomposable regularizers," in *Proc. NIPS*, 2009, pp. 1348–1356.
- [36] Y. Fan, R. He, J. Liang, and B. Hu, "Self-paced learning: An implicit regularization perspective," in *Proc. 31st AAAI Conf. Artif. Intell.*, 2017, pp. 1877–1883.
- [37] R. J. Little and D. B. Rubin, *Statistical Analysis With Missing Data*. Hoboken, NJ, USA: Wiley, 2014.
- [38] S. Ganan and D. McClure, "Bayesian image analysis: An application to single photon emission tomography," *Amer. Stat. Assoc.*, pp. 12–18, 1985.
- [39] S. Geman and D. E. McClure, "Statistical methods for tomographic image reconstruction," in *Proc. Session ICI Bull. (ICI)*, vol. 4, 1987, pp. 1–68.
- [40] I. Daubechies, R. DeVore, M. Fornasier, and C. S. Güntürk, "Iteratively reweighted least squares minimization for sparse recovery," *Commun. Pure Appl. Math.*, vol. 63, no. 1, pp. 1–38, Jan. 2010.
- [41] X. He, D. Cai, and P. Niyogi, "Laplacian score for feature selection," in *Proc. NIPS*, 2005, pp. 507–514.
- [42] Y. Zhang, Z. Sun, R. He, and T. Tan, "Robust subspace clustering via half-quadratic minimization," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2013, pp. 3096–3103.
- [43] R. He, W.-S. Zheng, T. Tan, and Z. Sun, "Half-quadratic-based iterative minimization for robust sparse representation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 2, pp. 261–275, Feb. 2014.

**Heng Tao Shen** (Senior Member, IEEE) received the B.Sc. (Hons.) and Ph.D. degrees from the Department of Computer Science, National University of Singapore, Singapore, in 2000 and 2004, respectively.

He then joined The University of Queensland, Brisbane, QLD, Australia, where he became a Professor in late 2011. He is currently a Professor and the Dean of the School of Computer Science and Engineering, the Executive Dean of the AI Research Institute, and the Director of the Centre for Future Media, University of Electronic Science and Technology of China, Chengdu, China. His research interests mainly include multimedia search, computer vision, artificial intelligence, and big data management. He has published over 280 peer-reviewed papers, including over 80 IEEE/ACM TRANSACTIONS.

Dr. Shen is a fellow of OSA and an ACM Distinguished Member. He received seven best paper awards from international conferences, including the Best Paper Award from the ACM Multimedia 2017 and the Best Paper Award–Honorable Mention from ACM SIGIR 2017. He is/was an Associate Editor of *ACM Transactions of Data Science*, the IEEE TRANSACTIONS ON IMAGE PROCESSING, the IEEE TRANSACTIONS ON MULTIMEDIA, and the IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING.

**Yonghua Zhu** is currently pursuing the master's degree with Guangxi University, Nanning, China.

His current research interests include data mining and machine learning.

**Wei Zheng** is currently pursuing the master's degree with Guangxi Normal University, Guilin, China.

His current research interests include data mining and pattern recognition.

**Xiaofeng Zhu** (Senior Member, IEEE) is currently a Faculty Member with the University of Electronic Science and Technology of China, Chengdu, China. His current research interests include large-scale multimedia retrieval, feature selection, sparse learning, data preprocess, and medical image analysis.