



# Multi-label feature selection based on logistic regression and manifold learning

Yao Zhang<sup>1,2</sup> · Yingcang Ma<sup>1</sup> · Xiaofei Yang<sup>1</sup>

Accepted: 10 November 2021 / Published online: 4 January 2022

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

## Abstract

Like traditional single-label learning, multi-label learning is also faced with the problem of dimensional disaster. Feature selection is an effective technique for dimensionality reduction and learning efficiency improvement of high-dimensional data. This paper combined logistic regression, manifold learning, and sparse regularization to construct a joint framework for multi-label feature selection (LMFS). Firstly, the sparsity of the feature weight matrix is constrained by the  $L_{2,1}$ -norm. Secondly, the feature manifold and label manifold can constrain the feature weight matrix to fit the information of data and label better. An iterative updating algorithm is designed, and the convergence of the algorithm is proved. Finally, the LMFS algorithm is compared with DRMFS, SCLS, and other algorithms on eight classical multi-label data sets. The experimental results show the effectiveness of the LMFS algorithm.

**Keywords** Feature selection · Manifold learning · Multi-label learning ·  $L_{2,1}$ -norm · Logistic regression

## 1 Introduction

With the rapid development of the Internet and digital acquisition equipment in recent years, the scale of data that needs to be analyzed and processed in classification problems has increased dramatically. These data may contain not only single-label data but also a large number of multi-label data. Each instance has only one label in the single-label data, and different labels are mutually independent. While in the multi-label data, a sample may belong to multiple labels simultaneously, and each label intersects with the other and is correlated. So far, the research of multi-label learning, which includes text classification, image annotation, video classification, biology, etc., has attracted the attention of many scholars.

In many practical applications mentioned above, multi-label data usually has thousands or even more features, which brings many problems to data analysis, decision-making, screening, and prediction [1]. For example, redundant and irrelevant features may affect the function of classifiers [2]. In order to solve these problems, we will select a subset of related and optimal features. The procedure is called feature selection. Feature selection has many advantages in learning algorithms, including reducing measurement cost and storage requirements, shortening training time, avoiding dimension disaster, reducing overfitting, and so on [3, 4]. Therefore, multi-label feature selection has become a research hot spot.

Based on label information and search strategy, feature selection methods are usually divided into two categories [5]. Based on the search strategy, feature selection can be divided into three categories: filter [6–8], wrapper [9, 10] and embedded [11, 12]. Among them, the embedded method combines the advantages of the filter method and wrapper method. They embed the feature selection process in the learning process. Because they do not evaluate the feature subset iteratively, they are more effective than the wrapper method [1].

Through the study of advanced models, it is found that most of them are based on linear mapping [13, 14] and information theory. With the development of research and the introduction of manifold learning [15–17], the multi-label feature selection system is constantly improved.

---

✉ Yingcang Ma  
mayincang@xpu.edu.cn

Yao Zhang  
ayunxiaobao@163.com

Xiaofei Yang  
yangxiaofei2002@163.com

<sup>1</sup> School of Science, Xi'an Polytechnic University, Xi'an, China

<sup>2</sup> Basic Education Department, Shandong Huayu University of Technology, Dezhou, China

A multi-label feature selection method based on mutual information and label correlation is proposed in [18]. A new multi-label feature selection based on label redundancy, called (LRFS), is proposed in [19]. It divides labels into independent labels and dependent labels and analyzes the differences between independent and dependent labels. Kernel alignment is introduced into multi-label learning to measure the consistency between feature space and label space. Moreover, a new multi-label feature selection method, which can automatically learn and deal with the importance of labels, is proposed in [20]. A new feature selection method of extended adaptive minimum absolute contraction selection operator (EALasso) is presented [21]. This method preserves the properties of determining the correct subset model and obtaining the optimal estimation accuracy, proposes an iterative optimization algorithm, and gives theoretical convergence proof. Some researchers put forward a multi-label feature selection method with multiple regularizations (MDFS) [22].

Moreover, they calculate the correlation between the feature and the local label and use the objective function that includes  $L_{21}$ -norm regularization. Through linear mapping and combining manifold learning and  $L_{21}$ -norm regularization, multi-label feature selection via feature manifold learning and sparsity regularization (MSSL) is proposed in [23]. A robust multi-label feature selection with dual-graph regularization (DRMFS) is proposed in [24]. In order to improve the robustness of the algorithm, the model uses  $L_{21}$ -norm mapping and combines label manifold with characteristic manifold to consider not only the correlation between features but also the correlation between labels. Finally, the  $L_{21}$ -norm is used to constrain the sparsity of the weight matrix.

To sum up, linear regression is often used in multi-label feature selection models. However, since labels are binary, it is not appropriate in most cases to assume a linear correlation between the sample space and the label space. Moreover, multi-label feature selection is used for multi-label classification. So from the perspective of classification, logistic regression is more suitable for the multi-label feature selection model. The reasons are as follows:

- 1) For multi-label data, label (dependent variable) is the discrete value (0 or 1), which is more suitable for logistic regression.
- 2) Logistic regression is a generalized linear model, which is equivalent to introducing nonlinearity into the model, which can improve the expression ability of the model and increase the fitting.
- 3) Linear regression directly analyzes the relationship between the dependent and independent variables, while logistic regression analyzes the relationship

between the probability of taking a particular value of the dependent variable and the independent variable.

From reasons 1 and 2, it can be seen that logistic regression is more suitable for data classification than linear regression, and it can be applied to a variety of data distribution including the distribution of the positive and the negative; from reason 3, it can be seen that logistic regression is more robust than linear regression.

Based on this problem, some scholars choose logistic regression to replace the least square regression in the model and improve the algorithm's function by improving the regular term. The author puts forward a correlation logistic regression model (CorrLog) for multi-label image classification, which extends the traditional logistic regression model to multi-label image classification [25]. A feature subset selection algorithm for mixed-integer optimal logistic regression is proposed in [26]. This paper presents a mixed-integer linear optimization problem, which can be solved using standard integer optimization software to approximate the logistic loss function piecewise. A robust logistic regression method based on the regularization of  $L_q$ -norm  $q \in [0, 1]$  is proposed in [27], which is a feasible and effective feature selection method.

However, the existing multi-label feature selection algorithm based on logistic regression ignores the feature manifold structure. The above multi-label feature selection algorithm [1, 3, 4, 19, 20] ignores the fitting of label information while paying attention to the feature manifold structure, and the above multi-label feature selection algorithm [21, 23] ignores the fitting of the feature manifold structure to label information while paying attention. Therefore, this study will combine the logistic regression model with the regularization of the feature map, label map, and  $L_{21}$ -norm sparse regularization to solve the problem of multi-label feature selection. The main contributions of this paper are as follows:

- 1) The assumption of linear correlation between sample space and label space is not applicable in most cases to solve the problem. Therefore, this paper uses logistic regression to construct a multi-label feature selection model.
- 2) Considering that feature selection should be based on sample and label matrices, few existing models consider both of these matrices. Therefore, the feature manifold and label manifold are combined into the multi-label feature selection model constructed by logistic regression. In addition, the  $L_{2,1}$ -norm sparse constraint is combined to construct "Multi-label feature selection based on logistic regression and manifold learning."

- 3) The model's optimal solution is realized, the optimal algorithm is designed, and the algorithm's convergence is proved.
- 4) A large number of experiments were designed and carried out on eight classical multi-label data sets, compared with five advanced multi-label feature selection algorithms (DRMFS, SCLS, etc.) and baseline, and the results prove the effectiveness of the LMFS algorithm.

The rest of this paper is organized as follows. Section 2 gives a multi-label feature selection model. In Section 3, the model is solved, an iterative algorithm for multi-label feature selection is proposed, and its time complexity is analyzed. In Section 4, the comprehensive experiment on six classical data sets shows that the algorithm proposed in this paper is superior to other algorithms. Finally, the conclusions and future work is presented in Section 5.

## 2 Problem description

### 2.1 Logistic regression model

Suppose a multi-label data set  $D = \{(d_i, y_i)\}_{i=1}^n$  consists of  $n$  independent samples with the same distribution. Let  $X = [x_1; x_2; \dots; x_n]$  be the augmented matrix of the data matrix,  $X \in R^{n \times d}$ , where  $x_i = [1, d_i]$ .  $Y = [y_1; y_2; \dots; y_n]$  be the label matrix,  $Y \in R^{n \times m}$ ,  $m$  is the number of classes. The value of  $y_{ij}$  is 0 or 1, indicating whether the  $i$ -th sample is associated with the  $j$ -th class. In the logistic regression, the posterior probability that sample  $x_i$  belongs to the  $j$ -th class is:

$$Pr(y_{ij} = 1|x_i) = g(x_i w_j) = \frac{\exp(x_i w_j)}{1 + \exp(x_i w_j)} \quad (1)$$

Thus the posterior probability that sample  $x_i$  does not belong to the  $j$ -th class is:

$$Pr(y_{ij} = 0|x_i) = 1 - g(x_i w_j) = \frac{1}{1 + \exp(x_i w_j)} \quad (2)$$

where  $W = [w_1, w_2, \dots, w_m]$  and  $W \in R^{d \times m}$ ;  $w_j$  is the  $j$ -th column vector of the coefficient matrix  $W$ .

If the maximum likelihood estimation method is used to estimate the coefficient matrix, then the likelihood function (joint probability distribution) of the logistic regression on the multi-label data set is:

$$P(W) = \prod_{j=1}^m \prod_{i=1}^n g(x_i w_j)^{y_{ij}} (1 - g(x_i w_j))^{1 - y_{ij}} \quad (3)$$

Since it is inconvenient to solve optimization  $\max P(W)$ , the minimum value of  $L(W)$  of negative log likelihood function for solving logistic regression is used to solve  $W$

$$L(W) = - \sum_{j=1}^m \sum_{i=1}^n [y_{ij} \ln(g(x_i w_j)) + (1 - y_{ij}) \ln(1 - g(x_i w_j))] \quad (4)$$

$$= - \sum_{j=1}^m \sum_{i=1}^n [y_{ij} x_i w_j + \ln(1 - g(x_i w_j))]$$

### 2.2 Sparse constraint

The logistic regression model may suffer from ill-posed problems, such as overfitting, multi-collinearity, and infinite solutions, which results in incorrect estimation of the coefficient matrix [28]. So in order to solve this problem, a widely used strategy is to introduce penalty terms into  $L(W)$ , which aims to achieve a stable and accurate logistic regression model in high-dimensional data. The so-called penalty function is usually expressed as follows, where  $\beta$  is the regularization parameter.

$$\min_W L(W) + \beta R(W) \quad (5)$$

For the  $i$ -th row vector  $W_i$  of the coefficient matrix  $W$ , it can be regarded as a vector that measures the importance of the  $i$ -th feature. Let  $f_i \in R^n$  be the  $i$ -th feature vector of the data matrix, and then the data matrix  $X$  can be expressed in the form of  $X = [f_1, f_2, \dots, f_d]$ .

The common term  $R(W)$  has various forms for different purposes. Take  $L_1$ -norm regular term and  $L_2$ -norm regular term. For example,  $L_1$ -norm is often used to guided sparsity;  $L_2$ -norm is often used to guided stability. Because  $\|W_i\|_2$  is generally used to measure the importance of feature  $f_i$ , in order to distinguish the importance of features better, here we take  $L_{21}$ -norm as the standard term  $R(W)$ , which not only guides the row sparsity of the sparse matrix but also is sensitive to singular values [29]. So the objective optimization problem can be written as:

$$\min_W L(W) + \frac{\beta}{2} \|W\|_{2,1} \quad (6)$$

where,  $\|W\|_{2,1} = \sum_{i=1}^d \|W_i\|_2$ .

### 2.3 Feature manifold learning

Considering that the parameter of each coefficient vector  $w_j$  in formula (6) is  $\beta$ , but according to the idea of binary conversion, the regularization parameter  $\beta$  may not apply to all coefficient vectors. In addition, the features are extracted from some manifolds called the feature manifold to [2, 15, 16]. This is an important technique that can obtain the structure of the feature weight feature manifold by exploring the geometry. According to the problem assumption, if the features  $f_i$  and  $f_j$  are closer, then their weight vectors  $W_i$  and  $W_j$  should also be closed. Therefore, a feature map regularization is constructed, which can adjust the regularization parameters of the coefficient vectors  $w_j$  according to the similarity between the features  $f_i$  and  $f_j$ .

Its expression is as follows:

$$\begin{aligned}
 & \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^d \|W_i - W_j\|_2^2 S_{ij} \\
 = & \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^d (W_i - W_j)(W_i - W_j)^T S_{ij} \\
 = & \sum_{i=1}^d W_i W_i^T M_{ii} - \sum_{i=1}^d \sum_{j=1}^d W_i W_j^T S_{ij} \tag{7} \\
 = & Tr(W^T (M - S)W) \\
 = & Tr(W^T L_S W)
 \end{aligned}$$

where,  $M \in R^{d \times d}$  is the diagonal matrix, and  $M_{ii} = \sum_{j=1}^d S_{ij}$  is the  $i$ -th diagonal element of  $M$ .  $L_S = M - S$  is the Laplacian matrix of the feature similarity matrix  $S$ , and  $S_{ij}$  is the  $i$ -th row and  $j$ -th element of the feature similarity matrix  $S$ , representing the similarity between features  $f_i$  and  $f_j$ . There are many ways to construct feature similarity matrix  $S$ , for example:

By using a kernel function, a feature association matrix  $S$  can be constructed, where  $t \in R$ :

$$S_{ij} = \begin{cases} \exp(-\frac{\|f_i - f_j\|_2^2}{t}), & \text{if } f_i \in N_K(f_j) \text{ or } f_j \in N_K(f_i) \\ 0, & \text{others} \end{cases} \tag{8}$$

where  $N_K(*)$  represents the  $k$ -nearest neighbor set of  $*$ . Through feature map regularization, the problem of feature selection is optimized:

$$\min_W L(W) + \frac{\lambda}{2} Tr(W^T L_S W) + \frac{\beta}{2} \|W\|_{2,1} \tag{9}$$

### 2.4 Label manifold learning

In order to better fit the label information while fitting the manifold structure. According to the problem, suppose: Let  $f(x_i W) = [g(x_i w_1), g(x_i w_2), \dots, g(x_i w_m)]$ ,  $f(x_i W) \in R^m$ , if the labels  $y_i$  and  $y_j$  are closer, then the probability  $f(x_i W)$  and  $f(x_j W)$  in the logistic regression model should also be closer, and according to the positive correlation between  $g(x_i w_j)$  and  $x_i w_j$ , the positive correlation between  $f(x_i W)$  and  $x_i W$  is deduced, thus  $x_i W$  should be closer to  $x_j W$ . Therefore, a regularization of the label graph is constructed, which can adjust the coefficient matrix  $W$  according to the similarity between the labels  $y_i$  and  $y_j$ , so that  $W$  can better fit the label information. Its expression is as follows:

$$\begin{aligned}
 & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \|x_i W - x_j W\|_2^2 A_{ij} \\
 = & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i W - x_j W)(x_i W - x_j W)^T A_{ij} \\
 = & \sum_{i=1}^n x_i W (x_i W)^T P_{ii} - \sum_{i=1}^n \sum_{j=1}^n x_i W (x_j W)^T A_{ij} \tag{10} \\
 = & Tr(W^T X^T (P - A) X W) \\
 = & Tr(W^T X^T L_A X W) \\
 = & Tr(W^T \overline{L_A} W)
 \end{aligned}$$

where,  $\overline{L_A} = X^T L_A X$  and  $P \in R^{n \times n}$  are diagonal matrices,  $P_{ii} = \sum_{j=1}^n A_{ij}$  is the  $i$ -th diagonal element of

$P$ .  $L_A = P - A$  is the Laplacian matrix of label similarity matrix  $A$ , and  $A_{ij}$  is the element of the  $i$ -th row and the  $j$ -th column of label similarity matrix  $A$ , representing the similarity between the labels  $y_i$  and  $y_j$ . The label similarity matrix  $A$  can be given by many methods, as follows:

By using a kernel function, a label association matrix  $A$  can be constructed, where  $t \in R$ :

$$A_{ij} = \begin{cases} \exp(-\frac{\|y_i - y_j\|_2^2}{t}), & \text{if } y_i \in N_K(y_j) \text{ or } y_j \in N_K(y_i) \\ 0, & \text{others} \end{cases} \tag{11}$$

As for several different calculation methods of feature similarity matrix  $S$  and label similarity matrix  $A$ . The impact on multi-label feature selection, we have made a simple analysis on the Image and Emotion data sets, set the parameter range as  $[0.001, 0.01, 0.1, 1, 10, 100, 1000]$  to search and get the best result. As shown in Fig. 1 below, and found that several methods are similar. So in the experiment part, we use kernel function to learn the feature similarity matrix and label similarity matrix.

Through label map regularization, the optimization feature selection problem is transformed into:

$$\min_W L(W) + \frac{\lambda}{2} Tr(W^T L_S W) + \frac{\beta}{2} \|W\|_{2,1} + \frac{\gamma}{2} Tr(W^T \overline{L_A} W) \tag{12}$$

## 3 Problem solving and proof of convergence

### 3.1 Problem solving

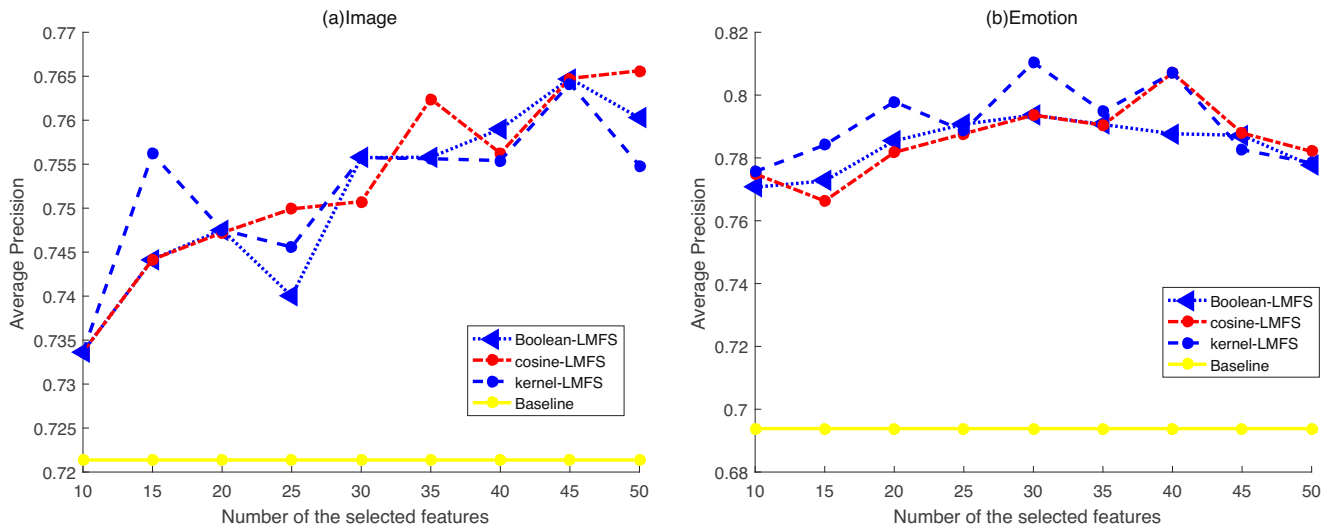
Due to the non-smoothness of  $L_{21}$ -norm, it is difficult to find the closed solution of the optimization problem in (12) directly. According to [29], this problem can be solved by another method. When  $W_i \neq 0 (i = 1, 2, \dots, d)$ , the derivative of  $\|W\|_{2,1}$  to  $W$  is:

$$\frac{\partial(\|W\|_{2,1})}{\partial W} = 2HW \tag{13}$$

where  $H \in R^{d \times d}$  is the diagonal matrix and the  $i$ -th diagonal element of  $H$  is:

$$H_{ii} = \frac{1}{2\|W_i\|_2} \tag{14}$$

Therefore, the derivative in  $L_{21}$ -norm can also be regarded as the derivative of  $Tr(W^T H W)$ . Since  $\|W\|_{2,1}$  is convex, the optimization problem of  $L_{21}$ -norm can be used



**Fig. 1** Average precision comparison of different similarity matrix methods when ML-KNN is used as the basic classifier. (the higher the result, the better)

to find the approximate solution of (12). Thus the objective function is transformed into:

$$obj(W) = L(W) + \frac{\lambda}{2}Tr(W^T L_S W) + \frac{\beta}{2}Tr(W^T H W) + \frac{\gamma}{2}Tr(W^T \overline{L}_A W) \quad (15)$$

For this problem, we can give an  $H$ , calculate  $W$  with the current  $H$ , and then update  $H$  based on the currently calculated  $W$ .

Since (15) is differentiable, it can be solved by the Newton–Raphson algorithm. The first derivative of (15) to  $W$  is:

$$\frac{\partial(obj(W))}{\partial W} = -X^T[Y - G(XW)] + \lambda L_S W + \beta H W + \gamma \overline{L}_A W \quad (16)$$

where  $G(XW) = [f(x_1 W); f(x_2 W); \dots; f(x_n W)]$ . The second derivative of (15) to  $W$  is:

$$\frac{\partial^2(obj(W))}{\partial W \partial W^T} = -X^T U X + \lambda L_S + \beta H + \gamma \overline{L}_A \quad (17)$$

Among them:

$$U = diag \sum_{j=1}^m [(1 - g(x_i w_j))g(x_i w_j)] \quad (18)$$

where  $i = 1, 2, \dots, d$ .

The updated formula for  $W$  is:

$$W^{t+1} = W^t - \left(\frac{\partial^2(obj(W))}{\partial W \partial W^T}\right)^{-1} \frac{\partial(obj(W))}{\partial W} \quad (19)$$

**Algorithm 1** LMFS.

Input: The data matrix  $X \in R^{n \times d}$ , the label matrix  $Y \in R^{n \times m}$ , three regularization parameters  $\lambda, \beta$  and  $\gamma$ , and the number of the selected features  $k$ .

1. Calculate the feature similarity matrix  $S$  according to (8), and calculate  $L = M - S$ .
2. Calculate the label similarity matrix  $A$  according to (9), and calculate  $L_A = P - A$  and  $\overline{L}_A = X^T L_A X$ .
3. Initialization matrix  $H$  as the unit array and coefficient matrix  $W$  as a matrix whose elements are all 0.
4. Repeat:
  - (a) Calculate the loss function value according to (12).
  - (b) Update  $U$  according to (18).
  - (c) Calculate the first derivative function and second derivative function of (15) according to (16) and (17).
  - (d) Update  $W$  according to (19).
  - (e) Update  $H$  according to (14).
5. Until convergence criterion has been satisfied.
6. Calculate and rank  $\|W_i\|_2 (i = 2, 3, \dots, d)$  to find the top  $k$  largest assignments to  $I$ .

Output: The feature selection result  $I$ .

In LMFS algorithm, the main purpose is to calculate the update  $U, \frac{\partial(obj(W))}{\partial W}, \frac{\partial^2(obj(W))}{\partial W \partial W^T}, W$  and  $H$ . In each iteration, the complexity of update  $U$  is  $O(mn^2)$ , the complexity

of update  $\frac{\partial(obj(W))}{\partial W}$  is  $O(d^2n)$ , the complexity of update  $\frac{\partial^2(obj(W))}{\partial W \partial W^T}$  is  $O(d^2)$ , the complexity of update  $W$  is  $O(d^2m)$ , and the complexity of update  $H$  is  $O(dm)$ . The LMFS algorithm has iterated  $t$  times in total. Therefore, the total complexity of the LMFS algorithm is  $O(t(d^2n + d^2m + d^2 + nm^2 + dm))$ , and the value of  $t$  is not large. Therefore, the running time of the LMFS algorithm processing data is greatly affected by the dimension  $d$  of data, the number of labels  $m$ , and the number of samples  $n$  in the data set.

### 3.2 Proof of convergence

In this section, we prove that the iterative procedure shown in Algorithm 1 is convergent. Therefore, in the  $t$ -th iteration, we know:

$$W^{t+1} = \underset{W}{\operatorname{argmin}} L(W) + \frac{\lambda}{2} \operatorname{Tr}(W^T L_S W) + \frac{\beta}{2} \operatorname{Tr}(W^T H^t W) + \frac{\gamma}{2} \operatorname{Tr}(W^T \overline{L}_A W) \tag{20}$$

where  $H_{ii}^t = \frac{1}{2\|W_i^t\|_2}$  ( $i = 1, 2, \dots, d$ ), so we have:

$$\begin{aligned} L(W^{t+1}) &+ \frac{\lambda}{2} \operatorname{Tr}((W^{t+1})^T L_S W^{t+1}) \\ &+ \frac{\beta}{2} \operatorname{Tr}((W^{t+1})^T H^t W^{t+1}) \\ &+ \frac{\gamma}{2} \operatorname{Tr}((W^{t+1})^T \overline{L}_A W^{t+1}) \leq L(W^t) \\ &+ \frac{\lambda}{2} \operatorname{Tr}((W^t)^T L_S W^t) \\ &+ \frac{\beta}{2} \operatorname{Tr}((W^t)^T H^t W^t) + \frac{\gamma}{2} \operatorname{Tr}((W^t)^T \overline{L}_A W^t) \end{aligned} \tag{21}$$

That is:

$$\begin{aligned} L(W^{t+1}) &+ \frac{\lambda}{2} \operatorname{Tr}((W^{t+1})^T L_S W^{t+1}) \\ &+ \frac{\gamma}{2} \operatorname{Tr}((W^{t+1})^T \overline{L}_A W^{t+1}) \\ &+ \frac{\beta}{2} \sum_{i=1}^d \frac{\|W_i^{t+1}\|_2^2}{2\|W_i^t\|_2} \leq L(W^t) \\ &+ \frac{\lambda}{2} \operatorname{Tr}((W^t)^T L_S W^t) \\ &+ \frac{\gamma}{2} \operatorname{Tr}((W^t)^T \overline{L}_A W^t) + \frac{\beta}{2} \sum_{i=1}^d \frac{\|W_i^t\|_2^2}{2\|W_i^t\|_2} \end{aligned} \tag{22}$$

It can be further transformed into:

$$\begin{aligned} L(W^{t+1}) &+ \frac{\lambda}{2} \operatorname{Tr}((W^{t+1})^T L_S W^{t+1}) \\ &+ \frac{\gamma}{2} \operatorname{Tr}((W^{t+1})^T \overline{L}_A W^{t+1}) \\ &+ \frac{\beta}{2} \|W^{t+1}\|_{2,1} - \frac{\beta}{2} (\|W^{t+1}\|_{2,1} \\ &- \sum_{i=1}^d \frac{\|W_i^{t+1}\|_2^2}{2\|W_i^t\|_2}) \leq L(W^t) \\ &+ \frac{\lambda}{2} \operatorname{Tr}((W^t)^T L_S W^t) + \frac{\gamma}{2} \operatorname{Tr}((W^t)^T \overline{L}_A W^t) \\ &+ \frac{\beta}{2} \|W^t\|_{2,1} - \frac{\beta}{2} (\|W^t\|_{2,1} - \sum_{i=1}^d \frac{\|W_i^t\|_2^2}{2\|W_i^t\|_2}) \end{aligned} \tag{23}$$

According to the inequality  $\sqrt{a} - \frac{a}{2\sqrt{b}} \leq \sqrt{b} - \frac{b}{2\sqrt{b}}$  for any positive numbers  $a$  and  $b$ , we have:

$$\|W_i^{t+1}\|_{2,1} - \frac{\|W_i^{t+1}\|_2^2}{2\|W_i^t\|_2} \leq \|W_i^t\|_{2,1} - \frac{\|W_i^t\|_2^2}{2\|W_i^t\|_2} \tag{24}$$

Which sums, we get:

$$\sum_{i=1}^d \left( \|W_i^{t+1}\|_{2,1} - \frac{\|W_i^{t+1}\|_2^2}{2\|W_i^t\|_2} \right) \leq \sum_{i=1}^d \left( \|W_i^t\|_{2,1} - \frac{\|W_i^t\|_2^2}{2\|W_i^t\|_2} \right) \tag{25}$$

Which implies:

$$\|W^{t+1}\|_{2,1} - \sum_{i=1}^d \frac{\|W_i^{t+1}\|_2^2}{2\|W_i^t\|_2} \leq \|W^t\|_{2,1} - \sum_{i=1}^d \frac{\|W_i^t\|_2^2}{2\|W_i^t\|_2} \tag{26}$$

In summary, the convergence of Algorithm 1 is proved.

## 4 Experiments and results

In order to verify the effectiveness of the LMFS algorithm, the experiment uses eight public data sets and compares its performance with some of the most advanced methods and baselines. At the same time, the experiment selects ML-KNN [30] as the representative of the multi-label classification algorithm for evaluation.

**Table 1** Dataset information

data	sample	features	label	training	test	species
Image	600	294	5	400	200	image
Emotion	593	72	6	391	202	music
Enron	1702	1001	53	1123	579	text
Business	5000	438	30	2000	3000	text
Computers	5000	681	33	2000	3000	text
Health	5000	612	32	2000	3000	text
Scene	2407	294	6	1211	1196	image
Yeast	2417	103	14	1500	917	biological

**Table 2** Hamming loss comparison of different algorithms under each data set

algorithms	LMFS	DRMFS	SCLS	MDMR	PMU	FIMF	Baseline
Business	<b>0.0264</b>	0.0280	0.0274	0.0273	0.0284	0.0274	<u>0.0269</u>
Image	<b>0.1950</b>	<u>0.2020</u>	0.2110	0.2240	0.2270	0.2340	0.2130
Emotion	<b>0.2063</b>	<u>0.2244</u>	0.2500	0.2409	0.2673	0.2252	0.2937
Health	<b>0.0370</b>	0.0391	<u>0.0389</u>	0.0391	0.0457	0.0407	0.0458
Computers	<b>0.0381</b>	<u>0.0393</u>	0.0398	0.0398	0.0416	0.0409	0.0412
Enron	<u>0.0486</u>	<b>0.0478</b>	0.0495	0.0505	0.0505	0.0501	0.0520
Scene	<u>0.1006</u>	0.1126	0.1073	0.1348	0.1137	0.1587	<b>0.0989</b>
Yeast	<u>0.1943</u>	<b>0.1938</b>	0.2006	0.1999	0.2006	0.2021	0.1980

#### 4.1 Dataset and experimental setup

The experiment uses eight public data sets from four different areas. The specific parameters of each data set are shown in Table 1:

In terms of experimental environment, all experimental related environments are: Microsoft Windows7 system, processor: Intel (R) Core (TM) i5-4210U CUP @ 1.70GHz 2.40GHz, memory: 4.00GB, programming software: Matlab R2016a.

To verify the effectiveness of the proposed feature selection method, the following most advanced state-of-the-art feature selection algorithms are compared:

- 1) Baseline: The results on various evaluation indicators after learning the data set directly with ML-KNN without any feature selection.
- 2) DRMFS [24]: A robust multi-label feature selection with dual-graph regularization was constructed by using feature graph and label graph to guide the sparsity between rows and within rows of the weight matrix, and  $L_{2,1}$  norm to guide its global properties and robust.

- 3) SCLS [31]: Multi-label feature selection method based on scalable standards.
- 4) MDMR [32]: Multi-label feature selection through an evaluation metric that combines mutual information with maximum dependency and minimum redundancy.
- 5) PMU [33]: A multi-label feature selection algorithm based on mutual information. Multi-label feature selection is performed by selecting the dependency between the selected feature and the label.
- 6) FIMF [34]: A fast multi-label feature selection method based on information theory feature ranking. Based on information theory, a scoring function that evaluates the importance of features is derived, and its calculation cost is analyzed.

In order to ensure the fairness of the experiment, in terms of parameter setting:

The number of nearest neighbors  $K$  for the multi-label classification algorithm ML-KNN is set to 10, and the value of smooth  $S$  is set to 1. For MDMR, PMU, and FIMF, we discretize the data set using the equal-width intervals [35]. For FIMF, we set  $Q = 10$  and  $b = 2$ . For DRMFS

**Table 3** Ranking loss comparison of different algorithms under each data set

algorithms	LMFS	DRMFS	SCLS	MDMR	PMU	FIMF	Baseline
Business	<u>0.0382</u>	0.0443	0.0405	0.0404	0.0444	0.0423	<b>0.0374</b>
Image	<b>0.2000</b>	0.2213	<u>0.2167</u>	0.2550	0.2483	0.2662	0.2333
Emotion	<b>0.1662</b>	<u>0.1687</u>	0.2056	0.1994	0.2570	0.2012	0.2829
Health	<b>0.0530</b>	0.0577	<u>0.0562</u>	0.0566	0.0679	0.0578	0.0605
Computers	<b>0.0844</b>	0.0934	0.0909	<u>0.0903</u>	0.0980	0.0955	0.0922
Enron	<b>0.0885</b>	<u>0.0896</u>	0.0921	0.0944	0.0949	0.0935	0.0938
Scene	<u>0.1014</u>	0.1087	0.1129	0.1444	0.1290	0.1994	<b>0.0931</b>
Yeast	<u>0.1677</u>	<b>0.1656</b>	0.1745	0.1710	0.1723	0.1747	0.1715

**Table 4** One-error comparison of different algorithms under each data set

algorithms	LMFS	DRMFS	SCLS	MDMR	PMU	FIMF	Baseline
Business	<b>0.1177</b>	0.1303	0.1240	0.1237	0.1320	0.1260	<u>0.1213</u>
Image	<b>0.3650</b>	<u>0.3950</u>	0.4000	0.4450	0.4700	0.5000	0.4350
Emotion	<b>0.2426</b>	<u>0.2772</u>	0.3614	0.3564	0.3614	0.3515	0.4059
Health	<b>0.3197</b>	<u>0.3333</u>	0.3410	0.3373	0.4403	0.3723	0.4207
Computers	<b>0.4150</b>	<u>0.4340</u>	0.4580	0.4543	0.4700	0.4627	0.4367
Enron	<u>0.2297</u>	<b>0.2124</b>	0.2470	0.2435	0.2694	0.2453	0.3040
Scene	<u>0.2651</u>	0.2968	0.2977	0.3905	0.3904	0.4983	<b>0.2425</b>
Yeast	<b>0.2094</b>	<u>0.2137</u>	0.2268	0.2366	0.2366	0.2366	0.2345

and LMFS, we use (8) to calculate the similarity matrix between features. The above settings are the default settings of the algorithms. In addition, For DRMFS and other comparative algorithms, the experiment adjusts the regularization parameters of all methods by the “grid search” strategy from [0.001, 0.01, 0.1, 1, 10, 100, 1000]. For the feature dimension, we set the number of selected features as [10, 15, 20, 25, 30, 35, 40, 45, 50]. The maximum number of iterations for alliterative algorithms is fixed as 50. At the same time, the size of neighborhood  $K$  is set as 5. For all multi-label feature selection algorithms, the experiments show the best results from the optimal parameters.

**4.2 Evaluation metrics**

The performance evaluation of the multi-label learning systems is different from the single-label learning systems. The evaluation criteria of the multi-label learning system are more complicated. The experiment uses five evaluation criteria: Hamming loss, Ranking loss, One-error, Coverage, and Average precision in ML-KNN. The specific contents of the five evaluation criteria are as follows:

Suppose there is a test data set  $D = \{(x_i, y_i)\}_{i=1}^n$ , where  $n$  is the number of test samples. Given test sample  $x_i$ , the binary label vector that is predicted by the multi-label classifier is denoted as  $h(x_i)$ , and the rank of the  $l$ -th label prediction is denoted as  $rank_l(l)$ .

1) Hamming loss: evaluates the percentage of mislabeled labels, i.e., a label belonging to the instance is not predicted or a label not belonging to the instance is predicted. The smaller the value, the better the performance.

$$HL(D) = \frac{1}{n} \sum_{i=1}^n \frac{1}{m} \|h(x_i) \Delta y_i\|_1 \tag{27}$$

where  $\Delta$  represents the symmetric difference between the two sets, and returns those values that appear only in one of the sets,  $HL(D) \in [0, 1]$ .

2) Ranking loss: evaluates the proportion of reverse-order label pairs, that is, the case where unrelated labels are more relevant than related labels. The smaller the value, the better the performance.

$$RL(D) = \frac{1}{n} \sum_{i=1}^n \frac{1}{1_m^T y_i 1_m^T \bar{y}_i} \sum_{l: y_l'=1} \sum_{l': y_{l'}'=0} (\delta(rank_l(l) \geq rank_{l'}(l'))) \tag{28}$$

where  $\bar{y}_i$  is the complement of  $y_i$  in  $Y$ .  $RL(D) \in [0, 1]$ .

3) One-error: evaluates the proportion of samples that “the most relevant label is not” in “real labels”. The smaller

**Table 5** Coverage comparison of different algorithms under each data set

algorithms	LMFS	DRMFS	SCLS	MDMR	PMU	FIMF	Baseline
Business	<u>2.1990</u>	2.4123	2.3050	2.2917	2.4020	2.3600	<b>2.1847</b>
Image	<b>1.0800</b>	<u>1.1600</u>	1.1650	1.3200	1.2900	1.3550	1.2150
Emotion	<b>1.8911</b>	<u>1.9158</u>	2.1139	2.0891	2.3614	2.0545	2.4901
Health	<b>3.0053</b>	3.2060	<u>3.1350</u>	3.1647	3.5690	3.1780	3.3047
Computers	<b>4.1187</b>	4.4987	4.3697	<u>4.3480</u>	4.6520	4.5580	4.4160
Enron	<b>12.6790</b>	<u>12.8450</u>	13.0415	13.1606	13.4128	13.2038	13.2055
Scene	<u>0.6104</u>	0.6421	0.6681	0.8253	0.7492	1.0953	<b>0.5686</b>
Yeast	<u>6.2955</u>	<b>6.2857</b>	6.4482	6.3642	6.3708	6.3740	6.4144



**Table 6** Average precision comparison of different algorithms under each data set

algorithms	LMFS	DRMFS	SCLS	MDMR	PMU	FIMF	Baseline
Business	<b>0.8809</b>	0.8689	0.8758	0.8757	0.8690	0.8730	<u>0.8798</u>
Image	<b>0.7642</b>	0.7434	<u>0.7437</u>	0.7058	0.7002	0.6791	0.7214
Emotion	<b>0.8104</b>	<u>0.7917</u>	0.7496	0.7551	0.7143	0.7510	0.6938
Health	<b>0.7429</b>	<u>0.7245</u>	0.7159	<u>0.7256</u>	0.6593	0.7090	0.6812
Computers	<b>0.6590</b>	<u>0.6344</u>	0.6317	0.6304	0.6093	0.6203	0.6334
Enron	<b>0.6742</b>	<u>0.6704</u>	0.6589	0.6566	0.6483	0.6548	0.6232
Scene	<u>0.8359</u>	0.8158	0.8163	0.7633	0.8034	0.6906	<b>0.8512</b>
Yeast	<b>0.7667</b>	<u>0.7653</u>	0.7563	0.7579	0.7562	0.7552	0.7585

the value, the better the performance.

$$OE(D) = \frac{1}{n} \sum_{i=1}^n \delta(y_i^{l_i} = 0) \quad (29)$$

where  $l_i = \operatorname{argmin}_{l \in [1, m]} \operatorname{rank}_i(l)$  and  $\delta$  are indicator functions,  $OE(D) \in [0, 1]$ .

4) Coverage: evaluates how many steps the "sorted label list" needs to move, on the average, to cover the true related label set. The smaller the value, the better the performance.

$$CV(D) = \frac{1}{n} \sum_{i=1}^n \operatorname{argmax}_{l: y_i^l = 1} \operatorname{rank}_i(l) - 1 \quad (30)$$

where  $CV(D) \in [1, m - 1]$ .

5) Average precision: evaluates the proportion of those labels that are more relevant than particular labels. The larger the value, the better the performance.

$$AP(D) = \frac{1}{n} \sum_{i=1}^n \frac{1}{1_m^T y_i} \sum_{l: y_i^l = 1} \frac{\operatorname{prec}_i(l)}{\operatorname{rank}_i(l)} \quad (31)$$

where  $\operatorname{prec}_i(l) = \sum_{l': y_i^{l'} = 1} \delta(\operatorname{rank}_i(l) \geq \operatorname{rank}_i(l'))$  and  $AP(D) \in [0, 1]$ .

### 4.3 Experimental results

The proposed multi-label feature selection algorithm has been tested in six public data sets with extensive experi-

ments. Comparing with several state-of-the-art algorithms, we consider evaluation metrics of Hamming loss, Ranking loss, One-error, Coverage, and Average precision to evaluate the performance of the above multi-label feature selection methods. Tables 2, 3, 4, 5 and 6 show the best results of all the feature selection methods from the optimal parameters. The best performance is indicated in bold font in these tables, and the second-best performance is underlined. Tables 2, 3, 4, 5 and 6 report Hamming loss, Ranking loss, One-error, Coverage, and Average precision comparison of different algorithms in each data set.

First of all, feature selection is adequate. It reduces the number of features, shortens the classifier's running time, and improves the performance of the classification algorithm. Secondly, as shown in Tables 2, 3, 4, 5 and 6, although the performance of the LMFS algorithm on data sets Business, Scene and Yeast is slightly inadequate, the performance of the LMFS algorithm on data sets Business and Scene is second only to the Baseline. In addition, the LMFS algorithm has the best performance on other data sets.

In order to visually show the relative performance of the LMFS algorithm and other comparing algorithms, Figs. 2, 3, 4, 5 and 6 shows the performance of all multi-label feature selection algorithms. As the number of selected features changes, the value of the metrics will also change. Therefore the  $x$ -axis represents the number of features selected by each feature selection algorithm, and the  $y$ -axis represents the performance of the evaluation metrics after classification feature selection. These results show that the proposed LMFS algorithm is superior to the previous ones among almost all data sets in most cases.

Specifically, Figs. 2 to 6 show the Hamming loss, Ranking loss, One-error, Coverage, and Average precision comparison of different feature selection methods when

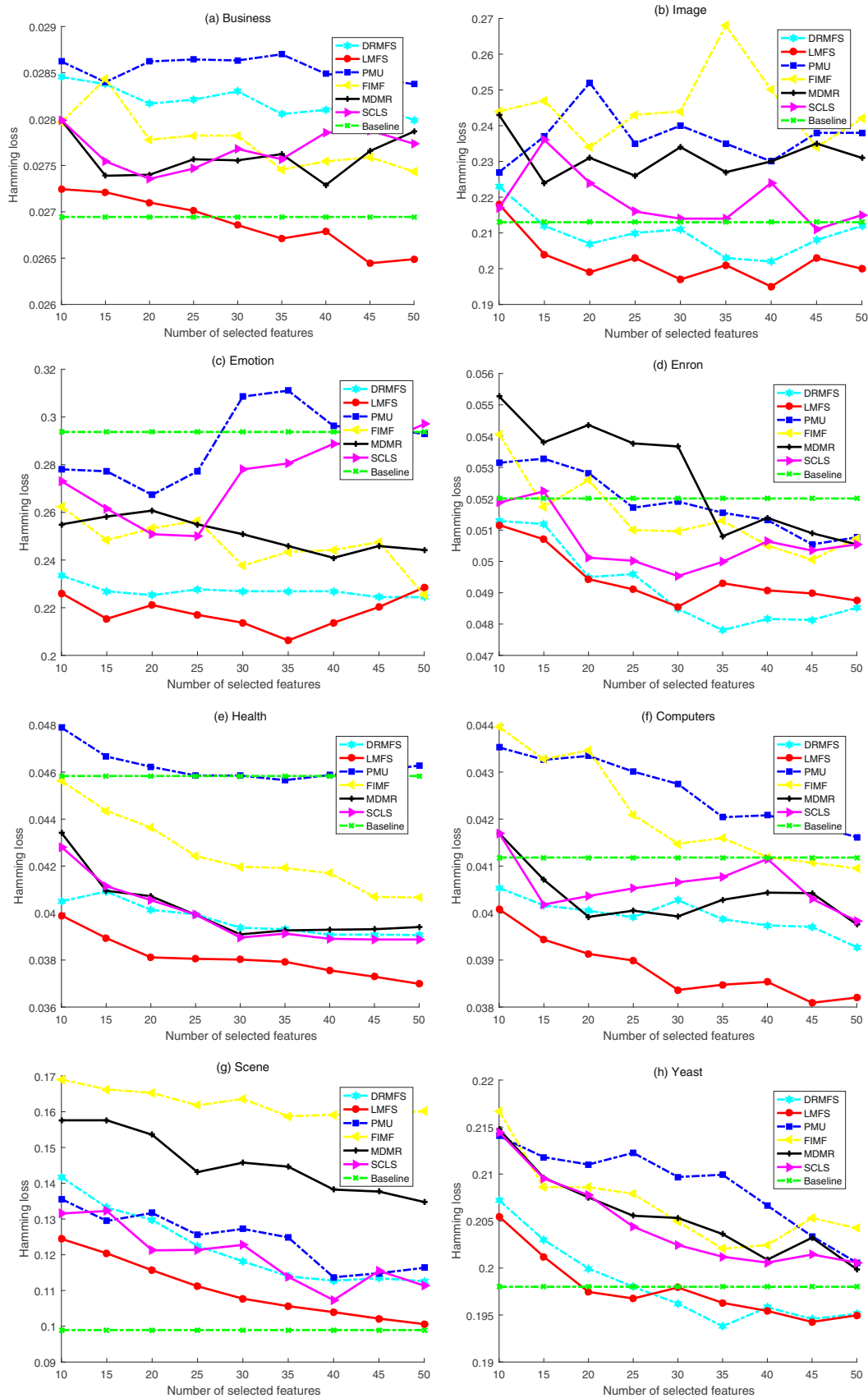
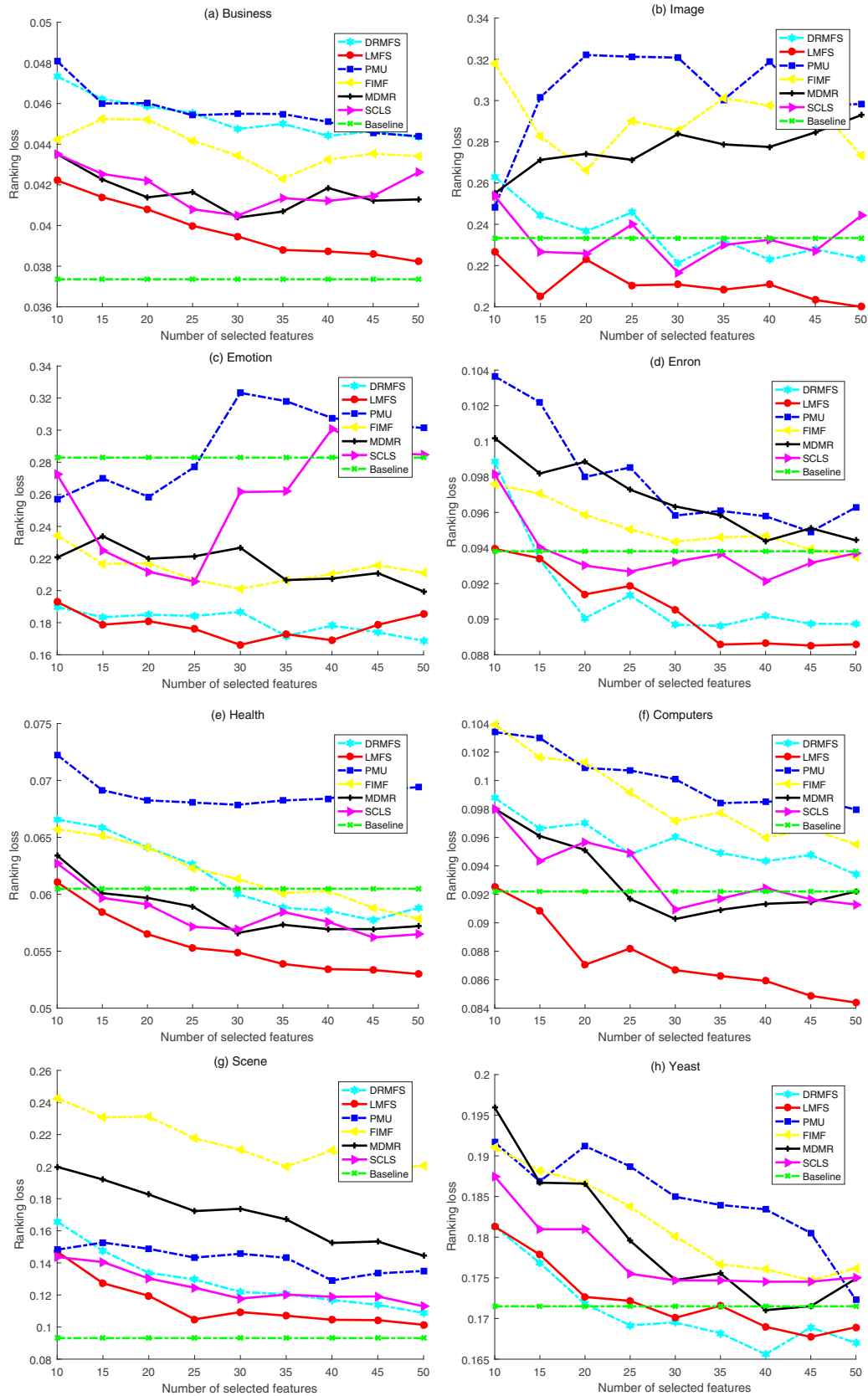


Fig. 2 Hamming loss comparison of different feature selection methods when ML-KNN is used as the basic classifier. (the lower the result, the better)



**Fig. 3** Ranking loss comparison of different feature selection methods when ML-KNN is used as the basic classifier. (the lower the result, the better)

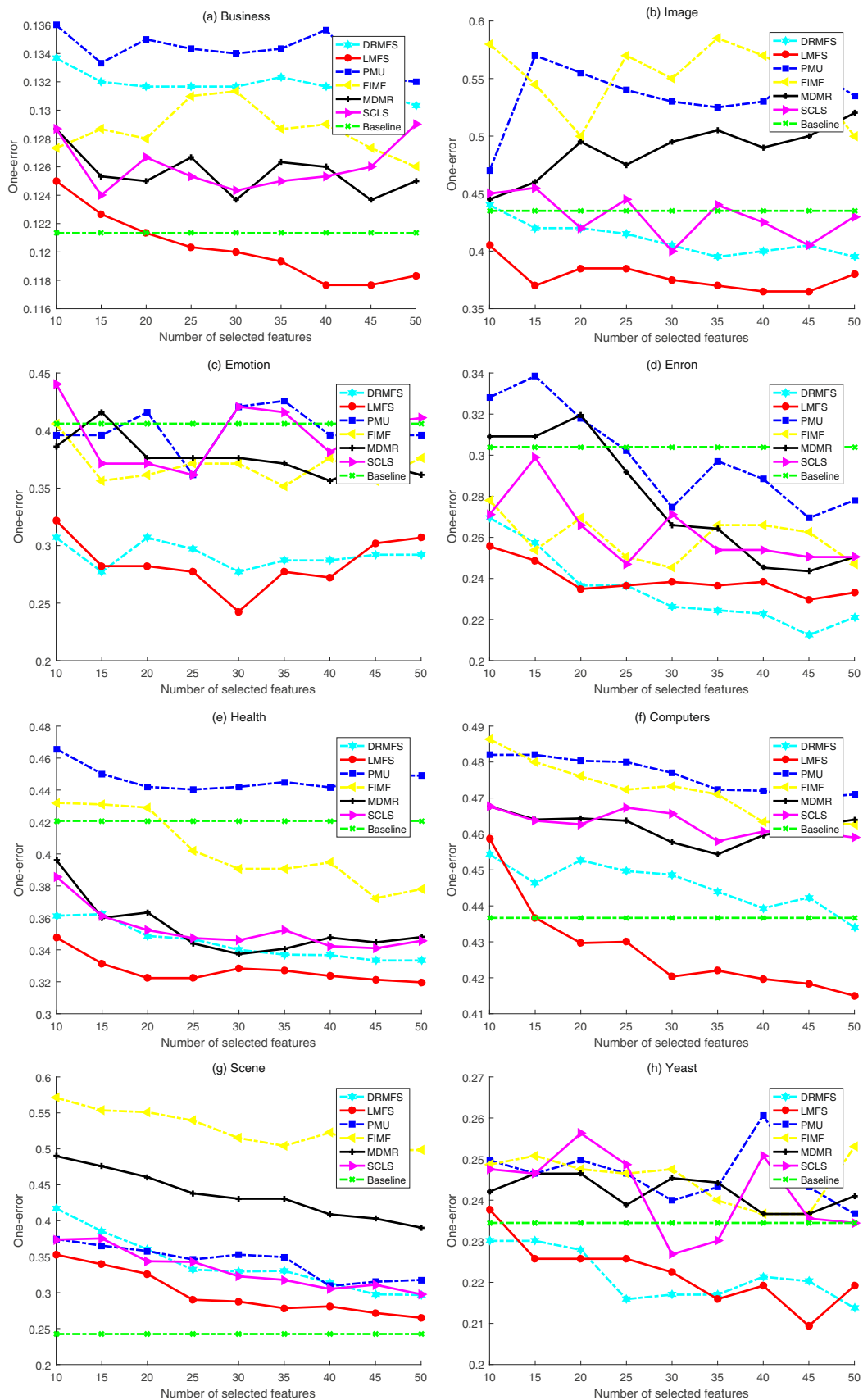


Fig. 4 One-error comparison of different feature selection methods when ML-KNN is used as the basic classifier. (the lower the result, the better)

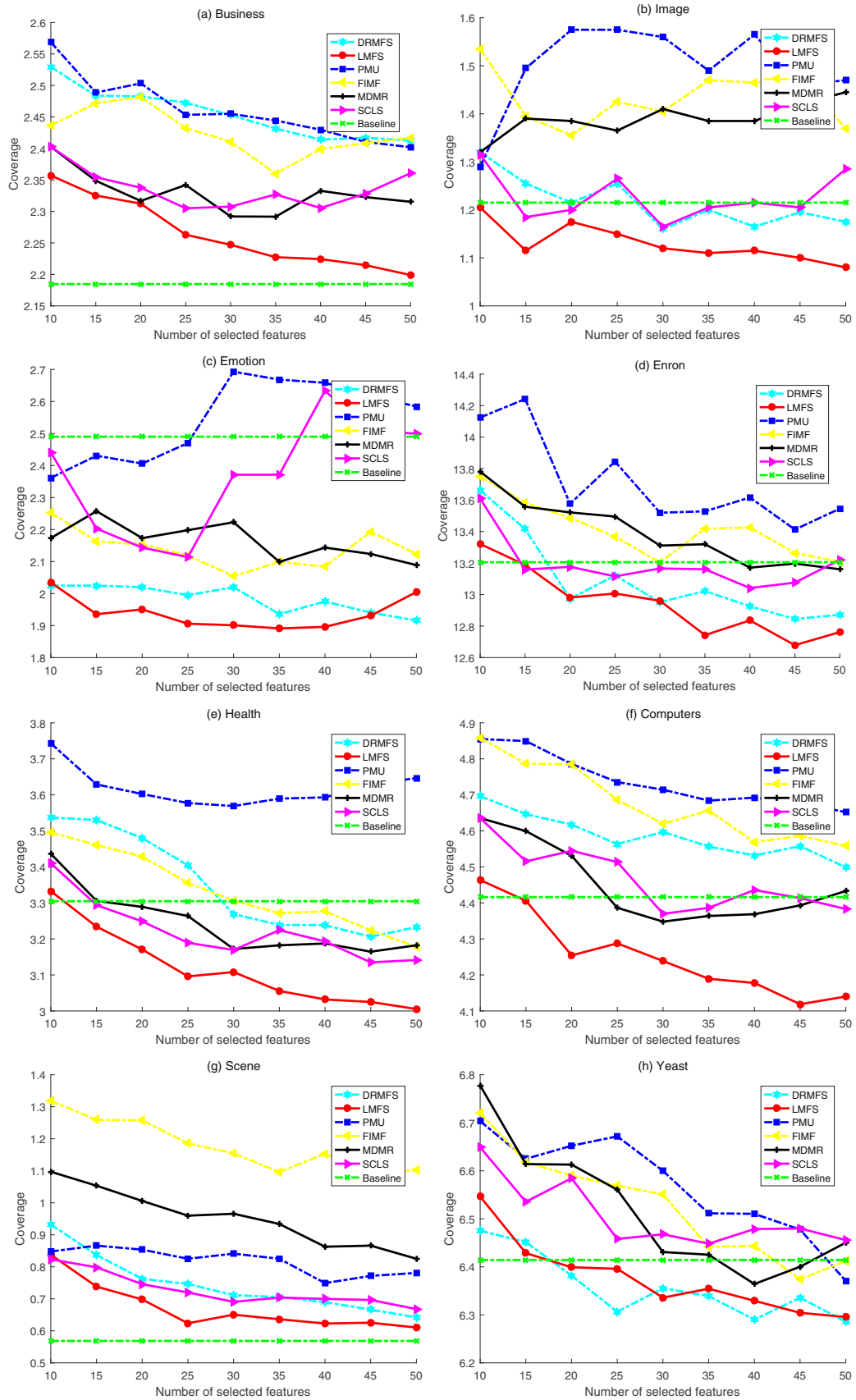
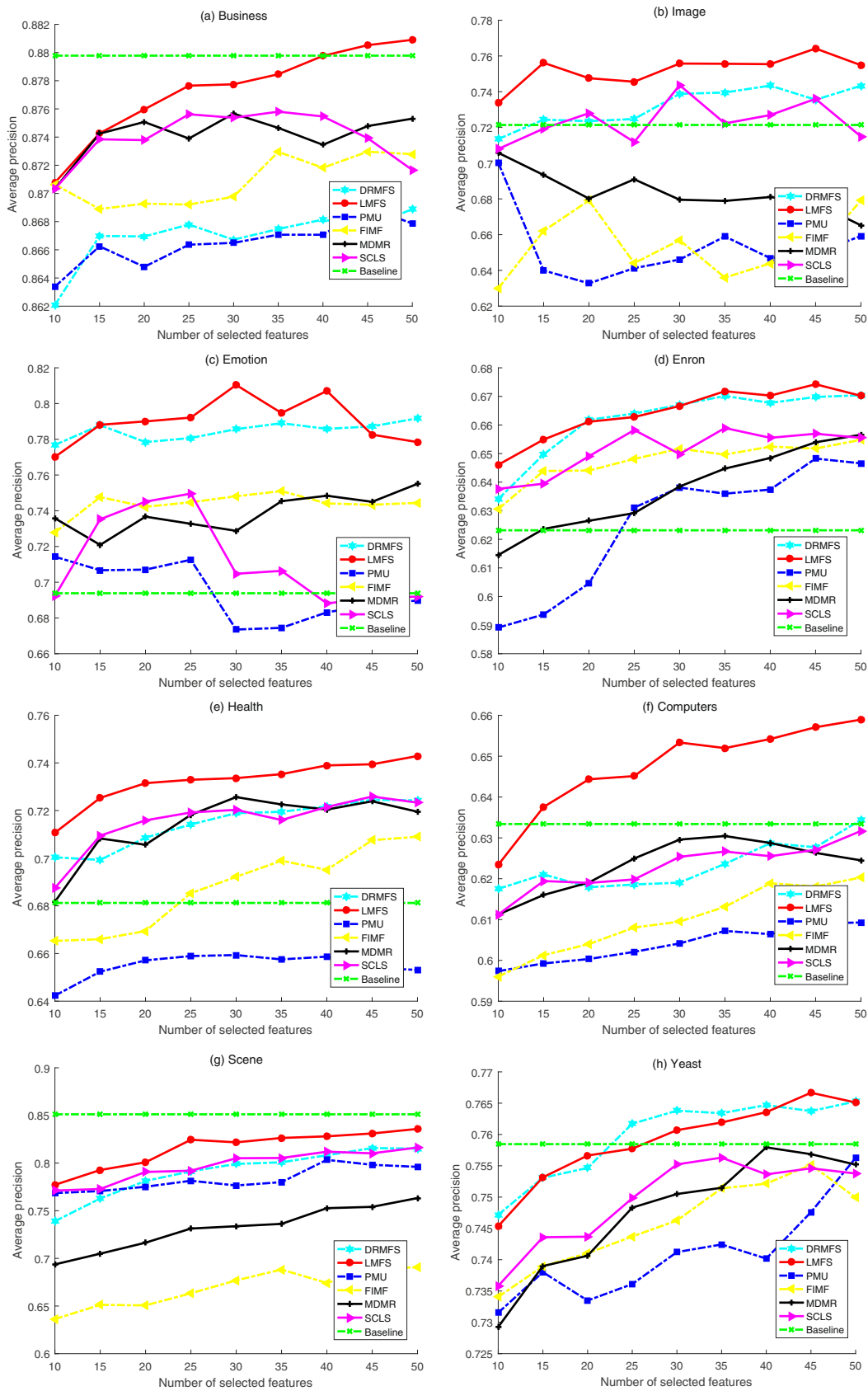
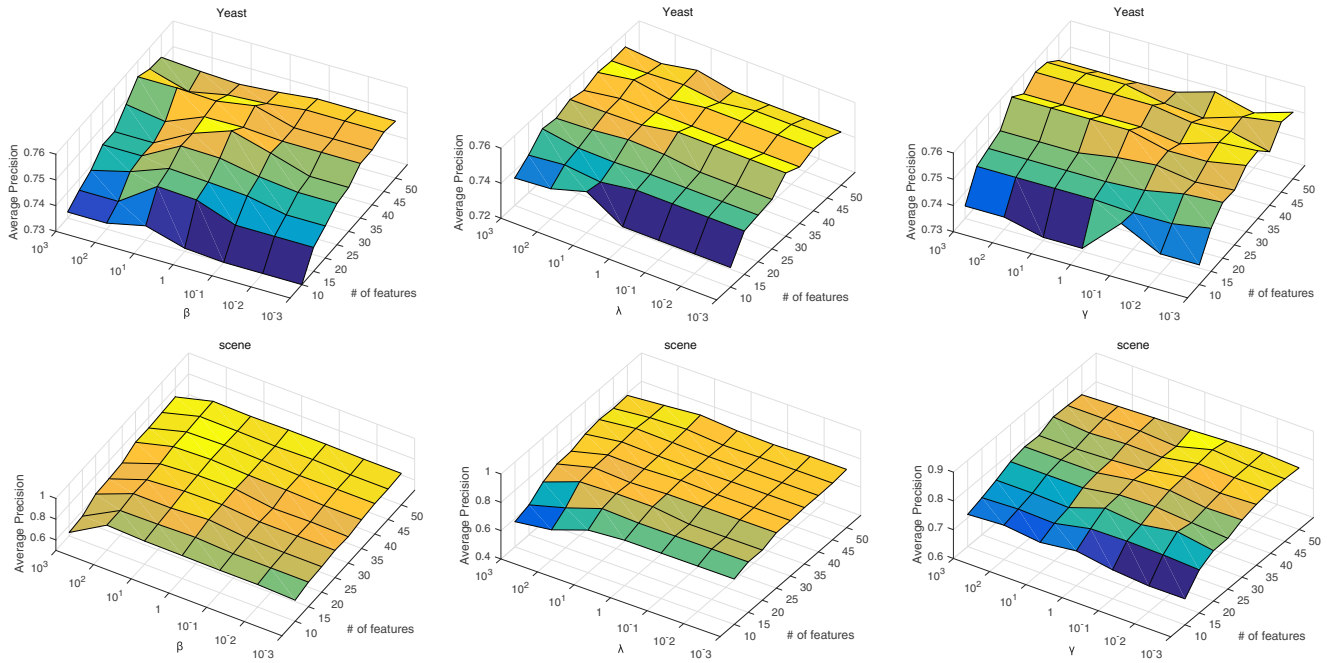


Fig. 5 Coverage comparison of different feature selection methods when ML-KNN is used as the basic classifier. (the lower the result, the better)



**Fig. 6** Average precision comparison of different feature selection methods when ML-KNN is used as the basic classifier. (the higher the result, the better)

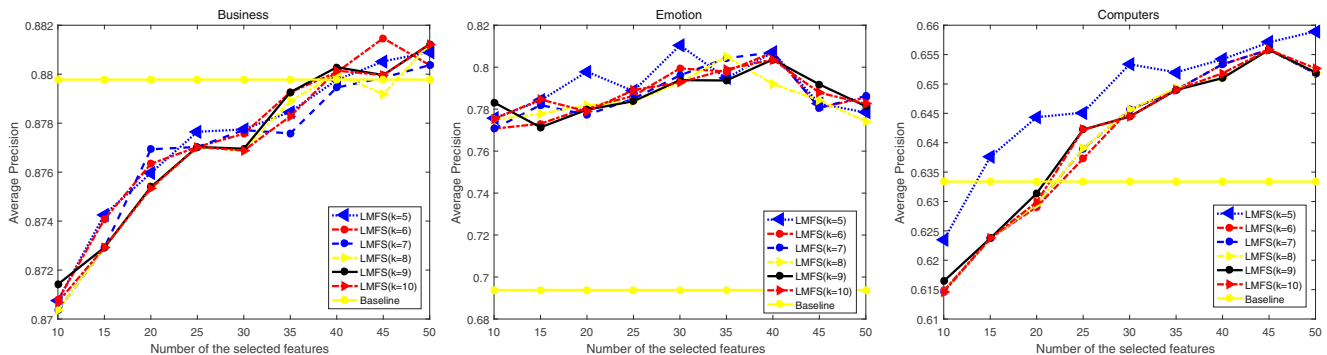


**Fig. 7** The change of Average precision with parameters in Enron and Scene data set

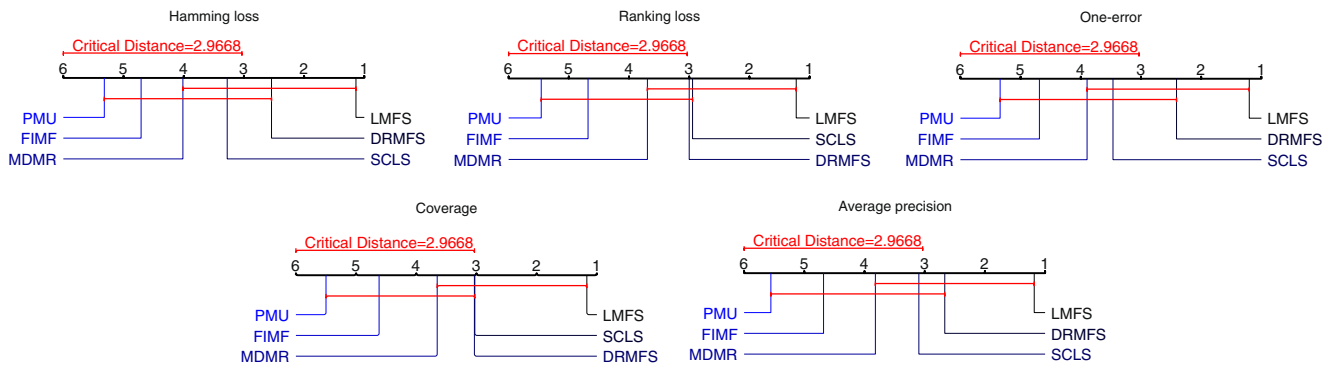
using ML-KNN as the primary classifier. From a, b, c, e, f and g in Fig. 2a, b, e, f and g in Fig. 3a, b, e, f and g in Fig. 4; and a, b, c, d, e, f and g in Fig. 5, we can see that the curve of LMFS algorithm is lower than that of all comparison algorithms, even for the other images in Figs. 2 to 5. The curves of the LMFS algorithm are significantly lower than those of the SCLS algorithm, MDMR algorithm, PMU algorithm, and FIMF algorithm. Even in the subfigure of Figs. 2 and 4, only the LMFS algorithm’s curves are below the baseline. From a, b, c, e, f, and g in Fig. 6, we can see that the curve of LMFS is significantly higher than that of all comparison algorithms, even in a and f, only the LMFS algorithm’s curves were above the baseline. Thus, it can be seen that the proposed LMFS algorithm can reduce irrelevant or redundant features.

In addition, to explore the influence of parameters on the performance of the LMFS algorithm, we choose two different kinds of data sets: music data set Scene and biological data set Yeast. For parameters  $\lambda$ ,  $\beta$  and  $\gamma$ , we fix two of them as 1. The influence of another parameter on the performance of the LMFS algorithm is discussed under the selection of a different number of features. The parameter range is set as [0.001, 0.01, 0.1, 1, 10, 100, 1000], the number of feature selection is set to [10, 15, 20, 25, 30, 35, 40, 45, 50]. The experimental results are shown in Fig. 7.

Specifically, the performance of the algorithm will change with the change of parameters. As shown in Fig. 7, for different data sets, the optimal range of parameters is different. For example, on the Scene data set the optimal range of parameter  $\beta$  is [10, 100], the optimal range of parameter  $\lambda$  is [10, 1000], and the optimal range



**Fig. 8** The influence of the nearest neighbor parameter k on the performance of the algorithm



**Fig. 9** Bonferroni-Dunn test results in the form of average rank diagrams. Groups of feature selection algorithms that are not significantly different (at  $p = 0.1$ ) are connected

of parameter  $\gamma$  is  $[0.01, 0.1]$ . Due to the different basic structures of different data sets, the parameters in the Yeast data set are more sensitive, but the parameters in the Scene data set are not sensitive. At the same time, in order to explore the influence of the nearest neighbor parameter  $K$  on the performance of the algorithm, let  $K = [5, 6, 7, 8, 9, 10]$ , Experiments were carried out on three data sets, Business, Emotion and Computers, the experimental results are shown in Fig. 8, we can see that the performance of the algorithm is sensitive to  $K$ .

As shown in Fig. 9, the horizontal axis represents the sorting of multi-label feature selection algorithms under each index; from left to right, the algorithm's performance is getting better and better, the best performing algorithm is on the far right side. At the same time, we report the results of the Bonferroni-Dunn test in the form of an average rank graph, the algorithm groups with no significant difference ( $P < 0.1$ ) were connected, if the difference of average ranking reaches the critical value of the difference (CD), then there is significant difference [36]. Although the LMFS algorithm has no significant difference with the DRMFS algorithm, SCLS algorithm, and MDMR algorithm in all indicators, the LMFS algorithm always has a significant difference with PMU and FIMF algorithm, and the LMFS algorithm consistently ranks on the right side. Therefore, compared with other methods, the LMFS algorithm shows better performance.

## 5 Conclusions and outlook

In this paper, logistic regression is combined with feature manifold learning, sparse regularization, and label map regularization to study the multi-label feature selection problem. Sparseness has been widely used in regression-based feature selection methods. In order to overcome the shortcomings that when dealing with the regression

coefficient of features, the existing feature selection method based on logistic regression fails to consider the geometric structure of the feature manifold; and that the existing linear regression feature selection method fails to consider the relevant label information between the geometric structure of the feature manifold and the manifold feature coefficient, we embed feature map regularization method and label map regularization method into the multi-label feature selection problem based on logistic regression to obtain the regression coefficients, so that the regression coefficients are smooth relative to the feature manifold without losing the relevant label information. We also design an iterative update algorithm to prove the convergence of the LMFS algorithm. Another direction in the future is to extend this method to study the semi-supervised feature selection.

**Acknowledgments** This work was supported by the Natural Science Foundation of China (61976130), the Key Research and Development Project of Shaanxi Province (2018KW-021), the Natural Science Foundation of Shaanxi Province (2020JQ-923).

## Declarations

**Conflict of Interests** The authors declared that they have no conflicts of interest to this work. We declare that we do not have any commercial or associative interest that represents a conflict of interest in connection with the work submitted.

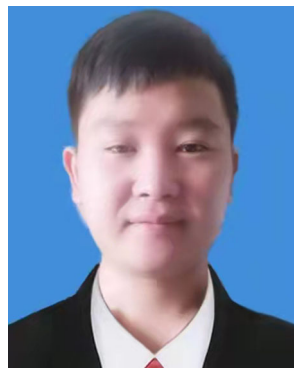
## References

1. Cai J, Luo JW, Wang SL et al (2018) Feature selection in machine learning: a new perspective[J]. *Neurocomputing* 300:70–79
2. Tang C, Liu XW, Zhu XZ et al (2019) Feature selective projection with low-rank embedding and dual laplacian regularization[J]. *IEEE Trans Knowl Data Eng* 32(9):1747–1760
3. Bermingham ML, Pong-Wong R, Spiliopoulou A et al (2015) Application of high-dimensional feature selection: evaluation for genomic prediction in man[J]. *Scientific Reports* 5:10312



4. Sun X, Liu YH, Li J et al (2012) Using cooperative game theory to optimize the feature selection problem[J]. *Neurocomputing* 97:86–93
5. Zhang R, Nie FP, Li XL et al (2019) Feature selection with multi-view data: a survey[J]. *Information Fusion* 50:158–167
6. Ding CC, Zhao M, Lin J et al (2019) Multi-objective iterative optimization algorithm based optimal wavelet filter selection for multi-fault diagnosis of rolling element bearings[J]. *ISA Transactions* 82:199–215
7. Labani M, Moradi P, Ahmadizar F et al (2018) A novel multivariate filter method for feature selection in text classification problems[J]. *Engineering Applications of Artificial Intelligence* 70:25–37
8. Yao C, Liu YF, Jiang B et al (2017) LLE score: a new filter-based unsupervised feature selection method based on nonlinear manifold embedding and its application to image recognition[J]. *IEEE Transactions on Image Processing* 26(11):5257–5269
9. Gonzalez J, Ortega J, Damas M et al (2019) A new multi-objective wrapper method for feature selection - accuracy and stability analysis for BCI[J]. *Neurocomputing* 333:407–418
10. Swati J, Hongmei H, Karl J (2018) Information gain directed genetic algorithm wrapper feature selection for credit rating[J]. *Applied Soft Computing* 69:541–553
11. Maldonado S, López J (2018) Dealing with high-dimensional class-imbalanced datasets: embedded feature selection for SVM classification[J]. *Applied Soft Computing* 67:94–105
12. Kong YC, Yu TW (2018) A graph-embedded deep feedforward network for disease outcome classification and feature selection using gene expression data[J]. *Bioinformatics* 34(21):3727–3737
13. Yang XH, Jiang XY, Tian CX et al (2020) Inverse projection group sparse representation for tumor classification: a low rank variation dictionary approach[J]. *Knowledge-Based Systems* 196:105768
14. Deng TQ, Ye DS, Ma R et al (2020) Low-rank local tangent space embedding for subspace clustering[J]. *Inf Sci* 508:1–21
15. Xiao Q, Dai JH, Luo JW et al (2019) Multi-view manifold regularized learning-based method for prioritizing candidate disease miRNAs[J]. *Knowledge-Based Systems* 175:118–129
16. Tang C, Zheng X, Liu XW et al (2021) Cross-view locality preserved diversity and consensus learning for multi-view unsupervised feature selection[J]. *IEEE Transactions on Knowledge and Data Engineering* 99:1–1
17. Tang C, Liu XW, Li MM et al (2018) Robust unsupervised feature selection via dual self-representation and manifold regularization[J]. *Knowledge-Based Systems* 145:109–120
18. Sun ZQ, Zhang J, Dai L et al (2019) Mutual information based multi-label feature selection via constrained convex optimization[J]. *Neurocomputing* 329:447–456
19. Zhang P, Liu GX, Gao WF (2019) Distinguishing two types of labels for multi-label feature selection[J]. *Pattern Recogn* 95:72–82
20. Chen LL, Chen DG (2019) Alignment based feature selection for multi-label learning[J]. *Neural Processing Letters* 50(7):28–36
21. Chen SB, Zhang YM, Ding CHQ et al (2019) Extended adaptive lasso for multi-class and multi-label feature selection[J]. *Knowledge-Based Systems* 173:28–36
22. Zhang J, Luo ZM, Li CD et al (2019) Manifold regularized discriminative feature selection for multi-label learning[J]. *Pattern Recognition* 95:136–150
23. Cai ZL, Zhu W (2018) Multi-label feature selection via feature manifold learning and sparsity regularization[J]. *Int J Machine Learning Cybern* 9(8):1321–1334
24. Hu JC, Li YH, Gao WF et al (2020) Robust multi-label feature selection with dual-graph regularization[J]. *Knowledge-Based Systems* 203:106126
25. Li Q, Xie B, You J et al (2016) Correlated logistic model with elastic net regularization for multilabel image classification[J]. *IEEE Transactions on Image Processing* 25(8):3801–3813
26. Sato T, Takano Y, Miyashiro R et al (2016) Feature subset selection for logistic regression via mixed integer optimization[J]. *Computational Optimization and Applications* 64(3):865–880
27. Yang ZY, Liang Y, Zhang H et al (2018) Robust sparse logistic regression with the  $L_q$  ( $0 < q < 1$ ) regularization for feature selection using gene expression data[J]. *IEEE Access* 6:68586–68595
28. Pan XL, Xu YT (2021) A safe feature elimination rule for L1-regularized logistic regression[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. <https://doi.org/10.1109/TPAMI.2021.3071138>
29. Zhang R, Nie FP, Li X (2017) Self-weighted supervised discriminative feature selection[J]. *IEEE Trans Neural Netw Learn Syst* 29(8):3913–3918
30. Zhang ML, Zhou ZH (2007) ML-KNN: a lazy learning approach to multi-label learning[J]. *Pattern Recognition* 40(7):2038–2048
31. Lee J, Kim DW (2017) SCLS: multi-label feature selection based on scalable criterion for large label set[J]. *Pattern Recognition* 66:342–352
32. Lin YJ, Hu QH, Liu JH et al (2015) Multi-label feature selection based on max-dependency and min-redundancy[J]. *Neurocomputing* 168:92–103
33. Lee J, Lim H, Kim DW (2012) Approximating mutual information for multi-label feature selection[J]. *Electron Lett* 48(15):929–930
34. Lee J, Kim DW (2015) Fast multi-label feature selection based on information-theoretic feature ranking[J]. *Pattern Recognition* 48(9):2761–2771
35. Dougherty J, Kohavi R, Sahami M et al (1995) Supervised and unsupervised discretization of continuous features[J]. In: *Machine learning: proceedings of the 12th international conference*, vol 2, pp 194–202
36. Demiar J, Schuurmans D (2006) Statistical comparisons of classifiers over multiple data sets[J]. *Journal of Machine Learning Research* 7(1):1–30

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Yao Zhang** received the B.S. degree in applied Mathematics from Yancheng Teachers University, China, 2018, the M.S. degree from Xi'an Polytechnic University, China, 2021. His research interests include multi-label learning, feature selection and clustering.



**Yingcang Ma** received the M.S. degree in Basic Mathematics from Shaanxi Normal University, China, in 2002, and the Ph.D. degree in Computer science and Software from Northwestern Polytechnical University, China, in 2006. He is currently a professor in the school of science, Xi'an Polytechnic University, China. His research interests include machine learning, fuzzy logic and neutrosophic theory.



**Xiaofei Yang** received the B.S. degree in applied mathematics from Luoyang Normal University, China, 2006, the M.S. degree in pure mathematics from Shaanxi Normal University, China, 2009, and the Ph.D. degree in pure mathematics from Shaanxi Normal University, China, 2012. His research interests include data mining, image processing and fuzzy mathematics.