

# Compositional Metric Learning for Multi-Label Classification

Yan-Ping Sun<sup>1,2,3</sup>, Min-Ling Zhang (✉)<sup>1,2,3</sup>

<sup>1</sup> School of Computer Science and Engineering, Southeast University, Nanjing 210096, China

<sup>2</sup> Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, China  
<sup>3</sup> Collaborative Innovation Center for Wireless Communications Technology, Nanjing 211100, China

© Higher Education Press and Springer-Verlag Berlin Heidelberg 2019

**Abstract** Multi-label classification aims to assign a set of proper labels for each instance, where distance metric learning can help improve the generalization ability of instance-based multi-label classification models. Existing multi-label metric learning techniques work by utilizing pairwise constraints to enforce that examples with similar label assignments should have close distance in the embedded feature space. In this paper, a novel distance metric learning approach for multi-label classification is proposed by modeling structural interactions between instance space and label space. On one hand, compositional distance metric is employed which adopts the representation of a weighted sum of rank-1 PSD matrices based on component bases. On the other hand, compositional weights are optimized by exploiting triplet similarity constraints derived from both instance and label spaces. Due to the compositional nature of employed distance metric, the resulting problem admits quadratic programming formulation with linear optimization complexity w.r.t. the number of training examples. We also derive the generalization bound for the proposed approach based on algorithmic robustness analysis of the compositional metric. Extensive experiments on sixteen benchmark data sets clearly validate the usefulness of compositional metric in yielding effective distance metric for multi-label classification.

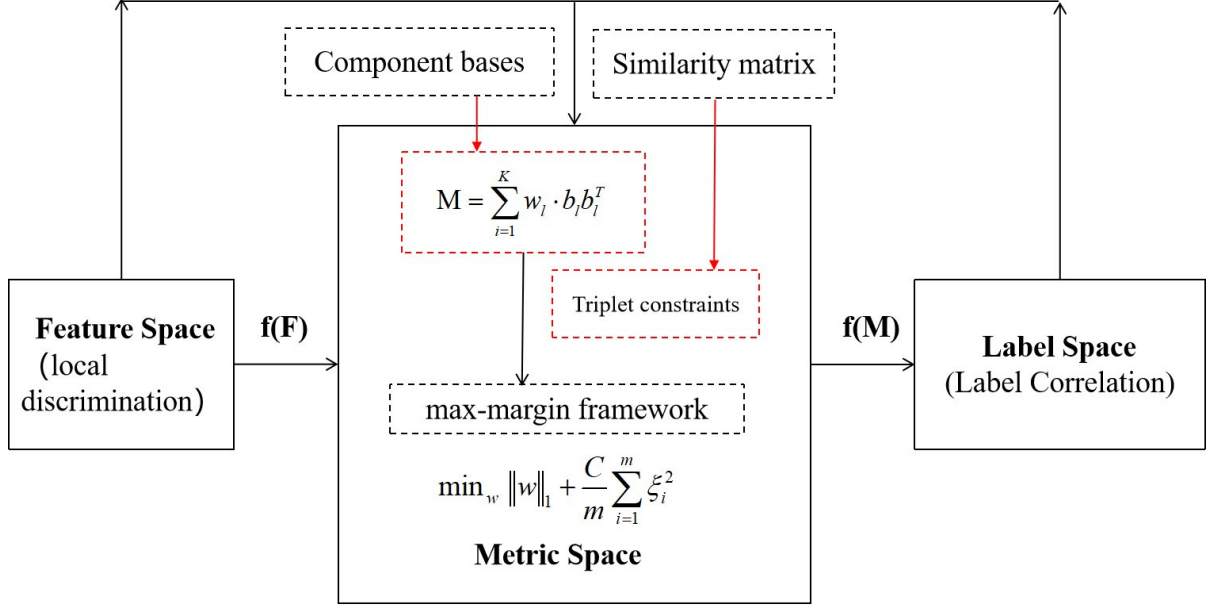
**Keywords** machine learning, multi-label learning, metric learning, compositional metric, positive semidefinite matrix decomposition

## 1 Introduction

In multi-label classification, each instance is associated with multiple class labels simultaneously and the task is to learn a predictive model mapping from instance to the set of proper labels [1, 2]. In recent years, multi-label classification techniques have been widely applied to learn from real-world objects with rich semantics [3–7].

Distance metric learning serves as a popular strategy to facilitate supervised learning, where a positive semi-definite (PSD) matrix  $\mathbf{M} \geq 0$  is usually learned to parameterize the distance in embedded feature space [8, 9]. Some recent attempts show promising results of learning distance metric to build multi-label classification models with stronger generalization performance [10, 11]. Specifically, given training examples  $(\mathbf{x}_i, \mathbf{y}_i)$  and  $(\mathbf{x}_j, \mathbf{y}_j)$ , their distance in the embedded feature space  $d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^{\top} \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j)}$  should move closer if  $\mathbf{y}_i$  is similar to  $\mathbf{y}_j$  in the label space. This strategy can be instantiated in different ways such as large margin output coding [10, 12, 13] or pairwise similarity preservation [11, 14, 15].

In this paper, a novel multi-label distance metric learning approach named COMMU, i.e. *COMpositional METric for Multi-label classification*, is proposed. Compared to existing approaches for multi-label metric learning, COMMU considers a more advanced strategy by modeling structural interactions between instance space and label space. In Figure 1, the general framework of COMMU for multi-label distance metric learning is illustrated. Specifically, the multi-label distance metric is assumed to adopt the compositional representation



**Fig. 1** The multi-label distance metric learning framework of COMMU. The original feature space is mapped into the distance metric feature space based on the compositional distance metric, based on which the prediction on multi-label output space will be made. Specifically, each component of the distance metric is generated by employing triplet constraints derived from similarity relationships in both instance and label spaces.

with a weighted sum of rank-1 PSD matrices. Here, the rank-1 PSD matrix corresponds to the outer product of component bases generated by encoding discriminative information of class labels. Furthermore, the weights forming the compositional distance metric are optimized by exploiting triplet constraints derived from similarity relationships in both instance and label spaces. Experimental studies across sixteen benchmark multi-label data sets show that COMMU is capable of significantly improving the generalization performance of instance-based multi-label classification models with the learned compositional distance metric.

The rest of this paper is organized as follows. Section 2 presents technical details of the proposed approach. Section 3 provides the corresponding theoretical analysis. Section 4 reports experimental results of comparative studies. Section 5 briefly discusses related works. Finally, Section 6 concludes this paper.

## 2 The COMMU Approach

Formally, let  $\mathcal{X} = \mathbb{R}^d$  be the instance space and  $\mathcal{Y} = \{\lambda_1, \lambda_2, \dots, \lambda_q\}$  be the label space with  $q$  class labels. Multi-label classification aims to learn a predictive function  $h: \mathcal{X} \mapsto 2^{\mathcal{Y}}$  from the training set  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq m\}$ , where  $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iq})^T$  is the labeling vector associated with  $\mathbf{x}_i$  such that  $y_{il} = 1$  if  $\lambda_l$  is a relevant label for  $\mathbf{x}_i$  and  $y_{il} = 0$

otherwise.

To model the distance metric  $\mathbf{M}$  with enriched structural information, COMMU chooses to adopt the compositional representation with a weighted sum of  $K$  rank-1 PSD matrices [16, 17]:

$$\mathbf{M} = \sum_{l=1}^K w_l \cdot \mathbf{b}_l \mathbf{b}_l^T \quad (1)$$

Here,  $\mathbf{b}_l \in \mathbb{R}^d$  is the  $d$ -dimensional component base and  $w_l \geq 0$  is the corresponding nonnegative compositional weight. In this way, one can simplify the parameterization complexity of the distance metric from  $O(d^2)$  to  $O(K)$  with  $\mathbf{w} = [w_1, w_2, \dots, w_K]^T$ . More importantly, the compositional decomposition naturally enables the encoding of discriminative information into the distance metric. Specifically, COMMU generates one component base for each class label in the label space (i.e.  $K = q$ ).

For the  $l$ -th class label  $\lambda_l \in \mathcal{Y}$  ( $1 \leq l \leq q$ ), COMMU considers the difference between the mean of positive examples and negative examples w.r.t.  $\lambda_l$ :

$$\mathbf{b}_l = \frac{\sum_{\mathbf{x}_i \in \mathcal{P}_l} \mathbf{u}}{|\mathcal{P}_l|} - \frac{\sum_{\mathbf{x}_i \in \mathcal{N}_l} \mathbf{v}}{|\mathcal{N}_l|} \quad (2)$$

Here,  $\mathcal{P}_l = \{\mathbf{x}_i \mid y_{il} = 1, (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}\}$  and  $\mathcal{N}_l = \{\mathbf{x}_i \mid y_{il} = 0, (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}\}$  correspond to the set of positive examples and negative examples w.r.t.  $\lambda_l$  respectively. Conceptually, the statistic in Eq.(2) is used to reflect holistic labeling distribution

of class label, which has shown to be beneficial for encoding discriminative information in the feature space [18–20].

To optimize the parameters  $\mathbf{w}$  for distance metric  $\mathbf{M}$ , a set of constraints are specified to characterize the properties which  $\mathbf{M}$  are expected to possess. Given the multi-label training example  $(\mathbf{x}_i, \mathbf{y}_i)$  and other two reference examples  $\{(\mathbf{x}_j, \mathbf{y}_j), (\mathbf{x}_k, \mathbf{y}_k)\}$ , it is desirable that  $d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j)$  should be smaller than  $d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_k)$  if  $\mathbf{x}_i$  is semantically more similar to  $\mathbf{x}_j$  than  $\mathbf{x}_k$ . Under traditional single-label scenario, the semantic similarity can be easily measured by considering whether two examples have the same class label [8, 21]. However, under multi-label scenario, it is impractical to measure semantic similarity by considering exact labeling equivalence due to the combinatorial nature of multiple class labels. For COMMU, the semantic similarity matrix  $\mathbf{S} = [s_{ij}]_{m \times m}$  is calculated by synergizing discriminative information from both input space and label space:

$$s_{ij} = \mathbf{y}_i^{\top} \mathbf{G} \mathbf{y}_j \quad (3)$$

where  $\mathbf{G} = (\alpha \mathbf{A} + (1 - \alpha) \mathbf{C})$

$$\mathbf{A} = [a_{lh}]_{q \times q} \quad \text{with} \quad a_{lh} = \frac{\sum_{i=1}^m y_{il} \cdot y_{ih}}{\sum_{i=1}^m y_{il}}$$

$$\mathbf{C} = [c_{lh}]_{q \times q} \quad \text{with} \quad c_{lh} = \frac{\mathbf{b}_l^{\top} \mathbf{b}_h}{\|\mathbf{b}_l\| \cdot \|\mathbf{b}_h\|}$$

Here,  $a_{lh}$  corresponds to the fraction of examples with label  $y_l$  which also have label  $y_h$ . It is noteworthy that  $a_{lh} = a_{hl}$  does not necessarily hold here to reflect the fact that correlations among class labels are usually *asymmetric* [22, 23]. Furthermore,  $c_{lh}$  corresponds to the cosine similarity between compositional bases. The coefficient  $\alpha$  balances relative contributions from label space (i.e.  $\mathbf{A}$ ) and instance space (i.e.  $\mathbf{C}$ ) in calculating the semantic similarity.

Thereafter, the set of “similar” and “dissimilar” examples for training instance  $\mathbf{x}_i$  are determined as:

$$\begin{aligned} \mathcal{Z}_i &= \{\mathbf{x}_j \mid s_{ij} \geq \theta, j \neq i, 1 \leq j \leq m\} \\ \tilde{\mathcal{Z}}_i &= \{\mathbf{x}_k \mid s_{ik} < \theta, k \neq i, 1 \leq k \leq m\} \end{aligned} \quad (4)$$

Here,  $\theta$  is used as the thresholding parameter for measuring semantic similarity. Accordingly, the following set of triplets are generated by utilizing subset  $\mathcal{K}_i \subseteq \mathcal{Z}_i$  ( $\tilde{\mathcal{K}}_i \subseteq \tilde{\mathcal{Z}}_i$ ) which consists of top  $k$  ( $\tilde{k}$ ) instances with highest semantic similarity in  $\mathcal{Z}_i$  ( $\tilde{\mathcal{Z}}_i$ ):

$$\mathcal{R} = \{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \mid 1 \leq i \leq m, \mathbf{x}_j \in \mathcal{K}_i, \mathbf{x}_k \in \tilde{\mathcal{K}}_i\} \quad (5)$$

Here,  $\mathcal{R}$  contains a total of  $m \cdot k \cdot \tilde{k}$  triplets. Based on Eq.(5), COMMU learns the compositional distance metric by solving

**Table 1** The pseudo-code of COMMU.

**Inputs:**

- $\mathcal{D}$ : multi-label training set  $\{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq m\}$
- $\alpha$ : balancing parameter in Eq.(3) with  $\alpha \in (0, 1)$
- $C$ : cost parameter in Eq.(6) with  $C > 0$
- $\theta$ : thresholding parameter in Eq.(4)

**Outputs:**

- $\mathbf{w}$ : compositional weight vector for the distance metric

**Process:**

- 1: **for**  $l = 1$  to  $q$  **do**
- 2:   Generate component base  $\mathbf{b}_l$  according to Eq.(1);
- 3: **end for**
- 4: Calculate the similarity matrix  $\mathbf{S}$  according to Eq.(3);
- 5: Form the set of triplets  $\mathcal{R}$  according to Eq.(5);
- 6: Initialize FISTA procedure with  $\mathbf{w}_0 = \mathbf{w}_1 = \frac{1}{q} \cdot \mathbf{1}_{q \times 1}$ ,  $\tau_0 = \tau_1 = 0.01mC$ ,  $\eta = 0.4$ , and  $t_0 = t_1 = 1$ ;
- 7: Set  $r = 1$  and  $\tilde{\mathbf{w}}_1 = \mathbf{w}_1$ ;
- 8: **repeat**
- 9:   Set  $L = \eta \cdot \tau_r$ ;
- 10: **repeat**
- 11:   Calculate  $\mathbf{a}^* = \tilde{\mathbf{w}}_r - \frac{1}{L}(\nabla f(\tilde{\mathbf{w}}_r) + \mathbf{1}_{q \times 1})$ ;
- 12:   **if**  $F(\Pi_+(\mathbf{a}^*)) \leq Q_L(\Pi_+(\mathbf{a}^*), \tilde{\mathbf{w}}_r)$  **then**
- 13:      $\tau_{r+1} = L$ ;
- 14:     **go to step 19**;
- 15:   **else**
- 16:      $L = \frac{1}{\eta}L$ ;
- 17:   **end if**
- 18: **until** false
- 19:    $\mathbf{w}_{r+1} = \Pi_+(\mathbf{a}^*)$ ;
- 20:    $t_{r+1} = \frac{1 + \sqrt{1 + 4t_r^2}}{2}$ ;
- 21:    $r = r + 1$ ;
- 22:    $\tilde{\mathbf{w}}_r = \mathbf{w}_{r-1} + \frac{t_{r-1}-1}{t_r} \cdot (\mathbf{w}_{r-1} - \mathbf{w}_{r-2})$ ;
- 23: **until** convergence
- 24: **Return**  $\mathbf{w} = \mathbf{w}_r$ ;

the following optimization problem with triplet constraints:

$$\min_{\mathbf{w}} \|\mathbf{w}\|_1 + \frac{C}{m} \sum_{i=1}^m \xi_i^2 \quad \text{s.t. :} \quad (6)$$

$$d_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_k) - d_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_j) \geq \Delta(\mathbf{y}_j, \mathbf{y}_k) - \xi_i$$

$$(\forall (\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in \mathcal{R})$$

$$\xi_i \geq 0, \quad w_i \geq 0 \quad (1 \leq i \leq m)$$

Here,  $d_{\mathbf{M}}^2(\mathbf{x}, \mathbf{x}') = \sum_{l=1}^q w_l (\mathbf{x} - \mathbf{x}')^{\top} \mathbf{b}_l \mathbf{b}_l^{\top} (\mathbf{x} - \mathbf{x}')$  corresponds to the distance between two instances in the embedded feature space and  $\Delta(\mathbf{y}_j, \mathbf{y}_k) = \mathbf{y}_j^{\top} (\mathbf{1}_{q \times q} - \mathbf{G}) \mathbf{y}_k$  corresponds to the dissimilarity between two labeling vectors. Accordingly, the slack variable  $\xi_i$  corresponds to:

$$\xi_i = \max \left( 0, \max_{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in \mathcal{R}} (\Delta(\mathbf{y}_j, \mathbf{y}_k) - (d_{\mathbf{M}}^2(\mathbf{x}_{i,k}) - d_{\mathbf{M}}^2(\mathbf{x}_{i,j})) \right) \quad (7)$$

**Table 2** Characteristics of the benchmark multi-label data sets.

Data set	$ \mathcal{D} $	$\dim(\mathcal{D})$	$CL(\mathcal{D})$	$LCard(\mathcal{D})$	$LDen(\mathcal{D})$	$DL(\mathcal{D})$	$PDL(\mathcal{D})$	Domain
genbase	662	1186	27	1.252	0.046	32	0.048	biology
Society	2000	636	27	1.692	0.063	329	0.165	text
Social	2000	1047	39	1.283	0.033	137	0.069	text
Reference	2000	793	33	1.169	0.035	132	0.066	text
Health	2000	612	32	1.662	0.052	164	0.082	text
Education	2000	550	33	1.461	0.044	200	0.1	text
Computers	2000	681	33	1.508	0.046	148	0.074	text
Business	2000	438	30	1.588	0.053	96	0.048	text
Arts	2000	462	26	1.636	0.063	254	0.127	text
yeast	2417	103	14	4.237	0.303	198	0.082	biology
corel5k	5000	499	374	3.522	0.009	3175	0.635	images
rcv1-subset1	6000	944	101	2.88	0.029	1028	0.171	text
corel16k001	13766	500	153	2.859	0.019	4937	0.359	images
eurlex-dc	19348	100	412	1.292	0.003	1615	0.083	text
eurlex-sm	19348	100	201	2.213	0.011	2504	0.129	text
eurlex	19314	1854	815	4.273	0.0052	14763	0.764	text

Therefore, the solution to Eq.(6) can be obtained by optimizing the following equivalent problem:

$$\min_{\mathbf{w}} F(\mathbf{w}) \equiv f(\mathbf{w}) + g(\mathbf{w}) \quad (8)$$

Here,  $f(\mathbf{w}) = \frac{c}{m} \sum_{i=1}^m \xi_i^2$  whose gradient  $\nabla f$  is Lipschitz continuous w.r.t.  $\mathbf{w}$  [10, 24] and  $g(\mathbf{w}) = \|\mathbf{w}\|_1$  is convex. For optimization problem admitting such decomposition, its solution can be obtained by employing the *FISTA (Fast Iterative Shrinkage-Thresholding Algorithm)* procedure [25, 26]. Specifically, given the current solution  $\mathbf{w}$ , the solution at next iteration is solved by minimizing the following quadratic programming problem:

$$Q_L(\mathbf{a}, \mathbf{w}) = f(\mathbf{w}) + \langle \nabla f(\mathbf{w}), \mathbf{a} - \mathbf{w} \rangle + \frac{L}{2} \|\mathbf{a} - \mathbf{w}\|^2 + g(\mathbf{a}) \quad (9)$$

where  $L > 0$  is the Lipschitz constant for  $f(\mathbf{w})$ . By setting the gradient of Eq.(9) to zero, one can obtain the minimizer  $\mathbf{a}^* = \mathbf{w} - \frac{1}{L}(\nabla f(\mathbf{w}) + \mathbf{1}_{q \times 1})$ . To ensure nonnegativity of compositional weights for the distance metric  $\mathbf{M}$ , the iterative solution  $\mathbf{a}^*$  will be mapped to  $\Pi_+(\mathbf{a}^*)$  by setting negative elements in  $\mathbf{a}^*$  to zero.

Table 1 summarizes the complete procedure of COMMU. Firstly, a set of compositional bases are generated by discriminative information encoding (steps 1 to 3). After that, the set of triplet constraints are specified by considering semantic similarity among training examples (steps 4-5). Thirdly, the compositional weights are learned by invoking the FISTA

iterative optimization procedure (steps 6-24).<sup>1)</sup>

### 3 Theoretical Analysis

In this section, we provide a theoretical analysis of our approach in the form of a generalization bound based on algorithmic robustness analysis for metric learning [27].

Given a multi-label dataset  $\mathcal{S} = \{\mathbf{z} = (\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  drawn i.i.d. from a distribution  $P$  over the labelled space  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ , where the label vector  $\mathbf{y}_i$  simultaneously contains multiple labels. Assume that  $\|\mathbf{x}\| \leq R$  (for some convenient norm),  $\forall \mathbf{x} \in \mathcal{X}$ . Different from the single-label setting, COMMU defines the multi-label semantic similarity matrix to construct the triplet  $(\mathbf{z}, \mathbf{z}', \mathbf{z}'')$ , where  $\mathbf{y}$  is similar to  $\mathbf{y}'$  and dissimilar to  $\mathbf{y}''$ . Let  $\mathcal{S}_R$  be the set of all admissible triplets built from  $\mathcal{S}$  and  $L(\mathbf{w}, \mathbf{z}, \mathbf{z}', \mathbf{z}'') = [\Delta(\mathbf{y}', \mathbf{y}'') + d_w(\mathbf{x}, \mathbf{x}') - d_w(\mathbf{x}, \mathbf{x}'')]_+$  denote the multi-label triple loss function in Eq.(6), which is uniformly upper-bounded by a constant  $U$ .

The empirical loss  $R_{emp}^{S_R}(\mathbf{w})$  of  $\mathbf{w}$  on  $\mathcal{S}_R$  is defined as

$$R_{emp}^{S_R}(\mathbf{w}) = \frac{1}{|\mathcal{S}_R|} \sum_{(\mathbf{z}, \mathbf{z}', \mathbf{z}'') \in \mathcal{S}_R} L(\mathbf{w}, \mathbf{z}, \mathbf{z}', \mathbf{z}''),$$

and its expected loss  $R(\mathbf{w})$  over distribution  $P$  as

$$R(\mathbf{w}) = \mathbb{E}_{\mathbf{z}, \mathbf{z}', \mathbf{z}'' \sim P} L(\mathbf{w}, \mathbf{z}, \mathbf{z}', \mathbf{z}'').$$

<sup>1)</sup> The FISTA procedure terminates when the value of the objective function in Eq.(8) does not significantly decrease for two consecutive solutions  $\mathbf{w}_r$  and  $\mathbf{w}_{r+1}$ , i.e.  $F(\mathbf{w}_r) - F(\mathbf{w}_{r+1}) \leq \epsilon \cdot F(\mathbf{w}_r)$  with  $\epsilon = 0.001$ .

**Table 3** Predictive performance (mean  $\pm$  std. deviation) of each distance metric learning approach in terms of *ranking loss*.

Data set	kNN-				MLKNN-			
	COMMU	LM	NJE	ORIGINAL	COMMU	LM	NJE	ORIGINAL
genbase	0.008 $\pm$ 0.008	0.146 $\pm$ 0.071	<b>0.006<math>\pm</math>0.007</b>	0.007 $\pm$ 0.008	<b>0.006<math>\pm</math>0.007</b>	0.009 $\pm$ 0.007	0.008 $\pm$ 0.009	0.007 $\pm$ 0.007
Society	0.261 $\pm$ 0.017	0.306 $\pm$ 0.019	<b>0.215<math>\pm</math>0.021</b>	0.263 $\pm$ 0.02	0.142 $\pm$ 0.007	0.164 $\pm$ 0.013	0.202 $\pm$ 0.016	<b>0.138<math>\pm</math>0.006</b>
Social	0.165 $\pm$ 0.016	0.260 $\pm$ 0.023	<b>0.137<math>\pm</math>0.023</b>	0.160 $\pm$ 0.02	<b>0.065<math>\pm</math>0.008</b>	0.084 $\pm$ 0.007	0.099 $\pm$ 0.013	0.065 $\pm$ 0.007
Reference	0.182 $\pm$ 0.032	0.303 $\pm$ 0.034	<b>0.143<math>\pm</math>0.017</b>	0.246 $\pm$ 0.032	<b>0.063<math>\pm</math>0.009</b>	0.103 $\pm$ 0.010	0.093 $\pm$ 0.010	0.083 $\pm$ 0.011
Health	0.154 $\pm$ 0.031	0.184 $\pm$ 0.036	<b>0.115<math>\pm</math>0.015</b>	0.199 $\pm$ 0.03	<b>0.057<math>\pm</math>0.010</b>	0.062 $\pm$ 0.010	0.075 $\pm$ 0.010	0.063 $\pm$ 0.009
Education	0.210 $\pm$ 0.023	0.299 $\pm$ 0.02	<b>0.168<math>\pm</math>0.020</b>	0.209 $\pm$ 0.024	<b>0.087<math>\pm</math>0.006</b>	0.105 $\pm$ 0.011	0.131 $\pm$ 0.014	0.087 $\pm$ 0.006
Computers	0.186 $\pm$ 0.029	0.285 $\pm$ 0.017	<b>0.132<math>\pm</math>0.013</b>	0.192 $\pm$ 0.041	<b>0.082<math>\pm</math>0.008</b>	0.093 $\pm$ 0.011	0.114 $\pm$ 0.011	0.082 $\pm$ 0.006
Business	<b>0.090<math>\pm</math>0.012</b>	0.122 $\pm$ 0.015	0.089 $\pm$ 0.015	0.092 $\pm$ 0.012	<b>0.037<math>\pm</math>0.005</b>	0.043 $\pm$ 0.007	0.075 $\pm$ 0.013	0.038 $\pm$ 0.006
Arts	0.270 $\pm$ 0.028	0.315 $\pm$ 0.026	<b>0.202<math>\pm</math>0.025</b>	0.302 $\pm$ 0.031	<b>0.149<math>\pm</math>0.013</b>	0.168 $\pm$ 0.012	0.181 $\pm$ 0.023	0.153 $\pm$ 0.015
yeast	0.197 $\pm$ 0.006	0.322 $\pm$ 0.015	<b>0.188<math>\pm</math>0.014</b>	0.195 $\pm$ 0.009	0.176 $\pm$ 0.008	0.182 $\pm$ 0.009	0.197 $\pm$ 0.016	<b>0.175<math>\pm</math>0.008</b>
core15k	<b>0.456<math>\pm</math>0.016</b>	0.675 $\pm$ 0.011	-	0.578 $\pm$ 0.027	<b>0.120<math>\pm</math>0.006</b>	0.129 $\pm$ 0.006	-	0.132 $\pm$ 0.006
rcv1-subset1	<b>0.191<math>\pm</math>0.030</b>	0.307 $\pm$ 0.013	-	0.226 $\pm$ 0.010	<b>0.073<math>\pm</math>0.008</b>	0.098 $\pm$ 0.007	-	0.080 $\pm$ 0.004
core116k	<b>0.495<math>\pm</math>0.008</b>	0.680 $\pm$ 0.004	-	0.537 $\pm$ 0.017	<b>0.169<math>\pm</math>0.002</b>	0.173 $\pm$ 0.002	-	0.175 $\pm$ 0.001
eurlex-dc	<b>0.376<math>\pm</math>0.034</b>	0.637 $\pm$ 0.023	-	0.376 $\pm$ 0.035	<b>0.094<math>\pm</math>0.009</b>	0.125 $\pm$ 0.008	-	0.094 $\pm$ 0.009
eurlex-sm	<b>0.191<math>\pm</math>0.003</b>	0.402 $\pm$ 0.006	-	0.191 $\pm$ 0.003	<b>0.051<math>\pm</math>0.001</b>	0.073 $\pm$ 0.001	-	0.051 $\pm$ 0.001
eurlex	0.923 $\pm$ 0.000	<b>0.922<math>\pm</math>0.000</b>	-	0.985 $\pm$ 0.000	0.320 $\pm$ 0.000	0.326 $\pm$ 0.000	-	<b>0.316<math>\pm</math>0.000</b>

**Table 4** Predictive performance (mean  $\pm$  std. deviation) of each distance metric learning approach in terms of *coverage*.

Data set	kNN-				MLKNN-			
	COMMU	LM	NJE	ORIGINAL	COMMU	LM	NJE	ORIGINAL
genbase	0.019 $\pm$ 0.011	0.106 $\pm$ 0.053	0.020 $\pm$ 0.017	<b>0.014<math>\pm</math>0.005</b>	<b>0.021<math>\pm</math>0.012</b>	0.025 $\pm$ 0.011	0.230 $\pm$ 0.019	0.021 $\pm$ 0.013
Society	<b>0.272<math>\pm</math>0.025</b>	0.284 $\pm$ 0.008	0.307 $\pm$ 0.030	0.277 $\pm$ 0.001	0.208 $\pm$ 0.015	0.228 $\pm$ 0.017	0.289 $\pm$ 0.029	<b>0.203<math>\pm</math>0.001</b>
Social	<b>0.120<math>\pm</math>0.015</b>	0.144 $\pm$ 0.053	0.175 $\pm$ 0.029	0.120 $\pm$ 0.001	<b>0.087<math>\pm</math>0.013</b>	0.107 $\pm$ 0.012	0.129 $\pm$ 0.017	0.088 $\pm$ 0.001
Reference	0.153 $\pm$ 0.014	<b>0.138<math>\pm</math>0.011</b>	0.165 $\pm$ 0.017	0.150 $\pm$ 0.001	0.104 $\pm$ 0.011	0.117 $\pm$ 0.012	0.108 $\pm$ 0.009	<b>0.097<math>\pm</math>0.001</b>
Health	0.142 $\pm$ 0.018	<b>0.138<math>\pm</math>0.011</b>	0.198 $\pm$ 0.022	0.169 $\pm$ 0.001	<b>0.098<math>\pm</math>0.013</b>	0.106 $\pm$ 0.012	0.131 $\pm$ 0.015	0.104 $\pm$ 0.001
Education	<b>0.167<math>\pm</math>0.011</b>	0.192 $\pm$ 0.016	0.228 $\pm$ 0.026	0.168 $\pm$ 0.001	<b>0.116<math>\pm</math>0.009</b>	0.135 $\pm$ 0.013	0.173 $\pm$ 0.013	0.116 $\pm$ 0.001
Computers	<b>0.153<math>\pm</math>0.019</b>	0.167 $\pm$ 0.003	0.182 $\pm$ 0.021	0.160 $\pm$ 0.001	<b>0.118<math>\pm</math>0.013</b>	0.131 $\pm$ 0.016	0.156 $\pm$ 0.017	0.118 $\pm$ 0.001
Business	<b>0.090<math>\pm</math>0.012</b>	0.098 $\pm$ 0.013	0.149 $\pm$ 0.016	0.095 $\pm$ 0.001	<b>0.071<math>\pm</math>0.007</b>	0.080 $\pm$ 0.010	0.130 $\pm$ 0.020	0.073 $\pm$ 0.001
Arts	<b>0.281<math>\pm</math>0.028</b>	0.295 $\pm$ 0.004	0.287 $\pm$ 0.031	0.304 $\pm$ 0.001	<b>0.209<math>\pm</math>0.020</b>	0.226 $\pm$ 0.016	0.253 $\pm$ 0.032	0.213 $\pm$ 0.001
yeast	<b>0.462<math>\pm</math>0.010</b>	0.560 $\pm$ 0.011	0.470 $\pm$ 0.020	0.473 $\pm$ 0.011	0.456 $\pm$ 0.010	0.469 $\pm$ 0.007	0.483 $\pm$ 0.022	<b>0.454<math>\pm</math>0.012</b>
core15k	<b>0.359<math>\pm</math>0.017</b>	0.792 $\pm$ 0.011	-	0.739 $\pm$ 0.019	<b>0.280<math>\pm</math>0.013</b>	0.294 $\pm$ 0.012	-	0.302 $\pm$ 0.014
rcv1-subset1	<b>0.244<math>\pm</math>0.023</b>	0.290 $\pm$ 0.016	-	0.267 $\pm$ 0.011	<b>0.165<math>\pm</math>0.016</b>	0.202 $\pm$ 0.012	-	0.178 $\pm$ 0.010
core116k	<b>0.511<math>\pm</math>0.006</b>	0.586 $\pm$ 0.003	-	0.548 $\pm$ 0.009	<b>0.328<math>\pm</math>0.003</b>	0.335 $\pm$ 0.003	-	0.339 $\pm$ 0.002
eurlex-dc	<b>0.276<math>\pm</math>0.026</b>	0.449 $\pm$ 0.013	-	0.276 $\pm$ 0.026	<b>0.114<math>\pm</math>0.010</b>	0.150 $\pm$ 0.009	-	0.114 $\pm$ 0.010
eurlex-sm	<b>0.233<math>\pm</math>0.003</b>	0.429 $\pm$ 0.004	-	0.233 $\pm$ 0.004	<b>0.092<math>\pm</math>0.001</b>	0.126 $\pm$ 0.002	-	0.092 $\pm$ 0.001
eurlex	0.642 $\pm$ 0.000	<b>0.640<math>\pm</math>0.000</b>	-	0.640 $\pm$ 0.000	<b>0.600<math>\pm</math>0.000</b>	0.609 $\pm$ 0.000	-	0.600 $\pm$ 0.000

The goal of the theoretical analysis is to bound the deviation between  $R(w)$  and  $R_{emp}^{S_R}(w)$ , where  $w$  is the metric coefficient to learn.

**Theorem 1.** Let  $w^*$  be the optimal solution to COMMU with  $K$  basis elements,  $C > 0$  and the triplet  $S_R$  constructed from  $S = \{z = (x_i, y_i)\}_{i=1}^n$ . Let  $K^* \leq K$  be the number of nonzero entries in  $w^*$ . Assume the norm of any instance bounded by some constant  $R$  and the loss  $L$  uniformly upper-bounded by some constant  $U$ . Then for any  $\delta > 0$ , with probability at least  $1 - \delta$  we have:

$$\left| R(w^*) - R_{emp}^{S_R}(w^*) \right| \leq 16\gamma R K^* \Delta_0 C + \theta + 3U \sqrt{\frac{2^{q+1} \mathcal{N}(\gamma, \mathcal{X}, \|\cdot\|_1) \ln 2 + 2 \ln \frac{1}{\delta}}{n}},$$

where  $2^q \mathcal{N}(\gamma, \mathcal{X}, \|\cdot\|_1)$  is the size of an  $\gamma$ -cover of  $\mathcal{Z}$ ,  $\Delta_0$  is the maximum distance of the dissimilarity set and  $\theta$  is the threshold of the triplet construction in Eq.(4). This bound has a standard  $O(1/\sqrt{n})$  asymptotic convergence rate.<sup>2)</sup> The detailed proofs can be found in the Appendixes.

<sup>2)</sup> In robustness bounds, the cover radius  $\gamma$  can be made arbitrarily close to zero at the expense of increasing  $\mathcal{N}(\gamma, \mathcal{Z}, \rho)$ . Since  $\mathcal{N}(\gamma, \mathcal{Z}, \rho) = 2^q \mathcal{N}(\gamma, \mathcal{X}, \|\cdot\|_1)$  appears in the second term, the right hand side of the bound indeed goes to zero when  $n \rightarrow \infty$ . This is in accordance with other similar learning bounds.

**Table 5** Predictive performance (mean  $\pm$  std. deviation) of each distance metric learning approach in terms of *average precision*.

Data set	kNN-				MLKNN-			
	COMMU	LM	NJE	ORIGINAL	COMMU	LM	NJE	ORIGINAL
genbase	<b>0.991<math>\pm</math>0.007</b>	0.841 $\pm$ 0.067	0.921 $\pm$ 0.015	0.986 $\pm$ 0.010	<b>0.988<math>\pm</math>0.011</b>	0.983 $\pm$ 0.008	0.933 $\pm$ 0.016	0.984 $\pm$ 0.012
Society	0.544 $\pm$ 0.016	<b>0.555<math>\pm</math>0.026</b>	0.131 $\pm$ 0.013	0.539 $\pm$ 0.018	0.584 $\pm$ 0.011	0.527 $\pm$ 0.031	0.148 $\pm$ 0.020	<b>0.586<math>\pm</math>0.014</b>
Social	<b>0.120<math>\pm</math>0.015</b>	<b>0.662<math>\pm</math>0.025</b>	0.068 $\pm$ 0.005	0.664 $\pm$ 0.027	0.705 $\pm$ 0.018	0.674 $\pm$ 0.017	0.139 $\pm$ 0.117	<b>0.706<math>\pm</math>0.021</b>
Reference	0.579 $\pm$ 0.032	<b>0.592<math>\pm</math>0.024</b>	0.078 $\pm$ 0.012	0.435 $\pm$ 0.025	<b>0.659<math>\pm</math>0.040</b>	0.585 $\pm$ 0.020	0.087 $\pm$ 0.028	0.627 $\pm$ 0.031
Health	0.588 $\pm$ 0.066	<b>0.703<math>\pm</math>0.041</b>	0.106 $\pm$ 0.010	0.474 $\pm$ 0.031	0.690 $\pm$ 0.046	<b>0.726<math>\pm</math>0.036</b>	0.090 $\pm$ 0.013	0.659 $\pm$ 0.040
Education	0.516 $\pm$ 0.022	0.514 $\pm$ 0.020	0.124 $\pm$ 0.009	<b>0.520<math>\pm</math>0.027</b>	0.571 $\pm$ 0.015	0.534 $\pm$ 0.034	0.126 $\pm$ 0.013	<b>0.572<math>\pm</math>0.019</b>
Computers	<b>0.625<math>\pm</math>0.030</b>	0.600 $\pm$ 0.027	0.092 $\pm$ 0.006	0.613 $\pm$ 0.024	<b>0.653<math>\pm</math>0.020</b>	0.643 $\pm$ 0.022	0.102 $\pm$ 0.035	0.652 $\pm$ 0.017
Business	<b>0.872<math>\pm</math>0.016</b>	0.846 $\pm$ 0.019	0.088 $\pm$ 0.005	0.864 $\pm$ 0.018	<b>0.882<math>\pm</math>0.015</b>	0.870 $\pm$ 0.017	0.077 $\pm$ 0.008	0.878 $\pm$ 0.019
Arts	0.473 $\pm$ 0.033	<b>0.526<math>\pm</math>0.028</b>	0.173 $\pm$ 0.016	0.422 $\pm$ 0.029	<b>0.533<math>\pm</math>0.018</b>	0.494 $\pm$ 0.026	0.157 $\pm$ 0.017	0.513 $\pm$ 0.029
yeast	<b>0.752<math>\pm</math>0.009</b>	0.659 $\pm$ 0.011	0.381 $\pm$ 0.012	0.746 $\pm$ 0.011	<b>0.754<math>\pm</math>0.010</b>	0.746 $\pm$ 0.010	0.457 $\pm$ 0.051	0.753 $\pm$ 0.009
corel5k	<b>0.251<math>\pm</math>0.011</b>	0.191 $\pm$ 0.012	-	0.154 $\pm$ 0.013	<b>0.303<math>\pm</math>0.013</b>	0.288 $\pm$ 0.009	-	0.252 $\pm$ 0.013
rcv1-subset1	0.522 $\pm$ 0.028	<b>0.527<math>\pm</math>0.017</b>	-	0.488 $\pm$ 0.014	<b>0.554<math>\pm</math>0.017</b>	0.449 $\pm$ 0.011	-	0.539 $\pm$ 0.013
corel16k	<b>0.212<math>\pm</math>0.004</b>	0.185 $\pm$ 0.002	-	0.184 $\pm$ 0.006	0.293 $\pm$ 0.004	<b>0.303<math>\pm</math>0.002</b>	-	0.279 $\pm$ 0.002
eurlex-dc	<b>0.440<math>\pm</math>0.027</b>	0.310 $\pm$ 0.018	-	0.440 $\pm$ 0.027	<b>0.464<math>\pm</math>0.027</b>	0.371 $\pm$ 0.018	-	0.464 $\pm$ 0.027
eurlex-sm	<b>0.609<math>\pm</math>0.004</b>	0.510 $\pm$ 0.006	-	0.609 $\pm$ 0.004	<b>0.652<math>\pm</math>0.004</b>	0.560 $\pm$ 0.003	-	0.652 $\pm$ 0.004
eurlex	0.023 $\pm$ 0.000	<b>0.030<math>\pm</math>0.000</b>	-	0.011 $\pm$ 0.000	0.032 $\pm$ 0.000	<b>0.040<math>\pm</math>0.000</b>	-	0.033 $\pm$ 0.000

**Table 6** Predictive performance (mean  $\pm$  std. deviation) of each distance metric learning approach in terms of *micro-F1*.

Data set	kNN-				MLKNN-			
	COMMU	LM	NJE	ORIGINAL	COMMU	LM	NJE	ORIGINAL
genbase	<b>0.957<math>\pm</math>0.020</b>	0.848 $\pm$ 0.066	0.806 $\pm$ 0.245	0.950 $\pm$ 0.025	0.942 $\pm$ 0.029	<b>0.951<math>\pm</math>0.030</b>	0.843 $\pm$ 0.125	0.945 $\pm$ 0.031
Society	0.377 $\pm$ 0.011	0.404 $\pm$ 0.032	<b>0.415<math>\pm</math>0.028</b>	0.381 $\pm$ 0.019	0.318 $\pm$ 0.020	0.345 $\pm$ 0.020	<b>0.410<math>\pm</math>0.020</b>	0.312 $\pm$ 0.020
Social	0.506 $\pm$ 0.035	0.556 $\pm$ 0.032	<b>0.587<math>\pm</math>0.026</b>	0.507 $\pm$ 0.031	0.517 $\pm$ 0.020	0.502 $\pm$ 0.019	<b>0.583<math>\pm</math>0.135</b>	0.519 $\pm$ 0.018
Reference	0.373 $\pm$ 0.041	0.493 $\pm$ 0.028	<b>0.530<math>\pm</math>0.033</b>	0.295 $\pm$ 0.032	0.388 $\pm$ 0.021	0.441 $\pm$ 0.015	<b>0.514<math>\pm</math>0.015</b>	0.385 $\pm$ 0.011
Health	0.451 $\pm$ 0.058	0.588 $\pm$ 0.043	<b>0.625<math>\pm</math>0.029</b>	0.367 $\pm$ 0.033	0.439 $\pm$ 0.004	0.562 $\pm$ 0.011	<b>0.618<math>\pm</math>0.020</b>	0.388 $\pm$ 0.004
Education	0.367 $\pm$ 0.025	0.384 $\pm$ 0.027	<b>0.430<math>\pm</math>0.027</b>	0.376 $\pm$ 0.031	0.252 $\pm$ 0.016	0.355 $\pm$ 0.008	<b>0.426<math>\pm</math>0.135</b>	0.243 $\pm$ 0.010
Computers	0.470 $\pm$ 0.020	<b>0.475<math>\pm</math>0.032</b>	0.471 $\pm$ 0.026	0.453 $\pm$ 0.020	0.406 $\pm$ 0.002	0.474 $\pm$ 0.005	<b>0.478<math>\pm</math>0.015</b>	0.376 $\pm$ 0.001
Business	0.715 $\pm$ 0.018	<b>0.726<math>\pm</math>0.021</b>	0.339 $\pm$ 0.307	0.674 $\pm$ 0.017	<b>0.696<math>\pm</math>0.029</b>	0.718 $\pm$ 0.024	0.603 $\pm$ 0.135	0.693 $\pm$ 0.029
Arts	0.322 $\pm$ 0.031	0.376 $\pm$ 0.027	<b>0.413<math>\pm</math>0.040</b>	0.273 $\pm$ 0.032	0.195 $\pm$ 0.004	0.304 $\pm$ 0.008	<b>0.404<math>\pm</math>0.015</b>	0.143 $\pm$ 0.004
yeast	0.610 $\pm$ 0.013	0.572 $\pm$ 0.016	0.408 $\pm$ 0.176	<b>0.641<math>\pm</math>0.012</b>	0.634 $\pm$ 0.011	0.623 $\pm$ 0.015	0.417 $\pm$ 0.174	<b>0.635<math>\pm</math>0.011</b>
corel5k	<b>0.229<math>\pm</math>0.008</b>	0.193 $\pm$ 0.020	-	0.122 $\pm$ 0.015	0.067 $\pm$ 0.012	<b>0.113<math>\pm</math>0.011</b>	-	0.030 $\pm$ 0.007
rcv1-subset1	0.434 $\pm$ 0.024	<b>0.446<math>\pm</math>0.016</b>	-	0.398 $\pm$ 0.009	0.301 $\pm$ 0.018	<b>0.321<math>\pm</math>0.011</b>	-	0.285 $\pm$ 0.012
corel16k	<b>0.168<math>\pm</math>0.005</b>	0.047 $\pm$ 0.002	-	0.021 $\pm$ 0.006	0.015 $\pm$ 0.003	<b>0.046<math>\pm</math>0.007</b>	-	0.009 $\pm$ 0.002
eurlex-dc	<b>0.363<math>\pm</math>0.027</b>	0.287 $\pm$ 0.023	-	0.363 $\pm$ 0.027	<b>0.324<math>\pm</math>0.032</b>	0.254 $\pm$ 0.030	-	0.324 $\pm$ 0.032
eurlex-sm	<b>0.537<math>\pm</math>0.004</b>	0.476 $\pm$ 0.005	-	0.506 $\pm$ 0.003	<b>0.554<math>\pm</math>0.004</b>	0.475 $\pm$ 0.007	-	0.554 $\pm$ 0.004
eurlex	<b>0.009<math>\pm</math>0.000</b>	0.002 $\pm$ 0.000	-	0.002 $\pm$ 0.000	0.002 $\pm$ 0.000	<b>0.013<math>\pm</math>0.000</b>	-	0.002 $\pm$ 0.000

## 4 Experiments

### 4.1 Experimental Setup

#### 4.1.1 Data Sets

To comprehensively evaluate the effectiveness of COMMU, sixteen benchmark multi-label data sets are collected for experimental studies.<sup>3)</sup> Given a multi-label data set  $S$ , we use  $|S|$ ,  $dim(S)$  and  $L(S)$  to represent its number of examples, num-

ber of features and number of class labels respectively. In addition, properties of  $S$  are further characterized by several useful multi-label statistics [28], including label cardinality  $LCard(S)$ , label density  $LDen(S)$ , distinct labelsets  $DL(S)$  and proportion of distinct label sets  $PDL(S)$ .

Table 2 summarizes characteristics of the benchmark multi-label data sets, which are roughly ordered according to  $|S|$ . As shown in Table 2, these data sets serve as a solid basis for comparative studies which exhibit diversified properties in terms of different multi-label statistics.

<sup>3)</sup> Publicly available at <http://mulan.sourceforge.net/datasets.html>, <http://waikato.github.io/meka/datasets/> and <http://manikvarma.org/downloads/XC/XMLRepository>.

**Table 7** Predictive performance (mean  $\pm$  std. deviation) of each distance metric learning approach in terms of *macro-F1*.

Data set	<i>k</i> NN-				MLKNN-			
	COMMU	LM	NJE	ORIGINAL	COMMU	LM	NJE	ORIGINAL
genbase	<b>0.896<math>\pm</math>0.055</b>	0.848 $\pm$ 0.069	0.439 $\pm$ 0.156	0.864 $\pm$ 0.062	0.820 $\pm$ 0.081	<b>0.836<math>\pm</math>0.078</b>	0.485 $\pm$ 0.123	0.836 $\pm$ 0.083
Society	0.255 $\pm$ 0.040	0.150 $\pm$ 0.024	0.131 $\pm$ 0.014	<b>0.268<math>\pm</math>0.031</b>	0.200 $\pm$ 0.085	<b>0.266<math>\pm</math>0.033</b>	0.110 $\pm$ 0.012	0.203 $\pm$ 0.034
Social	<b>0.403<math>\pm</math>0.066</b>	0.127 $\pm$ 0.028	0.089 $\pm$ 0.014	0.401 $\pm$ 0.068	<b>0.384<math>\pm</math>0.085</b>	0.365 $\pm$ 0.056	0.082 $\pm$ 0.017	0.382 $\pm$ 0.061
Reference	<b>0.538<math>\pm</math>0.052</b>	0.133 $\pm$ 0.020	0.134 $\pm$ 0.017	0.475 $\pm$ 0.061	0.486 $\pm$ 0.038	<b>0.500<math>\pm</math>0.058</b>	0.122 $\pm$ 0.017	0.466 $\pm$ 0.067
Health	<b>0.546<math>\pm</math>0.049</b>	0.209 $\pm$ 0.039	0.216 $\pm$ 0.026	0.536 $\pm$ 0.042	0.520 $\pm$ 0.014	<b>0.565<math>\pm</math>0.048</b>	0.201 $\pm$ 0.025	0.486 $\pm$ 0.038
Education	0.454 $\pm$ 0.052	0.108 $\pm$ 0.013	0.133 $\pm$ 0.018	<b>0.461<math>\pm</math>0.052</b>	0.408 $\pm$ 0.030	<b>0.448<math>\pm</math>0.062</b>	0.118 $\pm$ 0.015	0.409 $\pm$ 0.067
Computers	<b>0.347<math>\pm</math>0.052</b>	0.137 $\pm$ 0.022	0.095 $\pm$ 0.015	0.327 $\pm$ 0.051	0.291 $\pm$ 0.002	<b>0.339<math>\pm</math>0.044</b>	0.095 $\pm$ 0.011	0.294 $\pm$ 0.043
Business	0.376 $\pm$ 0.042	0.167 $\pm$ 0.023	0.094 $\pm$ 0.027	<b>0.398<math>\pm</math>0.044</b>	0.374 $\pm$ 0.025	<b>0.402<math>\pm</math>0.050</b>	0.135 $\pm$ 0.018	0.364 $\pm$ 0.051
Arts	<b>0.255<math>\pm</math>0.038</b>	0.150 $\pm$ 0.018	0.157 $\pm$ 0.021	0.240 $\pm$ 0.035	0.187 $\pm$ 0.013	<b>0.233<math>\pm</math>0.039</b>	0.146 $\pm$ 0.020	0.173 $\pm$ 0.045
yeast	0.473 $\pm$ 0.014	0.394 $\pm$ 0.022	0.438 $\pm$ 0.024	<b>0.468<math>\pm</math>0.017</b>	0.381 $\pm$ 0.024	0.355 $\pm$ 0.031	<b>0.429<math>\pm</math>0.019</b>	0.381 $\pm$ 0.023
corel5k	0.237 $\pm$ 0.009	<b>0.339<math>\pm</math>0.014</b>	-	0.325 $\pm$ 0.013	0.328 $\pm$ 0.014	<b>0.329<math>\pm</math>0.014</b>	-	0.321 $\pm$ 0.013
rcv1-subset1	0.307 $\pm$ 0.030	<b>0.324<math>\pm</math>0.032</b>	-	0.286 $\pm$ 0.020	0.213 $\pm$ 0.028	<b>0.223<math>\pm</math>0.020</b>	-	0.205 $\pm$ 0.022
corel16k	<b>0.055<math>\pm</math>0.003</b>	0.017 $\pm$ 0.002	-	0.005 $\pm$ 0.001	0.011 $\pm$ 0.002	<b>0.022<math>\pm</math>0.002</b>	-	0.008 $\pm$ 0.002
eurlex-dc	<b>0.325<math>\pm</math>0.025</b>	0.278 $\pm$ 0.027	-	0.325 $\pm$ 0.025	<b>0.296<math>\pm</math>0.025</b>	0.268 $\pm$ 0.023	-	0.296 $\pm$ 0.025
eurlex-sm	<b>0.252<math>\pm</math>0.012</b>	0.205 $\pm$ 0.011	-	0.182 $\pm$ 0.009	<b>0.224<math>\pm</math>0.013</b>	0.145 $\pm$ 0.009	-	0.224 $\pm$ 0.013
eurlex	<b>8.784e-4<math>\pm</math>0.000</b>	2.993e4 $\pm$ 0.000	-	2.914e-4 $\pm$ 0.000	3.400e-4 $\pm$ 0.000	<b>0.002<math>\pm</math>0.000</b>	-	1.831e-4 $\pm$ 0.000

**Table 8** Win/tie/loss counts (pairwise *t*-test at 0.05 significance level) between COMMU and the comparing approaches.

	<i>k</i> NN-COMMU against			MLKNN-COMMU against		
	<i>k</i> NN-LM	<i>k</i> NN-NJE	<i>k</i> NN-ORIGINAL	MLKNN-LM	MLKNN-NJE	MLKNN-ORIGINAL
<i>ranking loss</i>	14/1/0	1/2/7	5/10/0	12/3/0	8/2/0	3/12/0
<i>coverage</i>	9/6/0	7/3/0	11/4/0	12/3/0	8/2/0	11/4/0
<i>average precision</i>	7/5/3	10/0/0	6/9/0	9/5/1	10/0/0	3/12/0
<i>micro-F1</i>	6/4/5	2/2/6	7/8/0	2/4/9	3/7/0	5/10/0
<i>macro-F1</i>	12/3/0	9/1/0	2/13/0	2/9/4	9/1/0	1/14/0
<b>In Total</b>	<b>48/19/8</b>	<b>29/8/13</b>	<b>31/44/0</b>	<b>37/24/14</b>	<b>38/12/0</b>	<b>23/52/0</b>

#### 4.1.2 Comparing Algorithms

Based on the learned distance metric, it is desirable to show whether the performance of instance-based multi-label classification models can be improved along with the distance measure in embedded feature space. Accordingly, the vanilla *k*NN method and the MLKNN method [29] are utilized as two natural choices for instance-based multi-label classification models. In this paper, the effectiveness of COMMU is compared against two state-of-the-art multi-label metric learning approaches:

- LM [10]: Based on the maximum margin output coding formulation [12], LM learns the distance metric by maximizing the margin of embedded feature vectors and labeling vectors.
- NJE [11]: Based on the Jaccard distance between labeling vectors, NJE learns the distance metric by preserving the similarity of instances in the embedded feature space w.r.t. the labeling Jaccard distance.

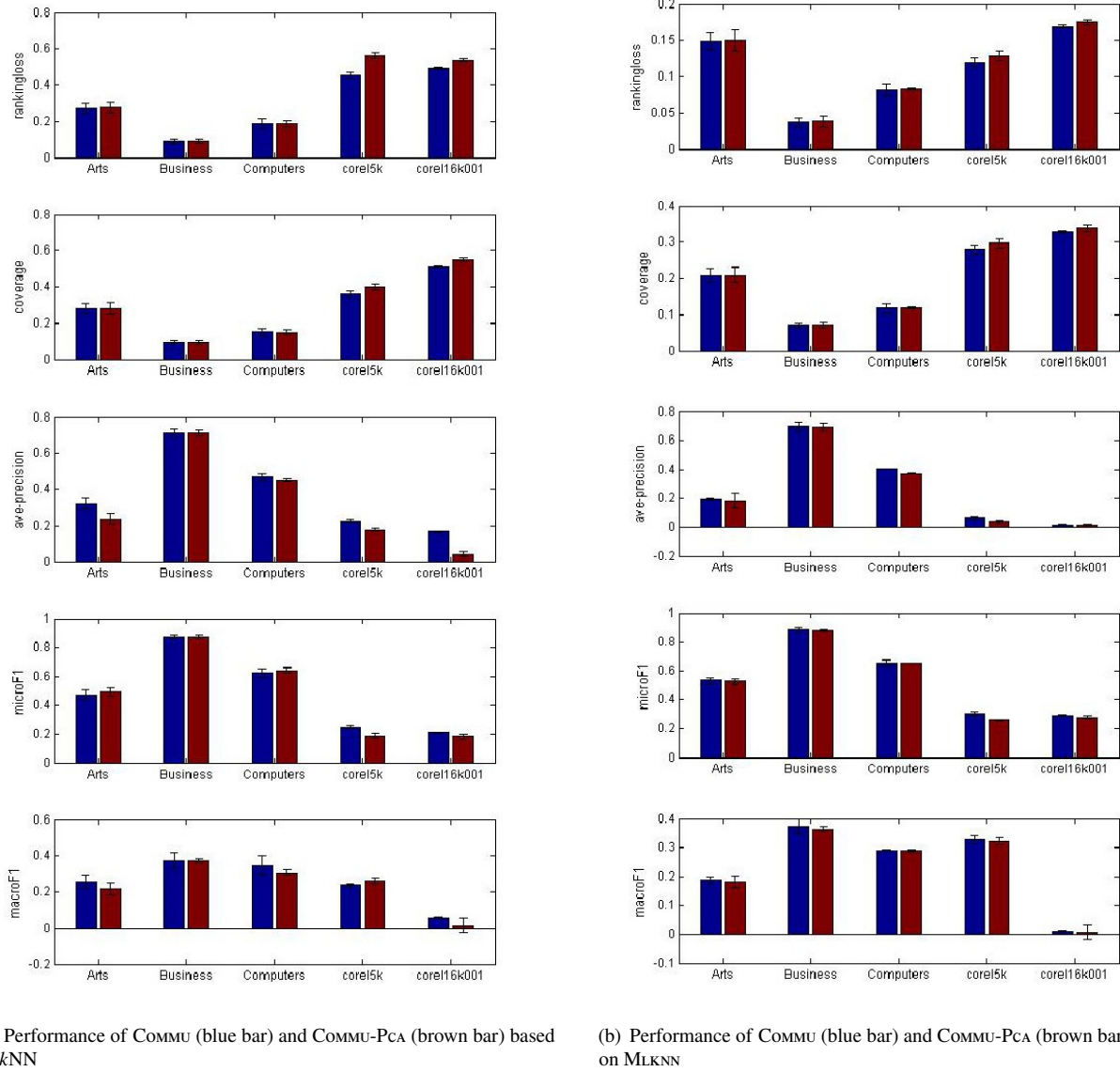
Given the multi-label classification model  $\mathcal{A}$  ( $\mathcal{A} \in \{kNN,$

MLKNN}), its improved version by incorporating the learned distance metric is denoted as  $\mathcal{A}$ -COMMU,  $\mathcal{A}$ -LM and  $\mathcal{A}$ -NJE respectively.

Parameters suggested in the literatures are used to instantiate LM and NJE. As shown in Table 1, the balancing parameter  $\alpha$ , cost parameter  $C$  and thresholding parameter  $\theta$  for COMMU are chosen among  $\{0.1, 0.2, \dots, 1\}$ ,  $\{1, 2, \dots, 10\}$  and  $\{0.1, 0.2, \dots, 1\}$  with cross-validation on the training set. In addition, the number of nearest neighbors used by *k*NN and MLKNN are set to be 10.

#### 4.2 Experimental Results

In this paper, the classification performance is evaluated in terms of five popular multi-label evaluation criteria including *ranking loss*, *coverage*, *average precision*, *micro-F1* and *macro-F1* [1, 2]. For *ranking loss* and *coverage*, the smaller the criterion value the better the performance. For *average precision*, *micro-F1* and *macro-F1*, the greater the criterion value the better the performance. Tables 3-7 report the de-



**Fig. 2** Performance of COMMU and COMMU-PCA based on  $k$ NN and MLKNN in terms of *ranking loss*, *coverage*, *average precision*, *macro-F1* and *micro-F1* (top to bottom) on five data sets.

tailed experimental results of each comparing approach in terms of *ranking loss*, *coverage*, *average precision*, *micro-F1* and *macro-F1* when the learned distance metric is incorporated with  $k$ NN and MLKNN for multi-label prediction. On each data set, ten-fold cross-validation is performed where the mean criterion value as well as the standard deviation are recorded.<sup>4)</sup>

Given the experimental data set and evaluation criterion, pairwise  $t$ -test at 0.05 significance level is conducted to show whether the performance of COMMU is significantly differ-

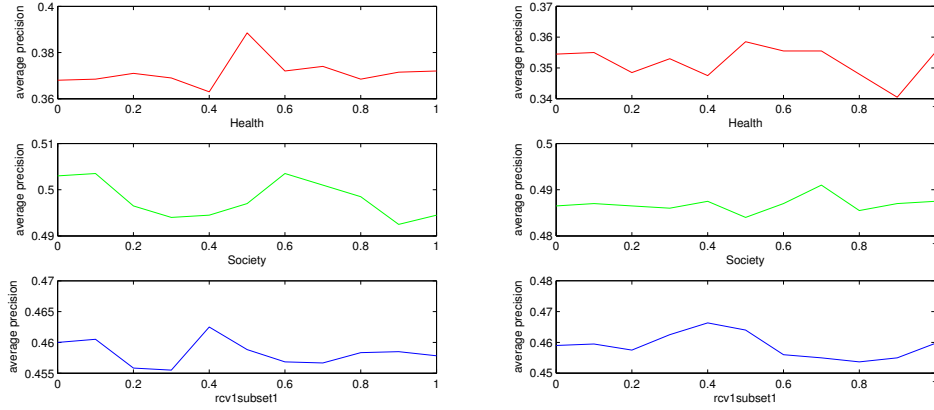
ent to the comparing approaches. Table 8 summarizes the win/tie/loss counts between COMMU and the comparing approaches in terms of each evaluation criterion.

Overall, the following observations can be made based on the reported experimental results:

- Across all evaluation metrics,  $k$ NN-COMMU ranks 1st in 53.8% cases and ranks 2nd in 28.8% cases while MLKNN-COMMU ranks 1st in 52.5% cases and ranks 2nd in 31.3% cases. It is impressive that whenever  $k$ NN or MLKNN are utilized to make multi-label prediction, their counterpart versions ( $k$ NN-COMMU or MLKNN-COMMU) always achieve significantly better or at least comparable performance

<sup>4)</sup> One exception is the *eurllex* dataset from extreme multi-label classification repository, where the predefined training and testing split are used for performance evaluation.





**Fig. 3** Performance of COMMU (in terms of *average precision*) changes with varying value of parameter  $\alpha$  based on  $k$ NN (left column) and MLKNN (right column) on three data sets.

after employing the learned distance metric (Table 8).

- As shown in Table 8, compared to LM, COMMU can lead to superior performance in 64.0% cases for  $k$ NN ( $k$ NN-COMMU against  $k$ NN-LM) and 49.3% cases for MLKNN (MLKNN-COMMU against MLKNN-LM). Compared to NJE, COMMU can lead to superior performance in 58.0% cases for  $k$ NN ( $k$ NN-COMMU against  $k$ NN-NJE) and 76.0% cases for MLKNN (MLKNN-COMMU against MLKNN-NJE).
- As shown in Tables 3, 4, 5, 6,7, the performance advantage of COMMU is more pronounced than the comparing approaches on data sets with larger number of class labels (i.e. *core15k*, *rcv1-subset1*, *core116k*, *eurllex-dc* and *eurllex-sm*). This desirable merit might be attributed to the compositional nature of the distance metric employed by COMMU, where the  $L_1$  regularization term  $\|\mathbf{w}\|_1$  in Eq.(6) can help identify component bases whose corresponding class labels do bring beneficial information for distance metric learning.
- As shown in Tables 3-7, the performance of  $k$ NN-NJE and MLKNN-NJE are not reported on data sets at large scale ( $|\mathcal{S}| \geq 5,000$ ) due to its quadratic training complexity w.r.t. the number of training examples. Specifically, let  $m$ ,  $q$  and  $d$  denote the number of training examples, number of class labels and number of features, the training complexities for COMMU, LM and NJE correspond to  $O((d + q^2)m)$ ,  $O(q^3 + mdq^2)$  and  $O(m^2q + qdm \log(m))$  respectively. For LM, the main computation is the SVD operation of the PSD matrix in each iteration. For NJE, firstly the target vectors is solved at  $O(m^2t)$  ( $t$  is the dimension of the target vector which is unfixed) and then the embedder is learned at  $O(tdm \log(m))$ . For COMMU, to achieve an  $\epsilon$ -solution, the number of iterations needed

by *FISTA* update is  $O(\frac{1}{\sqrt{\epsilon}})$ . At each iteration, projections onto the positive semi-definite cone are performed to solve the coefficient vector  $\mathbf{w}$ . Therefore, the training stage complexity for each iteration is  $O((d + q^2)m\tilde{k}\bar{k})$  with  $k, \tilde{k}$  being the values specified in Eq.(5).

### 4.3 Further Analysis

#### 4.3.1 Effectiveness of Component Bases Generation

We further investigate the effectiveness of COMMU's strategy in generating component bases by encoding discriminative information in label space (Eq.(2)). Specifically, we derive a variant of COMMU (COMMU-PCA) by setting the component bases to the principal components yielded with top  $q$  eigenvalues by conducting PCA over the training instances. Figure 2 compares the performance of COMMU and COMMU-PCA based on  $k$ NN and MLKNN in terms of *ranking loss*, *coverage*, *average precision*, *macro-F1* and *micro-F1* on five data sets, which clearly show the benefits of exploiting discriminative information in generating component bases for COMMU.

#### 4.3.2 Parameter Sensitivity

The parameter  $\alpha$  in Eq.(3) represents the relative contributions from label space and instance space in calculating the semantic similarity. In Figure 3, the performance of COMMU (in terms of *average precision*) on three data sets are illustrated as the parameter  $\alpha$  increases from 0 to 1 with stepsize 0.1 (left column:  $k$ NN; right column: MLKNN). It is obvious that the parameter setting of  $\alpha$  has significant influence on classification performance of the COMMU approach. Therefore, the value of  $\alpha$  are chosen among  $\{0.1, 0.2, \dots, 1\}$

**Table 9** Training time of comparing algorithms on five data sets (in seconds).

Training time	$k$ NN-				MLKNN-			
	COMMU	LM	NJE	ORIGINAL	COMMU	LM	NJE	ORIGINAL
Arts	960.583	106.042	383.512	18.069	1005.790	176.966	347.050	84.309
Business	650.475	123.689	377.070	17.122	690.933	188.033	378.474	76.942
Computers	2699.429	211.573	435.263	36.847	2872.315	478.411	646.539	298.919
yeast	93.605	50.947	364.206	5.519	102.203	71.077	365.362	18.366
genbase	409.308	16.15	182.503	9.158	484.188	130.709	295.765	122.345
corel16k	25341.622	23128.438	-	18585.445	16032.433	8695.786	-	4111.488

with cross-validation on the training set in the experimental studies.

#### 4.3.3 Cost of Training Time

Table 9 reports the training time of comparing algorithms on five data sets. For COMMU, the cost of training time is generally higher than ORIGINAL and comparable to LM and NJE.

the embedded feature vector  $\mathbf{V}^T \mathbf{x}$  naturally follows from the mapping induced by the learned PSD matrix  $\mathbf{M} = \mathbf{V}\mathbf{V}^T$ . Correspondingly, dimensionality reduction serves as the most popular techniques for manipulating multi-label features [34, 35]. There are some other strategies to manipulate the feature space for multi-label learning such as label-specific features [19, 20, 36, 37], meta-level features [38, 39] and multi-view features [40–43].

## 5 Related Works

In Section 4, the performance of COMMU is compared against LM and NJE which to the best of our knowledge are the only two available works on multi-label metric learning. LM [10] adapts the maximum margin output coding formulation [12] for distance metric learning, where the encoding projections are optimized by maximizing the margin of embedded feature vectors and labeling vectors. NJE [11] learns the distance metric by preserving pairwise similarity of labeling vectors in the embedded feature space, where the Jaccard distance is utilized for similarity measurement. Other than the single instance representation, there have been some works on distance metric learning for multi-instance multi-label data [18, 30, 31].

Exploitation of label correlations plays a key role for the success of multi-label classification, where numerous multi-label learning techniques have been proposed by considering different orders of label correlations [1, 2, 32]. Full-order label correlations are considered by LM via linear projection of the labeling vector, while first-order label correlations are considered by NJE via bitwise Jaccard distance measurement. For the proposed COMMU approach, label correlations are brought into the compositional structure of distance metric with label-dependent component bases.

Distance metric learning plays an important role in real-world applications (such as Person Re-ID [33]) in measuring similarity between objects. Generally, distance metric learning can be viewed as feature manipulation techniques where

## 6 Conclusion

In this paper, the problem of distance metric learning for multi-label classification is studied. A novel multi-label metric learning approach named COMMU is proposed, which assumes compositional representation for distance metric. Specifically, component bases as well as triplet constraints are generated by exploiting semantic similarity in label space, and the resulting optimization problem is iteratively solved with linear complexity w.r.t. the number of training examples. Theoretical analysis as well as extensive experiments clearly validate the effectiveness of the proposed compositional distance metric for multi-label classification.

In the future, it is interesting to leverage auxiliary information such as domain knowledge [44] to facilitate multi-label distance metric learning. Furthermore, it is worthwhile to investigate strategies of combining distance metric learning with other popular mechanisms such as feature selection [45–47] for multi-label classification.

## Appendix

Given a multi-label dataset  $\mathcal{S} = \{z = (\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  drawn i.i.d. from a distribution  $P$  over the labelled space  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ , where the label vector  $\mathbf{y}_i$  simultaneously contains multiple labels. Assume that  $\|\mathbf{x}\| \leq R$  (for some convenient norm),  $\forall \mathbf{x} \in \mathcal{X}$ . Different from the single-label setting, COMMU defines the

multi-label semantic similarity matrix to construct the triplet  $(z, z', z'')$ , where  $y$  is similar to  $y'$  and dissimilar to  $y''$ . Let  $\mathcal{S}_R$  be the set of all admissible triplets built from instances in  $\mathcal{S}$ .

Let  $L(h, z, z', z'')$  be the loss suffered by some hypothesis  $h$  on triplet  $(z, z', z'')$  with the convention that  $L$  returns 0 for non-admissible triplets. Assume  $L$  to be uniformly upper-bounded by a constant  $U$ . The empirical loss  $R_{emp}^{S_R}(h)$  of  $h$  on  $\mathcal{S}_R$  is defined as

$$R_{emp}^{S_R}(h) = \frac{1}{|\mathcal{S}_R|} \sum_{(z, z', z'') \in \mathcal{S}_R} L(h, z, z', z''),$$

and its expected loss  $R(h)$  over distribution  $P$  as

$$R(h) = \mathbb{E}_{z, z', z'' \sim P} L(h, z, z', z'').$$

The goal of the theoretical analysis is to bound the deviation between  $R(\mathcal{A}_{S_R})$  and  $R_{emp}^{S_R}(\mathcal{A}_{S_R})$ , where  $\mathcal{A}_{S_R}$  is the hypothesis learned by algorithm  $\mathcal{A}$  on  $\mathcal{S}_R$ .

### Theoretical Basis

To derive the generalization bounds of COMMU, we use the recent framework of algorithmic robustness in metric learning [16, 27]. Algorithmic robustness is the ability of an algorithm to perform “similarly” on a training example and on a test example that are “close”. The proximity of points is based on a partitioning of the space  $\mathcal{Z}$ : two examples are close to each other if they lie in the same region. The partition is based on the notion of covering number.

**Definition 1** (Covering number). For a metric space  $(\mathcal{M}, \rho)$  and  $\nu \subset \mathcal{M}$ , we say that  $\hat{\nu} \subset \nu$  is a  $\gamma$ -cover of  $\nu$  if  $\forall t \in \nu, \exists \hat{t} \in \hat{\nu}$  such that  $\rho(t, \hat{t}) \leq \gamma$ . The  $\gamma$ -covering number of  $\nu$  is

$$N(\gamma, \mathcal{X}, \rho) = \min \{|\hat{\nu}| : \hat{\nu} \text{ is a } \gamma\text{-cover of } \nu\}.$$

In particular, when  $\mathcal{X}$  is compact,  $N(\gamma, \mathcal{X}, \rho)$  is finite, leading to a finite cover. Then,  $\mathcal{Z}$  can be partitioned into  $|\mathcal{Y}| N(\gamma, \mathcal{X}, \rho)$  subsets such that if two examples  $z = (\mathbf{x}, \mathbf{y})$  and  $z' = (\mathbf{x}', \mathbf{y}')$  belong to the same subset, then  $\mathbf{y} = \mathbf{y}'$  and  $\rho(\mathbf{x}, \mathbf{x}') \leq \gamma$ . The definition of robustness for tripletwise loss functions is as follows.

**Definition 2** (Robustness for metric learning) [27]. An algorithm  $\mathcal{A}$  is  $(N, \epsilon(\cdot))$  robust for  $N \in \mathbb{N}$  and  $\epsilon(\cdot) : (\mathcal{Z} \times \mathcal{Z})^n \rightarrow \mathbb{R}$  if  $\mathcal{Z}$  can be partitioned into  $N$  disjoint sets, denoted by  $\{Q_i\}_{i=1}^N$ , such that the following holds for all  $\mathcal{S} \in \mathcal{Z}^n$ :  $\forall (z_1, z_2, z_3) \in \mathcal{S}_R, \forall z, z', z'' \in \mathcal{Z}, \forall i, j, k \in [N] : \text{if } z_1, z \in Q_i, z_2, z' \in Q_j, z_3, z'' \in Q_k$  then

$$|L(\mathcal{A}_{S_R}, z_1, z_2, z_3) - L(\mathcal{A}_{S_R}, z, z', z'')| \leq \epsilon(\mathcal{S}_R),$$

where  $\mathcal{A}_{S_R}$  is the hypothesis learned by  $\mathcal{A}$  on  $\mathcal{S}_R$ .

$N$  and  $\epsilon(\cdot)$  quantify the robustness of the algorithm and depend on the training data. The work [27] showed that a metric learning algorithm that satisfies Definition 2 has the following generalization guarantees.

**Theorem 2.** If a learning algorithm  $\mathcal{A}$  is  $(N, \epsilon(\cdot))$ -robust and the training data consists of the triplets  $\mathcal{S}_R$  obtained from a sample  $\mathcal{S}$  generated by  $n$  i.i.d draws from  $P$ , then for any  $\delta > 0$ , with probability at least  $1 - \delta$  we have :

$$|L(\mathcal{A}_{S_R}, z_1, z_2, z_3) - L(\mathcal{A}_{S_R}, z, z', z'')| \leq \epsilon(\mathcal{S}_R) + 3U \sqrt{\frac{2N \ln 2 + 2 \ln \frac{1}{\delta}}{n}}.$$

Additionally, as shown in [27], the following theorem, which basically says that if a metric learning algorithm has approximately the same loss for triplets that are close to each other and then it is robust, can be used to determine the robustness of the algorithm more conveniently.

**Theorem 3.** Fix  $\gamma > 0$  and a metric  $\rho$  of  $\mathcal{Z}$ . Suppose that  $\forall z_1, z_2, z_3, z, z', z'' : (z_1, z_2, z_3) \in \mathcal{S}_R, \rho(z_1, z) \leq \gamma, \rho(z_2, z') \leq \gamma, \rho(z_3, z'') \leq \gamma, \mathcal{A}$  satisfies

$$|L(\mathcal{A}_{S_R}, z_1, z_2, z_3) - L(\mathcal{A}_{S_R}, z, z', z'')| \leq \epsilon(\mathcal{S}_R),$$

and  $N(\frac{\gamma}{2}, \mathcal{Z}, \rho) < \infty$ . Then the algorithm  $\mathcal{A}$  is  $(N(\frac{\gamma}{2}, \mathcal{Z}, \rho), \epsilon(\mathcal{S}_R))$ -robust.

### Generalization Bounds for COMMU

To derive the generalization bound of COMMU, the main work is to prove its robustness, which contains the computation for  $N$  and  $\epsilon(\mathcal{S}_R)$ .

The loss function of COMMU is defined as:

$$L(\mathbf{w}, z, z', z'') = [\Delta(\mathbf{y}', \mathbf{y}'') + d_w(\mathbf{x}, \mathbf{x}') - d_w(\mathbf{x}, \mathbf{x}'')]_+.$$

Let  $\mathbf{w}^*$  be the optimal solution to COMMU. By optimality of  $\mathbf{w}^*$  we have:

$$L(\mathbf{w}^*, z, z', z'') + \frac{1}{C} \|\mathbf{w}^*\|_1 \leq L(\mathbf{0}, z, z', z'') + \frac{1}{C} \|\mathbf{0}\|_1 = \Delta(\mathbf{y}', \mathbf{y}'') = \mathbf{y}'^T \mathbf{y}'' - \mathbf{y}'^T \mathbf{G} \mathbf{y}'',$$

where the second item  $\mathbf{y}'^T \mathbf{G} \mathbf{y}'' \leq \theta$  because of the dissimilarity between  $\mathbf{y}'$  and  $\mathbf{y}''$ . Let  $\Delta_0 = \mathbf{y}'^T \mathbf{y}''$ , thus  $\Delta_0 - \theta \leq \Delta(\mathbf{y}', \mathbf{y}'') \leq \Delta_0$  and  $\|\mathbf{w}^*\|_1 \leq \Delta_0 C$ .

$M^* = \sum_{i=1}^K w_i^* \mathbf{b}_i \mathbf{b}_i^T$  is the corresponding metric. The norm of the basis element  $\mathbf{b}_i$  is bounded by 1. Based on Holder's inequality and the bound on  $\mathbf{w}^*$  and  $\mathbf{b}_i$ 's, the bound for  $M^*$  is derived.

$$\|M^*\|_1 = \left\| \sum_{i=1}^K w_i^* \mathbf{b}_i \mathbf{b}_i^T \right\|_1 = \left\| \sum_{i: w_i \neq 0} w_i^* \mathbf{b}_i \mathbf{b}_i^T \right\|_1 \leq \|\mathbf{w}^*\|_1 \sum_{i: w_i \neq 0} \|\mathbf{b}_i\|_\infty \|\mathbf{b}_i\|_\infty \leq K^* \Delta_0 C,$$

$$\begin{aligned}
& \left| [\Delta(\mathbf{y}_2, \mathbf{y}_3) + d_w^*(\mathbf{x}_1, \mathbf{x}_2) - d_w^*(\mathbf{x}_1, \mathbf{x}_3)]_+ - [\Delta(\mathbf{y}'_2, \mathbf{y}'_3) + d_w^*(\mathbf{x}'_1, \mathbf{x}'_2) - d_w^*(\mathbf{x}'_1, \mathbf{x}'_3)]_+ \right| \\
& \leq \left| \Delta(\mathbf{y}_2, \mathbf{y}_3) - \Delta(\mathbf{y}'_2, \mathbf{y}'_3) + d_w^*(\mathbf{x}_1, \mathbf{x}_2) - d_w^*(\mathbf{x}'_1, \mathbf{x}'_2) + d_w^*(\mathbf{x}'_1, \mathbf{x}'_3) - d_w^*(\mathbf{x}_1, \mathbf{x}_3) \right| \\
& \leq \left| \Delta(\mathbf{y}_2, \mathbf{y}_3) - \Delta(\mathbf{y}'_2, \mathbf{y}'_3) \right| + \\
& \quad \left| (\mathbf{x}_1 - \mathbf{x}_2)^T M^*(\mathbf{x}_1 - \mathbf{x}_2) + (\mathbf{x}_1 - \mathbf{x}_2)^T M^*(\mathbf{x}'_1 - \mathbf{x}'_2) - (\mathbf{x}_1 - \mathbf{x}_2)^T M^*(\mathbf{x}'_1 - \mathbf{x}'_2) - (\mathbf{x}'_1 - \mathbf{x}'_2)^T M^*(\mathbf{x}'_1 - \mathbf{x}'_2) \right| + \\
& \quad \left| (\mathbf{x}'_1 - \mathbf{x}'_3)^T M^*(\mathbf{x}'_1 - \mathbf{x}'_3) - (\mathbf{x}'_1 - \mathbf{x}'_3)^T M^*(\mathbf{x}_1 - \mathbf{x}_3) + (\mathbf{x}'_1 - \mathbf{x}'_3)^T M^*(\mathbf{x}_1 - \mathbf{x}_3) - (\mathbf{x}_1 - \mathbf{x}_3)^T M^*(\mathbf{x}_1 - \mathbf{x}_3) \right| \\
& = \left| \Delta(\mathbf{y}_2, \mathbf{y}_3) - \Delta(\mathbf{y}'_2, \mathbf{y}'_3) \right| + \left| (\mathbf{x}_1 - \mathbf{x}_2)^T M^*(\mathbf{x}_1 - \mathbf{x}_2 - (\mathbf{x}'_1 - \mathbf{x}'_2)) + (\mathbf{x}_1 - \mathbf{x}_2 - (\mathbf{x}'_1 - \mathbf{x}'_2))^T M^*(\mathbf{x}'_1 - \mathbf{x}'_2) \right| + \\
& \quad \left| (\mathbf{x}'_1 - \mathbf{x}'_3)^T M^*(\mathbf{x}'_1 - \mathbf{x}'_3 - (\mathbf{x}_1 - \mathbf{x}_3)) + (\mathbf{x}'_1 - \mathbf{x}'_3 - (\mathbf{x}_1 - \mathbf{x}_3))^T M^*(\mathbf{x}_1 - \mathbf{x}_3) \right| \\
& \leq \left| \Delta(\mathbf{y}_2, \mathbf{y}_3) - \Delta(\mathbf{y}'_2, \mathbf{y}'_3) \right| + \left| (\mathbf{x}_1 - \mathbf{x}_2)^T M^*(\mathbf{x}_1 - \mathbf{x}'_1) \right| + \left| (\mathbf{x}_1 - \mathbf{x}_2)^T M^*(\mathbf{x}'_2 - \mathbf{x}_2) \right| + \\
& \quad \left| (\mathbf{x}_1 - \mathbf{x}'_1)^T M^*(\mathbf{x}'_1 - \mathbf{x}'_2) \right| + \left| (\mathbf{x}'_2 - \mathbf{x}_2)^T M^*(\mathbf{x}'_1 - \mathbf{x}'_2) \right| + \left| (\mathbf{x}'_1 - \mathbf{x}'_3)^T M^*(\mathbf{x}'_1 - \mathbf{x}_1) \right| + \\
& \quad \left| (\mathbf{x}'_1 - \mathbf{x}'_3)^T M^*(\mathbf{x}_3 - \mathbf{x}'_3) \right| + \left| (\mathbf{x}'_1 - \mathbf{x}_1)^T M^*(\mathbf{x}_1 - \mathbf{x}_3) \right| + \left| (\mathbf{x}_3 - \mathbf{x}'_3)^T M^*(\mathbf{x}_1 - \mathbf{x}_3) \right| \\
& \leq \left| \Delta(\mathbf{y}_2, \mathbf{y}_3) - \Delta(\mathbf{y}'_2, \mathbf{y}'_3) \right| + \|\mathbf{x}_1 - \mathbf{x}_2\|_\infty \|M^*\|_1 \|\mathbf{x}_1 - \mathbf{x}'_1\|_1 + \|\mathbf{x}_1 - \mathbf{x}_2\|_\infty \|M^*\|_1 \|\mathbf{x}'_2 - \mathbf{x}_2\|_1 + \\
& \quad \|\mathbf{x}_1 - \mathbf{x}'_1\|_1 \|M^*\|_1 \|\mathbf{x}'_1 - \mathbf{x}'_2\|_\infty + \|\mathbf{x}'_2 - \mathbf{x}_2\|_1 \|M^*\|_1 \|\mathbf{x}'_1 - \mathbf{x}'_2\|_\infty + \|\mathbf{x}'_1 - \mathbf{x}'_3\|_\infty \|M^*\|_1 \|\mathbf{x}'_1 - \mathbf{x}_1\|_1 + \\
& \quad \|\mathbf{x}'_1 - \mathbf{x}'_3\|_\infty \|M^*\|_1 \|\mathbf{x}_3 - \mathbf{x}'_3\|_1 + \|\mathbf{x}_3 - \mathbf{x}'_3\|_1 \|M^*\|_1 \|\mathbf{x}_1 - \mathbf{x}_3\|_\infty + \|\mathbf{x}'_1 - \mathbf{x}_1\|_1 \|M^*\|_1 \|\mathbf{x}_1 - \mathbf{x}_3\|_\infty \\
& \leq 16\gamma RK^* \Delta_0 C + \theta
\end{aligned}$$

where  $K^* \leq K$  is the number of nonzero entries in  $w^*$ .

According to Definition 1,  $\mathcal{Z}$  can be partitioned into  $2^q \mathcal{N}(\gamma, \mathcal{X}, \rho)$  subsets, where  $q$  is the size of label vector. COMMU constructs the triplet  $(z, z', z'')$  by the multi-label semantic similarity matrix. For  $z_1, z_2, z_3, z'_1, z'_2, z'_3 \in \mathcal{Z}$ , if  $y_1$  is similar to  $y'_1$ ,  $\|x_1 - x'_1\|_1 \leq \gamma$ ,  $y_2$  is similar to  $y'_2$ ,  $\|x_2 - x'_2\|_1 \leq \gamma$ ,  $y_3$  is similar to  $y'_3$ ,  $\|x_3 - x'_3\|_1 \leq \gamma$ , then  $(z_1, z_2, z_3)$  and  $(z'_1, z'_2, z'_3)$  are either both admissible or non-admissible triplets.

In the non-admissible case, it can be seen from definition that their respective loss is 0 and so is the deviation between the losses. In the admissible case we have the above result, by the property that the hinge loss is 1-Lipschitz (the first  $\leq$ ), Holder's inequality (the 2th-4th  $\leq$ ) and  $\|x_i - x_j\|_\infty \leq 2R$  ( $\|x\| \leq R, \forall x \in \mathcal{X}$ ),  $|\Delta(\mathbf{y}_2, \mathbf{y}_3) - \Delta(\mathbf{y}'_2, \mathbf{y}'_3)| \leq \theta$  ( $\Delta_0 - \theta \leq \Delta(\mathbf{y}', \mathbf{y}'') \leq \Delta_0$ ) in the last  $\leq$ . Thus COMMU is  $(2^q \mathcal{N}(\gamma, \mathcal{X}, \|\cdot\|_1), 16\gamma RK^* \Delta_0 C + \theta)$ -robust and the generalization bound follows.

$$\begin{aligned}
& \left| R(\mathcal{A}_{S_R}) - R_{emp}^{S_R}(\mathcal{A}_{S_R}) \right| \leq 16\gamma RK^* \Delta_0 C + \theta + \\
& \quad 3U \sqrt{\frac{2^{q+1} \mathcal{N}(\gamma, \mathcal{X}, \|\cdot\|_1) \ln 2 + 2 \ln \frac{1}{\delta}}{n}}.
\end{aligned}$$

## References

1. M.-L. Zhang and Z.-H. Zhou. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 26(8):1819–1837, 2014.
2. E. Gibaja and S. Ventura. A tutorial on multilabel learning. *ACM Computing Surveys*, 47(3):Article 52, 2015.
3. F. Briggs, B. Lakshminarayanan, L. Neal, X. Z. Fern, R. Raich, S. J. Hadley, A. S. Hadley, and M. G. Betts. Acoustic classification of multiple simultaneous bird species: A multi-instance multi-label approach. *Journal of the Acoustical Society of America*, 131(6):4640–4650, 2012.
4. R. Cabral, F. De la Torre, J. P. Costeira, and A. Bernardino. Matrix completion for weakly-supervised multi-label image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(1):121–135, 2015.
5. J. Liu, W.-C. Chang, Y. Wu, and Y. Yang. Deep learning for extreme multi-label text classification. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 115–124, Tokyo, Japan, 2017.
6. X. Pan, Y.-X. Fan, J. Jia, and H.-B. Shen. Identifying rna-binding proteins using multi-label deep learning. *Science China Information Sciences*, page 62:19103, 2019.
7. L. Sun, H. Ge, and W. Kang. Non-negative matrix factorization based modeling and training algorithm for multi-label learning. *Frontiers of Computer Science*, 13(6):1243–1254, 2019.
8. A. Bellet, A. Habrard, and M. Sebban. Metric learning. *Synthesis*

- Lectures on Artificial Intelligence and Machine Learning*, 9(1):1–151, 2015.
9. F. Wang and J. Sun. Survey on distance metric learning and dimensionality reduction in data mining. *Data Mining and Knowledge Discovery*, 29(2):534–564, 2015.
  10. W. Liu and I. W. Tsang. Large margin metric learning for multi-label prediction. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence*, pages 2800–2806, Austin, TX, 2015.
  11. H. Gouk and B. Pfahringer and M. Cree. Learning distance metrics for multi-label classification. In *Proceedings of the 8th Asian Conference on Machine Learning*, pages 318–333, Hamilton, New Zealand, 2016.
  12. Y. Zhang and J. Schneider. Maximum margin output coding. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1575–1582, Edinburgh, UK, 2012.
  13. Y. Verma and C. V. Jawahar. Image annotation by propagating labels from semantic neighbourhoods. *International Journal of Computer Vision*, 121(1):126–148, 2017.
  14. H. Gouk, B. Pfahringer, and M. Cree. Learning similarity metrics by factorising adjacency matrices. *CoRR*, abs/1511.06442, 2015.
  15. J. Ni, J. Liu, C. Zhang, D. Ye, and Z. Ma. Fine-grained patient similarity measuring using deep metric learning. In *Proceedings of the 26th ACM International Conference on Information and Knowledge Management*, pages 1189–1198, Singapore, 2017.
  16. Y. Shi, A. Bellet, and F. Sha. Sparse compositional metric learning. In *Proceedings of the 28th AAAI Conference on Artificial Intelligence*, pages 2078–2084, Québec City, Canada, 2014.
  17. J. St.Amand and J. Huan. Sparse compositional local metric learning. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1097–1104, Halifax, Canada, 2017.
  18. Z.-H. Zhou, M.-L. Zhang, S.-J. Huang, and Y.-F. Li. Multi-instance multi-label learning. *Artificial Intelligence*, 176(1):2291–2320, 2012.
  19. M.-L. Zhang and L. Wu. LIFT: Multi-label learning with label-specific features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(1):107–120, 2015.
  20. J. Huang, G. Li, Q. Huang, and X. Wu. Learning label-specific features and class-dependent labels for multi-label classification. *IEEE Transactions on Knowledge and Data Engineering*, 28(12):3309–3323, 2016.
  21. K. Q. Weinberger and L. K. Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10:207–244, 2009.
  22. S.-J. Huang and Z.-H. Zhou. Multi-label learning by exploiting label correlations locally. In *Proceedings of the 26th AAAI Conference on Artificial Intelligence*, pages 949–955, Toronto, Canada, 2012.
  23. Y. Zhu, J. Kwok, and Z.-H. Zhou. Multi-label learning with global and local correlation. *IEEE Transactions on Knowledge and Data Engineering*, 30(6):1081–1094, 2018.
  24. G.-X. Yuan, C.-H. Ho, and C.-J. Lin. An improved GLMNET for L1-regularized logistic regression. *Journal of Machine Learning Research*, 13:1999–2030, 2012.
  25. A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *Siam Journal on Imaging Sciences*, 2(1):183–202, 2009.
  26. K.-C. Toh and S. Yun. An accelerated proximal gradient algorithm for nuclear norm regularized least squares problems. *Pacific Journal of Optimization*, 6(3):615–640, 2010.
  27. A. Bellet and A. Habrard. Robustness and generalization for metric learning. *Neurocomputing*, 151(14):259–267, 2015.
  28. J. Read, B. Pfahringer, G. Holmes, and E. Frank. Classifier chains for multi-label classification. *Machine Learning*, 85(3):333–359, 2011.
  29. M.-L. Zhang and Z.-H. Zhou. ML-kNN: A lazy learning approach to multi-label learning. *Pattern Recognition*, 40(7):2038–2048, 2007.
  30. J. Rong, S. Wang, and Z.-H. Zhou. Learning a distance metric from multi-instance multi-label data. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 896–902, Miami, FL, 2009.
  31. Y. Verma and C. V. Jawahar. A robust distance with correlated metric learning for multi-instance multi-label data. In *Proceedings of the 24th ACM International Conference on Multimedia*, pages 441–445, Amsterdam, The Netherlands, 2016.
  32. M.-L. Zhang, Y.-K. Li, Y.-Y. Liu, and X. Geng. Binary relevance for multi-label learning: An overview. *Frontiers of Computer Science*, 12(2):191–202, 2018.
  33. Y. Wu, Y. Lin, X. Dong, Y. Yan, W. Bian, and Y. Yang. Progressive learning for person re-identification with one example. *IEEE Transactions on Image Processing*, 28(6):2872–2881, 2019.
  34. L. Sun, S. Ji, and J. Ye. *Multi-label Dimensionality Reduction*. Chapman and Hall/CRC, Boca Ration, FL, 2013.
  35. R. B. Pereira, A. Plastino, B. Zadrozny, and L. H. C. Merschmann. Categorizing feature selection methods for multi-label classification. *Artificial Intelligence Review*, 49(1):57–78, 2018.
  36. J. Zhang, C. Li, D. Cao, Y. Lin, S. Su, L. Dai, and S. Li. Multi-label learning with label-specific features by resolving label correlations. *Knowledge-Based Systems*, 159:148–157, 2018.
  37. Z.-S. Chen and M.-L. Zhang. Multi-label learning with regularization enriched label-specific features. In *Proceedings of the 11th Asian Conference on Machine Learning*, pages 411–424, Nagoya, Japan, 2019.
  38. Y. Yang and S. Gopal. Multilabel classification with meta-level features in a learning-to-rank framework. *Machine Learning*, 88(1-2):47–68, 2012.
  39. S. Canuto, M. A. Gonçalves, and F. Benevenuto. Exploiting new sentiment-based meta-level features for effective sentiment analysis. In *Proceedings of the 9th ACM International Conference on Web Search and Data Mining*, pages 53–62, San Francisco, CA, 2016.
  40. X. Zhu, X. Li, and S. Zhang. Block-row sparse multiview multilabel learning for image classification. *IEEE Transactions on Cybernetics*, 46(2):450–461, 2016.
  41. C. Zhang, Z. Yu, Q. Hu, P. Zhu, X. Liu, and X. Wang. Latent semantic aware multi-view multi-label classification. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, pages 4414–4421, New Orleans, LA, 2018.
  42. X. Wu, Q.-G. Chen, Y. Hu, D.-B. Wang, X. Chang, X. Wang, and M.-L.

- Zhang. Multi-view multi-label learning with view-specific information extraction. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 3884–3890, Macau, China, 2019.
43. R. Zhang, F. Nie, X. Li, and X. Wei. Feature selection with multi-view data: A survey. *Information Fusion*, 50:158–167, 2019.
  44. Z.-H. Zhou. Abductive learning: Towards bridging machine learning and logical reasoning. *Science China Information Sciences*, page 62:076101, 2019.
  45. Y. Yang, Z. Ma, A. G. Hauptmann, and N. Sebe. Feature selection for multimedia analysis by sharing information among multiple tasks. *IEEE Transactions on Multimedia*, 15(3):661–669, 2013.
  46. R. Zhang, F. Nie, and X. Li. Self-weighted supervised discriminative feature selection. *IEEE Transactions on Neural Networks and Learning Systems*, 29(8):3913–3918, 2018.
  47. R. Zhang, F. Nie, Y. Wang, and X. Li. Unsupervised feature selection via adaptive multimeasure fusion. *IEEE Transactions on Neural Networks and Learning Systems*, 30(9):2886–2892, 2019.