

# Pattern Recognition

## Multi-Label Feature Selection with Global and Local label correlation

--Manuscript Draft--

<b>Manuscript Number:</b>	PR-D-23-01769
<b>Article Type:</b>	Full Length Article
<b>Section/Category:</b>	Features
<b>Keywords:</b>	feature selection; multi-label learning; label correlation; nonnegative matrix factorization
<b>Corresponding Author:</b>	Fardin Akhlaghian Tab, Ph.D. University of Kurdistan Sanandaj, IRAN, ISLAMIC REPUBLIC OF
<b>First Author:</b>	Mohammad Faraji
<b>Order of Authors:</b>	Mohammad Faraji  Seyed Amjad Seyed  Fardin Akhlaghian Tab, Ph.D.  Reza Mahmoodi
<b>Abstract:</b>	In various application domains, high-dimensional multi-label data has become more prevalent, presenting two significant challenges: instances with high-dimensional features and a large number of labels. In the context of multi-label feature selection, the objective is to choose a subset of features from a given set that is highly pertinent for predicting multiple labels or categories associated with each instance. However, certain characteristics of multi-label classification, such as label dependencies and imbalanced label distribution, have often been overlooked although they hold valuable insights for designing effective multi-label feature selection algorithms. In this paper, we propose a feature selection model which exploits explicit global and local label correlations to select discriminative features across multiple labels. In addition, by representing the feature matrix and label matrix in a shared latent space, the model aims to capture the underlying correlations between features and labels. The shared representation can reveal common patterns or relationships that exist across multiple labels and features. An objective function involving $L_{2,1}$ -norm regularization is formulated, and an alternating optimization-based iterative algorithm is designed to obtain the sparse coefficients for multi-label feature selection. The proposed method was evaluated on twelve real-world multi-label datasets using six evaluation metrics, through comprehensive experiments. The results indicate its effectiveness, surpassing that of several representative methods.
<b>Suggested Reviewers:</b>	Ronghua Shang, PhD Professor, Xidian University rhshang@mail.xidian.edu.cn Prof. Shang recently has published a related paper in the Pattern Recognition journal.  <a href="https://doi.org/10.1016/j.patcog.2021.107873">https://doi.org/10.1016/j.patcog.2021.107873</a>  Rui Huang, PhD Professor, Shanghai University huangr@shu.edu.cn Prof. Huang recently has published a related paper in the Pattern Recognition journal.  <a href="https://doi.org/10.1016/j.patcog.2021.108149">https://doi.org/10.1016/j.patcog.2021.108149</a>

# Multi-Label Feature Selection with Global and Local Label Correlation

Mohammad Faraji, Seyed Amjad Seyed, Fardin Akhlaghian Tab\*, Reza Mahmoodi  
Department of Computer Engineering, University of Kurdistan, Sanandaj, Iran

---

\*Corresponding Author.

Email addresses: [mohammad.faraji@uok.ac.ir](mailto:mohammad.faraji@uok.ac.ir) (Mohammad Faraji), [amjadseyedi@uok.ac.ir](mailto:amjadseyedi@uok.ac.ir) (Seyed Amjad Seyed), [f.akhlaghian@uok.ac.ir](mailto:f.akhlaghian@uok.ac.ir) (Fardin Akhlaghian Tab), [reza.mahmoodi@uok.ac.ir](mailto:reza.mahmoodi@uok.ac.ir) (Reza Mahmoodi).

## Highlights

- The proposed joint model learns implicit label correlation from multi-label data.
- To extract relevant features, the global and local label correlation are utilized.
- The data manifold regularization is employed to preserve the local structure.
- A unified optimization approach is provided to solve the proposed objective function.
- Results demonstrate the effectiveness of the method in multi-label feature selection.

# Multi-Label Feature Selection with Global and Local label correlation

Mohammad Faraji, Seyed Amjad Seyedi, Fardin Akhlaghian Tab\*, Reza Mahmoodi

*Department of Computer Engineering, University of Kurdistan, Sanandaj, Iran*

---

## Abstract

In various application domains, high-dimensional multi-label data has become more prevalent, presenting two significant challenges: instances with high-dimensional features and a large number of labels. In the context of multi-label feature selection, the objective is to choose a subset of features from a given set that is highly pertinent for predicting multiple labels or categories associated with each instance. However, certain characteristics of multi-label classification, such as label dependencies and imbalanced label distribution, have often been overlooked although they hold valuable insights for designing effective multi-label feature selection algorithms. In this paper, we propose a feature selection model which exploits explicit global and local label correlations to select discriminative features across multiple labels. In addition, by representing the feature matrix and label matrix in a shared latent space, the model aims to capture the underlying correlations between features and labels. The shared representation can reveal common patterns or relationships that exist across multiple labels and features. An objective function involving  $L_{2,1}$ -norm regularization is formulated, and an alternating optimization-based iterative algorithm is designed to obtain the sparse coefficients for multi-label feature selection. The proposed method was evaluated on twelve real-world multi-label datasets using six evaluation metrics, through comprehensive experiments. The results indicate its effectiveness, surpassing that of several representative methods.

*Keywords:* feature selection, multi-label learning, label correlation, nonnegative matrix factorization

---

\*Corresponding author

Email addresses: [mohammad.faraji@uok.ac.ir](mailto:mohammad.faraji@uok.ac.ir) (Mohammad Faraji), [amjadseyedi@uok.ac.ir](mailto:amjadseyedi@uok.ac.ir) (Seyed Amjad Seyedi), [f.akhlaghian@uok.ac.ir](mailto:f.akhlaghian@uok.ac.ir) (Fardin Akhlaghian Tab), [reza.mahmoodi@uok.ac.ir](mailto:reza.mahmoodi@uok.ac.ir) (Reza Mahmoodi)

<sup>1</sup> **1. Introduction**

<sup>2</sup> In recent years, feature selection has become increasingly important due to the detrimental  
<sup>3</sup> effects of the high dimensionality curse. These effects include elevated learning complexity [1, 2],  
<sup>4</sup> heightened space allocation [3], and reduced classification performance [4]. The number of features  
<sup>5</sup> can be decreased and classification accuracy can be increased by choosing only pertinent features  
<sup>6</sup> and eliminating redundant and irrelevant ones [5, 6]. The use of feature selection algorithms  
<sup>7</sup> has been widespread in several disciplines, including recommender system [7, 8, 9], picture and  
<sup>8</sup> video annotation [10], microarray data processing, and genomics [11]. Since objects in the real  
<sup>9</sup> world can have numerous labels at once, substantial research has been done on feature selection  
<sup>10</sup> in these domains. For instance, a single gene in genomics can perform several tasks, including  
<sup>11</sup> photosynthesis, protein breakdown, and signal transmission [12]. In music analysis, one piece of  
<sup>12</sup> music can have several different emotions at the same time, like sadness, joy, and scariness [13],  
<sup>13</sup> and how a newswire story can be classified in news categorization.

<sup>14</sup> Multi-label learning typically makes the correlation between the labels an assumption, unlike  
<sup>15</sup> multi-class problems. Because of this, extracting and using the correlation of labels, which has  
<sup>16</sup> made it into an NP-hard problem, is the main challenge in multi-label learning in addition to  
<sup>17</sup> the dimension of label space. Currently used techniques for selecting multi-label features empha-  
<sup>18</sup> size the extraction of label correlation, label-feature relevance, and feature correlation. Existing  
<sup>19</sup> approaches for multi-label feature selection are classified into two groups in order to handle multi-  
<sup>20</sup> label data: problem transformation models and algorithm adaption models [14]. In problem  
<sup>21</sup> transformation, multi-label data is converted to single-label data using problem transformation  
<sup>22</sup> techniques, then single-label feature selection methods are applied to the single-label data. To  
<sup>23</sup> directly handle multi-label data, researchers have also offered a number of algorithm adaption  
<sup>24</sup> models for multi-label feature selection [15, 16]. Label correlations are primarily extracted by  
<sup>25</sup> algorithm adaptation techniques to direct the feature selection procedure. Adaptation methods  
<sup>26</sup> outperform problem transformation methods that ignore label correlations in terms of classifi-  
<sup>27</sup> cation performance because they take label correlation into account. The three approaches for  
<sup>28</sup> algorithm adoption are first-order, second-order, and high-order. The correlation between labels  
<sup>29</sup> is not taken into account in the first-order approach because labels are thought of as independent  
<sup>30</sup> of one another. The second-order approach takes into account the correlation between label pairs,  
<sup>31</sup> and in practice, we will have a ranking between labels that are related and those that are not.  
<sup>32</sup> The correlation between a subset or the entire set of labels is taken into account in the high-order  
<sup>33</sup> approach.

34 In multi-label learning, the presence of label correlations can yield significant insights. For  
35 instance, if the labels "Ski" and "Snow" exist, it is highly likely that the label "Winter Sport" will  
36 also be there. Likewise, if the labels "hot" and "sunny" are both present, it is highly unlikely  
37 that the label "snow" will appear. Incorporating label correlations of various degrees is a goal  
38 of recent studies on multi-label learning [17]. Some studies [18] [19] [20] concentrate mainly on  
39 global label correlation that apply to all instances. However, some label correlation [21], [22]  
40 are unique to a local data subset. For instance, in internet scope, the word "Amazon" refers to  
41 "Amazon shop", whereas in Nature scope, the word "Amazon" refers to "Amazon Forest". Global  
42 or local label correlations have been the focus of prior research. But it's clearly better and more  
43 desirable to take into account both of them in the Multi-Label Feature Selection problem.

44 In light of the above analysis, this paper presents a novel feature selection method named  
45 Multi-Label Feature Selection with Global and Local label correlation (MLFS-GLOCAL). In this  
46 method, the label information is mapped into a low-dimensional reduced space that captures the  
47 implicit correlations among multiple labels. In the light of this assumption that correlated features  
48 must share similar labels, we map feature information into the same low-dimensional reduced  
49 space. Therefore, this method embeds the label and feature correlations into a shared space.  
50 Though this low-rank structure can be regarded as implicitly exploiting label correlations, there  
51 is still a potential capacity to consider label correlations explicitly. Therefore, we extract both  
52 global and local label correlations from training label information and it encourages the prediction  
53 to be similar on highly correlated labels. In this way, to find relevant features across multiple  
54 labels, the proposed method incorporates implicit and explicit label correlations and alleviates  
55 the negative influences of imperfect label information. Additionally, a manifold regularization  
56 is imposed to preserve the local feature structure in the latent space. Finally, to enhance the  
57 interpretability and to guide feature selection, we introduce sparse nonnegative matrix regression  
58 with  $L_{2,1}$ -norm. We highlight the following contributions made by this paper:

- 59 • Using global and local label correlation for selecting discriminative features.
- 60 • Fusing label and feature information into a shared low-dimensional space for extracting  
61 label-feature relevance.
- 62 • Introducing a Nonnegative Matrix Factorization (NMF) [23] model that has inherent clus-  
63 tering and interpretability properties for selecting the most discriminative features.
- 64 • Adopting a graph regularization to guarantee the consistency between the original feature  
65 space and the latent space.

- 66 • Developing an effective optimization scheme to solve the MLFS-GLOCAL method.

67 The remainder of the paper is structured as shown below. Section 2 presents foundational  
 68 concepts and related works. The details of the proposed model are presented in Section 3. Section  
 69 4 presents the experimental results. Finally, the conclusion and future works are provided in  
 70 Section 5.

71 **2. Preliminaries**

72 In this paper, we use bold uppercase letters to signify matrices, such as matrix  $\mathbf{A}$ . When  
 73 matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{a}_i$  and  $\mathbf{a}^{(j)}$  represent the  $i$ -th column and the  $j$ -th row of  $\mathbf{A}$ , respectively.  
 74 Additionally, scalars are signified by lowercase letters like  $b$ , whereas vectors are represented by  
 75 bold italicized lowercase letters like  $\mathbf{b}$ .  $\mathbf{A}^\top$  and  $\text{Tr}(\mathbf{A})$  stand in for the transpose and trace of  
 76  $\mathbf{A}$ , respectively.  $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d A_{i,j}^2}$  and  $\|\mathbf{A}\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^d A_{i,j}^2}$  demonstrate the  
 77 Frobenius norm and  $L_{2,1}$ -norm of matrix  $\mathbf{A}$ , respectively, where  $A_{ij}$  represents the  $(i, j)$ -th entry  
 78 in matrix  $\mathbf{A}$ . We define  $\mathbf{X} \in \mathbb{R}^{n \times d}$  as a  $d$  feature space with  $n$  instances, and  $\mathbf{Y} \in \mathbb{R}^{n \times l}$  indicates  
 79 the same instances in a  $l$  label space. If the  $i$ -th sample has the  $j$ -label, then  $Y_{i,j} = 1$ ; otherwise,  
 80  $Y_{i,j} = 0$ .

81 *2.1. Related Work*

82 In recent years, there has been an increase in the prevalence of feature selection techniques  
 83 for managing multi-label data, where each instance is associated with multiple labels. The re-  
 84 search on multi-label feature selection has advanced quickly as more multi-label data have been  
 85 studied. Through the research on the most advanced multi-label feature selection, the existing  
 86 multi-label feature selection methods are mainly based on information-theoretic and embedded-  
 87 based methods. Information-theoretic techniques use mutual information or conditional mutual  
 88 information to extract the correlation between each candidate feature and each class label. As an  
 89 illustration, the Max-Dependence and Min-Redundancy (MDMR) [24] approach selects the best  
 90 feature subset by increasing the feature dependence between features and labels while reducing  
 91 feature redundancy. The following objective function shows how MDMR actually works:

$$J(\mathbf{f}_k) = \sum_{\mathbf{l}_i \in L} I(\mathbf{f}_k; \mathbf{l}_i) - \frac{1}{|S|} \sum_{\mathbf{f}_j \in S} \left\{ I(\mathbf{f}_k; \mathbf{f}_j) - \sum_{\mathbf{l}_i \in L} I(\mathbf{f}_k; \mathbf{l}_i | \mathbf{f}_j) \right\}, \quad (1)$$

92 where  $\mathbf{f}_k$ ,  $\mathbf{f}_j$ , and  $S$ , stand for candidate feature, already-selected feature, and the already-selected  
 93 feature subset, respectively. In contrast,  $I(\mathbf{f}_k; \mathbf{l}_i)$  quantifies the feature dependence. The feature  
 94 redundancy is measured by  $I(\mathbf{f}_k; \mathbf{f}_j) - \sum_{\mathbf{l}_i \in L} I(\mathbf{f}_k; \mathbf{l}_i | \mathbf{f}_j)$ . Similar to this, the SCLS feature  
 95 selection approach proposed by *Lee et al.* [25] consists of the feature relevance term and scalable  
 96 relevance evaluation.

$$J(\mathbf{f}_k) = \sum_{\mathbf{l}_i \in L} I(\mathbf{f}_k; \mathbf{l}_i) - \sum_{\mathbf{f}_j \in S} \frac{I(\mathbf{f}_k; \mathbf{f}_j)}{H(\mathbf{f}_k)} \sum_{\mathbf{l}_i \in L} I(\mathbf{f}_k; \mathbf{l}_i), \quad (2)$$

97 where the  $k$ -th feature's entropy is denoted by  $H(\mathbf{f}_k)$ . Recently, the LRFS technique for multi-  
 98 label feature selection was suggested [26]. It is based on redundant labels. Labels are divided  
 99 into independent and dependent categories by LRFS. Following is a presentation of the model:

$$J(\mathbf{f}_k) = LR(\mathbf{f}_k; L) - \frac{1}{|S|} \sum_{\mathbf{f}_j \in S} I(\mathbf{f}_k; \mathbf{f}_j) = \sum_{\mathbf{l}_i \in L} \left\{ \sum_{\mathbf{l}_i \neq \mathbf{l}_j, \mathbf{l}_j \in L} I(\mathbf{f}_k; \mathbf{l}_j | \mathbf{l}_i) - \frac{1}{|S|} \sum_{\mathbf{f}_j \in S} I(\mathbf{f}_k; \mathbf{f}_j) \right\}. \quad (3)$$

100 Methods based on information theory disregard high-order interaction connections between  
 101 features and labels. The significance of each individual feature or label is therefore a key factor  
 102 in how effective these strategies are. Instead, multi-label feature selection embedded-based ap-  
 103 proaches emphasize the use of label correlations, using label correlations to choose the compact  
 104 feature subset. In recent years, there have been several different sparse embedded-based feature  
 105 selection techniques [16, 27, 28]. As one of the basic methods, *Nie et al.* [27] introduced the  
 106 effective and Robust Feature Selection via the joint  $L_{2,1}$ -norm minimization (RFS) which can be  
 107 formulated as follows:

$$\min_{\mathbf{W}} \sum_{i=1}^n \|\mathbf{x}_i \mathbf{W} - \mathbf{y}_i\|_2 + \gamma \|\mathbf{W}\|_{2,1} = \min_{\mathbf{W}} \|\mathbf{X} \mathbf{W} - \mathbf{Y}\|_{2,1} + \gamma \|\mathbf{W}\|_{2,1}, \quad (4)$$

108 where  $\mathbf{W} \in \mathbb{R}^{d \times c}$  is a coefficient matrix, and the most discriminating features are chosen using  
 109  $\|\mathbf{W}\|_{2,1}$ .

110 Although RFS is a multi-class feature selection method, its framework is suitable for multi-  
 111 label ones, and it is employed in numerous multi-label feature selection methods. For instance,  
 112 Zhu et al. [6] developed the RFS model for Missing Label Multi-Label Feature Selection (MLMLFS).  
 113 This robust linear regression utilized a graph regularization based on the assumption that simi-

lar instances have similar labels. Additionally, a subset of discriminant features was selected by imposing an  $L_{2,1}$ ,  $p$ -norm constraint with  $0 < p \leq 1$ . Since learning the multi-label regression models on a binary label matrix is challenging, recent researches attempt to utilize alternatives for label matrix. These studies can be categorized into two methods either using a latent label matrix or a pseudo-label matrix as the regression goal. The latent-based method known as MIFS stands for Multi-label Informed Feature Selection [16] and takes advantage of implicit label correlations to choose the most discriminating features. Additionally, MIFS takes into account the reduced low-dimensional label matrix to prevent exponential growth in the total number of features and labels. The objective function of MIFS is as follows:

$$\min_{\mathbf{W}, \mathbf{V}, \mathbf{B}} \|\mathbf{X}\mathbf{W} - \mathbf{V}\|_F^2 + \alpha \|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2 + \beta \text{Tr}(\mathbf{V}^\top \mathbf{L}\mathbf{V}) + \gamma \|\mathbf{W}\|_{2,1}, \quad (5)$$

where  $\mathbf{W}$  is a coefficient matrix of features and  $\mathbf{X}$  is a feature matrix. The label matrix  $\mathbf{Y}$ 's coefficient matrix is denoted by the  $\mathbf{V}$ , while its latent semantics matrix is denoted by the  $\mathbf{B}$ . In other words, MIFS extracts common latent information from the feature and label matrices. This pioneer concept has been developed in some state-of-the-art multi-label feature selection methods. The MIFS method decomposes the multi-label matrix into two-factor matrices, which contain entries of mixed signs. As a result, interpreting them can be challenging. *Brayteey et al.* [29] presented CMFS, a multi-label feature selection approach that maps feature and label matrices into a shared low-dimensional space by a joint tri-factorization and local structure preservation. Also, this model attempts to maximize the dependence between latent correlations that are extracted from feature and label matrices.

$$\begin{aligned} \min_{\mathbf{V}, \mathbf{L}, \mathbf{Q}, \mathbf{P}, \mathbf{B}} & \|\mathbf{X} - \mathbf{V}\mathbf{L}\mathbf{Q}\|_F^2 + \alpha \|\mathbf{Y} - \mathbf{V}\mathbf{P}\mathbf{B}\|_F^2 + \beta \|\mathbf{L} - \mathbf{P}\|_F^2 + \epsilon \text{Tr}(\mathbf{R}(\mathbf{V}\mathbf{P}\mathbf{B})^\top \mathbf{V}\mathbf{P}\mathbf{B}) + \gamma \|\mathbf{Q}\|_{2,1} \\ \text{s.t. } & \mathbf{V}, \mathbf{L}, \mathbf{Q}, \mathbf{P}, \mathbf{B} \geq 0. \end{aligned} \quad (6)$$

Shared Common Mode Feature Selection (SCMFS) [30] is another extension of MIFS that similarly recovers the shared latent information between feature space and label space in an NMF-based model. To verify the matrix regression term, SCMFS adds a decoder factorization term on the feature matrix.

$$\min_{\mathbf{V}, \mathbf{B}, \mathbf{W}} \|\mathbf{X}\mathbf{W} - \mathbf{V}\|_F^2 + \alpha \|\mathbf{X} - \mathbf{V}\mathbf{Q}\|_F^2 + \beta \|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2 + \gamma \|\mathbf{W}\|_{2,1}, \text{ s.t. } \mathbf{W}, \mathbf{V}, \mathbf{Q}, \mathbf{B} \geq 0. \quad (7)$$

137 More recently, *Gao et al.* [31] proposed the Shared Structure Feature Selection (SSFS) that  
138 changed the regression term (encoder structure) in the MIFS to a factorization term (decoder  
139 structure), and similar to the SCMFS, optimized its cost function in an NMF-based model.

$$\min_{\mathbf{V}, \mathbf{Q}, \mathbf{M}} \|\mathbf{X} - \mathbf{V}\mathbf{Q}^\top\|_F^2 + \alpha \|\mathbf{Y} - \mathbf{V}\mathbf{M}\|_F^2 + \beta \text{Tr}(\mathbf{V}^\top \mathbf{L} \mathbf{V}) + \gamma \|\mathbf{Q}\|_{2,1}, \text{ s.t. } \mathbf{V}, \mathbf{M}, \mathbf{Q} \geq 0. \quad (8)$$

140 There are various pseudo-label-based methods in the MLFS literature that make a connection  
141 between the label matrix and its alternative in different ways. For example, *Huang et al.* [32]  
142 proposed an MLFS method with manifold regularization and dependence maximization (MRDM).  
143 This method replaces the label space with a manifold embedding. This embedding is constrained  
144 through manifold regularization. In addition, they use Hilbert Schmidt Independence Criterion  
145 (HSIC) as a regularization to maximize the dependence between the manifold embedding and the  
146 label matrix.

$$\min_{\mathbf{W}, \mathbf{Z}^\top \mathbf{Z} = \mathbf{I}} \|\mathbf{X}\mathbf{W} - \mathbf{Z}\|_F^2 + \alpha \text{Tr}(\mathbf{Z}^\top \mathbf{L} \mathbf{Z}) - \beta \text{Tr}(\mathbf{H} \mathbf{Z} \mathbf{Z}^\top \mathbf{H} \mathbf{Y} \mathbf{Y}^\top) + \gamma \|\mathbf{W}\|_{2,1}. \quad (9)$$

147 *Fan et al.* [33] proposed a dual manifold regularized framework that embeds the feature matrix  
148 into two latent spaces (label space and cluster space). Also, this model utilizes a regularization  
149 to maximize the dependence between the manifold embedding and the label matrix. The implicit  
150 label correlation is exploited by preserving the global and local structural information.

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{V}, \mathbf{F}, \mathbf{Q}} & \|\mathbf{X}^\top \mathbf{W} - \mathbf{V}\|_{2,1} + \alpha \|\mathbf{W}\|_{2,1} + \beta \|\mathbf{W} - \mathbf{F}\mathbf{Q}\|_F^2 + \gamma \text{Tr}(\mathbf{F}^\top \mathbf{X} \mathbf{L} \mathbf{X}^\top \mathbf{F}) \\ & + \lambda (\text{Tr}[(\mathbf{V} - \mathbf{Y})^\top \mathbf{E}(\mathbf{V} - \mathbf{Y})] + \text{Tr}(\mathbf{V}^\top \mathbf{L} \mathbf{V})), \text{ s.t. } \mathbf{V} \geq 0, \mathbf{W}^\top \mathbf{W} = \mathbf{I}, \mathbf{F}^\top \mathbf{F} = \mathbf{I}. \end{aligned} \quad (10)$$

151 Similarly, the authors [34] proposed a robust method that integrates multi-label feature se-  
152 lection and the local discriminant model. It clusters the feature weight matrix by incorporating  
153 discriminative information to exploit implicit label correlation.

$$\begin{aligned} & \min_{\mathbf{W}, \mathbf{b}, \mathbf{L}, \mathbf{P}} \|\mathbf{X}\mathbf{W} + \mathbf{1}_n \mathbf{b}^\top - \mathbf{Y}\|_{2,1} + \alpha \|\mathbf{W} - \mathbf{L}\mathbf{P}\|_F^2 \\ & + \beta \text{Tr}(\mathbf{L}^\top \mathbf{M} \mathbf{L}) + \gamma \|\mathbf{W}\|_{2,1}, \quad \text{s.t.} \quad \mathbf{L}^\top \mathbf{L} = \mathbf{I}. \end{aligned} \quad (11)$$

154 *Zhang et al.* [15] presented an NMDG or nonnegative multi-label feature selection with dy-  
155 namic graph constraints. In NMDG, the pseudo-label matrix is trained using label manifolds and  
156 linear regression. Additionally, the feature manifold is merged with the pseudo-label to create the  
157 dynamic graph Laplacian matrix, which is then utilized to constrain the learning of the feature  
158 weight matrix.

$$\begin{aligned} & \min_{\mathbf{W}, \mathbf{b}, \mathbf{F}} \|\mathbf{X}\mathbf{W} + \mathbf{1}_n \mathbf{b}^\top - \mathbf{F}\|_F^2 + \alpha \text{Tr}(\mathbf{F}^\top \mathbf{L}_Y \mathbf{F}) + \beta \text{Tr}(\mathbf{W} \mathbf{L}_{\mathbf{F}^\top} \mathbf{W}^\top) \\ & + \gamma \text{Tr}(\mathbf{W}^\top \mathbf{L}_{\mathbf{X}^\top} \mathbf{W}), \quad \text{s.t.} \quad (\mathbf{W}, \mathbf{F}) \geq 0. \end{aligned} \quad (12)$$

159 Besides the multi-label feature selection methods mentioned, there are some specific multi-  
160 label learning methods that explore high-order label correlations. As a pioneering work, *Zhu et*  
161 *al.* [35] presented the GLOCAL multi-label correlation learning strategy, which at the same time  
162 retrieves missing labels, trains the classifier, and utilizes global and local label correlations by  
163 optimizing the label manifolds and learning a latent label representation. *Zhao et al.* [36] intro-  
164 duced a multi-label learning method called LSGL. Considering the suppositions of global label  
165 consistency and local label smoothness, this method learns a label correlation matrix. LSGL at-  
166 tempts to extract label correlation from the global and local perspectives in a self-representation  
167 model and a local data structure framework, respectively. *Kumar et al.* [37] presented the Trans-  
168 formation of Low-Rank Label Subspace for Multi-label Learning with Missing Labels (LRMML).  
169 The framework extracts global label correlation by a self-representation model and models the  
170 transformation of the low-rank label subspace, which is used to restore missing labels and train  
171 the classifier.

### 172 3. Proposed Method

173 This section introduces the Multi-Label Feature Selection with Global and Local Label Cor-  
174 relation (MLFS-GLOCAL), which uses both global and local label correlations for selecting the  
175 relevant and non-redundant features. The method's success is attributed to four primary factors:

176 (1) To provide a more insightful feature-label representation, it extracts the low-rank structure  
 177 from the label and feature matrices, which also provides an implicit label correlation (Section  
 178 3.1); (2) It imposes a penalty that preserves the smoothness of local mapping in the shared latent  
 179 space. (Section 3.2); (3) It can leverage information from all labels by taking into account both  
 180 global and local label correlations. (Section 3.3); (Section 3.4) It combines the aforementioned  
 181 into a single joint learning problem and employs an effective alternating minimization approach  
 182 for optimization. Figure 1 provides a schematic representation of the MLFS-GLOCAL model.

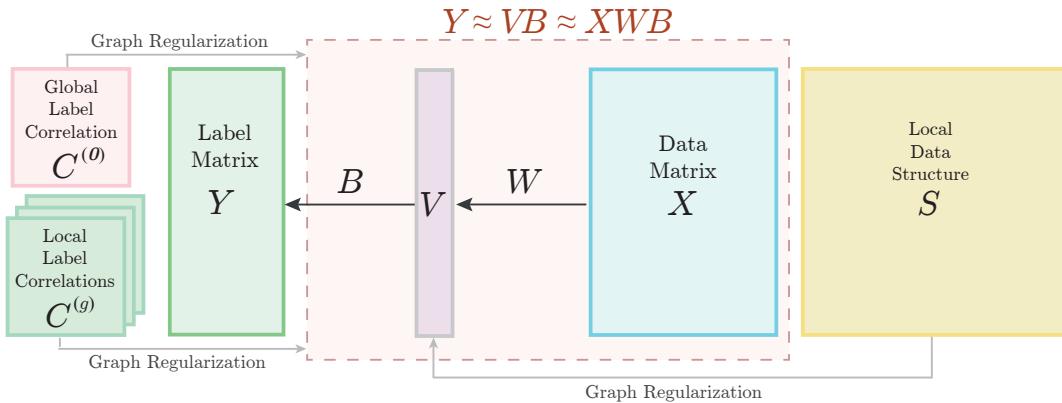


Figure 1: Schematic representation of the proposed multi-label feature selection model (MLFS-GLOCAL).

### 183 3.1. Shared latent space

184 We define the feature matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , and demonstrate the ground-truth label matrix  
 185  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_l\} \subseteq \{0, 1\} \in \mathbb{R}^{n \times l}$ , where  $Y_{ij} = 1$  if  $i$ -th instance would have the  $j$ -th label, and  
 186 otherwise  $Y_{ij} = 0$ . In multi-label problems, the labels are associated, hence it is common to  
 187 assume that the label matrix is low-rank. The label matrix  $\mathbf{Y}$  is sparse, binary-valued, and high  
 188 dimensional, therefore learning from this matrix is difficult. Let the rank of  $\mathbf{Y}$  be  $k < \min(n, l)$ .  
 189  $\mathbf{Y}$  can be factorized into two smaller matrices as follows:

$$\mathbf{Y} \simeq \mathbf{V}\mathbf{B}, \quad (13)$$

190 where  $\mathbf{V} \in \mathbb{R}^{n \times k}$  and  $\mathbf{B} \in \mathbb{R}^{k \times l}$ . Intuitively,  $\mathbf{V}$  reflects the latent labels, which are more compact  
 191 and semantically abstract than the original labels. while  $\mathbf{B}$  reflects how the original labels are  
 192 correlated to the latent labels. The significant information of the original label matrix  $\mathbf{Y}$  is  
 193 coded in the low-dimensional matrix  $\mathbf{V}$ , which also reduces the label matrix's unfavorable data.  
 194 Matrices  $\mathbf{V}$  and  $\mathbf{B}$  can be obtained by minimizing the label representation error  $\|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2$ . The

195 latent label representation, denoted as  $\mathbf{V}$ , is a low-dimensional, real-valued matrix that contains  
 196 dense information. It is comparatively easier to learn a continuous mapping from the feature  
 197 space to the latent label space than to the original label space [35]. Similar features plan to  
 198 have similar labels, which is one important assumption regarding correlations between features  
 199 and labels. As a result, the shared information between the feature space and the label space  
 200 should be consistent. We consider  $\mathbf{V}$  to be a shared factor matrix between the feature matrix  
 201 and the label matrix. We learn a matrix  $\mathbf{W} \in \mathbb{R}^{k \times l}$  to map instances to the latent space. By  
 202 decreasing the feature representation error  $\|\mathbf{V} - \mathbf{X}\mathbf{W}\|_F^2$ , the feature weight matrix  $\mathbf{W}$  would  
 203 be learnt. By integrating label representation and feature representation losses in a Nonnegative  
 204 Matrix Factorization framework [23], we learn shared latent space through the following objective  
 205 function:

$$\min_{\mathbf{V}, \mathbf{B}, \mathbf{W}} \|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2 + \|\mathbf{V} - \mathbf{X}\mathbf{W}\|_F^2, \quad \text{s.t. } \mathbf{V}, \mathbf{B}, \mathbf{W} \geq 0. \quad (14)$$

206 *3.2. Local structure preservation*

207 In many machine learning applications, the feature space can be high-dimensional and noisy,  
 208 making it challenging to extract meaningful information from the data. One popular approach is  
 209 to transform the original feature space into a lower-dimensional latent space, where the underlying  
 210 structure of the data can be more easily captured. However, simply projecting the data into  
 211 a lower-dimensional space can result in a loss of information and may not accurately capture  
 212 the complex relationships among the features. In order to solve this issue, we use the graph  
 213 regularization technique [38] to provide a suitable transformation matrix  $\mathbf{V}$  that ensures the  
 214 coherence between the initial feature space and the latent structure space. According to the  
 215 underlying principle of this regularization, the closer correlation between two instances in the  
 216 feature matrix  $\mathbf{X}$  denotes the closer correlation between the two corresponding latent feature  
 217 variables  $\mathbf{v}^{(i)}$  and  $\mathbf{v}^{(j)}$  in the latent structure. Specifically, we utilize a general graph regularization  
 218 term to encourage the locality preservation in the latent space, ensuring that the latent variables  
 219 accurately reflect the structure of the original features. The graph regularization term can be  
 220 expressed as:

$$\min_{\mathbf{V}} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{v}^{(i)} - \mathbf{v}^{(j)}\|^2 S_{i,j} = \text{Tr}(\mathbf{V}^\top \mathbf{D} \mathbf{V}) - \text{Tr}(\mathbf{V}^\top \mathbf{S} \mathbf{V}) = \text{Tr}(\mathbf{V}^\top \mathbf{L} \mathbf{V}), \quad (15)$$

221 where  $\mathbf{D}$  is a matrix diagonal,  $\mathbf{S}$  indicates a symmetric affinity matrix, and  $\mathbf{L} = \mathbf{D} - \mathbf{S}$  is the  
 222 graph Laplacian matrix. Integrating the above terms into our model, we achieve the following  
 223 function:

$$\min_{\mathbf{V}, \mathbf{B}, \mathbf{W}} \|\mathbf{Y} - \mathbf{VB}\|_F^2 + \|\mathbf{V} - \mathbf{XW}\|_F^2 + \lambda_1 \text{Tr}(\mathbf{V}^\top \mathbf{LV}), \quad \text{s.t. } \mathbf{V}, \mathbf{B}, \mathbf{W} \geq 0, \quad (16)$$

224 where  $\lambda_1$  is the local structure parameter.

225 *3.3. Global and Local label correlation*

226 To effectively utilize the information of multiple labels, label correlations must be incorporated.  
 227 In the proposed method, we regularize the proposed model using high-order label correlation. In  
 228 this direction, the coexistence of global and local label correlations should be noted. To consider  
 229 both of them, we introduce label manifold regularizers in this part. The concept behind the global  
 230 manifold regularizer is derived from the instance-level manifold regularizer, as outlined in (15).  
 231 In particular, if two labels are highly correlated, their corresponding classifier outputs should be  
 232 more similar, and conversely, less correlated labels should produce less similar classifier outputs.  
 233 In other words, label correlations lead to similar classifier outputs. This paper employs cosine  
 234 similarity to quantify the global label correlation, denoted as  $\mathbf{C} \in \mathbb{R}^{l \times l}$ , which is calculated as  
 235 follows:

$$C_{ij} = \frac{\mathbf{y}_i^\top \mathbf{y}_j}{\|\mathbf{y}_i\| \|\mathbf{y}_j\|}, \quad (17)$$

236 where  $\mathbf{y}_i$  and  $\mathbf{y}_j$  indicate the  $i$ -th and  $j$ -th label vectors for all instances.

237 In the basic model (16), the label predicted for sample  $\mathbf{x}$  is  $\mathbf{f}(\mathbf{x})$ , where  $\mathbf{f}(\mathbf{x}) = \mathbf{xWB}$ . Let  
 238  $\mathbf{f} = \{f_1, \dots, f_l\}$ , where  $f_j(\mathbf{x})$  is the  $j$ -th anticipated label for sample  $\mathbf{x}$ . As a result, predictions  
 239 for all  $n$  instances are recorded in the prediction matrix  $\mathbf{F} \in \mathbb{R}^{n \times l}$ , where  $\mathbf{F} = \mathbf{XWB}$  contains  
 240 predictions. If the  $i$ -th and  $j$ -th labels are more correlated, their corresponding predict labels  
 241  $\mathbf{f}_i$  and  $\mathbf{f}_j$  should be more similar to each other. The label manifold regularization is defined  
 242 similarly to the instance-level manifold regularization (15) as follows:

$$\min_{\mathbf{F}} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \|\mathbf{f}_i - \mathbf{f}_j\|^2 C_{i,j} = \text{Tr}(\mathbf{F} \mathbf{A} \mathbf{F}^\top) - \text{Tr}(\mathbf{F} \mathbf{C} \mathbf{F}^\top) = \text{Tr}(\mathbf{F} \mathbf{P} \mathbf{F}^\top), \quad (18)$$

243 where  $\mathbf{A}$  is a diagonal matrix as  $A_{i,i} = \sum_{j=1}^l C_{i,j}$  and  $\mathbf{P} = \mathbf{A} - \mathbf{C}$ . By minimizing (18),  $\|\mathbf{f}_i - \mathbf{f}_j\|_2^2$   
 244 will be small. The manifold regularizer in (18) written as  $\text{Tr}(\mathbf{F}\mathbf{P}\mathbf{F}^\top)$ , where  $\mathbf{P} = \mathbf{A} - \mathbf{C}$  is  
 245 the label Laplacian matrix of  $S$  and  $\mathbf{F} = \mathbf{X}\mathbf{W}\mathbf{B}$  is the classifier output matrix. Since label  
 246 correlations can differ across different local regions, we propose the local manifold regularizer.  
 247 We assume that the dataset  $\mathbf{X}$  is divided into  $g$  groups  $\{\mathbf{X}_1, \dots, \mathbf{X}_g\}$ , where  $\mathbf{X}_m \in \mathbb{R}^{n_m \times d}$   
 248 matrix has  $n_m$  instances. By using clustering or domain knowledge, such as networks and gene  
 249 pathways [39, 40] in bioinformatics applications, it is possible to acquire this partitioning. Assume  
 250  $\mathbf{C}_m \in \mathbb{R}^{l \times l}$  is the local label correlation matrix of a group  $m$  and that  $\mathbf{Y}_m$  is the label submatrix  
 251 in  $\mathbf{Y}$  corresponding to  $\mathbf{X}_m$ . The same as (17), we calculate the local label correlations as follows:

$$C_{i,j}^{(m)} = \frac{\mathbf{y}_i^{(m)\top} \mathbf{y}_j^{(m)}}{\|\mathbf{y}_i^{(m)}\| \|\mathbf{y}_j^{(m)}\|}, \quad m \in \{1, \dots, g\}. \quad (19)$$

252 Analogous to the global label correlations, we motivate the classifier outputs to be similar on  
 253 the correlated labels as follows:

$$\begin{aligned} & \min_{\mathbf{F}} \sum_{m=1}^g \frac{n_m}{n} \sum_{i=1}^l \sum_{j=1}^l \|\mathbf{f}_i^{(m)} - \mathbf{f}_j^{(m)}\|^2 C_{i,j}^{(m)} \\ &= \sum_{m=1}^g \frac{n_m}{n} [\text{Tr}(\mathbf{F}_m \mathbf{A}_m \mathbf{F}_m^\top) - \text{Tr}(\mathbf{F}_m \mathbf{C}_m \mathbf{F}_m^\top)] = \sum_{m=1}^g \frac{n_m}{n} \text{Tr}(\mathbf{F}_m \mathbf{P}_m \mathbf{F}_m^\top), \end{aligned} \quad (20)$$

254 where  $\mathbf{F}_m = \mathbf{X}_m \mathbf{W} \mathbf{B}$  is the classifier output matrix for group  $m$  and  $\mathbf{P}_m$  is the Laplacian matrix  
 255 of  $\mathbf{C}_m$ . To cover the cluster imbalance, we scale each local label correlation regularization by a  
 256 coefficient  $n_m/n$ . Problem (16) now has the following optimization problem after the addition of  
 257 global and local manifold regularizers (18) and (20):

$$\begin{aligned} & \min_{\mathbf{V}, \mathbf{B}, \mathbf{W}} \|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2 + \|\mathbf{V} - \mathbf{X}\mathbf{W}\|_F^2 + \lambda_1 \text{Tr}(\mathbf{V}^\top \mathbf{L}\mathbf{V}) \\ &+ \lambda_2 [\text{Tr}(\mathbf{F}\mathbf{P}\mathbf{F}^\top) + \sum_{m=1}^g \frac{n_m}{n} \text{Tr}(\mathbf{F}_m \mathbf{P}_m \mathbf{F}_m^\top)], \quad \text{s.t. } \mathbf{V}, \mathbf{B}, \mathbf{W} \geq 0, \end{aligned} \quad (21)$$

258 where  $\lambda_2$  indicates the impact of global and local label correlation on the objective function.  
 259 Finally, our function uses the  $L_{2,1}$ -norm, which has been shown to be beneficial for feature  
 260 selection. As a result, the objective function is set up as follows:

$$\begin{aligned} & \min_{\mathbf{V}, \mathbf{B}, \mathbf{W}} \|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2 + \|\mathbf{V} - \mathbf{X}\mathbf{W}\|_F^2 + \lambda_1 \text{Tr}(\mathbf{V}^\top \mathbf{L}\mathbf{V}) + \lambda_2 [\text{Tr}(\mathbf{F}\mathbf{P}\mathbf{F}^\top) \\ & \quad + \sum_{m=1}^g \frac{n_m}{n} \text{Tr}(\mathbf{F}_m \mathbf{P}_m \mathbf{F}_m^\top)] + \lambda_3 \|\mathbf{W}\|_{2,1}, \quad \text{s.t. } \mathbf{V}, \mathbf{B}, \mathbf{W} \geq 0, \end{aligned} \quad (22)$$

where the sparsity of  $\mathbf{W}$  on rows is ensured by using the  $L_{2,1}$ -norm. Adjusting the objective function's sparsity is done using the  $\lambda_3$  parameter.

### 3.4. Optimization

Function (22) contains the  $L_{2,1}$ -norm regularization term. Since it is not smooth, a direct solution is not possible. Additionally, when variables  $\mathbf{V}$ ,  $\mathbf{B}$ , and  $\mathbf{W}$  are considered together, it is non-convex. In other words, the Hessian matrix created from the partial derivatives of the objective function at the second degree is not a matrix that is positively semi-definite. The objective function must be optimized such that it is convex for each iteration by fixing any three of the variables and updating the remaining one. The phrase  $\|\mathbf{W}\|_{2,1}$  is further relaxed using  $\text{Tr}(\mathbf{W}^\top \mathbf{D}\mathbf{W})$ ,  $\mathbf{D}$  is a diagonal matrix in this case [30]. The algorithm's iterative updating feature is available. The  $\mathbf{D}$  element is  $D_{ii} = 1/(\|\mathbf{w}_i\| + \epsilon)$ , ( $\epsilon \leftarrow 0$ ), where  $\epsilon$  stops the non-differentiable problem's disturbance. As a result, the objective function (22) can be rewritten as follows:

$$\begin{aligned} & \min_{\mathbf{V}, \mathbf{B}, \mathbf{W}} \|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2 + \|\mathbf{V} - \mathbf{X}\mathbf{W}\|_F^2 + \lambda_1 \text{Tr}(\mathbf{V}^\top \mathbf{L}\mathbf{V}) + \lambda_2 [\text{Tr}(\mathbf{F}\mathbf{P}\mathbf{F}^\top) \\ & \quad + \sum_{m=1}^g \frac{n_m}{n} \text{Tr}(\mathbf{F}_m \mathbf{P}_m \mathbf{F}_m^\top)] + \lambda_3 (\mathbf{W}^\top \mathbf{D}\mathbf{W}), \quad \text{s.t. } \mathbf{V}, \mathbf{B}, \mathbf{W} \geq 0. \end{aligned} \quad (23)$$

We introduce Lagrangian multipliers to incorporate nonnegative constraint conditions into the function.  $\Phi$ ,  $\Psi$ , and  $\Omega$  to restrict  $\mathbf{V}$ ,  $\mathbf{B}$ , and  $\mathbf{W}$  respectively, where  $\Phi \in \mathbb{R}^{n \times k}$ ,  $\Psi \in \mathbb{R}^{k \times l}$ ,  $\Omega \in \mathbb{R}^{d \times k}$ . As a result, the function (23) is equivalent to the following function:

$$\begin{aligned} & \min_{\mathbf{V}, \mathbf{B}, \mathbf{W}} \|\mathbf{Y} - \mathbf{V}\mathbf{B}\|_F^2 + \|\mathbf{V} - \mathbf{X}\mathbf{W}\|_F^2 + \lambda_1 \text{Tr}(\mathbf{V}^\top \mathbf{L}\mathbf{V}) + \lambda_2 [\text{Tr}(\mathbf{F}\mathbf{P}\mathbf{F}^\top) \\ & \quad + \sum_{m=1}^g \frac{n_m}{n} \text{Tr}(\mathbf{F}_m \mathbf{P}_m \mathbf{F}_m^\top)] + \lambda_3 \text{Tr}(\mathbf{W}^\top \mathbf{D}\mathbf{W}) - \text{Tr}(\Phi \mathbf{V}^\top) - \text{Tr}(\Psi \mathbf{B}^\top) - \text{Tr}(\Omega \mathbf{W}^\top). \end{aligned} \quad (24)$$

Given an arbitrary matrix  $\mathbf{A}$ ,  $\|\mathbf{A}\|_F^2 = \text{Tr}(\mathbf{A}^\top \mathbf{A})$ . Therefore, the function (24) becomes:

$$\begin{aligned}
C = & \text{Tr}[(\mathbf{Y} - \mathbf{V}\mathbf{B})^\top(\mathbf{Y} - \mathbf{V}\mathbf{B})] + \text{Tr}[(\mathbf{V} - \mathbf{X}\mathbf{W})^\top(\mathbf{V} - \mathbf{X}\mathbf{W})] \\
& + \lambda_1 \text{Tr}(\mathbf{V}^\top \mathbf{L}\mathbf{V}) + \lambda_2 [\text{Tr}(\mathbf{F}\mathbf{P}\mathbf{F}^\top) + \sum_{m=1}^g \frac{n_m}{n} \text{Tr}(\mathbf{F}_m \mathbf{P}_m \mathbf{F}_m^\top)] \\
& + 2\lambda_3 \text{Tr}(\mathbf{W}^\top \mathbf{D}\mathbf{W}) - \text{Tr}(\boldsymbol{\Phi}\mathbf{V}^\top) - \text{Tr}(\boldsymbol{\Psi}\mathbf{B}^\top) - \text{Tr}(\boldsymbol{\Omega}\mathbf{W}^\top),
\end{aligned} \tag{25}$$

277 where  $\mathbf{F} = \mathbf{X}\mathbf{W}\mathbf{B}$  and  $\mathbf{F}_m = \mathbf{X}_m\mathbf{W}\mathbf{B}$ ,  $\forall m \in \{1, 2, \dots, g\}$ .

278 The partial derivatives of function (25) with respect to the variables  $\mathbf{V}$ ,  $\mathbf{B}$ , and  $\mathbf{W}$  are:

$$\begin{aligned}
\frac{\partial C}{\partial \mathbf{W}} = & -\mathbf{X}^\top \mathbf{V} + \mathbf{X}^\top \mathbf{X}\mathbf{W} + \lambda_2 [\mathbf{X}^\top \mathbf{X}\mathbf{W}\mathbf{B}\mathbf{P}\mathbf{B}^\top \\
& + \sum_{m=1}^g \frac{n_m}{n} \mathbf{X}_m^\top \mathbf{X}_m \mathbf{W}\mathbf{B}\mathbf{P}_m \mathbf{B}^\top] + 2\lambda_3 \mathbf{D}\mathbf{W} - \boldsymbol{\Omega},
\end{aligned} \tag{26}$$

$$\frac{\partial C}{\partial \mathbf{B}} = -\mathbf{V}^\top \mathbf{Y} + \mathbf{V}^\top \mathbf{V}\mathbf{B} + \lambda_1 [\mathbf{W}^\top \mathbf{X}^\top \mathbf{X}\mathbf{W}\mathbf{B}\mathbf{P} + \sum_{m=1}^g \frac{n_m}{n} \mathbf{W}^\top \mathbf{X}_m^\top \mathbf{X}_m \mathbf{W}\mathbf{B}\mathbf{P}_m] - \boldsymbol{\Phi}, \tag{27}$$

279 and

$$\frac{\partial C}{\partial \mathbf{V}} = \mathbf{V} - 2\mathbf{X}\mathbf{W} - 2\mathbf{Y}\mathbf{B}^\top + \mathbf{V}\mathbf{B}\mathbf{B}^\top + \lambda_2 \mathbf{L}\mathbf{V} - \boldsymbol{\Psi}. \tag{28}$$

280 By setting the partial derivatives (26), (27), and (28) to zero, we can find critical points of the  
281 function. In accordance with the Karush-Kuhn-Tucker condition, we set  $\mathbf{W} \odot \boldsymbol{\Omega} = \mathbf{0}$ ,  $\mathbf{B} \odot \boldsymbol{\Phi} = \mathbf{0}$ ,  
282 and  $\mathbf{V} \odot \boldsymbol{\Psi} = \mathbf{0}$ , that are fixed point equations that the solution must satisfy at convergence. By  
283 solving these equations, we derive the following updating rules for the  $\mathbf{W}$ ,  $\mathbf{B}$ , and  $\mathbf{V}$ :

$$\mathbf{W} \leftarrow \mathbf{W} \odot \frac{\mathbf{X}^\top \mathbf{V} + \lambda_1 [(\mathbf{X}^\top \mathbf{X}\mathbf{W}\mathbf{B}\mathbf{C}\mathbf{B}^\top) + \sum_{m=1}^g \frac{n_m}{n} (\mathbf{X}_m^\top \mathbf{X}_m \mathbf{W}\mathbf{B}\mathbf{C}_m \mathbf{B}^\top)]}{\mathbf{X}^\top \mathbf{X}\mathbf{W} + \lambda_1 (\mathbf{X}^\top \mathbf{X}\mathbf{W}\mathbf{B}\mathbf{A}\mathbf{B}^\top) + \sum_{m=1}^g \frac{n_m}{n} (\mathbf{X}_m^\top \mathbf{X}_m \mathbf{W}\mathbf{B}\mathbf{A}_m \mathbf{B}^\top) + \lambda_3 (\mathbf{D}\mathbf{W})}, \tag{29}$$

---

**Algorithm 1** Multi-label Feature Selection with Global and Local label correlation (MLFS-GLOCAL)

---

**Input:** Feature matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and Label matrix  $\mathbf{Y} \in \mathbb{R}^{n \times c}$ , Regularization parameters  $\lambda_1, \lambda_2$ , and  $\lambda_3$ , and latent factor  $k$ ;  
**Output:** Feature score  $s_i = \|\mathbf{w}_i\|, \forall i \in \{1, 2, \dots, d\}$ .

- 1: **Initialize**  $\mathbf{V}, \mathbf{W}, \mathbf{B}$  randomly;  $t = 0$ ;
- 2: **while**  $t < \text{MaxIteration}$  **do**
- 3:   Update  $D_{ii} \leftarrow \frac{1}{\|\mathbf{w}_i\| + \epsilon}$ ;
- 4:   Update  $\mathbf{W}$  by (29);
- 5:   Update  $\mathbf{B}$  by (30);
- 6:   Update  $\mathbf{V}$  by (31);
- 7:    $t = t + 1$ ;
- 8: **end while**
- 9: **Return**  $\mathbf{W}$ ;
- 10: Evaluate the feature score by  $s_i = \|\mathbf{w}_i\|$ .

---

284

$$\mathbf{B} \leftarrow \mathbf{B} \odot \frac{\mathbf{V}^\top \mathbf{Y} + \lambda_1 [(\mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W} \mathbf{B} \mathbf{C}) + \sum_{m=1}^g \frac{n_m}{n} (\mathbf{W}^\top \mathbf{X}_m^\top \mathbf{X}_m \mathbf{W} \mathbf{B} \mathbf{C}_m)]}{\mathbf{V}^\top \mathbf{V} \mathbf{B} + \lambda_1 [(\mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W} \mathbf{B} \mathbf{A}) + \sum_{m=1}^g \frac{n_m}{n} (\mathbf{W}^\top \mathbf{X}_m^\top \mathbf{X}_m \mathbf{W} \mathbf{B} \mathbf{A}_m)]}, \quad (30)$$

285 and

$$\mathbf{V} \leftarrow \mathbf{V} \odot \frac{\mathbf{X} \mathbf{W} + \mathbf{Y} \mathbf{B}^\top + \lambda_2 \mathbf{S} \mathbf{V}}{\mathbf{V} + \mathbf{V} \mathbf{B} \mathbf{B}^\top + \lambda_2 \mathbf{D} \mathbf{V}}. \quad (31)$$

286 MLFS-GLOCAL, which ranks all the features using  $\|\mathbf{w}_i\|_2, (i = 1, \dots, d)$  in descending order,  
287 allowing us to obtain the top-k features. Algorithm 1 outlines the detailed solution for the  
288 MLFS-GLOCAL model.

289 **4. Experimental study**

290 This section conducts a comprehensive evaluation to assess the MLFS-GLOCAL model on 12  
291 real-world multi-label benchmark datasets, using six diverse evaluation metrics. The proposed  
292 model is compared to nine well-known and state-of-the-art feature selection methods. All of our  
293 tests are run on an Intel Core (TM) i7-9700K with a 3.6 GHz processor and 32 GB RAM.

294 *4.1. Datasets*

295 In our studies, we employed 12 datasets from Mulan Library's multi-label text and image  
296 classification. The multi-label Yahoo datasets refer to a collection of datasets that are used for

Table 1: The detailed information of the real-world datasets

Dataset	#Instance	#Feature	#Label
Arts	5000	462	26
Business	5000	438	30
Computers	5000	681	33
corel5k	5000	499	374
Education	5000	550	33
Entertainment	5000	640	21
Health	5000	612	32
Recreation	5000	606	22
Reference	5000	793	33
Science	5000	743	40
Social	5000	1047	39
Society	5000	636	27

297 multi-label classification tasks. These datasets were originally released by Yahoo Labs and consist  
 298 of a large number of text documents that have been annotated with multiple labels. Each dataset  
 299 includes a training set and a test set that each comprises 2000 and 3000 documents, respectively  
 300 [41]. The Corel5k multi-label dataset is a collection of images used for multi-label classification  
 301 tasks. It consists of 5,000 images from 374 different categories, with each image having multiple  
 302 labels assigned to it. Both Yahoo and Corel5k datasets have been widely used in research for  
 303 developing and evaluating multi-label classification algorithms. Table 1 describes the specifics of  
 304 each benchmark dataset.

305 *4.2. Evaluation Metrics*

306 To examine the performance of all competition methods, the Multi-Label kNN (ML-kNN)  
 307 algorithm [42] is defined as a benchmark classifier. The ML-kNN is frequently used for classifi-  
 308 cation in multi-label feature selection approaches [43, 44, 16, 30] because of its interpretability  
 309 and simplicity. We set  $k = 10$  for the number of nearest neighbors. Furthermore, we use six  
 310 commonly used assessment criteria, including Micro-F1, Macro-F1, Average Precision, Ranking  
 311 Loss, Hamming Loss, and Coverage Error. These assessment metrics' definitions are as follows:

- 312 • **Macro-F1** and **Micro-F1** both are on the measure's foundation of F1-measure. The  
 313 evaluation metric directly uses F-measure averaging to rate the precision of the predictions  
 314 made by the classifier label set.

$$Micro - F1 = \frac{\sum_{i=1}^l 2TP_i}{\sum_{i=1}^l (2TP_i + FP_i + FN_i)}, \quad (32)$$

315 and

$$Macro - F1 = \sum_{i=1}^l \frac{2TP_i}{2TP_i + FP_i + FN_i}, \quad (33)$$

316 where T and F are True and False, respectively; P and N are Positive and Negative, re-  
317 spectively; so TP, TN, FP, and FN are the number of combinations of T, F, P, and N,  
318 respectively.

- 319 • **Average precision** determines the percentage of labels that are more relevant than specific  
320 labels.

$$AP(D) = \frac{1}{n} \sum_{i=1}^n \frac{1}{1_m^\top y_i} \sum_{l:y_i^l} \frac{prec_i(l)}{rank_i(l)}, \quad (34)$$

321 where  $prec_i(l) = \sum_{l:y_i^l=1} \delta(rank_i(l) \geq rank_i(l'))$ , and  $AP(D) \in [0, 1]$ .

- 322 • **Ranking loss** is the proportion of label pairings that are in reverse order, or when unrelated  
323 labels are more important than related labels.

$$RL(D) = \frac{1}{n} \sum_{i=1}^n \frac{1}{1_m^\top y_i 1_m^\top \bar{y}_i} \sum_{l:y_i^l=1} \sum_{l':y_i^{l'}=0} (\delta(rank_i(l) \geq rank_i(l'))), \quad (35)$$

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

- 325 • **Hamming loss** determines the percentage of labels that are incorrectly labeled, meaning  
326 that either a label that belongs to the instance or one that doesn't is predicted.

$$HL(D) = \frac{1}{n} \sum_{i=1}^n \frac{1}{m} \|h(x_i) \Delta y_i\|_1. \quad (36)$$

327        The  $\Delta$  is used to denote the symmetric difference between the two sets, which is the set of  
328        values that appear exclusively in either one of the two sets.

- 329        • **Coverage Error** indicates the number of steps required to cover all of the positive labels  
330        associated with the cases by moving down the anticipated label ranking.

$$CV(D) = \frac{1}{n} \sum_{i=1}^n \arg \max_{l:y_i^l=1} rank_i(l) - 1. \quad (37)$$

331        A smaller value indicates better classification performance in terms of Ranking loss, Hamming  
332        loss, and Coverage Error when the ideal value is 0. While a higher value reflects better classifi-  
333        cation performance in terms of Micro-F1, Macro-F1, and Average accuracy, the ideal value is  
334        1.

335        *4.3. Compared Methods*

336        We evaluate our method by comparing it to the latest multi-label feature selection methods.  
337        Each method is briefly described in the following

- 338        • **MDMR** [24] is a MLFS method based on information theory that chooses features by  
339        simultaneously maximizing the dependency and minimizing redundancy.
- 340        • **SCLS** [25] is another MLFS algorithm that uses information theory and evaluates condi-  
341        tional relevance using a scalable relevance assessment criterion.
- 342        • **LRFS** [26] is based on label redundancy. The conditional mutual information is used to  
343        build a new feature relevance term that evaluates feature information.
- 344        • **MIFS** [16] is a well-known MLFS method that makes use of latent semantics of the multi-  
345        labels to choose the most important features.
- 346        • **CMFS** [29] is a method for studying structured information that relies on feature-side and  
347        label-side correlations.
- 348        • **SCMFS** [30] employs coupled nonnegative matrix factorization to create the shared com-  
349        mon model.
- 350        • **SSFS** [31] presents a Latent Structure Shared (LSS) term that shares and maintains both  
351        the latent feature and label structures.

Table 2: Micro-F1 results on real-world datasets. The best result is highlighted in **bold** style, while underline style indicates the second-best.

Datasets	SSFS	SCMFS	MIFS	CMFS	MRDM	NMDG	SCLS	LRFS	MDMR	MLFS-GLOCAL
Arts	0.2263	<b>0.2674</b>	0.1655	0.1753	0.2496	0.1714	0.1258	0.0601	0.1418	<b>0.3230</b>
Business	0.6816	<u>0.6927</u>	0.6846	0.6786	0.6892	0.6728	0.6760	0.6698	0.6704	<b>0.6966</b>
Computers	0.4264	0.4254	0.0388	0.4105	0.4217	0.4064	<u>0.4310</u>	0.4088	0.4083	<b>0.4511</b>
Corel5k	0.0276	<u>0.0397</u>	0.0388	0.0338	0.0392	0.0361	0.0377	0.0141	0.0208	<b>0.0495</b>
Education	0.3061	<u>0.3131</u>	0.2095	0.2871	0.2306	0.1191	0.1383	0.0910	0.1591	<b>0.3594</b>
Entertainment	0.3314	<b>0.3640</b>	0.2796	0.3041	0.3597	0.2350	0.2665	0.1219	0.2828	<b>0.4377</b>
Health	0.4887	<u>0.5158</u>	0.4647	0.4783	0.4904	0.4649	0.4385	0.3671	0.4193	<b>0.5744</b>
Recreation	0.2158	0.2689	0.2252	0.1694	<u>0.2760</u>	0.2104	0.1432	0.0693	0.1766	<b>0.3371</b>
Reference	0.4371	0.4563	0.4291	0.4259	<u>0.4570</u>	0.3864	0.4083	0.3752	0.3722	<b>0.4838</b>
Science	0.2185	<u>0.2615</u>	0.1725	0.1975	0.2153	0.1530	0.1422	0.0871	0.1406	<b>0.2978</b>
Social	0.5306	<u>0.5485</u>	0.4995	0.5132	0.5398	0.4035	0.4520	0.3069	0.4430	<b>0.5883</b>
Society	0.3534	0.3536	0.3429	0.3234	<u>0.3625</u>	0.3071	0.2904	0.2862	0.2823	<b>0.3711</b>

- **MRDM** [32] utilizes HSIC as a criterion To increase the reliance between the Manifold embedding and the class labels.
- **NMDG** [15] uses the dynamic graph laplacian matrix constructed by pseudo-label in the feature selection process.

#### 356 4.4. Experimental Results

357 In this experiment, we selected the top 20% of features to determine the average performances  
 358 for each technique. Tables 2–7 present the results of these experiments on the six different  
 359 evaluation criteria. The best results for each dataset are indicated by bold fonts, where the  
 360 higher the values, the better the classification performance. To provide a more robust evaluation  
 361 of performance, each method is run 10 times and the average results are reported for all datasets.  
 362 The tables demonstrate that, across most datasets, the proposed model produces the best results.  
 363 In addition, the best Micro-F1 score was obtained by MLFS-GLOCAL, which had a significant  
 364 lead over the second-best method on the Arts, Corel5k, Entertainment, Health, Recreation, and  
 365 Social datasets. These tables show that the proposed model is first in 67 of the 72 comparison  
 366 cases and comes in second in the remaining ones. These results confirm that our method can  
 367 be applied to a broad spectrum of datasets, as opposed to other techniques. On average, we  
 368 observed significant improvements in the values of these metrics, with 0.0367 for Micro-F1, 0.0157  
 369 for Macro-F1, 0.0443 for Average Precision, 0.222 for Coverage Error, 0.0014 for Hamming Loss,  
 370 and 0.0011 for Ranking Loss.

Table 3: Macro-F1 results on real-world datasets. The best result is highlighted in **bold** style, while underline style indicates the second-best.

Datasets	SSFS	SCMFS	MIFS	CMFS	MRDM	NMDG	SCLS	LRFS	MDMR	MLFS-GLOCAL
Arts	0.1052	0.1299	0.0765	0.0790	<u>0.1301</u>	0.0702	0.0534	0.0203	0.0638	<b>0.1597</b>
Business	0.0789	<u>0.1066</u>	0.0984	0.0762	0.0951	0.0633	0.0679	0.0427	0.0527	<b>0.1224</b>
Computers	0.1184	0.1205	0.0747	0.1128	<u>0.1401</u>	0.0514	0.0981	0.0503	0.0830	<b>0.1659</b>
Corel5k	0.0022	0.0032	0.0027	0.0024	0.0023	0.0019	<u>0.0035</u>	0.0018	0.0028	<b>0.0044</b>
Education	0.0936	<u>0.1151</u>	0.0657	0.0970	0.0880	0.0365	0.0488	0.0285	0.0451	<b>0.1228</b>
Entertainment	0.1736	<u>0.1918</u>	0.1291	0.1488	0.1907	0.1170	0.1191	0.0399	0.1138	<b>0.2155</b>
Health	0.1908	<u>0.1991</u>	0.1599	0.1844	0.1967	0.1405	0.1321	0.0673	0.1030	<b>0.2283</b>
Recreation	0.1317	0.1576	0.1347	0.1087	<u>0.1691</u>	0.1188	0.0802	0.0490	0.0905	<b>0.1756</b>
Reference	0.0940	<u>0.1136</u>	0.0892	0.0935	0.1093	0.0663	0.0690	0.0322	0.0613	<b>0.1203</b>
Science	0.0864	<u>0.0947</u>	0.0660	0.0784	0.0915	0.0538	0.0531	0.0310	0.0495	<b>0.1090</b>
Social	0.1349	<u>0.1362</u>	0.0974	0.1195	0.1177	0.0558	0.0937	0.0250	0.0611	<b>0.1567</b>
<b>Society</b>	0.1024	<u>0.1114</u>	0.0718	0.0693	0.1050	0.0750	0.0415	0.0362	0.0432	<b>0.1195</b>

Table 4: Average Precision results on real-world datasets. The best result is highlighted in **bold** style, while underline style indicates the second-best.

Datasets	SSFS	SCMFS	MIFS	CMFS	MRDM	NMDG	SCLS	LRFS	MDMR	MLFS-GLOCAL
Arts	<u>0.0698</u>	0.0677	0.0670	0.0690	0.0694	0.0551	0.0658	0.0659	0.0651	<b>0.0707</b>
Business	0.0627	<u>0.0630</u>	0.0605	0.0623	0.0626	0.0584	0.0580	0.0553	0.0558	<b>0.0689</b>
Computers	<u>0.0795</u>	0.0688	0.0589	0.0717	0.0722	0.0551	0.0551	0.0524	0.0598	<b>0.0807</b>
Corel5k	0.0104	0.0103	0.0111	0.0104	0.0104	0.0107	<u>0.0115</u>	0.0100	0.0104	<b>0.0118</b>
Education	0.0581	0.0584	0.0561	0.0555	<b>0.0640</b>	0.0493	0.0534	0.0482	0.0493	<u>0.0624</u>
Entertainment	0.1056	0.1080	0.1039	0.1064	0.1085	<b>0.1168</b>	0.1124	0.0693	0.1041	<u>0.1150</u>
Health	0.0855	<u>0.0889</u>	0.0779	0.0860	0.0877	0.0724	0.0669	0.0715	0.0763	<b>0.1105</b>
Recreation	0.0788	0.0837	0.0837	0.0739	0.0898	<u>0.0901</u>	0.0690	0.0678	0.0708	<b>0.1028</b>
Reference	0.0486	<u>0.0496</u>	0.0414	0.0485	0.0486	0.0457	0.0429	0.0388	0.0433	<b>0.0501</b>
Science	0.0488	0.0495	0.0456	0.0515	<u>0.0530</u>	0.0399	0.0408	0.0403	0.0427	<b>0.0535</b>
Social	0.0608	<u>0.0634</u>	0.0553	0.0577	0.0613	0.0465	0.0494	0.0354	0.0497	<b>0.0679</b>
Society	0.0680	<u>0.0703</u>	0.1195	0.0656	0.0686	0.0667	0.0653	0.0654	0.0653	<b>0.0709</b>

Table 5: Ranking Loss results on real-world datasets. The best result is highlighted in **bold** style, while underline style indicates the second-best.

Datasets	SSFS	SCMFS	MIFS	CMFS	MRDM	NMDG	SCLS	LRFS	MDMR	MLFS-GLOCAL
Arts	0.2066	0.2077	0.2095	0.2056	0.2012	0.2190	<u>0.1998</u>	0.2117	0.2124	<b>0.1984</b>
Business	0.0528	<b>0.0488</b>	0.0522	0.0526	0.0494	0.0547	0.0580	0.0583	0.0592	<u>0.0490</u>
Computers	0.1207	<u>0.1172</u>	0.1218	0.1208	0.1194	0.1254	0.1250	0.1322	0.1226	<b>0.1148</b>
Corel5k	0.2085	<u>0.2078</u>	0.2171	0.2080	0.2088	0.2162	0.2172	0.2146	0.2161	<b>0.2063</b>
Education	0.1211	0.1216	0.1274	<u>0.1201</u>	0.1231	0.1356	0.1353	0.1313	0.1374	<b>0.1168</b>
Entertainment	0.1664	<u>0.1609</u>	0.1695	0.1685	0.1641	0.1697	0.1851	0.1874	0.1726	<b>0.1577</b>
Health	0.0804	0.0807	0.0836	<u>0.0803</u>	0.0790	0.0895	0.0979	0.1084	0.1025	<b>0.0789</b>
Recreation	0.2409	0.2325	0.2451	0.2465	<u>0.2297</u>	0.2521	0.2693	0.2605	0.2575	<b>0.2288</b>
Reference	0.1009	0.0986	0.1043	0.1006	<u>0.0982</u>	0.1074	0.1116	0.1203	0.1176	<b>0.0975</b>
Science	0.1635	<u>0.1583</u>	0.1672	0.1639	0.1586	0.1739	0.2030	0.2015	0.2016	<b>0.1559</b>
Social	0.0757	<u>0.0732</u>	0.0825	0.0774	0.0741	0.0867	0.0863	0.1053	0.0957	<b>0.0728</b>
Society	0.1800	0.1813	0.1852	0.1865	0.1811	0.1902	<b>0.1754</b>	0.2058	0.2226	<u>0.1758</u>

Table 6: Hamming Loss results on real-world datasets. The best result is highlighted in **bold** style, while underline style indicates the second-best.

Datasets	SSFS	SCMFS	MIFS	CMFS	MRDM	NMDG	SCLS	LRFS	MDMR	MLFS-GLOCAL
Arts	0.0615	<u>0.0588</u>	0.0616	0.0625	0.0597	0.0622	0.0631	0.0633	0.0627	<b>0.0568</b>
Business	0.0288	<u>0.0278</u>	0.0281	0.0290	0.0286	0.0287	0.0286	0.0288	0.0288	<b>0.0276</b>
Computers	<u>0.0393</u>	0.0394	0.0404	0.0402	0.3979	0.0411	0.0396	0.0422	0.0412	<b>0.0388</b>
Corel5k	0.095170	0.009506	<u>0.009495</u>	0.009506	0.009504	0.009506	0.009516	0.009512	0.009535	<b>0.009495</b>
Education	0.0414	<u>0.0410</u>	0.0437	0.0418	0.0436	0.0450	0.0443	0.0444	0.0442	<b>0.0392</b>
Entertainment	0.0615	<u>0.0594</u>	0.0641	0.0630	0.0613	0.0658	0.0634	0.0676	0.0629	<b>0.0555</b>
Health	0.0427	<u>0.0407</u>	0.0436	0.0433	0.0423	0.0430	0.0461	0.0506	0.0466	<b>0.0382</b>
Recreation	0.0617	0.0592	0.0620	0.0633	<u>0.0590</u>	0.0610	0.0644	0.0647	0.0628	<b>0.0571</b>
Reference	0.0296	<u>0.0287</u>	0.0297	0.0297	0.0291	0.0294	0.0323	0.0347	0.0322	<b>0.0275</b>
Science	0.0349	<u>0.0342</u>	0.0358	0.0351	0.0346	0.0352	0.0357	0.0358	0.0354	<b>0.0336</b>
Social	0.0243	<u>0.0236</u>	0.0253	0.0247	0.0242	0.0281	0.0275	0.0313	0.0262	<b>0.0221</b>
Society	0.0563	<u>0.0553</u>	0.0571	0.0574	0.0554	0.0576	0.0590	0.0590	0.0583	<b>0.0545</b>

Table 7: Coverage Error results on real-world datasets. The best result is highlighted in **bold** style, while underline style indicates the second-best.

Datats	SSFS	SCMFS	MIFS	CMFS	MRDM	NMDG	SCLS	LRFS	MDMR	MLFS-GLOCAL
Arts	7.820	7.852	7.930	7.820	7.710	8.217	<u>7.662</u>	8.045	8.059	<b>7.627</b>
Business	3.731	<b>3.584</b>	3.686	3.740	3.609	3.808	3.933	4.009	4.060	<u>3.588</u>
Computers	6.440	6.340	6.458	6.495	<u>6.324</u>	6.619	6.723	6.849	6.673	<b>6.237</b>
Corel5k	164.89	<u>164.36</u>	170.95	164.83	164.81	172.02	171.95	171.96	172.00	<b>163.63</b>
Education	5.870	5.882	6.093	<u>5.820</u>	6.019	6.453	6.505	6.412	6.621	<b>5.735</b>
Entertainment	5.190	5.075	5.269	5.254	5.149	<u>5.027</u>	5.610	5.694	5.359	<b>5.010</b>
Health	4.985	4.955	5.096	4.960	<u>4.932</u>	5.378	5.663	6.032	5.920	<b>4.903</b>
Recreation	7.125	6.944	7.231	7.259	<u>6.882</u>	7.406	7.824	7.678	7.511	<b>6.835</b>
Reference	4.820	4.742	4.946	4.811	<u>4.724</u>	5.039	5.194	5.470	5.411	<b>4.698</b>
Science	8.867	8.630	9.040	8.885	<u>8.613</u>	9.303	10.698	10.637	10.624	<b>8.525</b>
Social	4.820	<u>4.739</u>	5.131	4.911	4.761	5.341	5.339	6.116	5.841	<b>4.719</b>
Society	7.621	7.640	7.754	7.775	7.613	7.899	<u>7.577</u>	8.536	7.580	<b>7.470</b>

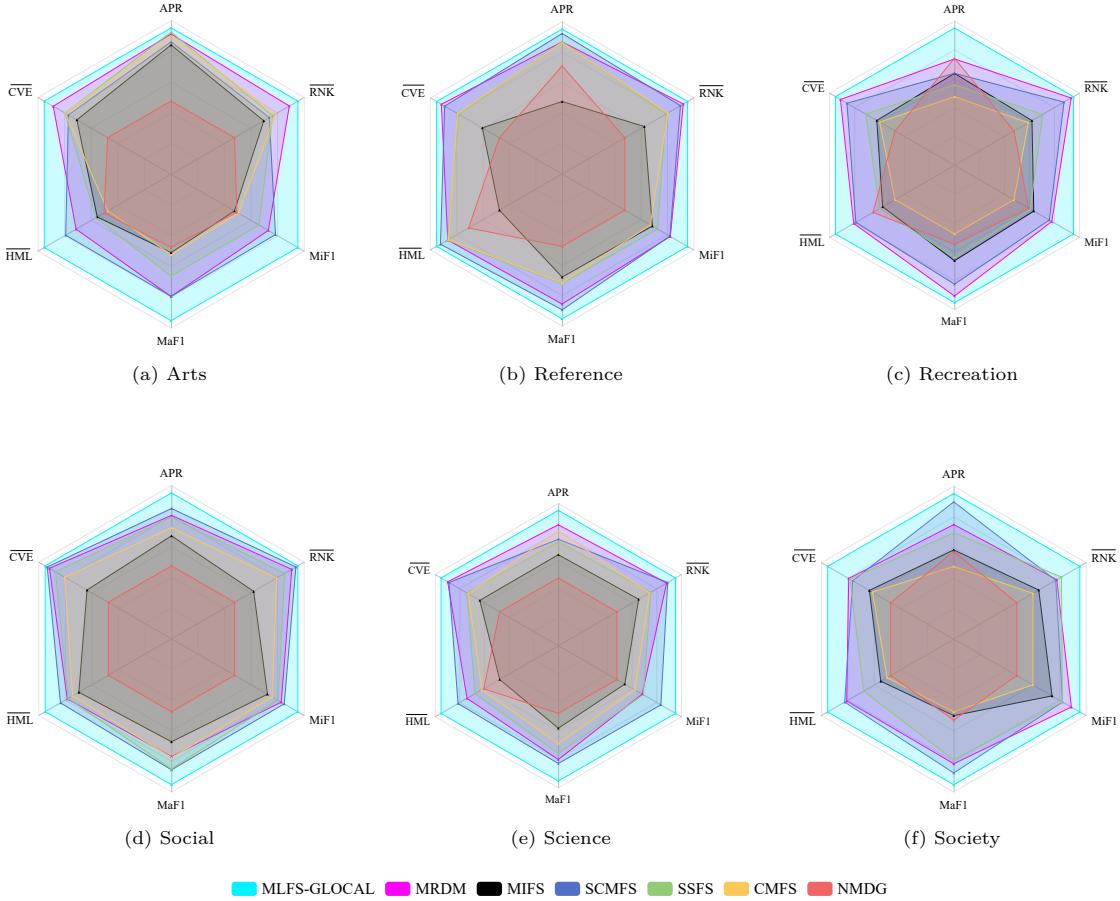


Figure 2: Universality analysis of the proposed method on the evaluation criteria for six different multi-label datasets is shown by the radar chart, where the APR=Average Precision, RANK=Ranking loss, MiF1=Micro-F1, MaF1=Macro-F1, HML=Hamming loss, and CVE=Coverage error.

In addition, we use radar diagrams to show the universality of our method on six different evaluation metrics in comparison with other methods. These metrics measure various aspects of the quality and effectiveness of methods [45]. It is worth mentioning that the Coverage error, Ranking loss, and hamming loss metrics are utilized inversely in order to maintain consistency with the other measures, so that a larger scale implies better performance for all measures. Also, to make the comparisons fair and clear, we normalize the data in Figure 2 so that all the values are between 0.5 and 1. The more area a method covers in the radar charts, the better it is in all the evaluation metrics. Figure 2 shows that our method has a larger area than the other methods, which means that it is more universal and superior in performance and evaluation criteria. This demonstrates that our method can handle different types of problems and situations better than the existing methods.

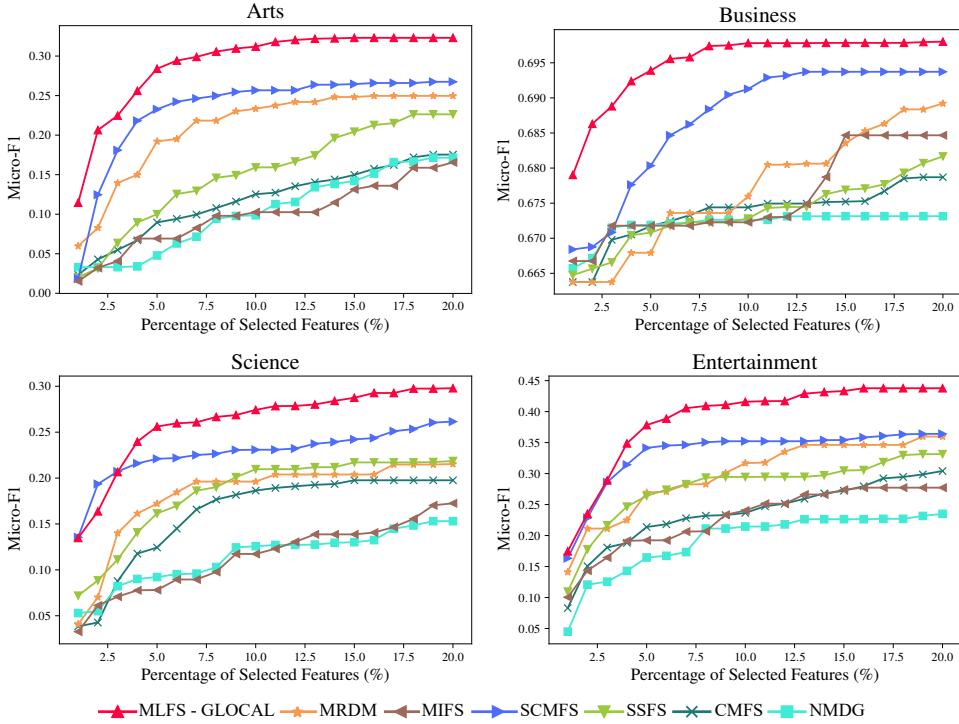


Figure 3: Robustness analysis on four benchmark datasets in terms of Micro-F1.

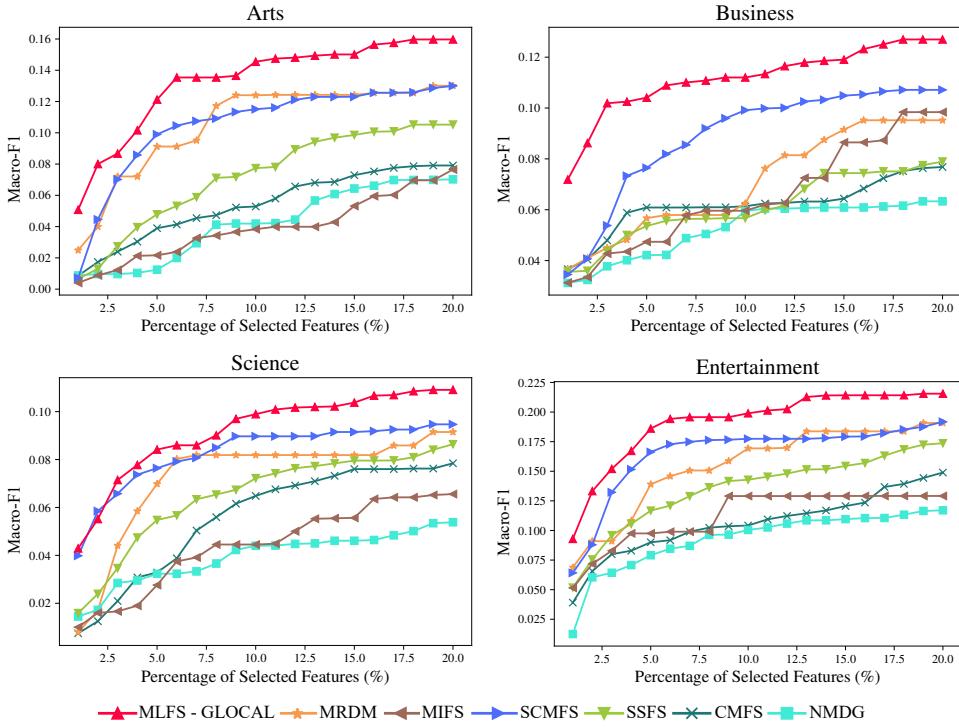


Figure 4: Robustness analysis on four benchmark datasets in terms of Macro-F1.

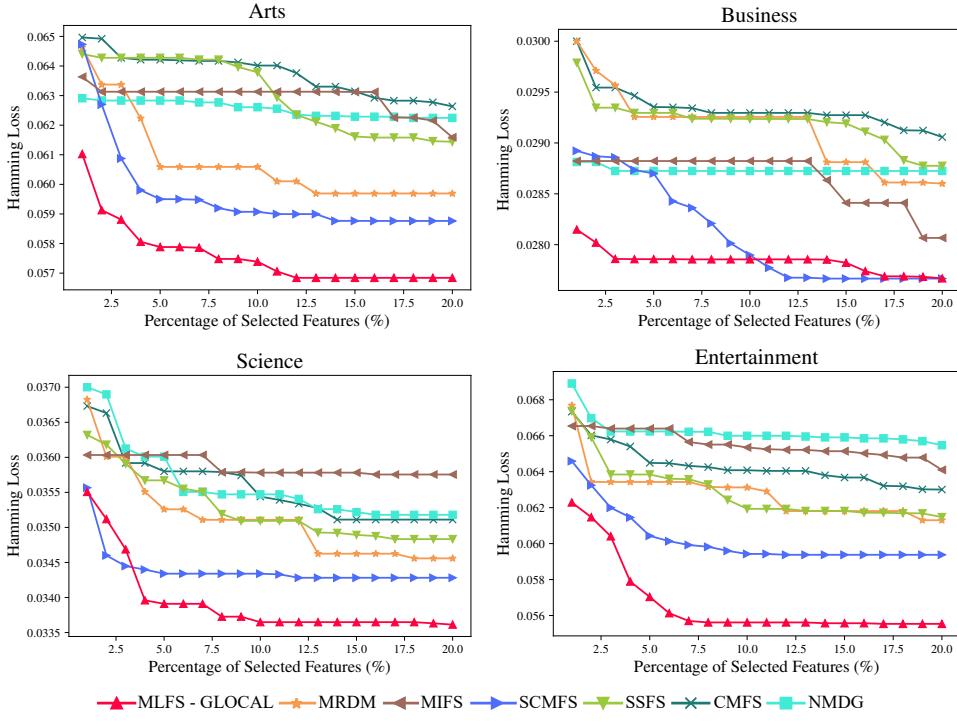


Figure 5: Robustness analysis on four benchmark datasets in terms of Hamming Loss.

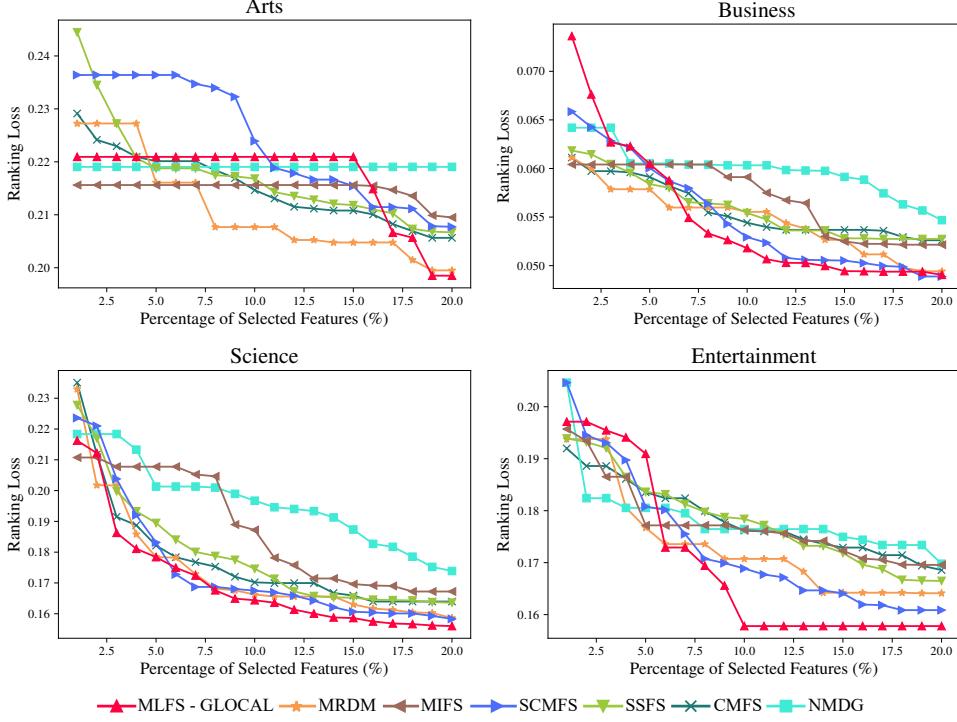


Figure 6: Robustness analysis on four benchmark datasets in terms of Ranking Loss.

382 The performance of MLFS-GLOCAL and other comparative approaches is illustrated visually  
383 using four datasets: Arts, Business, Corel5k, and Entertainment. The y-axis in Figures 3–6  
384 depicts the performance of various evaluation criteria, while the x-axis indicates the percentage  
385 of already-selected features. Our method is better at a low number of features and in Figures  
386 3–6, it is clear that each approach performs better as more features are selected. Moreover, we  
387 can compare the performance of each technique on the same dataset with the same metric. For  
388 example, Figure 3 shows the results for Arts, Business, Corel5k, and Entertainment datasets,  
389 where the Micro-F1 metric of our MLFS-GLOCAL algorithm is markedly better than that of  
390 several other algorithms.

391 *4.5. Analyzing the sensitivity of parameters*

392 In this section, we analyzed the influence of the hyperparameters, including the graph regu-  
393 larization parameter  $\lambda_1$ , the Global and Local label correlation parameter  $\lambda_2$ , and the sparsity  
394 parameter  $\lambda_3$ . Figures 7 and 8 illustrates the Micro-F1, Macro-F1, Hamming Loss, and Ranking  
395 Loss metrics of our method with various  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  on the six datasets. These Figures are  
396 plotted in 3D, i.e., three axes relate to  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  [46].

397 *4.5.1. Choosing the value of parameter  $\lambda_1$*

398 This parameter controls the manifold regularization term's effectiveness in the proposed  
399 method. In this parameter sensitivity analysis, the values in the set  $\{0, 0.01, 0.1, 1, 10, 100\}$  are  
400 selected for  $\lambda_1$  parameter across all datasets. From Figures 7 and 8, it can be deduced that the  
401 optimal value for this parameter is typically less than 1.

402 *4.5.2. Choosing the value of parameter  $\lambda_2$*

403 This parameter shows the effectiveness rate on global and local label correlation in our method  
404 where the analyzed values for  $\lambda_2$  parameter are  $\{0, 10^{-12}, 10^{-6}, 10^{-3}, 10^{-2}, 10^{-1}\}$ . As we can  
405 observe in Figures 7 and 8,  $\lambda_2$  with small values usually results a better performance in terms of  
406 Micro-F1 and Macro-F1 and with high values usually has better in terms of Hamming Loss and  
407 Coverage Error measures on the four datasets. Values near to zero or very large ones for this  
408 parameter may not perform relatively well.

409 *4.5.3. Choosing the value of parameter  $\lambda_3$*

410 The sparsity regularization term in the MLFS-GLOCAL method is controlled by  $\lambda_3$ . The  
411 range of  $\lambda_3$  is  $\{0, 0.001, 0.1, 0.5, 1, 10\}$ . The results show that  $\lambda_3$  is a delicate quantity that

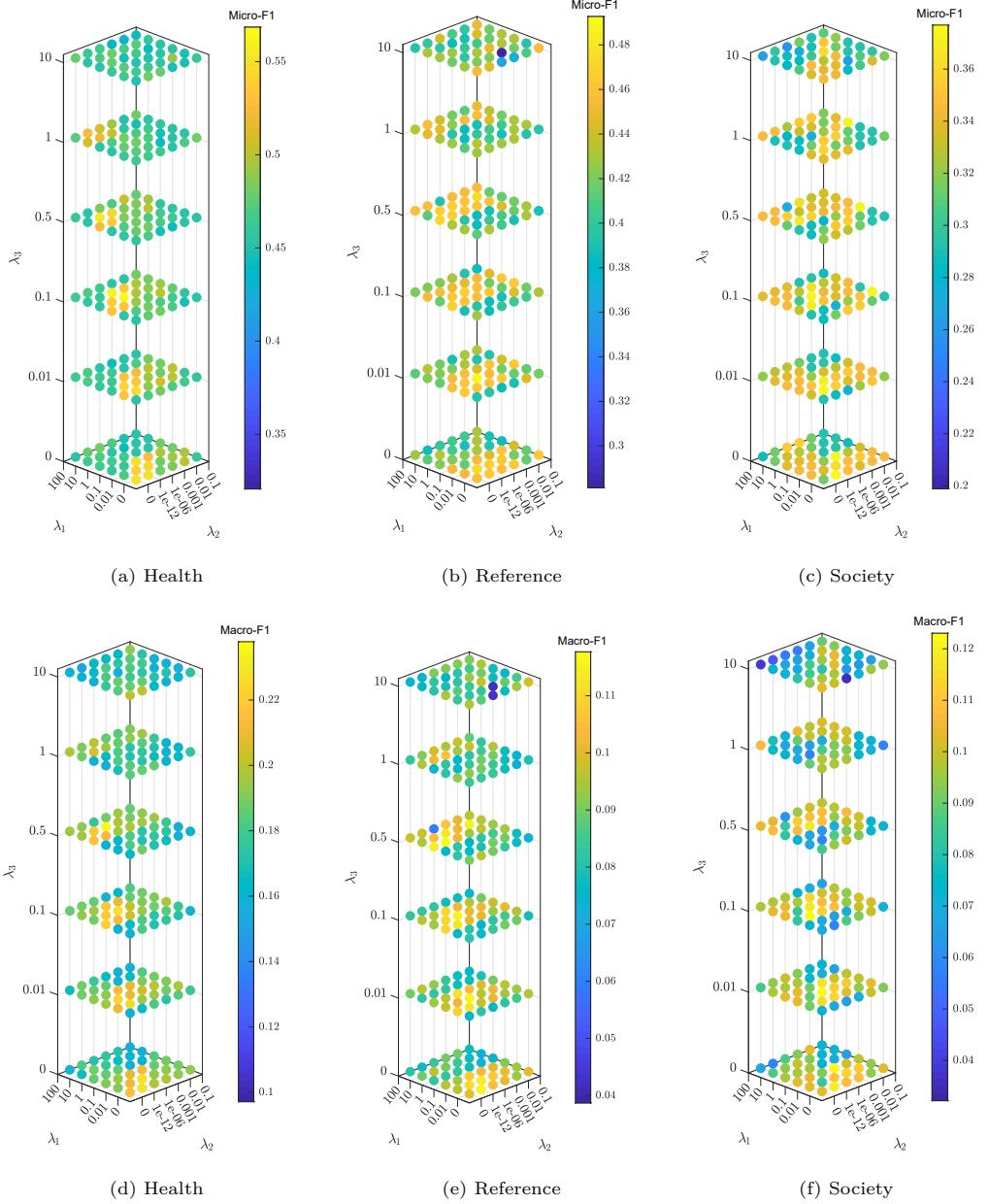


Figure 7: Parameter Analysis of our method with respect to the parameters  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  on six real-world datasets in terms of Micro-F1 and Macro-F1.

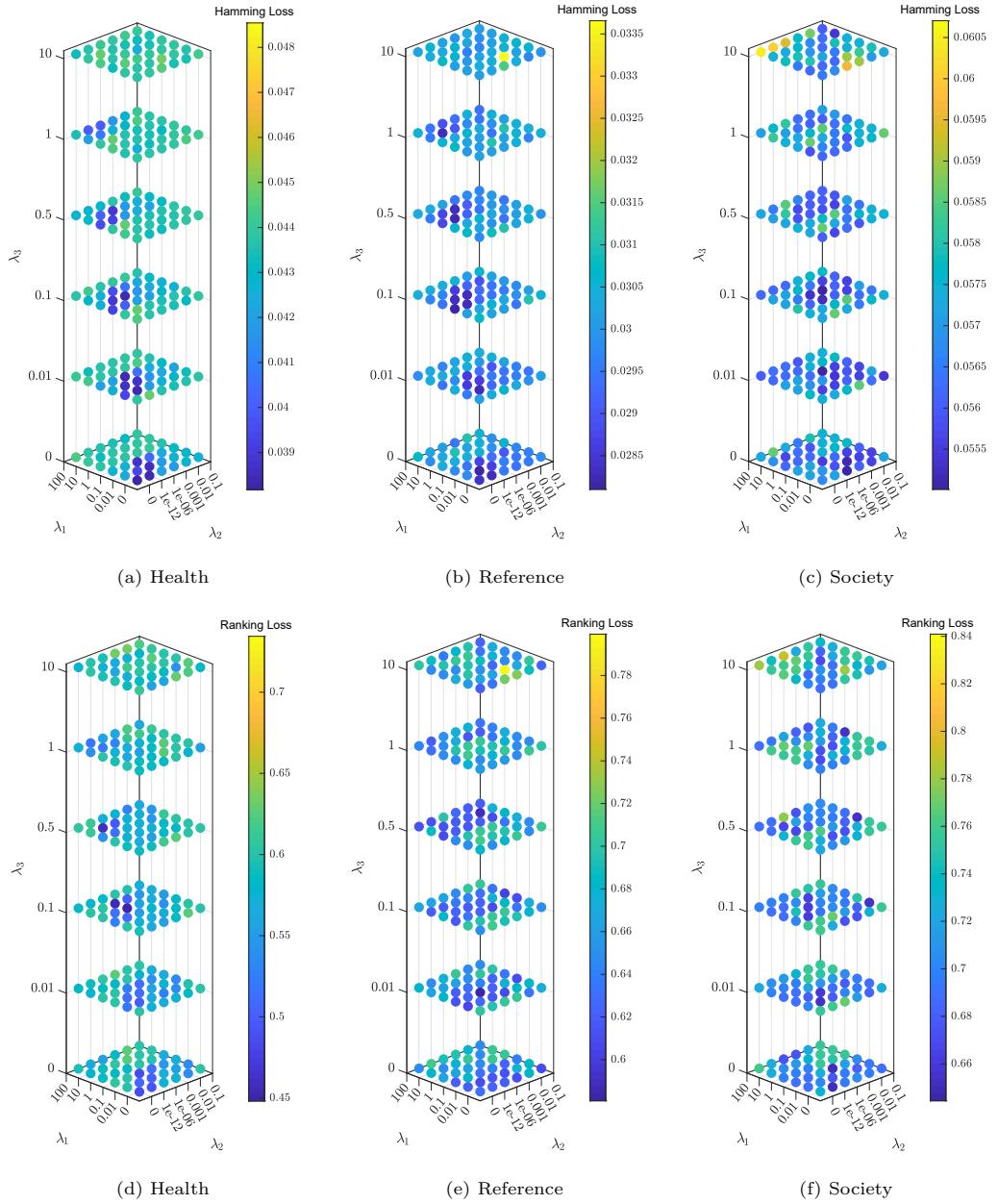


Figure 8: Parameter Analysis of our method with respect to the parameters  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  on six real-world datasets in terms of Hamming Loss and Ranking Loss.

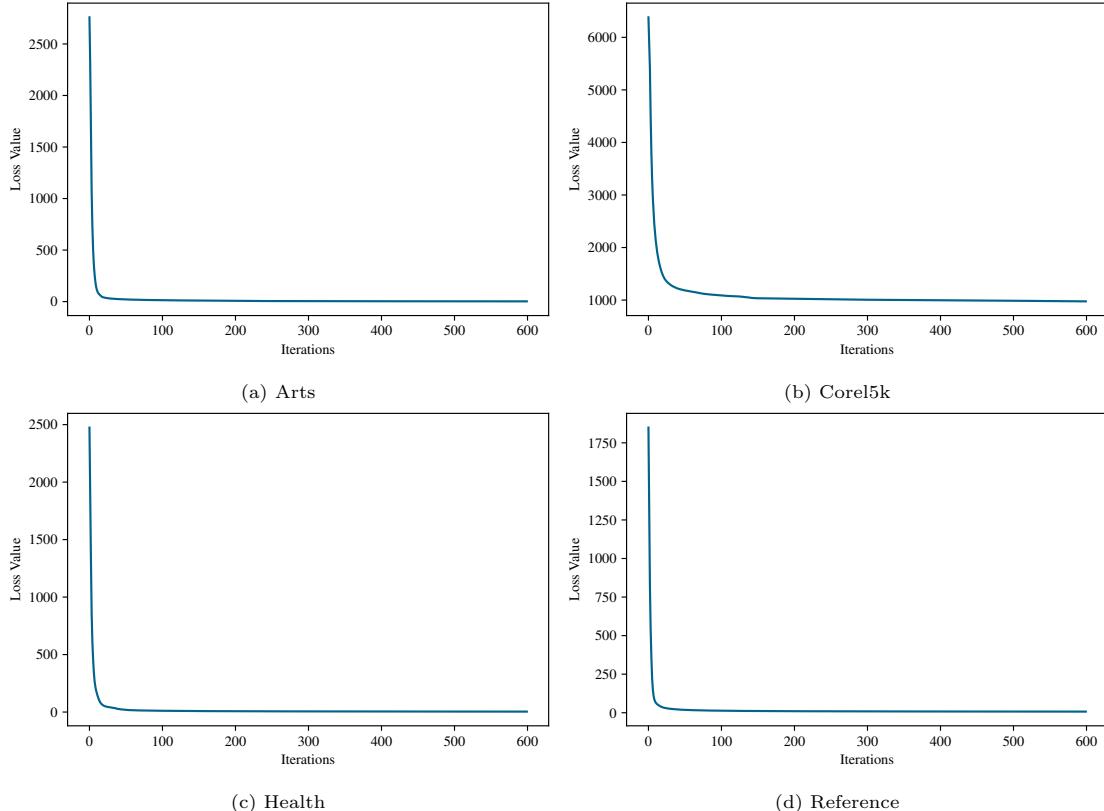


Figure 9: Convergence analysis of MLFS-GLOCAL model on the Arts, Corel5k, Health, and Reference datasets.

412 typically requires careful adjusting. The results indicate that choosing values less than 1 for this  
413 parameter is usually preferable.

414 4.6. Convergence Analysis

In this section, we evaluate the convergence behavior of our proposed MLFS-GLOCAL model  
 (22) by conducting experiments on four datasets: Arts, Corel5k, Health, and Reference. We apply  
 Algorithm 1 to each dataset and run it for 600 iterations. We plot the objective function value  
 against the number of iterations in Figure 9 to show how our method converges to a solution.  
 As can be seen from Figure 9, the objective function value decreases rapidly and steadily in the  
 initial iterations, indicating that our method quickly approaches an optimal solution. In the  
 later iterations, the objective function value changes very slightly, suggesting that our method  
 has reached a near-optimal solution. This demonstrates that our method has a fast and stable  
 convergence performance.

424    **5. Conclusion**

425    In this study, we propose a novel multi-label feature selection method by considering implicit  
426    and explicit label correlation. The global label correlation helps to exploit the underlying structure  
427    of the label space. It allows the model to learn relationships and dependencies between labels.  
428    Local label correlation, on the other hand, refers to the label associations that are specific to a local  
429    context. By considering the correlation between the labels and incorporating them into the feature  
430    representation, the model identifies the features that are most relevant to the labels. In addition,  
431    MLFS-GLOCAL learns the shared common mode between the feature matrix and label matrix  
432    to extract implicit label correlation and guide feature selection. This low-dimensional embedding  
433    is constrained through manifold regularization, which means that it has a similar local structure  
434    to the original data and retains the most valuable information in the data. An alternating  
435    optimization-based iterative algorithm is developed to solve the objective function with  $L_{2,1}$ -norm  
436    regularization. Finally, extensive experiments over variable multi-label datasets demonstrate  
437    the effectiveness of the proposed MLFS-GLOCAL against some state-of-the-art feature selection  
438    methods.

439    In future research, we should continue to learn the correlations of labels to obtain optimal  
440    feature subsets. Due to the expensive label information in the real world, we can focus on  
441    semi-supervised [47, 48] or missing multi-label learning methods [49] based on label correlations.  
442    Besides, it is also interesting to extend our work to a multi-label multi-view setting, where the  
443    data is represented by multiple feature sets or views, and each instance is associated with multiple  
444    labels.

445    **References**

- 446    [1] J. Liu, Y. Lin, Y. Li, W. Weng, S. Wu, Online multi-label streaming feature selection based  
447    on neighborhood rough set, Pattern Recognition 84 (2018) 273–287.
- 448    [2] F. Li, D. Miao, W. Pedrycz, Granular multi-label feature selection based on mutual informa-  
449    tion, Pattern Recognition 67 (2017) 410–423.
- 450    [3] Y. Lin, Q. Hu, J. Liu, J. Li, X. Wu, Streaming feature selection for multilabel learning based  
451    on fuzzy mutual information, IEEE Transactions on Fuzzy Systems 25 (6) (2017) 1491–1507.
- 452    [4] W. Gao, L. Hu, P. Zhang, Class-specific mutual information variation for feature selection,  
453    Pattern Recognition 79 (2018) 328–339.

- 454 [5] J. Huang, G. Li, Q. Huang, X. Wu, Joint feature selection and classification for multilabel  
455 learning, *IEEE transactions on cybernetics* 48 (3) (2017) 876–889.
- 456 [6] P. Zhu, Q. Xu, Q. Hu, C. Zhang, H. Zhao, Multi-label feature selection with missing labels,  
457 *Pattern Recognition* 74 (2018) 488–502.
- 458 [7] B. Tang, S. Kay, H. He, Toward optimal feature selection in naive bayes for text categoriza-  
459 tion, *IEEE transactions on knowledge and data engineering* 28 (9) (2016) 2508–2521.
- 460 [8] Z. Shajarian, S. A. Seyedi, P. Moradi, A clustering-based matrix factorization method to  
461 improve the accuracy of recommendation systems, in: 2017 Iranian Conference on Electrical  
462 Engineering (ICEE), 2017, pp. 2241–2246.
- 463 [9] R. Abdollahi, S. Amjad Seyedi, M. Reza Noorimehr, Asymmetric semi-nonnegative matrix  
464 factorization for directed graph clustering, in: 2020 10th International Conference on Com-  
465 puter and Knowledge Engineering (ICCKE), 2020, pp. 323–328.
- 466 [10] R. Hong, M. Wang, Y. Gao, D. Tao, X. Li, X. Wu, Image annotation by multiple-instance  
467 learning with discriminative feature mapping and selection, *IEEE transactions on cybernetics*  
468 44 (5) (2013) 669–680.
- 469 [11] J. Xie, M. Wang, S. Xu, Z. Huang, P. W. Grant, The unsupervised feature selection algo-  
470 rithms based on standard deviation and cosine similarity for genomic data analysis, *Frontiers*  
471 in *Genetics* 12 (2021) 684100.
- 472 [12] P. Wang, C. Domeniconi, Building semantic kernels for text classification using wikipedia,  
473 in: ACM SIGKDD international conference on Knowledge discovery and data mining, 2008,  
474 pp. 713–721.
- 475 [13] K. Trohidis, G. Tsoumakas, G. Kalliris, I. Vlahavas, Multi-label classification of music by  
476 emotion, *EURASIP Journal on Audio, Speech, and Music Processing* 2011 (1) (2011) 1–9.
- 477 [14] M.-L. Zhang, Z.-H. Zhou, A review on multi-label learning algorithms, *IEEE Transactions*  
478 on *Knowledge and Data Engineering* 26 (8) (2014) 1819–1837.
- 479 [15] Y. Zhang, Y. Ma, Non-negative multi-label feature selection with dynamic graph constraints,  
480 *Knowledge-Based Systems* 238 (2022) 107924.
- 481 [16] L. Jian, J. Li, K. Shu, H. Liu, Multi-label informed feature selection, in: Proceedings of the  
482 Twenty-Fifth International Joint Conference on Artificial Intelligence, 2016, pp. 1627–1633.

- 483 [17] Z.-H. Zhou, A brief introduction to weakly supervised learning, *National science review* 5 (1)  
484 (2018) 44–53.
- 485 [18] J. Fürnkranz, E. Hüllermeier, E. Loza Mencía, K. Brinker, Multilabel classification via cali-  
486 brated label ranking, *Machine learning* 73 (2008) 133–153.
- 487 [19] S. Ji, L. Tang, S. Yu, J. Ye, Extracting shared subspace for multi-label classification, in:  
488 Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery  
489 and data mining, 2008, pp. 381–389.
- 490 [20] J. Read, B. Pfahringer, G. Holmes, E. Frank, Classifier chains for multi-label classification,  
491 *Machine learning* 85 (2011) 333–359.
- 492 [21] S.-J. Huang, Z.-H. Zhou, Multi-label learning by exploiting label correlations locally, in:  
493 Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 26, 2012, pp. 949–955.
- 494 [22] W. Weng, Y. Lin, S. Wu, Y. Li, Y. Kang, Multi-label learning based on label-specific features  
495 and local pairwise label correlation, *Neurocomputing* 273 (2018) 385–394.
- 496 [23] D. D. Lee, H. S. Seung, Learning the parts of objects by non-negative matrix factorization,  
497 *Nature* 401 (6755) (1999) 788–791.
- 498 [24] Y. Lin, Q. Hu, J. Liu, J. Duan, Multi-label feature selection based on max-dependency and  
499 min-redundancy, *Neurocomputing* 168 (2015) 92–103.
- 500 [25] J. Lee, D.-W. Kim, Scls: Multi-label feature selection based on scalable criterion for large  
501 label set, *Pattern Recognition* 66 (2017) 342–352.
- 502 [26] P. Zhang, G. Liu, W. Gao, Distinguishing two types of labels for multi-label feature selection,  
503 *Pattern Recognition* 95 (2019) 72–82.
- 504 [27] F. Nie, H. Huang, X. Cai, C. Ding, Efficient and robust feature selection via joint l<sub>2,1</sub>-norms  
505 minimization, in: *Advances in Neural Information Processing Systems*, Vol. 23, 2010, pp.  
506 1813–1821.
- 507 [28] X. Cai, F. Nie, H. Huang, Exact top-k feature selection via l<sub>2,0</sub>-norm constraint, in: *PInter-  
508 national Joint Conference on Artificial Intelligence*, 2013, p. 12401246.
- 509 [29] A. Braytee, W. Liu, D. R. Catchpoole, P. J. Kennedy, Multi-label feature selection using  
510 correlation information, in: *Proceedings of the 2017 ACM on Conference on Information and  
511 Knowledge Management*, 2017, pp. 1649–1656.

- 512 [30] L. Hu, Y. Li, W. Gao, P. Zhang, J. Hu, Multi-label feature selection with shared common  
513 mode, *Pattern Recognition* 104 (2020) 107344.
- 514 [31] W. Gao, Y. Li, L. Hu, Multilabel feature selection with constrained latent structure shared  
515 term, *IEEE Transactions on Neural Networks and Learning Systems* 34 (3) (2023) 1253–1262.
- 516 [32] R. Huang, Z. Wu, Multi-label feature selection via manifold regularization and dependence  
517 maximization, *Pattern Recognition* 120 (2021) 108149.
- 518 [33] Y. Fan, J. Liu, P. Liu, Y. Du, W. Lan, S. Wu, Manifold learning with structured subspace  
519 for multi-label feature selection, *Pattern Recognition* 120 (2021) 108169.
- 520 [34] Y. Fan, J. Liu, W. Weng, B. Chen, Y. Chen, S. Wu, Multi-label feature selection with local  
521 discriminant model and label correlations, *Neurocomputing* 442 (2021) 98–115.
- 522 [35] Y. Zhu, J. T. Kwok, Z.-H. Zhou, Multi-label learning with global and local label correlation,  
523 *IEEE Transactions on Knowledge and Data Engineering* 30 (6) (2017) 1081–1094.
- 524 [36] D. Zhao, Q. Gao, Y. Lu, D. Sun, Learning multi-label label-specific features via global and  
525 local label correlations, *Soft Computing* 26 (5) (2022) 2225–2239.
- 526 [37] S. Kumar, R. Rastogi, Low rank label subspace transformation for multi-label learning with  
527 missing labels, *Information Sciences* 596 (2022) 53–72.
- 528 [38] D. Cai, X. He, J. Han, T. S. Huang, Graph regularized nonnegative matrix factorization for  
529 data representation, *IEEE transactions on pattern analysis and machine intelligence* 33 (8)  
530 (2010) 1548–1560.
- 531 [39] H.-Y. Chuang, E. Lee, Y.-T. Liu, D. Lee, T. Ideker, Network-based classification of breast  
532 cancer metastasis, *Molecular systems biology* 3 (1) (2007) 140.
- 533 [40] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert, M. A. Gillette,  
534 A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander, et al., Gene set enrichment analysis:  
535 a knowledge-based approach for interpreting genome-wide expression profiles, *Proceedings  
536 of the National Academy of Sciences* 102 (43) (2005) 15545–15550.
- 537 [41] G. Doquire, M. Verleysen, Mutual information-based feature selection for multilabel classifi-  
538 cation, *Neurocomputing* 122 (2013) 148–155.
- 539 [42] M.-L. Zhang, Z.-H. Zhou, Ml-knn: A lazy learning approach to multi-label learning, *Pattern  
540 recognition* 40 (7) (2007) 2038–2048.

- 541 [43] J. Liu, Y. Lin, S. Wu, C. Wang, Online multi-label group feature selection, Knowledge-Based  
542 Systems 143 (2018) 42–57.
- 543 [44] J. Zhang, Z. Luo, C. Li, C. Zhou, S. Li, Manifold regularized discriminative feature selection  
544 for multi-label learning, Pattern Recognition 95 (2019) 136–150.
- 545 [45] S. A. Seyedi, S. S. Ghodsi, F. Akhlaghian, M. Jalili, P. Moradi, Self-paced multi-label learning  
546 with diversity, in: Asian Conference on Machine Learning, PMLR, 2019, pp. 790–805.
- 547 [46] N. Salahian, F. A. Tab, S. A. Seyedi, J. Chavoshinejad, Deep autoencoder-like nmf with  
548 contrastive regularization and feature relationship preservation, Expert Systems with Appli-  
549 cations 214 (2023) 119051.
- 550 [47] J. Chavoshinejad, S. A. Seyedi, F. Akhlaghian Tab, N. Salahian, Self-supervised semi-  
551 supervised nonnegative matrix factorization for data clustering, Pattern Recognition 137  
552 (2023) 109282.
- 553 [48] S. A. Seyedi, P. Moradi, F. A. Tab, A weakly-supervised factorization method with dynamic  
554 graph embedding, in: 2017 Artificial Intelligence and Signal Processing Conference (AISP),  
555 2017, pp. 213–218.
- 556 [49] S. A. Seyedi, F. Akhlaghian Tab, A. Lotfi, N. Salahian, J. Chavoshinejad, Elastic adversarial  
557 deep nonnegative matrix factorization for matrix completion, Information Sciences 621 (2023)  
558 562–579.

**Mohammad Faraji**

Faraji has completed his master's degree in artificial intelligence and robotics in the department of Computer Engineering, University of Kurdistan. His research interests include machine learning, representation learning, matrix factorization, multi-label learning, and feature selection. Mohammad received his bachelor's degree in Software Engineering.

**Seyed Amjad Seyed**

Seyed is a graduate research assistant at the University of Kurdistan working on deep learning (from optimization and generalization aspects). He received his Master's in Artificial Intelligence from the Department of Computer Engineering at the University of Kurdistan in 2018. His work mainly focused on matrix factorization and low-rank approximation.

**Fardin Akhlaghian Tab**

Akhlaghian is the associate professor of Computer engineering at the University of Kurdistan, and his research focuses on machine learning and computer vision. He did his Ph.D. in Computer Vision at the University of Wollongong in 2005. He holds a master's degree from Tehran University of Tarbiat Modarres in 1992.

**Reza Mahmoodi**

Mahmoodi is a PhD candidate in Artificial intelligence at the University of Kurdistan. He has completed his master's degree in artificial intelligence and robotics in the department of Computer Engineering, University of Kurdistan. His research interests include machine learning, representation learning, deep learning, matrix factorization, and complex network analysis. Reza received his bachelor's degree in Software Engineering.

### **Declaration of Interest Statement**

Manuscript title:

#### **Multi-Label Feature Selection with Global and Local Label Correlation**

The authors whose names are listed immediately below certify that they have NO affiliations with or involvement in any organization or entity with any financial interest (such as honoraria; educational grants; participation in speakers' bureaus; membership, employment, consultancies, stock ownership, or other equity interest; and expert testimony or patent-licensing arrangements), or non-financial interest (such as personal or professional relationships, affiliations, knowledge or beliefs) in the subject matter or materials discussed in this manuscript.

**Authors' names:**

-

The authors whose names are listed immediately below report the following details of affiliation or involvement in an organization or entity with a financial or non-financial interest in the subject matter or materials discussed in this manuscript. Please specify the nature of the conflict on a separate sheet of paper if the space below is inadequate.

**Authors' names:**

Mohammad Faraji

Seyed Amjad Seyed

Fardin Akhlaghian Tab

Reza Mahmoodi

This statement is signed by all the authors to indicate agreement that the above information is true and correct (a photocopy of this form may be used if there are more than 10 authors):

<b>Author's name</b>	<b>Author's signature</b>	<b>Date</b>
Mohammad Faraji	M. Faraji	07/08/2023
Seyed Amjad Seyed	S. A. Seyed	07/08/2023
Fardin Akhlaghian Tab	F. Akhlaghian Tab	07/08/2023
Reza Mahmoodi	R. Mahmoodi	07/08/2023