# SPL-MLL: Selecting Predictable Landmarks for Multi-Label Learning

Junbing Li[1], Changqing Zhang[1,*], Pengfei Zhu[1],
Baoyuan Wu[2], Lei Chen[3], and Qinghua Hu[1]

[1] Tianjin Key Lab of Machine Learning, College of Intelligence and Computing,
Tianjin University, China
[2] The Chinese University of Hong Kong, Shenzhen, Tencent AI lab, China
[3] Nanjing University of Posts and Telecommunications, China
{lijunbing,zhangchangqing,zhupengfei,huqinghua}@tju.edu.cn
wubaoyuan1987@gmail.com, chenlei@njupt.edu.cn

**Abstract.** Although significant progress achieved, multi-label classification is still challenging due to the complexity of correlations among different labels. Furthermore, modeling the relationships between input and some (dull) classes further increases the difficulty of accurately predicting all possible labels. In this work, we propose to select a small subset of labels as landmarks which are easy to predict according to input (predictable) and can well recover the other possible labels (representative). Different from existing methods which separate the landmark selection and landmark prediction in the 2-step manner, the proposed algorithm, termed Selecting Predictable Landmarks for Multi-Label Learning (SPL-MLL), jointly conducts landmark selection, landmark prediction, and label recovery in a unified framework, to ensure both the representativeness and predictableness for selected landmarks. We employ the Alternating Direction Method (ADM) to solve our problem. Empirical studies on real-world datasets show that our method achieves superior classification performance over other state-of-the-art methods.

**Keywords:** Multi-label learning, predictable landmarks, a unified framework

## 1 Introduction

Multi-label classification jointly assigns one sample with multiple tags reflecting its semantic content, which has been widely used in many real-world applications. In document classification, there are multiple topics for one document; in computer vision, one image may contain multiple types of object; in emotion analysis, there may be combined types of emotions, *e.g.*, relaxing and quiet. Though plenty of multi-label classification methods [27,21,37,22,12,19,31,32] have been proposed, multi-label classification is still a recognized challenging task due to the complexity of label correlations, and the difficulty of predicting all labels.

---
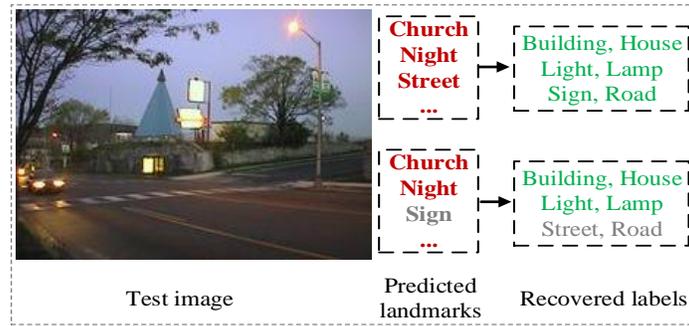
* Corresponding author: Changqing Zhang.

Fig. 1: Why landmarks should be predictable? With a test image, although the other related labels (*e.g.*, "street", "road") can be inferred from "sign", the landmark label "sign" itself is difficult to accurately predict (bottom row). While based on the image, the label "street" is more predictable, and accordingly, more related labels are correctly inferred (top row).

In real-world applications, labels are usually correlated, and simultaneously predicting all possible labels is usually rather difficult. Accordingly, there are techniques aiming to reduce the label space. The representative strategy is landmark based multi-label classification. Landmark based methods first select a small subset of representative labels as landmarks, where the landmark labels are able to establish interdependence with other labels. The algorithm [1] employ group-sparse learning to select a few labels as landmarks which can reconstruct the other labels. However, this method divides the landmark selection and landmark prediction as two separate processes. The method [3] performs label selection based on randomized sampling, where the sampling probability of each class label reflects its importance among all the labels. Another representative strategy is label embedding [13,24,7,39,40], which transforms label vectors into low-dimensional embeddings, where the correlations among labels can be implicitly encoded.

Although several methods have been proposed to reduce the dimensionality of label space, there are several limitations left behind these methods. First, existing landmark-based methods usually separate the landmark label selection and landmark prediction into two independent steps. Therefore, even the other labels could be easily recovered from the selected landmarks, but the landmarks themselves may be difficult to be accurately predicted with input (as shown in Fig. 1). Second, the label-embedding-based methods usually project the original label space into a low-dimensional embedding space, where embedded vectors can be used to recover the full-label set. Although the embedded vectors may be easy to predict, the embedding way (*i.e.*, dimensionality reduction) may cause information loss, and the label correlations are implicitly encoded thus lack interpretability. Considering the above issues, we jointly conduct landmarks selection, landmarks prediction and full-label recovery in a unified framework, and accord-

ingly, propose a novel multi-label learning method, termed *Selecting Predictable Landmarks for Multi-Label Learning* (**SPL-MLL**). The overview of SPL-MLL is shown in Fig. 2. The main advantages of the proposed algorithm include: (1) compared with existing landmark-based multi-label learning methods, SPL-MLL can select the landmarks which are both representative and predictable due to the unified objective; (2) compared with the embedding methods, SPL-MLL is more interpretable due to explicitly exploring the correlations with landmarks.

The contributions of this work are summarized as:

- We propose a novel landmark-based multi-label learning algorithm for complex correlations among labels. The landmarks bridge the intrinsic correlations among different labels, while also reduce the complexity of correlations and possible label noise.
- To the best of our knowledge, SPL-MLL is the first algorithm which simultaneously conducts landmark selection, landmark prediction, and full-label recovery in a unified objective, thus taking both representativeness and predictability for landmarks into account. This is quite different from the 2-step manner separating landmark selection and prediction.
- Extensive experiments on benchmark datasets are conducted, validating the effectiveness of the proposed method over state-of-the-arts.

## 2   Related Work

Generally, existing multi-label methods can be roughly categorized into three lines based on the order of label correlations [37]. The first-order strategy [4,36] tackles multi-label learning problem in the label-by-label manner, which ignores the co-existence of other labels. The second-order strategy [8,10,11] conducts multi-label learning problem by introducing the pairwise relations between different labels. For high-order strategy [14,22,28,30], multi-label learning problem is solved by establishing more complicated label relationships, which makes these approaches tend to be quite computationally expensive.

In order to reduce label space, there are approaches based on label embedding, which searches a low-dimensional subspace so that correlations among labels can be implicitly expressed [13,24,7,23,15,39,40,34,18,16]. Based on the low-dimensional latent label space, one can effectively reduce computation cost while performing multi-label prediction. The representative embedding based methods include: label embedding via random projections [13], principal label space transformation (PLST) [24] and its conditional version (CPLST) [7]. Beyond considering linear embedding functions, there are several approaches employing standard kernel functions (*e.g.*, low-degree polynomial kernels) for nonlinear label embedding. The work in [33] proposes a novel DNN architecture of Canonical-Correlated Autoencoder (C2AE), which is a DNN-based label embedding framework for multi-label classification, which is able to perform feature-aware label embedding and label-correlation aware prediction.

To explore label correlations, there are several landmark based multi-label classification models aiming to reduce the label space [1,3,40,5]. They usually

first select a small subset of labels as landmarks, which are supposed to be representative and able to establish interdependency with other labels. The work in [1] models landmark selection with group-sparsity technique. Following the assumption in [1], the method in [3] alleviates this problem of computation cost by proposing an valid label selection method based on randomized sampling, and utilizes the leverage score in the best rank-$k$ subspace of the label matrix to obtain the sampling probability of each label. It is noteworthy that these methods separate the landmark selection and landmark prediction in a 2-step manner, which can not simultaneously guarantee the representativeness and predictability of landmarks.
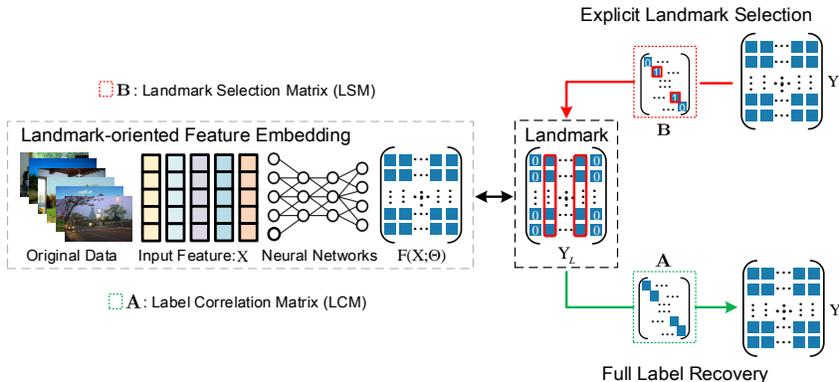


Fig. 2: Overview of SPL-MLL. The key component of our model is the landmark selection strategy, which induces the explicit landmark label matrix $\mathbf{Y}_L$. The matrix $\mathbf{B}$, termed as landmark selection matrix, is used to construct the landmark label matrix explicitly, while the matrix $\mathbf{A}$ is used to reconstruct all possible labels from landmarks. Benefitting from the explicit landmark label matrix $\mathbf{Y}_L$, the input is also able to be taken into account to ensure the predictable property for landmarks.

## 3   Our Algorithm: Selecting Predictable Landmarks for Multi-Label Learning

For clarification, we first provide the definitions for symbols and variables used through out this paper. Let $\mathcal{X} = \mathbb{R}^D$ and $\mathcal{Y} = \{0,1\}^C$ denote the feature space and label space, where $D$ and $C$ are the dimensionality of feature space and label space, respectively. Given training data with the form of instance-label pairs $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$, accordingly, the feature matrix can be represented as $\mathbf{X} \in \mathbb{R}^{N \times D}$, and the label matrix is represented as $\mathbf{Y} \in \mathbb{R}^{N \times C}$. The goal of multi-label learning is to learn a model $f : \mathcal{X} \to \mathcal{Y}$, to predict possible labels accurately for new coming instances. Motivated by the landmark strategy, we propose a novel algorithm for multi-label learning, termed SPL-MLL, *i.e.*, *Selecting Predictable*

*Landmarks for Multi-Label Learning.* SPL-MLL consists of two key components, *i.e.*, *Explicit Landmark Selection* and *Predictable Landmark Classification*.

### 3.1 Explicit Landmark Selection

Different from the 2-step manner [1] which only focuses on selecting landmarks that are most representative, our goal is to select landmarks which are both representative and predictable. There are two designed matrixes which are the keys to realize this goal. The first matrix is the **label correlation matrix (LCM) $\mathbf{A}$** used for recovering other labels with landmarks. In self-representation manner, the matrix $\mathbf{A} \in \mathbb{R}^{C \times C}$ is obtained which captures the correlation among labels and explores the interdependency between landmark labels and the others. In the work [1], the landmarks are selected implicitly. Specifically, the underlying assumption is $\mathbf{Y} = \mathbf{YA}$, where $\mathbf{A}$ is constrained by minimizing $||\mathbf{A}||_{2,1}$ to enforce the reconstruction of $\mathbf{Y}$ mainly base on a few labels, *i.e.*, landmarks. In our model, although the linear self-representation manner is also introduced in our model, we try to obtain the landmark label matrix explicitly in the objective function.

The second critical matrix is the **landmark selection matrix (LSM) $\mathbf{B}$**. Note that, in the work [1], there is no explicit landmark label matrix constructed and the selection result is implicitly encoded in $\mathbf{A}$ due to its sparsity in row. Both $\mathbf{YA}$ and $\mathbf{Y}$ ($\mathbf{YA} \approx \mathbf{Y}$) in [1] are full-label matrix. Different from [1], since we aim to jointly conduct landmark selection and learn a model to predict these selected landmarks instead of all labels, we need to explicitly derive a label matrix encoding the landmarks. To this end, we introduce the matrix $\mathbf{B} \in \mathbb{R}^{C \times C}$ which is a diagonal matrix, and each diagonal element is either 0 or 1, *i.e.*, $B_{ii} \in \{0, 1\}$. Then, we can obtain the explicit landmark label matrix $\mathbf{Y}_L$ with $\mathbf{Y}_L = \mathbf{YB}$. In this way, the columns corresponding to the landmarks in $\mathbf{Y}_L$ unchanged while the elements of other columns (corresponding to non-landmark labels) will be 0. It is noteworthy that $\mathbf{B}$ is learned in our model instead of being fixed in advance. Accordingly, the explicit landmark selection objective to minimize is induced as:

$$\begin{aligned}
\mathbf{\Gamma}(\mathbf{B}, \mathbf{A}) &= \|\mathbf{Y} - \mathbf{Y}_L \mathbf{A}\|_F^2 + \mathbf{\Omega}(\mathbf{B}) \\
&= \|\mathbf{Y} - \mathbf{YBA}\|_F^2 + \mathbf{\Omega}(\mathbf{B}), \\
&s.t.\ B_{ij} = 0,\ i \neq j;\ B_{ij} \in \{0, 1\},\ i = j.
\end{aligned} \tag{1}$$

Since it is difficult to strictly ensure the diagonal property for $\mathbf{B}$, a soft constraint $\mathbf{\Omega}(\mathbf{B})$ is introduced as follows:

$$\mathbf{\Omega}(\mathbf{B}) = \lambda_1 \|\mathbf{B} - \mathbf{I}\|_F^2 + \lambda_2 \|\mathbf{B}\|_{2,1}, \tag{2}$$

where the structure sparsity $||\mathbf{B}||_{2,1} = \sum_{i=1}^{C} \sqrt{\sum_{j=1}^{C} B_{ij}^2}$ is used to select a few landmarks, and the approximation to the identity matrix $\mathbf{I}$ ensures the labels corresponding to landmarks unchanged. The regularization parameter $\lambda_1$ and $\lambda_2$ control the degree of diagonal and sparsity property for $\mathbf{B}$, respectively. Notice

that the label correlation matrix $\mathbf{A}$ is learned automatically without constraint, the underlying assumption for the correlation is sparse (similar to the existing work [1]) which is jointly ensured by the sparse landmark selection matrix $\mathbf{B}$. Then, we can obtain the explicit landmark label matrix $\mathbf{Y}_L = \mathbf{YB}$, and train a prediction model exactly for the landmarks.

### 3.2   Predictable Landmark Classification

Now, we firstly consider learning the classification model for accurately predicting landmarks instead of all possible labels. Beyond label correlation, modeling $\mathcal{X} \to \mathcal{Y}$ is also critical in multi-label classification. However, the traditional landmark-based multi-label classification algorithms usually separate landmark selection and landmark prediction, which may result in unpromising classification accuracy because the selected landmarks may be representative but difficult to be predicted (see Fig. 1). Recall that the goal of our model is to recover full labels with landmark labels, so our classification model only focuses on predicting landmarks $\mathbf{Y}_L$ based on $\mathbf{X}$ instead of full labels $\mathbf{Y}$. Accordingly, our predictable landmark classification objective to minimize is as follows:

$$\begin{aligned}
\mathbf{\Phi}(\mathbf{B}, \mathbf{\Theta}) &= \|\boldsymbol{f}(\mathbf{X}; \mathbf{\Theta})\mathbf{B} - \mathbf{Y}_L\|_F^2 \\
&= \|(\boldsymbol{f}(\mathbf{X}; \mathbf{\Theta}) - \mathbf{Y})\mathbf{B}\|_F^2,
\end{aligned} \tag{3}$$

where $\boldsymbol{f}(\cdot; \mathbf{\Theta})$ is the neural networks (parameterized by $\mathbf{\Theta}$) used for feature embedding and conducting classification for landmarks, which is implemented by fully connected neural networks.

### 3.3   Objective Function

Based on above considerations, a novel landmark-based multi-label classification algorithm, *i.e.*, Selecting Predictable Landmarks for Multi-Label Learning (SPL-MLL), is induced, which jointly learns landmark selection matrix, label correlation matrix, and landmark-oriented feature embedding in a unified framework. Specifically, the objective function of SPL-MLL for us to minimize is as follows:

$$\begin{aligned}
\mathcal{L}(\mathbf{B}, \mathbf{\Theta}, \mathbf{A}) &= \mathbf{\Phi}(\mathbf{B}, \mathbf{\Theta}) + \mathbf{\Gamma}(\mathbf{B}, \mathbf{A}) \\
&= \|(\boldsymbol{f}(\mathbf{X}; \mathbf{\Theta}) - \mathbf{Y})\mathbf{B}\|_F^2 \\
&\quad + \|\mathbf{Y} - \mathbf{YBA}\|_F^2 + \lambda_1 \|\mathbf{B} - \mathbf{I}\|_F^2 + \lambda_2 \|\mathbf{B}\|_{2,1}.
\end{aligned} \tag{4}$$

It is noteworthy that the critical role of matrix $\mathbf{B}$, which bridges the landmark selection and landmark classification model. With this strategy, the proposed model jointly selects predictable landmark labels, captures the correlations among labels, and discovers the nonlinear correlations between features and landmarks, accordingly, promotes the performance of multi-label prediction.

---

**Algorithm 1:** Algorithm of SPL-MLL

---

**Input:** Feature matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$, label matrix $\mathbf{Y} \in \mathbb{R}^{N \times C}$, parameters $\lambda_1, \lambda_2$.
**Initialize:** $\mathbf{B} = \mathbf{I}$, initialize randomly $\mathbf{A}$.
**while** *not converged* **do**
  Update the parameters $\boldsymbol{\Theta}$ of $\boldsymbol{f}(\cdot; \boldsymbol{\Theta})$;
  Update $\mathbf{B}$ by Eq.(5);
  Update $\mathbf{A}$ by Eq.(6);
**end**
**Output:** $\boldsymbol{f}(\cdot; \boldsymbol{\Theta}), \mathbf{B}, \mathbf{A}$.

---

### 3.4 Optimization

Since the objective function of our SPL-MLL is not jointly convex for all the variables, we optimize our objective function by employing Alternating Direction Minimization(ADM) [17] strategy. To optimize the objective function in Eq. (4), we should solve three subproblems with respect to $\boldsymbol{\Theta}$, $\mathbf{B}$ and $\mathbf{A}$, respectively. The optimization is cycled over updating different blocks of variables. We apply the technique of stochastic gradient descent for updating $\boldsymbol{\Theta}$, $\mathbf{B}$ and $\mathbf{A}$. The details of optimization are demonstrated as follows:

• Update networks. The back-propagation algorithm is employed to update the network parameters.

• Update $\mathbf{B}$. The gradient of $\mathcal{L}$ with respect to $\mathbf{B}$ can be derived as:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \mathbf{B}} = {} & 2(\boldsymbol{f}(\mathbf{X}; \boldsymbol{\Theta}) - \mathbf{Y})^T (\boldsymbol{f}(\mathbf{X}; \boldsymbol{\Theta}) - \mathbf{Y}) \mathbf{B} \\
& - 2\mathbf{Y}^T(\mathbf{Y} - \mathbf{Y}\mathbf{B}\mathbf{A})\mathbf{A}^T + 2\lambda_1(\mathbf{B} - \mathbf{I}) + 2\lambda_2 \mathbf{D}\mathbf{B},
\end{aligned} \tag{5}$$

where $\mathbf{D}$ is a diagonal matrix with $D_{ii} = \frac{1}{2\|\mathbf{B}_i\|}$. Accordingly, gradient descent is employed based on Eq. (5).

• Update $\mathbf{A}$. The gradient of $\mathcal{L}$ with respect to $\mathbf{A}$ can be derived as:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}} = -2\mathbf{B}^T \mathbf{Y}^T (\mathbf{Y} - \mathbf{Y}\mathbf{B}\mathbf{A}). \tag{6}$$

then $\mathbf{A}$ is updated by applying gradient descent based on Eq. (6). The optimization procedure of SPL-MLL is summarized as Algorithm 1.

Once the model of SPL-MLL is obtained, it can be easily applied for predicting the labels of test samples. Specifically, given a test input $\mathbf{x}$, it will be first transformed into $\boldsymbol{f}(\mathbf{x}; \boldsymbol{\Theta})$, followed by utilizing the learned mappings $\mathbf{B}$ and $\mathbf{A}$ to predict its all possible labels with $\mathbf{y} = \boldsymbol{f}(\mathbf{x}; \boldsymbol{\Theta})\mathbf{B}\mathbf{A}$.

## 4 Experiments

### 4.1 Experiment Settings

We conduct experiments on the following benchmark multi-label datasets: emotions [26], yeast [8], tmc2007 [6], scene [4], espgame [29] and pascal VOC 2007 [9].

Specifically, emotions and yeast are used for music and gene functional classification, respectively; tmc2007 is a large-scale text dataset, while scene, espgame and pascal voc 2007 belong to the domain of image. The description of features for emotions, yeast, tmc2007 and scene could be referred in [26,8,6,4]. For espgame and pascal voc 2007, the local descriptor DenseSift [20] is used. These datasets can be found in Mulan [1] and LEAR websites [2]. The detailed statistics information of each dataset is listed in Table 1. We employ the standard partitions for training and testing sets [1,2].

For the proposed SPL-MLL, we utilize neural networks for feature embedding and classification. The networks consists of 2 layers: for the first and second fully connected layer, 512 and 64 neurons are deployed, respectively. A leaky ReLU activation function is employed with the batch size being 64. In addition, we initialize the matrix $\mathbf{B}$ with $\mathbf{B} = \mathbf{I}$ which captures the most sparse correlation among labels and is beneficial to landmark selection. The regularization parameters, i.e., $\lambda_1$ and $\lambda_2$ are both fixed as 0.1 for all datasets and promising performance is obtained. In our experiments, we set the constraint $\mathbf{B}_{ij} = 0, i \neq j$ in each iteration of optimization. This strictly guarantees the diagonal property and can provide clear interpretability for landmarks. The experimental results show that both convergence of our model and promising performance are achieved with this constraint. Five diverse metrics are employed for performance evaluation.

Table 1: Statistics of datasets.

| dataset | #instances | #features | #labels | cardinality | domain |
|---|---|---|---|---|---|
| emotions | 593 | 72 | 6 | 1.9 | music |
| scene | 2407 | 294 | 6 | 1.1 | image |
| yeast | 2417 | 103 | 14 | 4.2 | biology |
| tmc2007 | 28596 | 500 | 22 | 2.2 | text |
| espgame | 20770 | 1000 | 268 | 4.7 | image |
| pascal VOC 2007 | 9963 | 1000 | 20 | 1.5 | image |

For *Hamming loss* and *Ranking loss*, smaller value indicates better classification quality, while larger value of *Average precision*, *Macro-F1* and *Micro-F1* means better performance. These evaluation metrics evaluate the performance of multi-label predictor from various aspects, and details of these evaluation metrics can be found in [37]. 10-fold cross-validation is performed for each method, which randomly holds 1/10 of training data for validation during each fold. We repeat each experiment 10 times and report the averaged results with standard derivations.

### 4.2   Experimental Results

**Comparison with state-of-the-art multi-label classification methods**
We compare our algorithm with both baseline and state-of-the-art multi-label

---

[1] http://mulan.sourceforge.net/datasets-mlc.html
[2] http://lear.inrialpes.fr/people/guillaumin/data.php

Table 2: Comparing results (mean ± std.) of multi-label learning algorithms. ↓ (↑) indicates the smaller (larger), the better. The values in red and blue indicate the best and the second best performances, respectively. • indicates that ours is better than the compared algorithms.

| Datasets | Methods | Ranking Loss ↓ | Hamming Loss ↓ | Average Precision ↑ | Micro-F1 ↑ | Macro-F1 ↑ |
|---|---|---|---|---|---|---|
| emotions | BR [27] | 0.309±0.021• | 0.265±0.015• | 0.687±0.017• | 0.592±0.025• | 0.590±0.016• |
| | LP [4] | 0.345±0.022• | 0.277±0.010• | 0.661±0.018• | 0.533±0.016• | 0.504±0.019• |
| | ML-kNN [36] | 0.173±0.015• | 0.209±0.021• | 0.794±0.016• | 0.650±0.031• | 0.607±0.033• |
| | EPS [21] | 0.183±0.014• | 0.208±0.010• | 0.780±0.017• | 0.664±0.012• | 0.655±0.018• |
| | ECC [22] | 0.198±0.021• | 0.228±0.022• | 0.766±0.014• | 0.617±0.013• | 0.597±0.019• |
| | RAkEL [28] | 0.217±0.026• | 0.219±0.013• | 0.766±0.031• | 0.634±0.023• | 0.618±0.036• |
| | CLR [10] | 0.199±0.024• | 0.255±0.012• | 0.762±0.024• | 0.614±0.037• | 0.601±0.038• |
| | MLML [12] | 0.184±0.015• | 0.197±0.013• | 0.719±0.018• | 0.661±0.039• | 0.650±0.047• |
| | MLFE [38] | 0.181±0.012• | 0.217±0.020• | 0.782±0.013• | 0.674±0.026• | 0.663±0.021• |
| | HNOML [35] | 0.173±0.012• | 0.192±0.005• | 0.784±0.011• | 0.672±0.014• | 0.660±0.029• |
| | Ours (linear) | 0.172±0.006 | 0.184±0.015 | 0.798±0.011 | 0.686±0.013 | 0.675±0.031 |
| | Ours | 0.170±0.004 | 0.175±0.021 | 0.815±0.014 | 0.698±0.021 | 0.687±0.024 |
| yeast | BR [27] | 0.322±0.011• | 0.253±0.004• | 0.614±0.008• | 0.569±0.014• | 0.386±0.011• |
| | LP [4] | 0.408±0.008• | 0.282±0.005• | 0.566±0.008• | 0.519±0.023• | 0.361±0.025• |
| | ML-kNN [36] | 0.171±0.006 | 0.218±0.004• | 0.757±0.011• | 0.636±0.012• | 0.357±0.021• |
| | EPS [21] | 0.205±0.003• | 0.214±0.005• | 0.731±0.017• | 0.625±0.015• | 0.372±0.014• |
| | ECC [22] | 0.187±0.007• | 0.209±0.009• | 0.745±0.012• | 0.618±0.013• | 0.369±0.017• |
| | RAkEL [28] | 0.250±0.005• | 0.232±0.005• | 0.710±0.009• | 0.632±0.009• | 0.430±0.012• |
| | CLR [10] | 0.187±0.005• | 0.222±0.005• | 0.745±0.008• | 0.628±0.012• | 0.400±0.018• |
| | MLML [12] | 0.178±0.002• | 0.224±0.005• | 0.757±0.009• | 0.641±0.014• | 0.443±0.025• |
| | MLFE [38] | 0.169±0.021 | 0.227±0.010• | 0.754±0.012• | 0.646±0.013• | 0.415±0.011• |
| | HNOML [35] | 0.179±0.007• | 0.222±0.004• | 0.757±0.011• | 0.648±0.006• | 0.421±0.016• |
| | Ours (linear) | 0.172±0.003 | 0.210±0.008 | 0.769±0.006 | 0.659±0.012 | 0.443±0.016 |
| | Ours | 0.171±0.004 | 0.201±0.006 | 0.786±0.005 | 0.667±0.011 | 0.451±0.023 |
| scene | BR [27] | 0.236±0.017• | 0.136±0.004• | 0.715±0.011• | 0.609±0.014• | 0.616±0.025• |
| | LP [4] | 0.219±0.010• | 0.149±0.006• | 0.722±0.010• | 0.585±0.016• | 0.592±0.011• |
| | ML-kNN [36] | 0.093±0.009• | 0.095±0.008• | 0.851±0.016• | 0.718±0.015• | 0.719±0.024• |
| | EPS [21] | 0.113±0.007• | 0.103±0.017• | 0.825±0.013• | 0.686±0.018• | 0.688±0.018• |
| | ECC [22] | 0.103±0.010• | 0.104±0.012• | 0.832±0.015• | 0.668±0.017• | 0.671±0.016• |
| | RAkEL [28] | 0.106±0.005• | 0.106±0.005• | 0.829±0.007• | 0.636±0.023• | 0.644±0.019• |
| | CLR [10] | 0.106±0.003• | 0.138±0.003• | 0.817±0.006• | 0.612±0.026• | 0.620±0.025• |
| | MLML [12] | 0.079±0.004• | 0.098±0.013• | 0.862±0.010 | 0.728±0.029• | 0.729±0.029• |
| | MLFE [38] | 0.079±0.002• | 0.094±0.003• | 0.858±0.013• | 0.732±0.021• | 0.734±0.019• |
| | HNOML [35] | 0.103±0.005• | 0.110±0.003• | 0.832±0.108• | 0.733±0.011• | 0.736±0.013• |
| | Ours (linear) | 0.073±0.003 | 0.083±0.006 | 0.861±0.005 | 0.738±0.012 | 0.742±0.021 |
| | Ours | 0.067±0.003 | 0.074±0.004 | 0.884±0.005 | 0.746±0.016 | 0.753±0.024 |
| espgame | BR [27] | 0.266±0.003• | 0.019±0.002• | 0.221±0.001• | 0.205±0.004• | 0.116±0.001• |
| | LP [4] | 0.496±0.003• | 0.031±0.001• | 0.055±0.004• | 0.109±0.003• | 0.060±0.002• |
| | ML-kNN [36] | 0.238±0.001• | 0.017±0.002 | 0.255±0.003• | 0.039±0.002• | 0.020±0.001• |
| | EPS [21] | 0.380±0.001• | 0.017±0.001 | 0.200±0.003• | 0.083±0.002• | 0.065±0.001• |
| | ECC [22] | 0.230±0.001• | 0.020±0.002• | 0.282±0.001• | 0.245±0.004• | 0.123±0.001• |
| | RAkEL [28] | 0.343±0.001• | 0.019±0.001• | 0.211±0.003• | 0.150±0.003• | 0.059±0.001• |
| | CLR [10] | 0.196±0.001 | 0.019±0.001• | 0.305±0.003 | 0.266±0.004• | 0.143±0.001• |
| | MLML [12] | 0.317±0.000• | 0.019±0.003• | 0.086±0.002• | 0.103±0.003• | 0.060±0.002• |
| | MLFE [38] | 0.312±0.012• | 0.020±0.001• | 0.268±0.011• | 0.260±0.003• | 0.134±0.004• |
| | HNOML [35] | 0.221±0.001• | 0.019±0.003• | 0.271±0.003• | 0.263±0.006• | 0.132±0.004• |
| | Ours (linear) | 0.223±0.002 | 0.017±0.001 | 0.289±0.002 | 0.269±0.002 | 0.143±0.002 |
| | Ours | 0.220±0.003 | 0.016±0.001 | 0.291±0.002 | 0.276±0.004 | 0.149±0.001 |
| tmc2007 | BR [27] | 0.037±0.007• | 0.031±0.004• | 0.899±0.025• | 0.834±0.014• | 0.719±0.011• |
| | LP [4] | 0.324±0.018• | 0.041±0.006• | 0.594±0.012• | 0.791±0.008• | 0.721±0.004• |
| | ML-kNN [36] | 0.031±0.006• | 0.058±0.004• | 0.844±0.017• | 0.682±0.003• | 0.493±0.002• |
| | EPS [21] | 0.021±0.004• | 0.033±0.005• | 0.927±0.007• | 0.829±0.009• | 0.722±0.010• |
| | ECC [22] | 0.017±0.006• | 0.026±0.003• | 0.925±0.006• | 0.862±0.014• | 0.763±0.007• |
| | RAkEL [28] | 0.038±0.008• | 0.024±0.002• | 0.923±0.005• | 0.870±0.011• | 0.756±0.006• |
| | CLR [10] | 0.018±0.005• | 0.034±0.004• | 0.923±0.011• | 0.825±0.013• | 0.711±0.011• |
| | MLML [12] | 0.018±0.001• | 0.021±0.001• | 0.921±0.002• | 0.865±0.011• | 0.769±0.008• |
| | MLFE [38] | 0.021±0.002• | 0.022±0.001• | 0.924±0.013• | 0.873±0.015• | 0.771±0.011• |
| | HNOML [35] | 0.023±0.002• | 0.017±0.001• | 0.919±0.003• | 0.858±0.014• | 0.762±0.016• |
| | Ours (linear) | 0.015±0.003 | 0.013±0.002 | 0.937±0.007 | 0.912±0.008 | 0.781±0.005 |
| | Ours | 0.012±0.004 | 0.011±0.001 | 0.945±0.007 | 0.944±0.007 | 0.792±0.010 |

classification methods. The binary relevance (BR) [27] and label powerset (LP) [4] act as baselines. We also compare ours with two ensemble methods, *i.e.*, ensemble of pruned sets (EPS) [21] and ensemble of classifier chains (ECC) [22], second-order approach - calibrated label ranking (CLR) [10] and high-order approach - random k-labelsets (RAkEL) [28], the lazy multi-label methods based on k-nearest neighbors (ML-kNN)[36] and feature-aware approach - multi-label manifold learning (MLML) [12], labeling information enrichment approach - Multi-label Learning with Feature-induced labeling information Enrichment(MLFE) [38] and robust approach for data with hybrid noise - hybrid noise-oriented multilabel learning (HNOML) [35]. We try our best to tune the parameters of all the above compared methods to the best performance according to the suggested ways in their literatures.

Table 3: Performance comparisons with approaches based on label space reduction.

| Datasets | tmc2007 | | espgame | |
|---|---|---|---|---|
| Methods / Metrics | Micro-F1↑ | Macro-F1↑ | Micro-F1↑ | Macro-F1↑ |
| MOPLMS [1] | 0.556±0.012 | 0.421±0.013 | 0.032±0.006 | 0.025±0.005 |
| ML-CSSP [3] | 0.604±0.014 | 0.432±0.015 | 0.035±0.004 | 0.023±0.006 |
| PBR [7] | 0.602±0.034 | 0.422±0.025 | 0.021±0.008 | 0.014±0.003 |
| CPLST [7] | 0.643±0.027 | 0.437±0.031 | 0.042±0.005 | 0.023±0.004 |
| FAIE [18] | 0.605±0.011 | 0.458±0.015 | 0.072±0.008 | 0.026±0.003 |
| Deep CPLST | 0.786±0.021 | 0.601±0.031 | 0.074±0.004 | 0.016±0.002 |
| Deep FAIE | 0.604±0.016 | 0.435±0.029 | 0.121±0.011 | 0.024±0.003 |
| LEML [34] | 0.704±0.013 | 0.616±0.022 | 0.148±0.004 | 0.082±0.001 |
| SLEEC [2] | 0.607±0.031 | 0.586±0.011 | 0.226±0.016 | 0.108±0.009 |
| DC2AE [33] | 0.808±0.017 | 0.757±0.027 | 0.256±0.013 | 0.121±0.009 |
| **Ours** | **0.944±0.007** | **0.792±0.010** | **0.276±0.004** | **0.149±0.001** |

As shown in Table 2, we report the quantitative experimental results of different methods on the benchmark datasets. Because above comparison methods are not based on neural networks, for fair comparisons, we also report the results of our model using the linear projections instead of neural networks for feature embedding. For each algorithm, the averaged performance with standard deviation are reported in terms of different metrics. As for each metric, "↑" indicates the larger the better while "↓" indicates the smaller the better. The red number and blue number indicate the best and the second best performances, respectively. According to Table 2, several observations are obtained as follows: 1) Compared with other multi-label classification methods, our algorithm achieves competitive performance on all the five benchmark datasets. For example, on emotions , scene and tmc2007, our SPL-MLL ranks as the first in terms of all metrics. 2) Compared with BR and LP, our SPL-MLL obtains much better performance on all datasets. The reason may be that these methods lack of sufficient ability to explore complex correlations among labels. 3) Compared with the three ensemble methods EPS, ECC and RAkEL, our algorithm always performs better, which further verifies the effectiveness of our SPL-MLL. 4) We also note that the performances of ML-kNN, CLR and MLML are also competitive, and the

performances of CLR are slightly better than ours on espgame in terms of some metrics. However, the performances of ours are more stable and robust for different datasets. For example, CLR performs unpromising on emotions, yeast, and scene. 5) Furthermore, compared with the latest and most advanced approaches MLFE and HNOML, our model outperforms them on all datasets in terms of most metrics. In short, our proposed SPL-MLL achieves promising and stable performance compared with state-of-the-art multi-label classification methods.
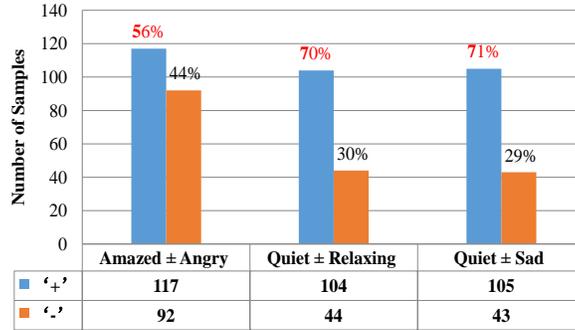


| | Amazed ± Angry | Quiet ± Relaxing | Quiet ± Sad |
|---|---|---|---|
| '+' | 117 | 104 | 105 |
| '-' | 92 | 44 | 43 |

Fig. 3: Visualization of the number of co-occurrence labels on emotions. '+' and '-' denote co-occurrence and no co-occurrence of two labels, respectively.
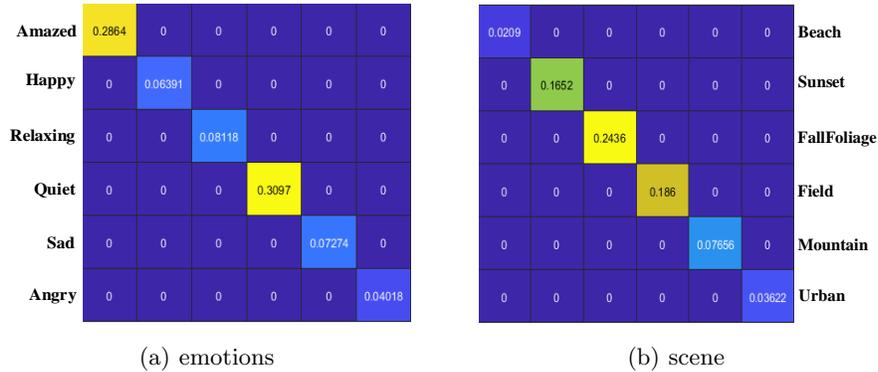


(a) emotions                    (b) scene

Fig. 4: Visualization of the landmark selection matrix **B**.

**Comparison with label space reduction methods** We compare our method with two typical landmark selection methods which conduct landmark selection with group-sparsity technique (MOPLMS) [1] and an efficient randomized sampling procedure (ML-CSSP) [3]. Moreover, we compare our method with the label embedding multi-label classification methods, which jointly reduce the label

space and explore the correlations among labels. Specifically, we conduct comparison with the following label embedding based methods: Conditional Principal Label Space Transformation (CPLST) [7], Feature-aware Implicit Label space Encoding (FaIE) [18], Low rank Empirical risk minimization for Multi-Label Learning (LEML) [34], Sparse Local Embeddings for Extreme Multi-label Classification (SLEEC) [2], and the baseline method of partial binary relevance (PBR) [7]. Furthermore, we replace the linear regressors in CPLST and FAIE with DNN regressors, and name them as Deep CPLST and Deep FAIE, respectively. The work in [33] proposes a novel DNN architecture of Canonical-Correlated Autoencoder (C2AE), which can exploit label correlation effectively. Since some methods (e.g., C2AE) reported the results in terms of Micro-F1 and Macro-F1 [25], we also provides results of different approaches in terms of these two metrics for convenient comparison as shown in Table 3. According to the results, it is observed that the performance of our model is much better than the landmark selection methods [1,3] which separate landmark selection and prediction in 2-step manner. Moreover, our SPL-MLL performs superiorly against these label embedding methods.

Table 4: Ablation studies for our model on different setting on pascal VOC 2007.

| Methods | Ranking Loss ↓ | Hamming Loss ↓ | Average Precision ↑ | Micro-F1 ↑ | Macro-F1 ↑ |
|---|---|---|---|---|---|
| MLFE [38] | 0.232±0.013 | 0.162±0.012 | 0.565±0.022 | 0.436±0.026 | 0.357±0.011 |
| HNOML [35] | 0.227±0.012 | 0.123±0.008 | 0.593±0.023 | 0.443±0.024 | 0.368±0.019 |
| NN-embeddings | 0.324±0.016 | 0.266±0.011 | 0.431±0.013 | 0.308±0.011 | 0.287±0.016 |
| Ours(NN + separated) | 0.243±0.014 | 0.194±0.015 | 0.521±0.024 | 0.384±0.009 | 0.311±0.017 |
| Ours(joint + linear) | 0.192±0.011 | 0.095±0.012 | 0.608±0.021 | 0.516±0.025 | 0.422±0.024 |
| Ours(joint + NN) | **0.184±0.012** | **0.083±0.013** | **0.616±0.018** | **0.586±0.018** | **0.495±0.017** |

**Ablation Studies** To investigate the advantage of our model on jointly conducting landmark selection, landmark prediction and label recovery in a unified framework, we further conduct comparison and ablation experiments on pascal VOC 2007. Specifically, we conduct ablation studies for our model under the following settings: (1) NN-embeddings: the features are directly encoded by neural network for full label recovery without landmark selection and landmark prediction; (2) Ours (NN + separated): our model is still based on the landmark selection strategy, but separates the landmark selection and landmark prediction in the 2-step manner like the work [1]; (3) Ours (joint + linear): our model employs the linear projections instead of neural networks for feature embedding. To further validate the performance improvement from our model, we also report the results of the latest and most advanced approaches MLFE [38] and HNOML [35]. The comparision results are shown in Table 4, which validates the superiority of conducting landmark selection, landmark prediction and label recovery in a unified framework.

**Insight for selected landmarks** To investigate the improvement of SPL-MLL, we visualize the landmark selection matrix **B** on emotions and scene. As

illustrated in Fig. 4, the values in yellow on the diagonal are much larger than the values in other colors, where the corresponding labels are selected landmarks. For emotions, "Amazed" and "Quiet" are most likely to be landmark labels, and "Amazed" is often accompanied by "Angry" in music, "Quiet" tends to occur simultaneously with "Relaxing" or "Sad". Thus, we can utilize the selected landmark labels to recover other related labels effectively. Similar, for scene, the label "FallFoliage" and "Field" are most likely to be landmark labels.

As shown in Fig. 3, we count the number of those samples with or without "Angry" when having "Amazed", which is represented as "Amazed ± Angry", and similarly we obtain "Quiet ± Relaxing" and "Quiet ± Sad". According to Fig. 3, it is observed that when the "Amazed" ("Quiet") emotion occurs, the probability that "Angry" ("Relaxing" and "Sad") occur simultaneously is 56% (70% and 71%). This statistics further support the reasonability of the selected landmark labels.



Fig. 5: Example predictions on espgame and scene.

**Result visualization & convergence experiment** For intuitive analysis, Fig. 5 shows some representative examples from espgame and scene. The correctly predicted landmark labels from our model are in red, while the labels in green, gray and black indicate the successfully predicted, missed predicted and wrongly predicted labels. Generally, although multi-label classification is rather challenging especially for the large label set, our model achieves competitive results. We find that a few labels of some samples are not correctly predicted, and the possible reasons are as follows. First, a few labels on some samples do not obviously correlate with other labels, which makes it difficult to accurately recover given selected landmark labels. Second, a few labels for some samples are associated with very small parts in images, making it difficult to predict accurately even taking the feature of images into account in our model.

For example, the image labeled with "Kitchen" and "Dining" as landmark labels has the following labels predicted correctly: "Table", "Chair", "Room", "Light", "Restaurant" and "Door". However, there are labels: "Sun" and "Flower" failed to be predicted. The main reasons is that the label "Sun" and "Flower" may be not strongly correlated with the selected landmark labels in the dataset.

There are a few landmarks failed to be predicted for some samples, even though our model aims to select predictable landmark labels. For example, for the rightmost picture in the bottom of Fig. 5, we predict successfully "Sky" and "Man" as landmarks while not able to obtain the more critical landmark label "Parachuting", which leads to failure prediction for "Falling", "Fall", "Sky-diving", "Parachuting". It can be seen that the parachute is rather difficult to predict due to the strong illumination.
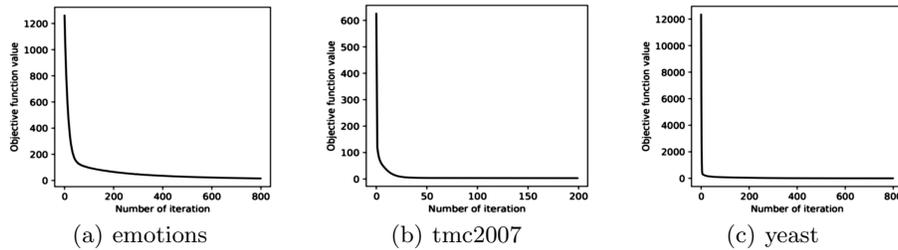


| (a) emotions | (b) tmc2007 | (c) yeast |

Fig. 6: Convergence experiment.

Fig. 6 gives the convergence experiments on emotion, yeast and tmc2007. Obviously, the results demonstrate that our method can converge within a small number of iterations.

## 5   Conclusions & Future Work

In this paper, we proposed a novel landmark-based multi-label classification algorithm, termed *SPL-MLL: Selecting Predictable Landmarks for Multi-Label Learning*. SPL-MLL jointly takes the representative and predictable properties for landmarks in a unified framework, avoiding separating landmark selection/prediction in the 2-step manner. Our key idea lies in selecting explicitly the landmarks which are both representative and predictable. The empirical experiments clearly demonstrate that our algorithm outperforms existing state-of-the-art methods. In the future, we will consider the end-to-end manner to extend our model for image annotation with large label set.

# References

1. Balasubramanian, K., Lebanon, G.: The landmark selection method for multiple output prediction. In: International Conference on Machine Learning (2012) 2, 3, 4, 5, 6, 10, 11, 12
2. Bhatia, K., Jain, H., Kar, P., Varma, M., Jain, P.: Sparse local embeddings for extreme multi-label classification. In: Advances in neural information processing systems. pp. 730–738 (2015) 10, 12
3. Bi, W., Kwok, J.: Efficient multi-label classification with many labels. In: International Conference on Machine Learning. pp. 405–413 (2013) 2, 3, 4, 10, 11, 12
4. Boutell, M.R., Luo, J., Shen, X., Brown, C.M.: Learning multi-label scene classification. Pattern recognition **37**(9), 1757–1771 (2004) 3, 7, 8, 9, 10
5. Boutsidis, C., Mahoney, M.W., Drineas, P.: An improved approximation algorithm for the column subset selection problem. In: Proceedings of the twentieth annual ACM-SIAM symposium on Discrete algorithms. pp. 968–977. SIAM (2009) 3
6. Charte, F., Rivera, A., del Jesus, M., Herrera, F.: Multilabel classification. problem analysis, metrics and techniques book repository 7, 8
7. Chen, Y.N., Lin, H.T.: Feature-aware label space dimension reduction for multi-label classification. In: Advances in Neural Information Processing Systems. pp. 1529–1537 (2012) 2, 3, 10, 12
8. Elisseeff, A., Weston, J.: A kernel method for multi-labelled classification. In: Advances in neural information processing systems. pp. 681–687 (2002) 3, 7, 8
9. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. International journal of computer vision **88**(2), 303–338 (2010) 7
10. Fürnkranz, J., Hüllermeier, E., Mencía, E.L., Brinker, K.: Multilabel classification via calibrated label ranking. Machine learning **73**(2), 133–153 (2008) 3, 9, 10
11. Ghamrawi, N., McCallum, A.: Collective multi-label classification. In: Proceedings of the 14th ACM international conference on Information and knowledge management. pp. 195–200. ACM (2005) 3
12. Hou, P., Geng, X., Zhang, M.L.: Multi-label manifold learning. In: Thirtieth AAAI Conference on Artificial Intelligence (2016) 1, 9, 10
13. Hsu, D.J., Kakade, S.M., Langford, J., Zhang, T.: Multi-label prediction via compressed sensing. In: Advances in neural information processing systems. pp. 772–780 (2009) 2, 3
14. Ji, S., Tang, L., Yu, S., Ye, J.: A shared-subspace learning framework for multi-label classification. ACM Transactions on Knowledge Discovery from Data (TKDD) **4**(2),  8 (2010) 3
15. Jia, X., Zheng, X., Li, W., Zhang, C., Li, Z.: Facial emotion distribution learning by exploiting low-rank label correlations locally. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 9841–9850 (2019) 3
16. Li, X., Guo, Y.: Multi-label classification with feature-aware non-linear label space transformation. In: Twenty-Fourth International Joint Conference on Artificial Intelligence (2015) 3
17. Lin, Z., Liu, R., Su, Z.: Linearized alternating direction method with adaptive penalty for low-rank representation. In: Advances in neural information processing systems. pp. 612–620 (2011) 7
18. Lin, Z., Ding, G., Hu, M., Wang, J.: Multi-label classification via feature-aware implicit label space encoding. In: International conference on machine learning. pp. 325–333 (2014) 3, 10, 12

19. Liu, J., Chang, W.C., Wu, Y., Yang, Y.: Deep learning for extreme multi-label text classification. In: Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 115–124. ACM (2017) 1

20. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. International journal of computer vision **60**(2), 91–110 (2004) 8

21. Read, J., Pfahringer, B., Holmes, G.: Multi-label classification using ensembles of pruned sets. In: 2008 Eighth IEEE International Conference on Data Mining. pp. 995–1000. IEEE (2008) 1, 9, 10

22. Read, J., Pfahringer, B., Holmes, G., Frank, E.: Classifier chains for multi-label classification. Machine learning **85**(3), 333 (2011) 1, 3, 9, 10

23. Ren, T., Jia, X., Li, W., Zhao, S.: Label distribution learning with label correlations via low-rank approximation. In: Proceedings of the 28th International Joint Conference on Artificial Intelligence. pp. 3325–3331. AAAI Press (2019) 3

24. Tai, F., Lin, H.T.: Multilabel classification with principal label space transformation. Neural Computation **24**(9), 2508–2542 (2012) 2, 3

25. Tang, L., Rajan, S., Narayanan, V.K.: Large scale multi-label classification via metalabeler. In: Proceedings of the 18th international conference on World wide web. pp. 211–220. ACM (2009) 12

26. Trohidis, K., Tsoumakas, G., Kalliris, G., Vlahavas, I.P.: Multi-label classification of music into emotions. In: ISMIR. vol. 8, pp. 325–330 (2008) 7, 8

27. Tsoumakas, G., Katakis, I.: Multi-label classification: An overview. International Journal of Data Warehousing and Mining (IJDWM) **3**(3), 1–13 (2007) 1, 9, 10

28. Tsoumakas, G., Katakis, I., Vlahavas, I.: Random k-labelsets for multilabel classification. IEEE Transactions on Knowledge and Data Engineering **23**(7), 1079–1089 (2011) 3, 9, 10

29. Von Ahn, L., Dabbish, L.: Labeling images with a computer game. In: Proceedings of the SIGCHI conference on Human factors in computing systems. pp. 319–326. ACM (2004) 7

30. Wu, B., Chen, W., Sun, P., Liu, W., Ghanem, B., Lyu, S.: Tagging like humans: Diverse and distinct image annotation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7967–7975 (2018) 3

31. Wu, B., Jia, F., Liu, W., Ghanem, B.: Diverse image annotation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2559–2567 (2017) 1

32. Wu, B., Jia, F., Liu, W., Ghanem, B., Lyu, S.: Multi-label learning with missing labels using mixed dependency graphs. International Journal of Computer Vision **126**(8), 875–896 (2018) 1

33. Yeh, C.K., Wu, W.C., Ko, W.J., Wang, Y.C.F.: Learning deep latent space for multi-label classification. In: Thirty-First AAAI Conference on Artificial Intelligence (2017) 3, 10, 12

34. Yu, H.F., Jain, P., Kar, P., Dhillon, I.: Large-scale multi-label learning with missing labels. In: International conference on machine learning. pp. 593–601 (2014) 3, 10, 12

35. Zhang, C., Yu, Z., Fu, H., Zhu, P., Chen, L., Hu, Q.: Hybrid noise-oriented multi-label learning. IEEE transactions on cybernetics (2019) 9, 10, 12

36. Zhang, M.L., Zhou, Z.H.: Ml-knn: A lazy learning approach to multi-label learning. Pattern recognition **40**(7), 2038–2048 (2007) 3, 9, 10

37. Zhang, M.L., Zhou, Z.H.: A review on multi-label learning algorithms. IEEE transactions on knowledge and data engineering **26**(8), 1819–1837 (2014) 1, 3, 8

38. Zhang, Q.W., Zhong, Y., Zhang, M.L.: Feature-induced labeling information enrichment for multi-label learning. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018) 9, 10, 12
39. Zhang, Y., Schneider, J.: Maximum margin output coding. arXiv preprint arXiv:1206.6478 (2012) 2, 3
40. Zhou, T., Tao, D., Wu, X.: Compressed labeling on distilled labelsets for multi-label learning. Machine Learning **88**(1-2), 69–126 (2012) 2, 3