

# Multi-Instance Multi-Label Learning

Zhi-Hua Zhou<sup>\*</sup>, Min-Ling Zhang, Sheng-Jun Huang, Yu-Feng Li

*National Key Laboratory for Novel Software Technology,*

*Nanjing University, Nanjing 210046, China*

---

## Abstract

In this paper, we propose the MIML (*Multi-Instance Multi-Label learning*) framework where an example is described by multiple instances and associated with multiple class labels. Compared to traditional learning frameworks, the MIML framework is more convenient and natural for representing complicated objects which have multiple semantic meanings. To learn from MIML examples, we propose the MIMLBOOST and MIMLSVM algorithms based on a simple degeneration strategy, and experiments show that solving problems involving complicated objects with multiple semantic meanings in the MIML framework can lead to good performance. Considering that the degeneration process may lose information, we propose the D-MIMLSVM algorithm which tackles MIML problems directly in a regularization framework. Moreover, we show that even when we do not have access to the real objects and thus cannot capture more information from real objects by using the MIML representation, MIML is still useful. We propose the INSDIF and SUBCOD algorithms. INSDIF works by transforming single-instances into the MIML representation for learning, while SUBCOD works by transforming single-label examples into the MIML representation for learning. Experiments show that in some tasks they are able to achieve better performance than learning the single-instances or single-label examples directly.

*Key words:* Machine Learning, Multi-Instance Multi-Label Learning, MIML, Multi-Label Learning, Multi-Instance Learning

---

<sup>\*</sup> Corresponding author. E-mail: zhouzh@lamda.nju.edu.cn

## 1 Introduction

In *traditional supervised learning*, an object is represented by an instance, i.e., a feature vector, and associated with a class label. Formally, let  $\mathcal{X}$  denote the instance space (or feature space) and  $\mathcal{Y}$  the set of class labels. The task is to learn a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  from a given data set  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ , where  $\mathbf{x}_i \in \mathcal{X}$  is an instance and  $y_i \in \mathcal{Y}$  is the known label of  $\mathbf{x}_i$ . Although this formalization is prevailing and successful, there are many real-world problems which do not fit in this framework well. In particular, each object in this framework belongs to only one concept and therefore the corresponding instance is associated with a single class label. However, many real-world objects are complicated, which may belong to multiple concepts simultaneously. For example, an image can belong to several classes simultaneously, e.g., *grasslands*, *lions*, *Africa*, etc.; a text document can be classified to several categories if it is viewed from different aspects, e.g., *scientific novel*, *Jules Verne's writing* or even *books on traveling*; a web page can be recognized as *news page*, *sports page*, *soccer page*, etc. In a specific real task, maybe only one of the multiple concepts is the right semantic meaning. For example, in image retrieval when a user is interested in an image with lions, s/he may be only interested in the concept *lions* instead of the other concepts *grasslands* and *Africa* associated with that image. The difficulty here is caused by those objects that involve multiple concepts. To choose the right semantic meaning for such objects for a specific scenario is the fundamental difficulty of many tasks. In contrast to starting from a large universe of all possible concepts involved in the task, it may be helpful to get the subset of concepts associated with the concerned object at first, and then make a choice in the small subset later. However, getting the subset of concepts, that is, assigning proper class labels to such objects, is still a challenging task.

We notice that as an alternative to representing an object by a single instance, in many cases it is possible to represent a complicated object using a set of instances. For example, multiple patches can be extracted from an image where each patch is described by an instance, and thus the image can be represented by a set of instances; multiple sections can be extracted from a document where each section is described by an instance, and thus the document can be represented by a set of instances; multiple links can be extracted from a web page where each link is

described by an instance, and thus the web page can be represented by a set of instances. Using multiple instances to represent those complicated objects may be helpful because some inherent patterns which are closely related to some labels may become explicit and clearer. In this paper, we propose the MIML (*Multi-Instance Multi-Label learning*) framework, where an example is described by multiple instances and associated with multiple class labels.

Compared to traditional learning frameworks, the MIML framework is more convenient and natural for representing complicated objects. To exploit the advantages of the MIML representation, new learning algorithms are needed. We propose the MIMLBOOST algorithm and the MIMLSVM algorithm based on a simple degeneration strategy, and experiments show that solving problems involving complicated objects with multiple semantic meanings under the MIML framework can lead to good performance. Considering that the degeneration process may lose information, we also propose the D-MIMLSVM (i.e., Direct MIMLSVM) algorithm which tackles MIML problems directly in a regularization framework. Experiments show that this “direct” algorithm outperforms the “indirect” MIMLSVM algorithm.

In some practical tasks we do not have access to the real objects themselves such as the real images and the real web pages; instead, we are given observational data where each real object has already been represented by a single instance. Thus, in such cases we cannot capture more information from the real objects using the MIML representation. Even in this situation, however, MIML is still useful. We propose the INSDIF (i.e., INSTance DIFferentiation) algorithm which transforms single-instances into MIML examples for learning. This algorithm is able to achieve a better performance than learning the single-instances directly in some tasks. This is not strange because for an object associated with multiple class labels, if it is described by only a single instance, the information corresponding to these labels are mixed and thus difficult for learning; if we can transform the single-instance into a set of instances in some proper ways, the mixed information might be detached to some extent and thus less difficult for learning.

MIML can also be helpful for learning single-label objects. We propose the SUBCOD (i.e., SUB-COncept Discovery) algorithm which works by discovering sub-concepts of the target concept at first and then transforming the data into MIML examples

for learning. This algorithm is able to achieve a better performance than learning the single-label examples directly in some tasks. This is also not strange because for a label corresponding to a high-level complicated concept, it may be quite difficult to learn this concept directly since many different lower-level concepts are mixed; if we can transform the single-label into a set of labels corresponding to some sub-concepts, which are relatively clearer and easier for learning, we can learn these labels at first and then derive the high-level complicated label based on them with a less difficulty.

The rest of this paper is organized as follows. In Section 2, we review some related work. In Section 3, we propose the MIML framework. In Section 4 we propose the MIMLBOOST and MIMLSVM algorithms, and apply them to tasks where the objects are represented as MIML examples. In Section 5 we present the D-MIMLSVM algorithm and compare it with the “indirect” MIMLSVM algorithm. In Sections 6 and 7, we study the usefulness of MIML when we do not have access to real objects. Concretely, in Section 6, we propose the INSDIF algorithm and show that using MIML can be better than learning single-instances directly; in Section 7 we propose the SUBCOD algorithm and show that using MIML can be better than learning single-label examples directly. Finally, we conclude the paper in Section 8.

## 2 Related Work

Much work has been devoted to the learning of multi-label examples under the umbrella of *multi-label learning*. Note that multi-label learning studies the problem where a real-world object described by one instance is associated with a number of class labels<sup>1</sup>, which is different from multi-class learning or multi-task learning [28]. In multi-class learning each object is only associated with a single label; while in multi-task learning different tasks may involve different domains and different data sets. Actually, traditional two-class and multi-class problems can both be cast into multi-label problems by restricting that each instance has only one label. The generality of multi-label problems, however, inevitably makes it more difficult to

---

<sup>1</sup> Most work on multi-label learning assumes that an instance can be associated with multiple valid labels, but there is also some work assuming that only one of the labels among those associated with an instance is correct [35].

address.

One famous approach to solving multi-label problems is Schapire and Singer’s ADABOOST.MH [56], which is an extension of ADABOOST and is the core of a successful multi-label learning system BOOSTEXTER [56]. This approach maintains a set of weights over both training examples and their labels in the training phase, where training examples and their corresponding labels that are hard (easy) to predict get incrementally higher (lower) weights. Later, De Comité et al. [22] used alternating decision trees [30] which are more powerful than decision stumps used in BOOSTEXTER to handle multi-label data and thus obtained the ADTBOOST.MH algorithm. Probabilistic generative models have been found useful in multi-label learning. McCallum [47] proposed a Bayesian approach for multi-label document classification, where a mixture probabilistic model (one mixture component per category) is assumed to generate each document and an EM algorithm is employed to learn the mixture weights and the word distributions in each mixture component. Ueda and Saito [65] presented another generative approach, which assumes that the multi-label text has a mixture of characteristic words appearing in single-label text belonging to each of the multi-labels. It is noteworthy that the generative models used in [47] and [65] are both based on learning text frequencies in documents, and are thus specific to text applications.

Many other multi-label learning algorithms have been developed, such as decision trees, neural networks,  $k$ -nearest neighbor classifiers, support vector machines, etc. Clare and King [21] developed a multi-label version of C4.5 decision trees through modifying the definition of entropy. Zhang and Zhou [79] presented multi-label neural network BP-MLL, which is derived from the Backpropagation algorithm by employing an error function to capture the fact that the labels belonging to an instance should be ranked higher than those not belonging to that instance. Zhang and Zhou [80] also proposed the ML- $k$ NN algorithm, which identifies the  $k$  nearest neighbors of the concerned instance and then assigns labels according to the maximum a posteriori principle. Elisseeff and Weston [27] proposed the RANKSVM algorithm for multi-label learning by defining a specific cost function and the corresponding margin for multi-label models. Other kinds of multi-label SVMs have been developed by Boutell et al. [11] and Godbole and Sarawagi [33]. In particular, by hierarchically approximating the Bayes optimal classifier for the H-loss,

Cesa-Bianchi et al. [15] proposed an algorithm which outperforms simple hierarchical SVMs. Recently, non-negative matrix factorization has also been applied to multi-label learning [43], and multi-label dimensionality reduction methods have been developed [74, 85].

Roughly speaking, earlier approaches to multi-label learning attempt to divide multi-label learning to a number of two-class classification problems [36, 72] or transform it into a label ranking problem [27, 56], while some later approaches try to exploit the correlation between the labels [43, 65, 85].

Most studies on multi-label learning focus on text categorization [22, 33, 39, 47, 56, 65, 74], and several studies aim to improve the performance of text categorization systems by exploiting additional information given by the hierarchical structure of classes [14, 15, 53] or unlabeled data [43]. In addition to text categorization, multi-label learning has also been found useful in many other tasks such as scene classification [11], image and video annotation [38, 48], bioinformatics [7, 12, 13, 21, 27], and even association rule mining [50, 63].

There is a lot of research on *multi-instance learning*, which studies the problem where a real-world object described by a number of instances is associated with a single class label. Here the training set is composed of many *bags* each containing multiple instances; a bag is labeled positively if it contains at least one positive instance and negatively otherwise. The goal is to label unseen bags correctly. Note that although the training bags are labeled, the labels of their instances are unknown. This learning framework was formalized by Dietterich et al. [24] when they were investigating drug activity prediction.

Long and Tan [44] studied the PAC-learnability of multi-instance learning and showed that if the instances in the bags are independently drawn from product distribution, the APR (Axis-Parallel Rectangle) proposed by Dietterich et al. [24] is PAC-learnable. Auer et al. [5] showed that if the instances in the bags are not independent then APR learning under the multi-instance learning framework is NP-hard. Moreover, they presented a theoretical algorithm that does not require product distribution, which was transformed into a practical algorithm named MULTINST [4]. Blum and Kalai [10] described a reduction from PAC-learning under the multi-instance learning framework to PAC-learning with one-sided random

classification noise. They also presented an algorithm with smaller sample complexity than that of the algorithm of Auer et al. [5].

Many multi-instance learning algorithms have been developed during the past decade. To name a few, DIVERSE DENSITY [45] and EM-DD [83],  $k$ -nearest neighbor algorithms CITATION- $k$ NN and BAYESIAN- $k$ NN [67], decision tree algorithms RELIC [54] and MITI [9], neural network algorithms BP-MIP and extensions [77,90] and RBF-MIP [78], rule learning algorithm RIPPER-MI [20], support vector machines and kernel methods MI-SVM and MI-SVM [3], DD-SVM [18], MISSSVM [88], MI-KERNEL [32], BAG-INSTANCE KERNEL [19], MARGINALIZED MI-KERNEL [42] and convex-hull method CH-FD [31], ensemble algorithms MI-ENSEMBLE [91], MI-BOOSTING [70] and MILBOOSTING [6], logistic regression algorithm MI-LR [51], etc. Actually almost all popular machine learning algorithms have their multi-instance versions. Most algorithms attempt to adapt single-instance supervised learning algorithms to the multi-instance representation, by shifting their focus from discrimination on instances to discrimination on bags [91]. Recently there is some proposal on adapting the multi-instance representation to single-instance algorithms by representation transformation [93].

It is worth mentioning that standard multi-instance learning [24] assumes that if a bag contains a positive instance then the bag is positive; this implies that there exists a *key instance* in a positive bag. Many algorithms were designed based on this assumption. For example, the point with maximal diverse density identified by the DIVERSE DENSITY algorithm [45] actually corresponds to a key instance; many SVM algorithms defined the margin of a positive bag by the margin of its *most* positive instance [3,19]. As the research of multi-instance learning goes on, however, some other assumptions have been introduced [29]. For example, in contrast to assuming that there is a key instance, some work has assumed that there is no key instance and every instance contributes to the bag label [17,70]. There is also an argument that the instances in the bags should not be treated independently [88]. All those assumptions have been put under the umbrella of multi-instance learning, and generally, in tackling real tasks it is difficult to know which assumption is the fittest. In other words, in different tasks multi-instance learning algorithms based on different assumptions may have different superiorities.

In the early years of the research of multi-instance learning, most work considered multi-instance classification with discrete-valued outputs. Later, multi-instance regression with real-valued outputs was studied [2, 52], and different versions of generalized multi-instance learning have been defined [58, 68]. The main difference between standard multi-instance learning and generalized multi-instance learning is that in standard multi-instance learning there is a single concept, and a bag is positive if it has an instance satisfying this concept; while in generalized multi-instance learning [58, 68] there are multiple concepts, and a bag is positive only when all concepts are satisfied (i.e., the bag contains instances from every concept). Recently, research on multi-instance clustering [82], multi-instance semi-supervised learning [49] and multi-instance active learning [60] have also been reported.

Multi-instance learning has also attracted the attention of the ILP community. It has been suggested that multi-instance problems could be regarded as a bias on inductive logic programming, and the multi-instance paradigm could be the key between the propositional and relational representations, being more expressive than the former, and much easier to learn than the latter [23]. Alphonse and Matwin [1] approximated a relational learning problem by a multi-instance problem, fed the resulting data to feature selection techniques adapted from propositional representations, and then transformed the filtered data back to relational representation for a relational learner. Thus, the expressive power of relational representation and the ease of feature selection on propositional representation are gracefully combined. This work confirms that multi-instance learning can really act as a bridge between propositional and relational learning.

Multi-instance learning techniques have already been applied to diverse applications including image categorization [17, 18], image retrieval [71, 84], text categorization [3, 60], web mining [86], spam detection [37], computer security [54], face detection [66, 76], computer-aided medical diagnosis [31], etc.

### 3 The MIML Framework

Let  $\mathcal{X}$  denote the instance space and  $\mathcal{Y}$  the set of class labels. Then, formally, the MIML task is defined as:

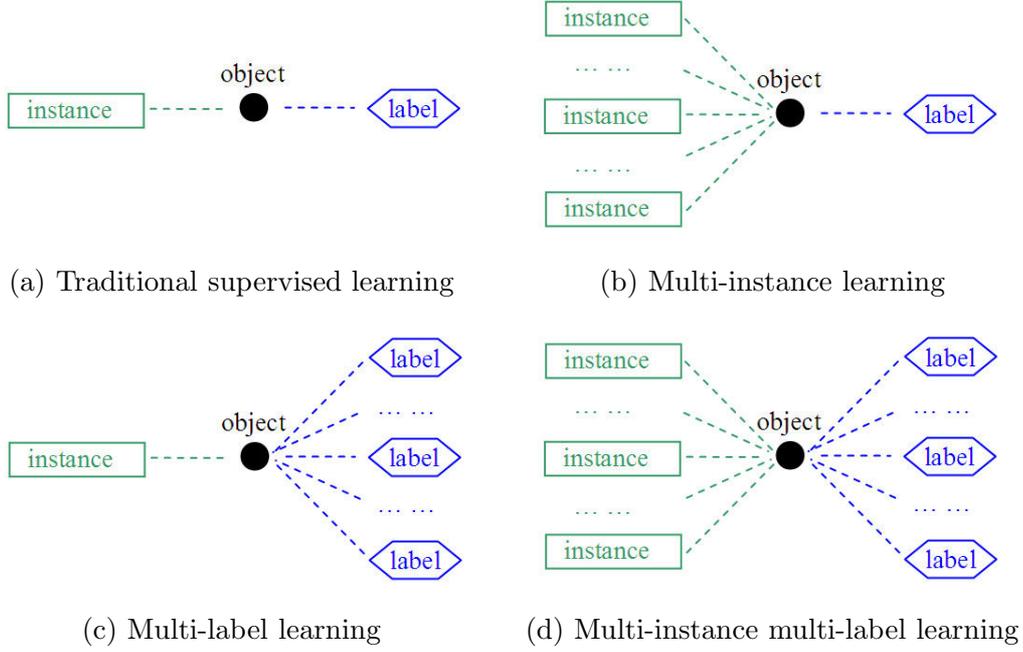


Fig. 1. Four different learning frameworks

- **MIML** (multi-instance multi-label learning): To learn a function  $f : 2^{\mathcal{X}} \rightarrow 2^{\mathcal{Y}}$  from a given data set  $\{(X_1, Y_1), (X_2, Y_2), \dots, (X_m, Y_m)\}$ , where  $X_i \subseteq \mathcal{X}$  is a set of instances  $\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{i, n_i}\}$ ,  $\mathbf{x}_{ij} \in \mathcal{X}$  ( $j = 1, 2, \dots, n_i$ ), and  $Y_i \subseteq \mathcal{Y}$  is a set of labels  $\{y_{i1}, y_{i2}, \dots, y_{i, l_i}\}$ ,  $y_{ik} \in \mathcal{Y}$  ( $k = 1, 2, \dots, l_i$ ). Here  $n_i$  denotes the number of instances in  $X_i$  and  $l_i$  the number of labels in  $Y_i$ .

It is interesting to compare MIML with the existing frameworks of traditional supervised learning, multi-instance learning, and multi-label learning.

- **Traditional supervised learning** (single-instance single-label learning): To learn a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  from a given data set  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ , where  $\mathbf{x}_i \in \mathcal{X}$  is an instance and  $y_i \in \mathcal{Y}$  is the known label of  $\mathbf{x}_i$ .
- **Multi-instance learning** (multi-instance single-label learning): To learn a function  $f : 2^{\mathcal{X}} \rightarrow \mathcal{Y}$  from a given data set  $\{(X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)\}$ , where  $X_i \subseteq \mathcal{X}$  is a set of instances  $\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{i, n_i}\}$ ,  $\mathbf{x}_{ij} \in \mathcal{X}$  ( $j = 1, 2, \dots, n_i$ ), and  $y_i \in \mathcal{Y}$  is the label of  $X_i$ .<sup>2</sup> Here  $n_i$  denotes the number of instances in  $X_i$ .

<sup>2</sup> According to notions used in multi-instance learning,  $(X_i, y_i)$  is a labeled *bag* while  $X_i$  an unlabeled bag.

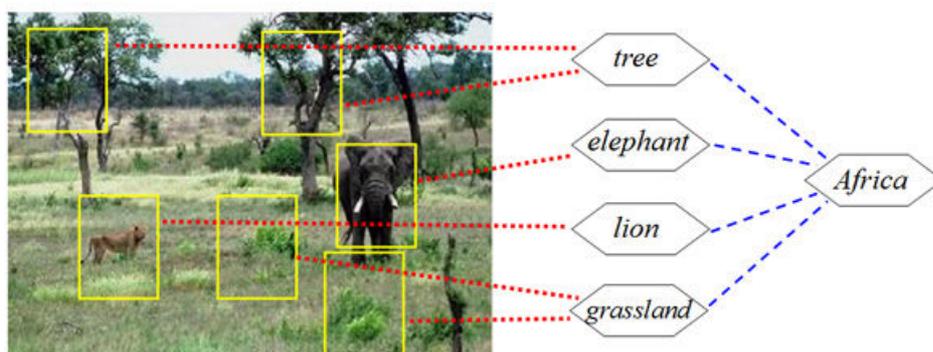
- **Multi-label learning** (single-instance multi-label learning): To learn a function  $f : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$  from a given data set  $\{(\mathbf{x}_1, Y_1), (\mathbf{x}_2, Y_2), \dots, (\mathbf{x}_m, Y_m)\}$ , where  $\mathbf{x}_i \in \mathcal{X}$  is an instance and  $Y_i \subseteq \mathcal{Y}$  is a set of labels  $\{y_{i1}, y_{i2}, \dots, y_{i, l_i}\}$ ,  $y_{ik} \in \mathcal{Y}$  ( $k = 1, 2, \dots, l_i$ ). Here  $l_i$  denotes the number of labels in  $Y_i$ .

From Fig. 1 we can see the differences among these learning frameworks. In fact, the *multi-* learning frameworks are resulted from the ambiguities in representing real-world objects. Multi-instance learning studies the ambiguity in the input space (or instance space), where an object has many alternative input descriptions, i.e., instances; multi-label learning studies the ambiguity in the output space (or label space), where an object has many alternative output descriptions, i.e., labels; while MIML considers the ambiguities in both the input and output spaces simultaneously. In solving real-world problems, having a good representation is often more important than having a strong learning algorithm, because a good representation may capture more meaningful information and make the learning task easier to tackle. Since many real objects are inherited with input ambiguity as well as output ambiguity, MIML is more natural and convenient for tasks involving such objects.

It is worth mentioning that MIML is more reasonable than (single-instance) multi-label learning in many cases. Suppose a multi-label object is described by one instance but associated with  $l$  number of class labels, namely  $\text{label}_1, \text{label}_2, \dots, \text{label}_l$ . If we represent the multi-label object using a set of  $n$  instances, namely  $\text{instance}_1, \text{instance}_2, \dots, \text{instance}_n$ , the underlying information in a single instance may become easier to exploit, and for each label the number of training instances can be significantly increased. So, transforming multi-label examples to MIML examples for learning may be beneficial in some tasks, which will be shown in Section 6. Moreover, when representing the multi-label object using a set of instances, the relation between the input patterns and the semantic meanings may become more easily discoverable. Note that in some cases, understanding why a particular object has a certain class label is even more important than simply making an accurate prediction, while MIML offers a possibility for this purpose. For example, under the MIML representation, we may discover that one object has  $\text{label}_1$  because it contains  $\text{instance}_n$ ; it has  $\text{label}_l$  because it contains  $\text{instance}_i$ ; while the occurrence of both  $\text{instance}_1$  and  $\text{instance}_i$  triggers  $\text{label}_j$ .



(a) *Africa* is a complicated high-level concept



(b) The concept *Africa* may become easier to learn through exploiting some sub-concepts

Fig. 2. MIML can be helpful in learning single-label examples involving complicated high-level concepts

MIML can also be helpful for learning single-label examples involving complicated high-level concepts. For example, as Fig. 2(a) shows, the concept *Africa* has a broad connotation and the images belonging to *Africa* have great variance, thus it is not easy to classify the top-left image in Fig. 2(a) into the *Africa* class correctly. However, if we can exploit some low-level sub-concepts that are less ambiguous and easier to learn, such as *tree*, *lions*, *elephant* and *grassland* shown in Fig. 2(b), it is possible to induce the concept *Africa* much easier than learning the concept *Africa* directly. The usefulness of MIML in this process will be shown in Section 7.

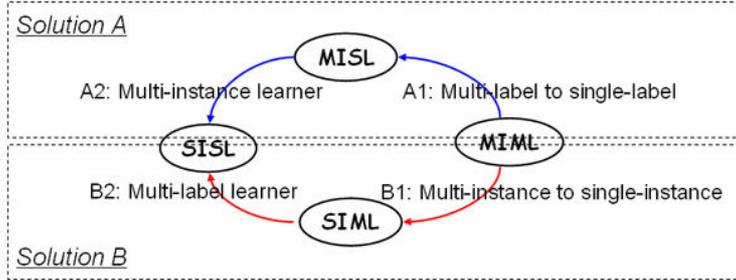


Fig. 3. The two general degeneration solutions.

#### 4 Solving MIML Problems by Degeneration

It is evident that traditional supervised learning is a degenerated version of multi-instance learning as well as a degenerated version of multi-label learning, while traditional supervised learning, multi-instance learning and multi-label learning are all degenerated versions of MIML. So, a simple idea to tackle MIML is to identify its equivalence in the traditional supervised learning framework, using multi-instance learning or multi-label learning as the bridge, as shown in Fig. 3.

- **Solution A:** Using multi-instance learning as the bridge:

The MIML learning task, i.e., to learn a function  $f : 2^{\mathcal{X}} \rightarrow 2^{\mathcal{Y}}$ , can be transformed into a multi-instance learning task, i.e., to learn a function  $f_{MIL} : 2^{\mathcal{X}} \times \mathcal{Y} \rightarrow \{-1, +1\}$ . For any  $y \in \mathcal{Y}$ ,  $f_{MIL}(X_i, y) = +1$  if  $y \in Y_i$  and  $-1$  otherwise. The proper labels for a new example  $X^*$  can be determined according to  $Y^* = \{y | \text{sign}[f_{MIL}(X^*, y)] = +1\}$ . This multi-instance learning task can be further transformed into a traditional supervised learning task, i.e., to learn a function  $f_{SISL} : \mathcal{X} \times \mathcal{Y} \rightarrow \{-1, +1\}$ , under a constraint specifying how to derive  $f_{MIL}(X_i, y)$  from  $f_{SISL}(\mathbf{x}_{ij}, y)$  ( $j = 1, 2, \dots, n_i$ ). For any  $y \in \mathcal{Y}$ ,  $f_{SISL}(\mathbf{x}_{ij}, y) = +1$  if  $y \in Y_i$  and  $-1$  otherwise. Here the constraint can be  $f_{MIL}(X_i, y) = \text{sign}[\sum_{j=1}^{n_i} f_{SISL}(\mathbf{x}_{ij}, y)]$  which has been used by Xu and Frank [70] in transforming multi-instance learning tasks into traditional supervised learning tasks. Note that other kinds of constraint can also be used here.

- **Solution B:** Using multi-label learning as the bridge:

The MIML learning task, i.e., to learn a function  $f : 2^{\mathcal{X}} \rightarrow 2^{\mathcal{Y}}$ , can be transformed into a multi-label learning task, i.e., to learn a function  $f_{MLL} : \mathcal{Z} \rightarrow 2^{\mathcal{Y}}$ . For any  $\mathbf{z}_i \in \mathcal{Z}$ ,  $f_{MLL}(\mathbf{z}_i) = f_{MIML}(X_i)$  if  $\mathbf{z}_i = \phi(X_i)$ ,  $\phi : 2^{\mathcal{X}} \rightarrow \mathcal{Z}$ . The proper

labels for a new example  $X^*$  can be determined according to  $Y^* = f_{MLL}(\phi(X^*))$ . This multi-label learning task can be further transformed into a traditional supervised learning task, i.e., to learn a function  $f_{SISL} : \mathcal{Z} \times \mathcal{Y} \rightarrow \{-1, +1\}$ . For any  $y \in \mathcal{Y}$ ,  $f_{SISL}(z_i, y) = +1$  if  $y \in Y_i$  and  $-1$  otherwise. That is,  $f_{MLL}(z_i) = \{y | f_{SISL}(z_i, y) = +1\}$ . Here the mapping  $\phi$  can be implemented with *constructive clustering* which was proposed by Zhou and Zhang [93] in transforming multi-instance bags into traditional single-instances. Note that other kinds of mappings can also be used here.

In the rest of this section we will propose two MIML algorithms, MIMLBOOST and MIMLSVM. MIMLBOOST is an illustration of Solution A, which uses *category-wise decomposition* for the A1 step in Fig. 3 and MIBOOSTING for A2; MIMLSVM is an illustration of Solution B, which uses *clustering-based representation transformation* for the B1 step and MLSVM for B2. Other MIML algorithms can be developed by taking alternative options. Both MIMLBOOST and MIMLSVM are quite simple. We will see that for dealing with complicated objects with multiple semantic meanings, good performance can be obtained under the MIML framework even by using such simple algorithms. This demonstrates that the MIML framework is very promising, and we expect better performance can be achieved in the future if researchers put forward more powerful MIML algorithms.

#### 4.1 MIMLBOOST

Now we propose the MIMLBOOST algorithm according to the first solution mentioned above, that is, identifying the equivalence in the traditional supervised learning framework using multi-instance learning as the bridge. Note that this strategy can also be used to derive other kinds of MIML algorithms.

Given any set  $\Omega$ , let  $|\Omega|$  denote its size, i.e., the number of elements in  $\Omega$ ; given any predicate  $\pi$ , let  $[[\pi]]$  be 1 if  $\pi$  holds and 0 otherwise; given  $(X_i, Y_i)$ , for any  $y \in \mathcal{Y}$ , let  $\Psi(X_i, y) = +1$  if  $y \in Y_i$  and  $-1$  otherwise, where  $\Psi$  is a function  $\Psi : 2^{\mathcal{X}} \times \mathcal{Y} \rightarrow \{-1, +1\}$  which judges whether a label  $y$  is a proper label of  $X_i$  or not. The basic assumption of MIMLBOOST is that the labels are independent

so that the MIML task can be decomposed into a series of multi-instance learning tasks to solve, by treating each label as a task. The pseudo-code of MIMLBOOST is summarized in Appendix A (Table A.1).

In the first step of MIMLBOOST, each MIML example  $(X_u, Y_u)$  ( $u = 1, 2, \dots, m$ ) is transformed into a set of  $|\mathcal{Y}|$  number of multi-instance bags, i.e.,  $\{[(X_u, y_1), \Psi(X_u, y_1)], [(X_u, y_2), \Psi(X_u, y_2)], \dots, [(X_u, y_{|\mathcal{Y}|}), \Psi(X_u, y_{|\mathcal{Y}|})]\}$ . Note that  $[(X_u, y_v), \Psi(X_u, y_v)]$  ( $v = 1, 2, \dots, |\mathcal{Y}|$ ) is a labeled multi-instance bag where  $(X_u, y_v)$  is a bag containing  $n_u$  number of instances, i.e.,  $\{(\mathbf{x}_{u1}, y_v), (\mathbf{x}_{u2}, y_v), \dots, (\mathbf{x}_{u, n_u}, y_v)\}$ , and  $\Psi(X_u, y_v) \in \{-1, +1\}$  is the label of this bag.

Thus, the original MIML data set is transformed into a multi-instance data set containing  $m \times |\mathcal{Y}|$  number of bags. We order them as  $[(X_1, y_1), \Psi(X_1, y_1)], \dots, [(X_1, y_{|\mathcal{Y}|}), \Psi(X_1, y_{|\mathcal{Y}|})], [(X_2, y_1), \Psi(X_2, y_1)], \dots, [(X_m, y_{|\mathcal{Y}|}), \Psi(X_m, y_{|\mathcal{Y}|})]$ , and let  $[(X^{(i)}, y^{(i)}), \Psi(X^{(i)}, y^{(i)})]$  denote the  $i$ -th of these  $m \times |\mathcal{Y}|$  number of bags which contains  $n_i$  number of instances.

Then, from the data set a multi-instance learning function  $f_{MIL}$  can be learned, which can accomplish the desired MIML function because  $f_{MIML}(X^*) = \{y | \text{sign}[f_{MIL}(X^*, y)] = +1\}$ . In this paper, the MIBOOSTING algorithm [70] is used to implement  $f_{MIL}$ . Note that by using MIBOOSTING, the MIMLBOOST algorithm assumes that all instances in a bag contribute independently in an equal way to the label of that bag.

For convenience, let  $(B, g)$  denote the bag  $[(X, y), \Psi(X, y)]$ ,  $B \in \mathcal{B}$ ,  $g \in \mathcal{G}$ , and  $E$  denotes the expectation. Then, here the goal is to learn a function  $\mathcal{F}(B)$  minimizing the bag-level exponential loss  $E_B E_{\mathcal{G}|B}[\exp(-g\mathcal{F}(B))]$ , which ultimately estimates the bag-level log-odds function  $\frac{1}{2} \log \frac{\text{Pr}(g=1|B)}{\text{Pr}(g=-1|B)}$  on the training set. In each boosting round, the aim is to expand  $\mathcal{F}(B)$  into  $\mathcal{F}(B) + cf(B)$ , i.e., adding a new weak classifier, so that the exponential loss is minimized. Assuming that all instances in a bag contribute equally and independently to the bag's label,  $f(B) = \frac{1}{n_B} \sum_j h(\mathbf{b}_j)$  can be derived, where  $h(\mathbf{b}_j) \in \{-1, +1\}$  is the prediction of the instance-level classifier  $h(\cdot)$  for the  $j$ -th instance of the bag  $B$ , and  $n_B$  is the number of instances in  $B$ .

It has been shown by [70] that the best  $f(B)$  to be added can be achieved by seek-

ing  $h(\cdot)$  which maximizes  $\sum_i \sum_{j=1}^{n_i} [\frac{1}{n_i} W^{(i)} g^{(i)} h(\mathbf{b}_j^{(i)})]$ , given the bag-level weights  $W = \exp(-g\mathcal{F}(B))$ . By assigning each instance the label of its bag and the corresponding weight  $W^{(i)}/n_i$ ,  $h(\cdot)$  can be learned by minimizing the weighted instance-level classification error. This actually corresponds to the Step 3a of MIMLBOOST. When  $f(B)$  is found, the best multiplier  $c > 0$  can be got by directly optimizing the exponential loss:

$$\begin{aligned} E_{\mathcal{B}} E_{\mathcal{G}|\mathcal{B}}[\exp(-g\mathcal{F}(B) + c(-gf(B)))] &= \sum_i W^{(i)} \exp \left[ c \left( -\frac{g^{(i)} \sum_j h(\mathbf{b}_j^{(i)})}{n_i} \right) \right] \\ &= \sum_i W^{(i)} \exp[(2e^{(i)} - 1)c], \end{aligned} \quad (1)$$

where  $e^{(i)} = \frac{1}{n_i} \sum_j \mathbb{I}[(h(\mathbf{b}_j^{(i)}) \neq g^{(i)})]$  (computed in Step 3b). Minimization of this expectation actually corresponds to Step 3d, where numeric optimization techniques such as quasi-Newton method can be used. Note that in Step 3c if  $e^{(i)} \geq 0.5$ , the Boosting process will stop [89]. Finally, the bag-level weights are updated in Step 3f according to the additive structure of  $\mathcal{F}(B)$ .

## 4.2 MIMLSVM

Now we propose the MIMLSVM algorithm according to the second solution mentioned before, that is, identifying the equivalence in the traditional supervised learning framework using multi-label learning as the bridge. Note that this strategy can also be used to derive other kinds of MIML algorithms.

Again, given any set  $\Omega$ , let  $|\Omega|$  denote its size, i.e., the number of elements in  $\Omega$ ; given  $(X_i, Y_i)$  and  $\mathbf{z}_i = \phi(X_i)$  where  $\phi: 2^{\mathcal{X}} \rightarrow \mathcal{Z}$ , for any  $y \in \mathcal{Y}$ , let  $\Phi(\mathbf{z}_i, y) = +1$  if  $y \in Y_i$  and  $-1$  otherwise, where  $\Phi$  is a function  $\Phi: \mathcal{Z} \times \mathcal{Y} \rightarrow \{-1, +1\}$ . The basic assumption of MIMLSVM is that the spatial distribution of the bags carries relevant information, and information helpful for label discrimination can be discovered by measuring the closeness between each bag and the representative bags identified through clustering. The pseudo-code of MIMLSVM is summarized in Appendix A (Table A.2).

In the first step of MIMLSVM, the  $X_u$  of each MIML example  $(X_u, Y_u)$  ( $u = 1, 2, \dots, m$ ) is collected and put into a data set  $\Gamma$ . Then, in the second step,  $k$ -medoids clustering is performed on  $\Gamma$ . Since each data item in  $\Gamma$ , i.e.  $X_u$ , is an

unlabeled multi-instance bag instead of a single instance, Hausdorff distance [26] is employed to measure the distance. The Hausdorff distance is a famous metric for measuring the distance between two bags of points, which has often been used in computer vision tasks; other techniques that can measure the distance between bags of points, such as the *set kernel* [32], can also be used here. In detail, given two bags  $A = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{n_A}\}$  and  $B = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{n_B}\}$ , the Hausdorff distance between  $A$  and  $B$  is defined as

$$d_H(A, B) = \max\left\{\max_{\mathbf{a} \in A} \min_{\mathbf{b} \in B} \|\mathbf{a} - \mathbf{b}\|, \max_{\mathbf{b} \in B} \min_{\mathbf{a} \in A} \|\mathbf{b} - \mathbf{a}\|\right\}, \quad (2)$$

where  $\|\mathbf{a} - \mathbf{b}\|$  measures the distance between the instances  $\mathbf{a}$  and  $\mathbf{b}$ , which takes the form of Euclidean distance here.

After the clustering process, the data set  $\Gamma$  is divided into  $k$  partitions, whose medoids are  $M_t$  ( $t = 1, 2, \dots, k$ ), respectively. With the help of these medoids, the original multi-instance example  $X_u$  is transformed into a  $k$ -dimensional numerical vector  $\mathbf{z}_u$ , where the  $i$ -th ( $i = 1, 2, \dots, k$ ) component of  $\mathbf{z}_u$  is the distance between  $X_u$  and  $M_i$ , that is,  $d_H(X_u, M_i)$ . In other words,  $\mathbf{z}_{ui}$  encodes some structure information of the data, that is, the relationship between  $X_u$  and the  $i$ -th partition of  $\Gamma$ . This process reassembles the *constructive clustering* process used by Zhou and Zhang [93] in transforming multi-instance examples into single-instance examples except that in [93] the clustering is executed at the instance level while here it is executed at the bag level. Thus, the original MIML examples  $(X_u, Y_u)$  ( $u = 1, 2, \dots, m$ ) have been transformed into multi-label examples  $(\mathbf{z}_u, Y_u)$  ( $u = 1, 2, \dots, m$ ), which corresponds to the Step 3 of MIMLSVM.

Then, from the data set a multi-label learning function  $f_{MLL}$  can be learned, which can accomplish the desired MIML function because  $f_{MIML}(X^*) = f_{MLL}(\mathbf{z}^*)$ . In this paper, the MLSVM algorithm [11] is used to implement  $f_{MLL}$ . Concretely, MLSVM decomposes the multi-label learning problem into multiple independent binary classification problems (one per class), where each example associated with the label set  $Y$  is regarded as a positive example when building SVM for any class  $y \in Y$ , while regarded as a negative example when building SVM for any class  $y \notin Y$ , as shown in the Step 4 of MIMLSVM. In making predictions, the *T-Criterion* [11] is used, which actually corresponds to the Step 5 of the MIMLSVM algorithm. That is, the test example is labeled by all the class labels with positive

SVM scores, except that when all the SVM scores are negative, the test example is labeled by the class label which is with the *top* (least negative) score.

### 4.3 Experiments

#### 4.3.1 Multi-Label Evaluation Criteria

In traditional supervised learning where each object has only one class label, *accuracy* is often used as the performance evaluation criterion. Typically, accuracy is defined as the percentage of test examples that are correctly classified. When learning with complicated objects associated with multiple labels simultaneously, however, accuracy becomes less meaningful. For example, if approach *A* missed one proper label while approach *B* missed four proper labels for a test example having five labels, it is obvious that *A* is better than *B*, but the accuracy of *A* and *B* may be identical because both of them incorrectly classified the test example.

Five criteria are often used for evaluating the performance of learning with multi-label examples [56,92]; they are *hamming loss*, *one-error*, *coverage*, *ranking loss* and *average precision*. Using the same denotation as that in Sections 3 and 4, given a test set  $S = \{(X_1, Y_1), (X_2, Y_2), \dots, (X_p, Y_p)\}$ , these five criteria are defined as below. Here,  $h(X_i)$  returns a set of proper labels of  $X_i$ ;  $h(X_i, y)$  returns a real-value indicating the confidence for  $y$  to be a proper label of  $X_i$ ;  $rank^h(X_i, y)$  returns the rank of  $y$  derived from  $h(X_i, y)$ .

- $hloss_S(h) = \frac{1}{p} \sum_{i=1}^p \frac{1}{|Y|} |h(X_i) \Delta Y_i|$ , where  $\Delta$  stands for the symmetric difference between two sets. The *hamming loss* evaluates how many times an object-label pair is misclassified, i.e., a proper label is missed or a wrong label is predicted. The performance is perfect when  $hloss_S(h) = 0$ ; the smaller the value of  $hloss_S(h)$ , the better the performance of  $h$ .
- $one-error_S(h) = \frac{1}{p} \sum_{i=1}^p \mathbb{1}[\arg \max_{y \in Y} h(X_i, y) \notin Y_i]$ . The *one-error* evaluates how many times the top-ranked label is not a proper label of the object. The performance is perfect when  $one-error_S(h) = 0$ ; the smaller the value of  $one-error_S(h)$ , the better the performance of  $h$ .

- $\text{coverage}_S(h) = \frac{1}{p} \sum_{i=1}^p \max_{y \in Y_i} \text{rank}^h(X_i, y) - 1$ . The *coverage* evaluates how far it is needed, on the average, to go down the list of labels in order to cover all the proper labels of the object. It is loosely related to precision at the level of perfect recall. The smaller the value of  $\text{coverage}_S(h)$ , the better the performance of  $h$ .
- $\text{rloss}_S(h) = \frac{1}{p} \sum_{i=1}^p \frac{1}{|Y_i| |\bar{Y}_i|} |\{(y_1, y_2) | h(X_i, y_1) \leq h(X_i, y_2), (y_1, y_2) \in Y_i \times \bar{Y}_i\}|$ , where  $\bar{Y}_i$  denotes the complementary set of  $Y_i$  in  $\mathcal{Y}$ . The *ranking loss* evaluates the average fraction of label pairs that are misordered for the object. The performance is perfect when  $\text{rloss}_S(h) = 0$ ; the smaller the value of  $\text{rloss}_S(h)$ , the better the performance of  $h$ .
- $\text{avgprec}_S(h) = \frac{1}{p} \sum_{i=1}^p \frac{1}{|Y_i|} \sum_{y \in Y_i} \frac{|\{y' | \text{rank}^h(X_i, y') \leq \text{rank}^h(X_i, y), y' \in Y_i\}|}{\text{rank}^h(X_i, y)}$ . The *average precision* evaluates the average fraction of proper labels ranked above a particular label  $y \in Y_i$ . The performance is perfect when  $\text{avgprec}_S(h) = 1$ ; the larger the value of  $\text{avgprec}_S(h)$ , the better the performance of  $h$ .

In addition to the above criteria, we design two new multi-label criteria, *average recall* and *average F1*, as below.

- $\text{avgrec}_S(h) = \frac{1}{p} \sum_{i=1}^p \frac{|\{y | \text{rank}^h(X_i, y) \leq |h(X_i)|, y \in Y_i\}|}{|Y_i|}$ . The *average recall* evaluates the average fraction of proper labels that have been predicted. The performance is perfect when  $\text{avgrec}_S(h) = 1$ ; the larger the value of  $\text{avgrec}_S(h)$ , the better the performance of  $h$ .
- $\text{avgF1}_S(h) = \frac{2 \times \text{avgprec}_S(h) \times \text{avgrec}_S(h)}{\text{avgprec}_S(h) + \text{avgrec}_S(h)}$ . The *average F1* expresses a tradeoff between the *average precision* and the *average recall*. The performance is perfect when  $\text{avgF1}_S(h) = 1$ ; the larger the value of  $\text{avgF1}_S(h)$ , the better the performance of  $h$ .

Note that since the above criteria measure the performance from different aspects, it is difficult for one algorithm to outperform another on every one of these criteria.

In the following we study the performance of MIML algorithms on two tasks involving complicated objects with multiple semantic meanings. We will show that for such tasks, MIML is a good choice, and good performance can be achieved even

by using simple MIML algorithms such as MIMLBOOST and MIMLSVM.

#### 4.3.2 Scene Classification

The scene classification data set consists of 2,000 natural scene images belonging to the classes *desert*, *mountains*, *sea*, *sunset* and *trees*. Over 22% of the images belong to multiple classes simultaneously. Each image has already been represented as a bag of nine instances generated by the SBN method [46], which uses a Gaussian filter to smooth the image and then subsamples the image to an  $8 \times 8$  matrix of *color blobs* where each blob is a  $2 \times 2$  set of pixels within the matrix. An instance corresponds to the combination of a single blob with its four neighboring blobs (up, down, left, right), which is described with 15 features. The first three features represent the mean R, G, B values of the central blob and the remaining twelve features express the differences in mean color values between the central blob and other four neighboring blobs respectively.<sup>3</sup>

We evaluate the performance of the MIML algorithms MIMLBOOST and MIMLSVM. Note that MIMLBOOST and MIMLSVM are merely proposed to illustrate the two general degeneration solutions to MIML problems shown in Fig. 3. We do not claim that they are the best algorithms that can be developed through the degeneration paths. There may exist other processes for transforming MIML examples into multi-instance single-label (MISL) examples or single-instance multi-label (SIML) examples. Even by using the same degeneration process as that used in MIMLBOOST and MIMLSVM, there are also many alternatives to realize the second step. For example, by using MI-SVM [3] to replace the MIBOOSTING used in MIMLBOOST and by using the two-layer neural network structure [81] to replace the MLSVM used in MIMLSVM, we get MIMLSVM<sub>mi</sub> and MIMLNN respectively. Their performance is also evaluated in our experiments.

We compare the MIML algorithms with several state-of-the-art algorithms for learning with multi-label examples, including ADTBOOST.MH [22], RANKSVM [27], MLSVM [11] and ML-*k*NN [80]; these algorithms have been introduced briefly in Section 2. Note that these are single-instance algorithms that regard each image as a 135-dimensional feature vector, which is obtained by concatenating the nine

---

<sup>3</sup> The data set is available at [http://lamda.nju.edu.cn/data\\_MIMLimage.ashx](http://lamda.nju.edu.cn/data_MIMLimage.ashx).

instances in the direction from upper-left to right-bottom.

The parameter configurations of RANKSVM, MLSVM and ML- $k$ NN are set by considering the strategies adopted in [27], [11] and [80] respectively. For RANKSVM, polynomial kernel is used where polynomial degrees of 2 to 9 are considered as in [27] and chosen by hold-out tests on training sets. For MLSVM, Gaussian kernel is used. For ML- $k$ NN, the number of nearest neighbors considered is set to 10.

The boosting rounds of ADTBOOST.MH and MIMLBOOST are set to 25 and 50, respectively; The performance of the two algorithms at different boosting rounds is shown in Appendix B (Fig. B.1), it can be observed that at those rounds the performance of the algorithms have become stable. Gaussian kernel LIBSVM [16] is used for the Step 3a of MIMLBOOST. The MIMLSVM and MIMLSVM<sub>mi</sub> are also realized with Gaussian kernels. The parameter  $k$  of MIMLSVM is set to be 20% of the number of training images; The performance of this algorithm with different  $k$  values is shown in Appendix B (Fig. B.2), it can be observed that the setting of  $k$  does not significantly affect the performance of MIMLSVM. Note that in Appendix B (Figs. B.1 and B.2) we plot  $1 - \text{average precision}$ ,  $1 - \text{average recall}$  and  $1 - \text{average } F1$  such that in all the figures, the lower the curve, the better the performance.

Here in the experiments, 1,500 images are used as training examples while the remaining 500 images are used for testing. Experiments are repeated for thirty runs by using random training/test partitions, and the average and standard deviation are summarized in Table 1,<sup>4</sup> where the best performance on each criterion has been highlighted in boldface.

Pairwise  $t$ -tests with 95% significance level disclose that all the MIML algorithms are significantly better than ADTBOOST.MH and MLSVM on all the seven evaluation criteria. This is impressive since as mentioned before, these evaluation criteria measure the learning performance from different aspects and one algorithm rarely outperforms another algorithm on all criteria. MIMLSVM and MIMLSVM<sub>mi</sub> are both significantly better than RANKSVM on all the evaluation criteria, while MIMLBOOST and MIMLNN are both significantly better than RANKSVM on the

---

<sup>4</sup> For the shared implementation of ADTBOOST.MH ([http://www.grappa.univ-lille3.fr/grappa/en\\_index.php3?info=software](http://www.grappa.univ-lille3.fr/grappa/en_index.php3?info=software)), *ranking loss*, *average recall* and *average F1* are not available in the program's outputs.

Table 1

Results (mean $\pm$ std.) on scene classification data set ( $\downarrow$  indicates ‘the smaller the better’;  $\uparrow$  indicates ‘the larger the better’)

Compared Algorithms	Evaluation Criteria						
	<i>hloss</i> $\downarrow$	<i>one-error</i> $\downarrow$	<i>coverage</i> $\downarrow$	<i>rloss</i> $\downarrow$	<i>aveprec</i> $\uparrow$	<i>avercl</i> $\uparrow$	<i>aveF1</i> $\uparrow$
MIMLBOOST	.193 $\pm$ .007	.347 $\pm$ .019	<b>.984<math>\pm</math>.049</b>	<b>.178<math>\pm</math>.011</b>	.779 $\pm$ .012	.433 $\pm$ .027	.556 $\pm$ .023
MIMLSVM	189 $\pm$ .009	.354 $\pm$ .022	1.087 $\pm$ .047	.201 $\pm$ .011	.765 $\pm$ .013	.556 $\pm$ .020	.644 $\pm$ .018
MIMLSVM <sub>mi</sub>	.195 $\pm$ .008	<b>.317<math>\pm</math>.018</b>	1.068 $\pm$ .052	.197 $\pm$ .011	<b>.783<math>\pm</math>.011</b>	<b>.587<math>\pm</math>.019</b>	<b>.671<math>\pm</math>.015</b>
MIMLNN	<b>.185<math>\pm</math>.008</b>	.351 $\pm$ .026	1.057 $\pm$ .054	.196 $\pm$ .013	.771 $\pm$ .015	.509 $\pm$ .022	.613 $\pm$ .020
ADTBOOST.MH	.211 $\pm$ .006	.436 $\pm$ .019	1.223 $\pm$ .050	N/A	.718 $\pm$ .012	N/A	N/A
RANKSVM	.210 $\pm$ .024	.395 $\pm$ .075	1.161 $\pm$ .154	.221 $\pm$ .040	.746 $\pm$ .044	.529 $\pm$ .068	.620 $\pm$ .059
MLSVM	.232 $\pm$ .004	.447 $\pm$ .023	1.217 $\pm$ .054	.233 $\pm$ .012	.712 $\pm$ .013	.073 $\pm$ .010	.132 $\pm$ .017
ML-kNN	.191 $\pm$ .006	.370 $\pm$ .017	1.085 $\pm$ .048	.203 $\pm$ .010	.759 $\pm$ .011	.407 $\pm$ .026	.529 $\pm$ .023

first five criteria. MIMLNN is significantly better than ML-kNN on all the evaluation criteria. Both MIMLBOOST and MIMLSVM<sub>mi</sub> are significantly better than ML-kNN on all criteria except *hamming loss*. MIMLSVM is significantly better than ML-kNN on *one-error*, *average precision*, *average recall* and *average F1*, while there are ties on the other criteria. Moreover, note that the best performance on all evaluation criteria are always attained by MIML algorithms. Overall, comparison on the scene classification task shows that the MIML algorithms can be significantly better than the non-MIML algorithms; this validates the powerfulness of the MIML framework.

#### 4.3.3 Text Categorization

The REUTERS-21578 data set is used in this experiment. The seven most frequent categories are considered. After removing documents that do not have labels or main texts, and randomly removing some documents that have only one label, a data set containing 2,000 documents is obtained, where over 14.9% documents have multiple labels. Each document is represented as a bag of instances according to the method used in [3]. Briefly, the instances are obtained by splitting each document into passages using overlapping windows of maximal 50 words each. As a result, there are 2,000 bags and the number of instances in each bag varies from 2 to 26 (3.6 on average). The instances are represented based on term frequency. The words with high frequencies are considered, excluding “function words” that

Table 2

Results (mean $\pm$ std.) on text categorization data set ( $\downarrow$  indicates ‘the smaller the better’;  $\uparrow$  indicates ‘the larger the better’)

Compared Algorithms	Evaluation Criteria						
	<i>hloss</i> $\downarrow$	<i>one-error</i> $\downarrow$	<i>coverage</i> $\downarrow$	<i>rloss</i> $\downarrow$	<i>aveprec</i> $\uparrow$	<i>averecl</i> $\uparrow$	<i>aveF1</i> $\uparrow$
MIMLBOOST	.053 $\pm$ .004	.094 $\pm$ .014	.387 $\pm$ .037	.035 $\pm$ .005	.937 $\pm$ .008	.792 $\pm$ .010	.858 $\pm$ .008
MIMLSVM	<b>.033<math>\pm</math>.003</b>	.066 $\pm$ .011	.313 $\pm$ .035	.023 $\pm$ .004	.956 $\pm$ .006	<b>.925<math>\pm</math>.010</b>	.940 $\pm$ .008
MIMLSVM <sub>mi</sub>	.041 $\pm$ .004	<b>.055<math>\pm</math>.009</b>	<b>.284<math>\pm</math>.030</b>	<b>.020<math>\pm</math>.003</b>	<b>.965<math>\pm</math>.005</b>	.921 $\pm$ .012	<b>.942<math>\pm</math>.007</b>
MIMLNN	.038 $\pm$ .002	.080 $\pm$ .010	.320 $\pm$ .030	.025 $\pm$ .003	.950 $\pm$ .006	.834 $\pm$ .011	.888 $\pm$ .008
ADTBOOST.MH	.055 $\pm$ .005	.120 $\pm$ .017	.409 $\pm$ .047	N/A	.926 $\pm$ .011	N/A	N/A
RANKSVM	.120 $\pm$ .013	.196 $\pm$ .126	.695 $\pm$ .466	.085 $\pm$ .077	.868 $\pm$ .092	.411 $\pm$ .059	.556 $\pm$ .068
MLSVM	.050 $\pm$ .003	.081 $\pm$ .011	.329 $\pm$ .029	.026 $\pm$ .003	.949 $\pm$ .006	.777 $\pm$ .016	.854 $\pm$ .011
ML- <i>k</i> NN	.049 $\pm$ .003	.126 $\pm$ .012	.440 $\pm$ .035	.045 $\pm$ .004	.920 $\pm$ .007	.821 $\pm$ .021	.867 $\pm$ .013

have been removed from the vocabulary using the SMART stop-list [55]. It has been found that based on document frequency, the dimensionality of the data set can be reduced to 1-10% without loss of effectiveness [73]. Thus, we use the top 2% frequent words, and therefore each instance is a 243-dimensional feature vector.<sup>5</sup>

The parameter configurations of RANKSVM, MLSVM and ML-*k*NN are set in the same way as in Section 4.3.2. The boosting rounds of ADTBOOST.MH and MIMLBOOST are set to 25 and 50, respectively. Linear kernels are used. The parameter *k* of MIMLSVM is set to be 20% of the number of training images. The single-instance algorithms regard each document as a 243-dimensional feature vector which is obtained by aggregating all the instances in the same bag; this is equivalent to represent the document using a sole term frequency feature vector.

Here in the experiments, 1,500 documents are used as training examples while the remaining 500 documents are used for testing. Experiments are repeated for thirty runs by using random training/test partitions, and the average and standard deviation are summarized in Table 2, where the best performance on each criterion has been highlighted in boldface.

Pairwise *t*-tests with 95% significance level disclose that, impressively, both MIMLSVM and MIMLSVM<sub>mi</sub> are significantly better than all the non-MIML algorithms. MIMLNN is significantly better than ADTBOOST.MH, RANKSVM, and ML-*k*NN on all the

<sup>5</sup> The data set is available at [http://lamda.nju.edu.cn/data\\_MIMLtext.ashx](http://lamda.nju.edu.cn/data_MIMLtext.ashx)

evaluation criteria; significantly better than MLSVM on *hamming loss*, *average recall* and *average F1* while there are ties on the other criteria. MIMLBOOST is significantly better than ADTBOOST.MH on all criteria except that there is a tie on *hamming loss*; significantly better than RANKSVM on all criteria; significantly better than MLSVM on *average recall* and there is a tie on *average F1*; significantly better than ML-*k*NN on *one-error*, *coverage*, *ranking loss* and *average precision*. Moreover, note that the best performance on all evaluation criteria are always attained by MIML algorithms. Overall, comparison on the text categorization task shows that the MIML algorithms are better than the non-MIML algorithms; this validates the powerfulness of the MIML framework.

## 5 Solving MIML Problems by Regularization

The degeneration methods presented in Section 4 may lose information during the degeneration process, and thus a “direct” MIML algorithm is desirable. In this section we propose a regularization method for MIML. In contrast to MIMLSVM and MIMLSVM<sub>*mi*</sub>, this method is developed from the regularization framework directly and so we call it D-MIMLSVM. The basic assumption of D-MIMLSVM is that the labels associated to the same example have some relatedness, and the performance of classifying the bags depends on the loss between the labels and the predictions on the bags as well as on the constituent instances. Moreover, considering that for any class label the number of positive examples is smaller than that of negative examples, this method incorporates a mechanism to deal with class imbalance. We employ the constrained concave-convex procedure (CCCP) which has well-studied convergence properties [62] to solve the resultant non-convex optimization problem. We also present a cutting plane algorithm that finds the solution efficiently.

### 5.1 The Loss Function

Given a set of MIML training examples  $\{(X_1, Y_1), (X_2, Y_2), \dots, (X_m, Y_m)\}$ , the goal of D-MIMLSVM is to learn a mapping  $\mathbf{f} : 2^{\mathcal{X}} \rightarrow 2^{\mathcal{Y}}$  where the proper label set for each bag  $X \subseteq \mathcal{X}$  corresponds to  $\mathbf{f}(X) \subseteq \mathcal{Y}$ . Specifically, D-MIMLSVM chooses to instantiate  $\mathbf{f}$  with  $T$  functions, i.e.  $\mathbf{f} = (f_1, f_2, \dots, f_T)$ , where  $T$  is the number of

labels in the label space  $\mathcal{Y} = \{l_1, l_2, \dots, l_T\}$ . Here, the  $t$ -th function  $f_t : 2^{\mathcal{X}} \rightarrow \mathcal{R}$  determines the belongingness of  $l_t$  for  $X$ , i.e.  $\mathbf{f}(X) = \{l_t \mid f_t(X) > 0, 1 \leq t \leq T\}$ . In addition, each single instance  $\mathbf{x} \in \mathcal{X}$  in a bag  $X$  can be viewed as a bag  $\{\mathbf{x}\}$  containing only one instance, such that  $\mathbf{f}(\{\mathbf{x}\}) = (f_1(\{\mathbf{x}\}), f_2(\{\mathbf{x}\}), \dots, f_T(\{\mathbf{x}\}))$  is also a well-defined function. For convenience,  $\mathbf{f}(\{\mathbf{x}\})$  and  $f_t(\{\mathbf{x}\})$  are simplified as  $\mathbf{f}(\mathbf{x})$  and  $f_t(\mathbf{x})$  in the rest of this section.

To train the component functions  $f_t (1 \leq t \leq T)$  in  $\mathbf{f}$ , D-MIMLSVM employs the following empirical loss function  $V$  involving two terms (balanced by  $\lambda$ ):

$$V(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \mathbf{f}) = V_1(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \mathbf{f}) + \lambda \cdot V_2(\{X_i\}_{i=1}^m, \mathbf{f}) \quad (3)$$

Here, the first term  $V_1$  considers the loss between the ground-truth label set of each training bag  $X_i$ , i.e.  $Y_i$ , to its predicted label set, i.e.  $\mathbf{f}(X_i)$ . Let  $y_{it} = 1$  if  $l_t \in Y_i$  holds ( $1 \leq i \leq m, 1 \leq t \leq T$ ). Otherwise,  $y_{it} = -1$ . Furthermore, let  $(z)_+ = \max(0, z)$  denote the hinge loss function. Accordingly, the first loss term  $V_1$  is defined as:

$$V_1(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \mathbf{f}) = \frac{1}{mT} \sum_{i=1}^m \sum_{t=1}^T (1 - y_{it} f_t(X_i))_+ \quad (4)$$

The second term  $V_2$  considers the loss between  $\mathbf{f}(X_i)$  and the predictions of  $X_i$ 's constituent instances, i.e.  $\{\mathbf{f}(\mathbf{x}_{ij}) \mid 1 \leq j \leq n_i\}$ , which reflects the relationships between the bag  $X_i$  and its instances  $\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{i, n_i}\}$ . Here, the common assumption in multi-instance learning is that the strength for  $X_i$  to hold a label is equal to the maximum strength for its instances to hold the label, i.e.  $f_t(X_i) = \max_{j=1, \dots, n_i} f_t(\mathbf{x}_{ij})$ .<sup>6</sup> Accordingly, the second loss term  $V_2$  is defined as:

$$V_2(\{X_i\}_{i=1}^m, \mathbf{f}) = \frac{1}{mT} \sum_{i=1}^m \sum_{t=1}^T l \left( f_t(X_i), \max_{j=1, \dots, n_i} f_t(\mathbf{x}_{ij}) \right) \quad (5)$$

Here,  $l(v_1, v_2)$  can be defined in various ways and is set to be the  $l_1$  loss in this paper, i.e.  $l(v_1, v_2) = |v_1 - v_2|$ . By combining Eq. 4 and Eq. 5, the empirical loss function  $V$  in Eq. 3 is then specified as:

---

<sup>6</sup> Note that this assumption may be restrictive to some extent. There are many cases where the label of the bag does not rely on the instance with the maximum predictions, as discussed in Section 2. In addition, in classification only the sign of prediction is important [19], i.e.  $\text{sign}(f_t(X_i)) = \text{sign}(\max_{j=1, \dots, n_i} f_t(\mathbf{x}_{ij}))$ . However, in this paper the above common assumption is still adopted due to its popularity and simplicity.

$$\begin{aligned}
V(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \mathbf{f}) &= \frac{1}{mT} \sum_{i=1}^m \sum_{t=1}^T (1 - y_{it} f_t(X_i))_+ \\
&\quad + \frac{\lambda}{mT} \sum_{i=1}^m \sum_{t=1}^T l\left(f_t(X_i), \max_{j=1, \dots, n_i} f_t(\mathbf{x}_{ij})\right)
\end{aligned} \tag{6}$$

## 5.2 Representer Theorem for MIML

For simplicity, we assume that each function  $f_t$  is a linear model, i.e.,  $f_t(\mathbf{x}) = \langle \mathbf{w}_t, \phi(\mathbf{x}) \rangle$  where  $\phi$  is the feature map induced by a kernel function  $k$  and  $\langle \cdot, \cdot \rangle$  denotes the standard inner product in the Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{H}$  induced by the kernel  $k$ . We recall that an instance can be regarded as a bag containing only one instance, so the kernel  $k$  can be any kernel defined on a set of instances, such as the *set kernel* [32]. In the case of classification, objects (bags or instances) are classified according to the sign of  $f_t$ .

D-MIMLSVM assumes that the labels associated with a bag should have some relatedness; otherwise they should not be associated with the bag simultaneously. To reflect this basic assumption, D-MIMLSVM regularizes the empirical loss function in Eq. 6 with an additional term  $\Omega(\mathbf{f})$ :

$$\Omega(\mathbf{f}) + \gamma \cdot V(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \mathbf{f}) \tag{7}$$

Here,  $\gamma$  is a regularization parameter balancing the model complexity  $\Omega(\mathbf{f})$  and the empirical risk  $V$ . Inspired by [28], we assume that the relatedness among the labels can be measured by the mean function  $\mathbf{w}_0$ ,

$$\mathbf{w}_0 = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t \tag{8}$$

The original idea in [28] is to minimize  $\sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}_0\|^2$  and meanwhile minimize  $\|\mathbf{w}_0\|^2$ , i.e. to set the regularizer as:

$$\Omega(\mathbf{f}) = \frac{1}{T} \sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}_0\|^2 + \eta \|\mathbf{w}_0\|^2 \tag{9}$$

According to Eq.8, the first term in the RHS of Eq. 9 can be rewritten as:

$$\frac{1}{T} \sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}_0\|^2 = \frac{1}{T} \sum_{t=1}^T \|\mathbf{w}_t\|^2 - \|\mathbf{w}_0\|^2 \tag{10}$$

Therefore, by substituting Eq. 10 into Eq. 9, the regularizer can be simplified as:

$$\Omega(\mathbf{f}) = \frac{1}{T} \sum_{t=1}^T \|\mathbf{w}_t\|^2 + \mu \|\mathbf{w}_0\|^2 \quad (11)$$

Further note that  $\|\mathbf{w}_t\|^2 = \|f_t\|_{\mathcal{H}}^2$  and  $\|\mathbf{w}_0\|^2 = \|\frac{\sum_{t=1}^T f_t}{T}\|_{\mathcal{H}}^2$ , by substituting Eq. 11 into Eq. 7, we have the regularization framework of D-MIMLSVM as follows:

$$\min_{\mathbf{f} \in \mathcal{H}} \frac{1}{T} \sum_{t=1}^T \|f_t\|_{\mathcal{H}}^2 + \mu \left\| \frac{\sum_{t=1}^T f_t}{T} \right\|_{\mathcal{H}}^2 + \gamma \cdot V(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \mathbf{f}) \quad (12)$$

Here,  $\mu$  is a parameter to trade off the discrepancy and commonness among the labels, that is, how similar or dissimilar the  $\mathbf{w}_t$ 's are. Refer to Eq. 10, we have  $\Omega(\mathbf{f}) = \frac{1}{T} \sum_{t=1}^T \|f_t\|_{\mathcal{H}}^2 + \mu \left\| \frac{\sum_{t=1}^T f_t}{T} \right\|_{\mathcal{H}}^2 = \frac{1}{T} \sum_{t=1}^T \|f_t - \frac{\sum_{t=1}^T f_t}{T}\|_{\mathcal{H}}^2 + (\mu + 1) \left\| \frac{\sum_{t=1}^T f_t}{T} \right\|_{\mathcal{H}}^2$ . Intuitively, when  $\mu + 1$  (or  $\mu$ ) is large, minimization of Eq. 12 will force  $\left\| \frac{\sum_{t=1}^T f_t}{T} \right\|_{\mathcal{H}}^2$  to tend to be zero and the discrepancy among the labels becomes more important; when  $\mu + 1$  (or  $\mu$ ) is small, minimization of Eq. 12 will force  $\|f_t - \frac{\sum_{t=1}^T f_t}{T}\|_{\mathcal{H}}^2$  to tend to be zero and the commonness among the labels becomes more important [28].

Given the above setup, we can prove the following representer theorem.

**Theorem 1** *The minimizer of the optimization problem 12 admits an expansion*

$$f_t(\mathbf{x}) = \sum_{i=1}^m \left( \alpha_{t,i0} k(\mathbf{x}, X_i) + \sum_{j=1}^{n_i} \alpha_{t,ij} k(\mathbf{x}, \mathbf{x}_{ij}) \right)$$

where all  $\alpha_{t,i0}, \alpha_{t,ij} \in \mathcal{R}$ .

**Proof.** Analogous to [28], we first introduce a combined feature map

$$\Psi(\mathbf{x}, t) = \left( \frac{\phi(\mathbf{x})}{\sqrt{r}}, \underbrace{\mathbf{0}, \dots, \mathbf{0}}_{t-1}, \phi(\mathbf{x}), \underbrace{\mathbf{0}, \dots, \mathbf{0}}_{T-t} \right)$$

and its decision function, i.e.,  $\hat{f}(\mathbf{x}, t) = \langle \hat{\mathbf{w}}, \Psi(\mathbf{x}, t) \rangle$  where

$$\hat{\mathbf{w}} = (\sqrt{r}\mathbf{w}_0, \mathbf{w}_1 - \mathbf{w}_0, \dots, \mathbf{w}_T - \mathbf{w}_0).$$

Here  $r = \mu T + T$ . Let  $\hat{k}$  denote the kernel function induced by  $\Psi$  and  $\hat{\mathcal{H}}$  is its corresponding RKHS. We have Eqs. 13 and 14.

$$\hat{f}(\mathbf{x}, t) = \langle \hat{\mathbf{w}}, \Psi(\mathbf{x}, t) \rangle = \langle (\mathbf{w}_0 + \mathbf{w}_t - \mathbf{w}_0), \phi(\mathbf{x}) \rangle = \langle \mathbf{w}_t, \phi(\mathbf{x}) \rangle = f_t(\mathbf{x}) \quad (13)$$

$$\|\hat{f}\|_{\hat{\mathcal{H}}}^2 = \|\hat{\mathbf{w}}\|^2 = \sum_{i=1}^T \|\mathbf{w}_t - \mathbf{w}_0\|^2 + r\|\mathbf{w}_0\|^2 = \sum_{i=1}^T \|\mathbf{w}_t\|^2 + \mu T\|\mathbf{w}_0\|^2 \quad (14)$$

Therefore, loss function in Eq.6 can be represented by  $\hat{V}(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \hat{f})$ , i.e.,

$$\begin{aligned} \hat{V}(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \hat{f}) &= \frac{1}{mT} \sum_{i=1}^m \sum_{t=1}^T \left(1 - y_{it} \hat{f}(X_i, t)\right)_+ \\ &+ \frac{\lambda}{mT} \sum_{i=1}^m \sum_{t=1}^T l\left(\hat{f}(X_i, t), \max_{j=1, \dots, n_i} \hat{f}(\mathbf{x}_{ij}, t)\right). \end{aligned} \quad (15)$$

Thus, Eq. 12 is equivalent to

$$\min_{\hat{f} \in \hat{\mathcal{H}}} \frac{1}{T} \|\hat{f}\|_{\hat{\mathcal{H}}}^2 + \gamma \hat{V}(\{X_i\}_{i=1}^m, \{Y_i\}_{i=1}^m, \hat{f}). \quad (16)$$

Note that  $\|\hat{f}\|_{\hat{\mathcal{H}}}^2 : [0, \infty) \rightarrow \mathcal{R}$  is a strictly monotonically increasing function. According to representer theorem (Theorem 4.2 in [57]), each minimizer  $\hat{f}$  of the functional risk in Eq. 16 admits a representation of the form

$$\hat{f}(\mathbf{x}, t) = \sum_{t=1}^T \sum_{i=1}^m \left( \beta_{t,i0} \hat{k}((X_i, t), (\mathbf{x}, t)) + \sum_{j=1}^{n_i} \beta_{t,ij} \hat{k}((\mathbf{x}_{ij}, t), (\mathbf{x}, t)) \right), \quad (17)$$

where  $\beta_{t,ij} \in \mathcal{R}$  and the corresponding weight vector  $\hat{\mathbf{w}}$  is represented as

$$\hat{\mathbf{w}} = \sum_{t=1}^T \sum_{i=1}^m \left( \beta_{t,i0} \Psi(X_i, t) + \sum_{j=1}^{n_i} \beta_{t,ij} \Psi(\mathbf{x}_{ij}, t) \right). \quad (18)$$

Finally, with Eqs. 13 and 18, we have

$$\begin{aligned} f_t(\mathbf{x}) &= \langle \mathbf{w}_t, \phi(\mathbf{x}) \rangle = \langle \mathbf{w}, \Psi(\mathbf{x}, t) \rangle \\ &= \sum_{i=1}^m \left( \alpha_{t,i0} k(\mathbf{x}, X_i) + \sum_{j=1}^{n_i} \alpha_{t,ij} k(\mathbf{x}, \mathbf{x}_{ij}) \right) \end{aligned} \quad (19)$$

where  $\alpha_{t,ij} = \frac{1}{\sqrt{r}}(\sum_t \beta_{t,ij}) + \beta_{t,ij}/r$ . □

Note that  $\mathbf{x}$  in Eq. 19 can be regarded not only as a bag  $X_i$  but also an instance  $\mathbf{x}_{ij}$ . In other words, both  $f_t(X_i)$  and  $f_t(\mathbf{x}_{ij})$  can be obtained by Eq. 19.

### 5.3 Optimization

Considering the use of  $l_1$  loss for  $l(v_1, v_2)$ , Eq.12 can be re-written as

$$\begin{aligned}
& \min_{\mathbf{f} \in \mathcal{H}, \boldsymbol{\xi}, \boldsymbol{\delta}} \frac{1}{T} \sum_{t=1}^T \|f_t\|_{\mathcal{H}}^2 + \mu \left\| \frac{\sum_{t=1}^T f_t}{T} \right\|_{\mathcal{H}}^2 + \frac{\gamma}{mT} \boldsymbol{\xi}' \mathbf{1} + \frac{\gamma\lambda}{mT} \boldsymbol{\delta}' \mathbf{1} \\
& \text{s.t. } y_{it} f_t(X_i) \geq 1 - \xi_{it}, \\
& \quad \boldsymbol{\xi} \geq \mathbf{0}, \\
& \quad -\delta_{it} \leq f_t(X_i) - \max_{j=1, \dots, n_i} f_t(\mathbf{x}_{ij}) \leq \delta_{it} \quad \forall i = 1, \dots, m, \quad t = 1, \dots, T \quad (20)
\end{aligned}$$

where  $\boldsymbol{\xi} = [\xi_{11}, \xi_{12}, \dots, \xi_{it}, \dots, \xi_{mT}]'$  are slack variables for the errors on the training bags for each label,  $\boldsymbol{\delta} = [\delta_{11}, \delta_{12}, \dots, \delta_{it}, \dots, \delta_{mT}]'$ , and  $\mathbf{0}$  and  $\mathbf{1}$  are all-zero and all-one vector, respectively.

Without loss of generality, assume that the bags and instances are ordered as  $(X_1, \dots, X_m, \mathbf{x}_{11}, \dots, \mathbf{x}_{1, n_1}, \dots, \mathbf{x}_{m,1}, \dots, \mathbf{x}_{m, n_m})$ . Thus, each object (bag or instance) in the training set can then be indexed by the following function  $\mathcal{I}$ , i.e.,

$$\begin{cases} \mathcal{I}(X_i) = i \\ \mathcal{I}(\mathbf{x}_{ij}) = m + \sum_{l=1}^{i-1} n_l + j \end{cases}$$

for  $j = 1, \dots, n_i$  and  $i = 1, \dots, m$ . With this ordering, we can obtain the  $(m+n) \times (m+n)$  kernel matrix  $\mathbf{K}$  defined on all objects in the training set, where  $n = \sum_{i=1}^m n_i$ . Denote the  $i$ -th column of  $\mathbf{K}$  by  $\mathbf{k}_i$ . According to theorem 1, we have  $f_t(X_i) = \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t + b_t$  and  $f_t(\mathbf{x}_{ij}) = \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t + b_t$ . Here, the bias  $b_t$  for each label is included.

According to definition of  $f_t$  in Eq. 19, Eq. 20 can be cast as the optimization problem

$$\begin{aligned}
& \min_{\mathbf{A}, \boldsymbol{\xi}, \boldsymbol{\delta}, \mathbf{b}} \frac{1}{2T} \sum_{t=1}^T \boldsymbol{\alpha}'_t \mathbf{K} \boldsymbol{\alpha}_t + \frac{\mu}{T^2} \mathbf{1}' \mathbf{A}' \mathbf{K} \mathbf{A} \mathbf{1} + \frac{\gamma}{mT} \boldsymbol{\xi}' \mathbf{1} + \frac{\gamma\lambda}{mT} \boldsymbol{\delta}' \mathbf{1} \quad (21) \\
& \text{s.t. } y_{it} (\mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t + b_t) \geq 1 - \xi_{it}, \\
& \quad \boldsymbol{\xi} \geq \mathbf{0}, \\
& \quad \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t - \delta_{it} \leq \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t, \\
& \quad \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t - \max_{j=1, \dots, n_i} \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t \leq \delta_{it},
\end{aligned}$$

where  $\mathbf{A} = [\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \dots, \boldsymbol{\alpha}_T]$  and  $\mathbf{b} = [b_1, b_2, \dots, b_T]'$ .

The above optimization problem is a non-convex optimization problem since the last constraint is non-convex. Note that this non-convex constraint is a difference between two convex functions, and thus the optimization problem can be solved

by CCCP [19, 62], which is one of the most standard techniques to solve such kind of non-convex optimization problems. CCCP is guaranteed to converge to a local minimum [75], and in many cases it can even converge to a global solution [25].

Here, for solving the optimization problem 21, CCCP works by solving a sequential convex quadratic problems. Concretely, given the initial subgradient  $\sum_{j=1}^{n_i} \rho_{ijt} \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t$  of  $\max_{j=1, \dots, n_i} \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t$ , we solve the following convex quadratic optimization (QP) problem

$$\begin{aligned} \min_{\mathbf{A}, \boldsymbol{\xi}, \delta, \mathbf{b}} \quad & \frac{1}{2T} \sum_{t=1}^T \boldsymbol{\alpha}'_t \mathbf{K} \boldsymbol{\alpha}_t + \frac{\mu}{T^2} \mathbf{1}' \mathbf{A}' \mathbf{K} \mathbf{A} \mathbf{1} + \frac{\gamma}{mT} \boldsymbol{\xi}' \mathbf{1} + \frac{\gamma \lambda}{mT} \boldsymbol{\delta}' \mathbf{1} \\ \text{s.t.} \quad & y_{it} (\mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t + b_t) \geq 1 - \xi_{it}, \\ & \boldsymbol{\xi} \geq \mathbf{0}, \\ & \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t - \delta_{it} \leq \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t, \\ & \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t - \sum_{j=1}^{n_i} \rho_{ijt} \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t \leq \delta_{it}. \end{aligned} \quad (22)$$

Then, in the next iteration we update  $\rho_{ijk}$  according to

$$\rho_{ijt} = \begin{cases} = 0, & \text{if } \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t \neq \max_{k=1, \dots, n_i} (\mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ik})} \boldsymbol{\alpha}_t), \\ = 1/n_d, & \text{otherwise,} \end{cases}$$

where  $n_d$  is the number of active  $\mathbf{x}_{ij}$ 's. It holds  $\sum_{j=1}^{n_i} \rho_{ijt} = 1$  for any  $t$ 's. The iteration continues and this procedure is guaranteed to converge to a local minimum.

#### 5.4 Handling Class-Imbalance

The above solution may be improved further if we explicitly take into account the instance-level class-imbalance, that is, for any class label the number of *positive* instances is smaller than the number of *negative* instances in MIML problems.

We can roughly estimate the *imbalance rate*, which is the ratio of the number of positive instances to that of negative instances, for each class label using the strategy adopted by [41]. In detail, for a specific label  $y \in \mathcal{Y}$ , we can divide the training bags  $\{(X_1, Y_1), (X_2, Y_2), \dots, (X_m, Y_m)\}$  into two subsets,  $A_1 = \{(X_i, Y_i) | y \in Y_i\}$  and  $A_2 = \{(X_i, Y_i) | y \notin Y_i\}$ . It is obvious that all the instances in  $A_2$  are negative to  $y$ . Then, for every  $(X_i, Y_i)$  in  $A_1$ , assuming that the instances of different labels is roughly equally distributed, the number of positive instances of  $y$  in  $(X_i, Y_i)$  is

roughly  $n_i \times \frac{1}{|Y_i|}$  where  $|Y_i|$  returns the number of labels in  $Y_i$ . Thus, the imbalance rate of  $y$  is:

$$ibr(y) = \sum_{\substack{i=1 \\ y \in Y_i}}^m \frac{n_i}{|Y_i|} \times \frac{1}{\sum_{i=1}^m n_i} = \sum_{\substack{i=1 \\ y \in Y_i}}^m \frac{n_i}{n \times |Y_i|}.$$

There are many class-imbalance learning methods [69]. One of the most popular and effective methods is *rescaling* [87], which can be incorporated into our framework easily. In short, after obtaining the estimated imbalance rate for every class label, we can use these rates to modulate the loss caused by different misclassifications.

In detail,  $\boldsymbol{\xi}$  in Eq. 22 is directly related to the hinge loss  $(1 - y_{it}f_t(X_i))_+$ . According to the rescaling method [87], without loss of generality, we can rewrite the loss function into Eq. 23.

$$\left( \frac{y_{it} + 1}{2} - y_{it} \times ibr(y_{it}) \right) (1 - y_{it}f_t(X_i)). \quad (23)$$

Let  $\boldsymbol{\tau} = [\tau_{11}, \tau_{12}, \dots, \tau_{it}, \dots, \tau_{mT}]$ , where  $\tau_{it} = \left( \frac{y_{it} + 1}{2} - y_{it} \times ibr(y_{it}) \right)$ . Then, to minimize the loss defined in Eq. 23, Eq. 22 becomes Eq. 24. Here  $\boldsymbol{\xi}'\boldsymbol{\tau}$  indicates the weighted loss after considering the instance-level class-imbalance. It is evident that the problem in Eq. 24 is still a standard QP problem.

$$\begin{aligned} \min_{\mathbf{A}, \boldsymbol{\xi}, \delta, \mathbf{b}} \quad & \frac{1}{2T} \sum_{t=1}^T \boldsymbol{\alpha}_t' \mathbf{K} \boldsymbol{\alpha}_t + \frac{\mu}{T^2} \mathbf{1}' \mathbf{A}' \mathbf{K} \mathbf{A} \mathbf{1} + \frac{\gamma}{mT} \boldsymbol{\xi}' \boldsymbol{\tau} + \frac{\gamma\lambda}{mT} \boldsymbol{\delta}' \mathbf{1} \\ \text{s.t.} \quad & y_{it}(\mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t + b_t) \geq 1 - \xi_{it}, \\ & \boldsymbol{\xi} \geq \mathbf{0}, \\ & \mathbf{k}'_{\mathcal{I}(x_{ij})} \boldsymbol{\alpha}_t - \delta_{it} \leq \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t, \\ & \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t - \sum_{j=1}^{n_i} \rho_{ijt} \mathbf{k}'_{\mathcal{I}(x_{ij})} \boldsymbol{\alpha}_t \leq \delta_{it}. \end{aligned} \quad (24)$$

### 5.5 Efficient Algorithm

Eq. 24 is a large-scale quadratic programming problem that involves many constraints and variables. To make it tractable and scalable, and observing that most of the constraints in Eq. 24 are redundant, we present an efficient algorithm which constructs a nested sequence of tighter relaxations of the original problem using the cutting plane method [40].

Similar to its use with structured prediction [64], we add a constraint (or a cut) that is most violated by the current solution, and then find the solution in the updated feasible region. Such a procedure will converge to an optimal (or  $\varepsilon$ -suboptimal) solution of the original problem. Moreover, Eq. 24 supports a natural problem decomposition since its constraint matrix is a block diagonal matrix, i.e., each block corresponds to one label.

The pseudo-code of the algorithm is summarized in Appendix A (Table A.3). We first initialize the working sets  $S_t$ 's as empty sets and the solutions as all zeros (Line 1). Then, instead of testing all the constraints, which is rather expensive when there are lots of constraints, we use the speedup heuristic as described in [61], i.e., we use  $p$  constraints to approximate the whole constraints (Line 4). Smola and Schölkopf [61] have shown that when  $p$  is larger than 59, the selected violated constraint is with probability 0.95 among the 5% most violated constraints among all constraints. The  $Loss_i$  (Line 5) is calculated as  $\max\{0, \mathbf{u}'\mathbf{x} - d\}$  where  $\mathbf{u}$  and  $d$  are the linear coefficients and bias of the  $i$ -th linear constraint, respectively. If the maximal  $Loss$  is lower than the given stopping criteria  $\varepsilon$  (we simply set  $\varepsilon$  as  $10^{-4}$  in our experiments), no update will be taken for the working set  $S_t$ ; otherwise the constraint with the maximal  $Loss$  will be added into  $S_t$  (lines 8 and 9). Once a new constraint is added, the solution will be re-computed with respect to  $S_t$  via solving a smaller quadratic program problem (line 10). The algorithm stops when there is no update for all  $S_t$ 's.

## 5.6 Experiments

The previous experiments in Section 4.3 have shown that different MIML algorithms have different advantages on different performance measures. In this section we propose the D-MIMLSVM algorithm. We do not claim that D-MIMLSVM is the best MIML algorithm. What we want to show is that, in contrast to heuristically solving the MIML problem by degeneration, developing algorithms from a regularization framework directly offers a better choice. So the most meaningful comparison is between the D-MIMLSVM, MIMLSVM and MIMLSVM<sub>*mi*</sub> algorithms, the latter two not being derived from the regularization framework directly.

To study the behavior of D-MIMLSVM, MIMLSVM and MIMLSVM<sub>*mi*</sub> under differ-

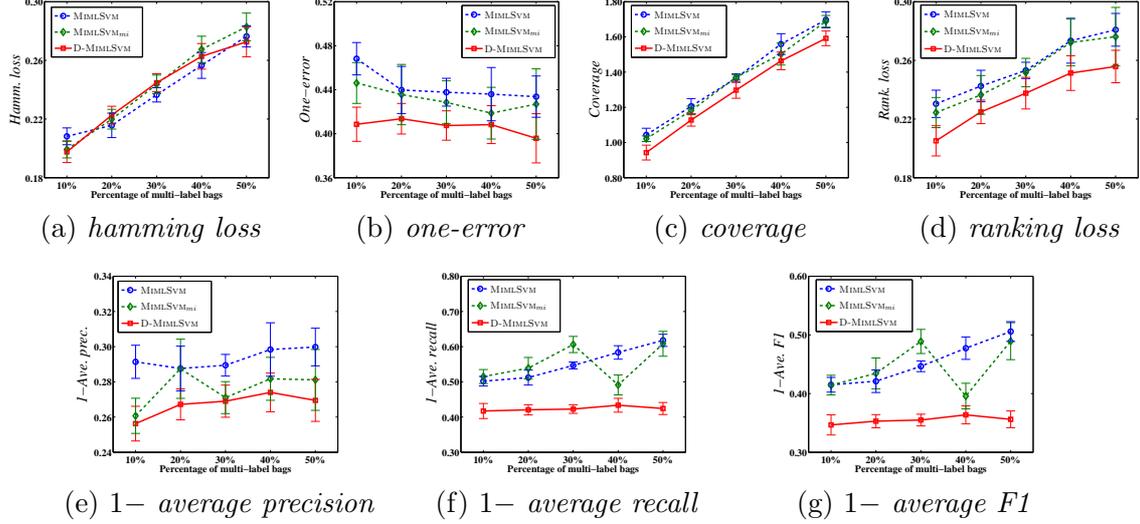


Fig. 4. Results on the scene classification data set with different percentage of multi-label data. The lower the curve, the better the performance.

ent amounts of multi-label data, we derive five data sets from the scene data used in Section 4.3.2. By randomly removing some single-label images, we obtain a data set where 30% (or 40%, or 50%) images belonging to multiple classes simultaneously; by randomly removing some multi-label images, we obtain a data set where 10% (or 20%) images belong to multiple classes simultaneously. A similar process is applied to the text data used in Section 4.3.3 to derive five data sets. On the derived data sets we use 25% data for training and the remaining 75% data for testing, and experiments are repeated for thirty runs with random training/test partitions. The parameters of D-MIMLSVM, MIMLSVM and MIMLSVM<sub>mi</sub> are all set by hold-out tests on training sets. Since D-MIMLSVM needs to solve a large optimization problem, although we have incorporated advanced mechanisms such as cutting-plane algorithm, the current D-MIMLSVM can only deal with moderate training set sizes.

The seven criteria introduced in Section 4.3.1 are used to evaluate the performance. The average and standard deviation are plotted in Figs. 4 and 5. Note that in the figures we plot 1–average precision, 1–average recall and 1–average F1 such that in all the figures, the lower the curve, the better the performance.

As shown in Figs. 4 and 5, the performance of D-MIMLSVM is better than those of MIMLSVM and MIMLSVM<sub>mi</sub> in most cases. Specifically, pairwise *t*-tests with

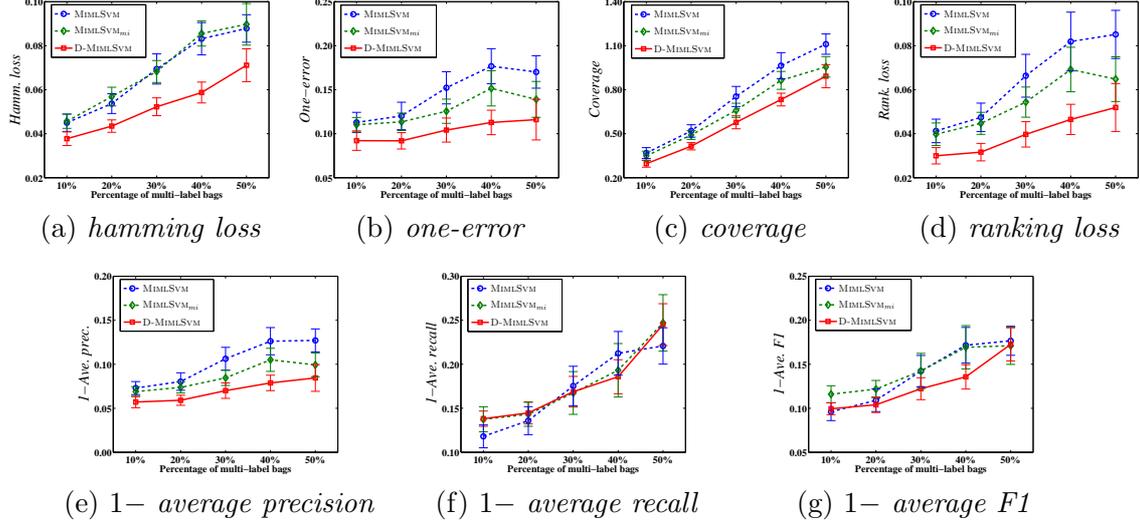


Fig. 5. Results on the text categorization data set with different percentage of multi-label data. The lower the curve, the better the performance.

95% significance level disclose that: a) On the scene classification task, among all the 35 configurations (7 evaluation criteria  $\times$  5 percentages of multi-label bags), the performance of D-MIMLSVM is superior to MIMLSVM and MIMLSVM<sub>mi</sub> in 88% and 80% cases, comparable to them in 6% and 20% cases, and inferior to them in only 6% and none cases; b) On the text categorization task, among all the 35 configurations, the performance of D-MIMLSVM is superior to MIMLSVM and MIMLSVM<sub>mi</sub> in 82% and 82% cases, comparable to them in 9% and 18% cases, and inferior to them in only 9% and none cases. The results suggest that D-MIMLSVM is a good choice for learning with moderate number of MIML examples.

### 5.7 Discussion

The regularization framework presented in this section has an important assumption, that is, all the class labels share some commonness, i.e., the  $\mathbf{w}_0$  in Eq. 8. This assumption makes the regularization easier to realize, however, it over-simplifies the real scenario. In fact, in real applications it is rare that all class labels share some commonness; it is more typical that some class labels share some commonness, but the commonness shared by different labels may be different. For example, class label  $y_1$  may share something with class label  $y_2$ , and  $y_2$  may share something with  $y_3$ , but maybe  $y_1$  shares nothing with  $y_3$ . So, a more reasonable assumption is

that different pairs of labels share different things (or even nothing). By considering this assumption, a more powerful method may be developed.

Actually, it is not difficult to modify the framework of Eq. 12 by replacing the role of  $\mathbf{w}_0$  by  $\mathbf{W}$  whose element  $\mathbf{W}_{ij}$  expresses the relatedness between the  $i$ -th and  $j$ -th class labels, that is,

$$\min \frac{1}{2T^2} \sum_{i,j} \|\mathbf{w}_i - \mathbf{W}_{ij}\|^2 + \frac{1}{T^2} \sum_{i,j} \mu_{ij} \|\mathbf{W}_{ij}\|^2 + \gamma \mathbf{V} . \quad (25)$$

Note that  $\mathbf{W}$  is a tensor and  $\mathbf{W}_{ij}$  is a vector.

To minimize Eq. 25, taking derivative to  $\mathbf{W}_{ij}$ , we have

$$-(\mathbf{w}_i - \mathbf{W}_{ij}) - (\mathbf{w}_j - \mathbf{W}_{ji}) + 2\mu_{ij}\mathbf{W}_{ij} + 2\mu_{ji}\mathbf{W}_{ji} = 0 .$$

Considering  $\mathbf{W}_{ij} = \mathbf{W}_{ji}$  and  $\mu_{ij} = \mu_{ji}$ , we have

$$-(\mathbf{w}_i - \mathbf{W}_{ij}) - (\mathbf{w}_j - \mathbf{W}_{ij}) + 4\mu_{ij}\mathbf{W}_{ij} = 0 ,$$

and so,

$$\mathbf{W}_{ij} = \frac{\mathbf{w}_i + \mathbf{w}_j}{4\mu_{ij} + 2} . \quad (26)$$

Put Eq. 26 into Eq. 25, we have

$$\min \frac{1}{2T^2} \sum_{i,j} \left\| \frac{(4\mu_{ij} + 1)\mathbf{w}_i - \mathbf{w}_j}{4\mu_{ij} + 2} \right\|^2 + \frac{1}{T^2} \sum_{i,j} \mu_{ij} \left\| \frac{\mathbf{w}_i + \mathbf{w}_j}{4\mu_{ij} + 2} \right\|^2 + \gamma \mathbf{V} . \quad (27)$$

After simplification, Eq. 25 becomes

$$\begin{aligned} \min \frac{1}{8T^2} \sum_{i,j} \left( \frac{16\mu_{ij}^2 + 10\mu_{ij} + 1}{(2\mu_{ij} + 1)^2} \|\mathbf{w}_i\|^2 + \frac{2\mu_{ij} + 1}{(2\mu_{ij} + 1)^2} \|\mathbf{w}_j\|^2 \right) \\ - \frac{1}{4T^2} \sum_{i,j} \frac{2\mu_{ij} + 1}{(2\mu_{ij} + 1)^2} \langle \mathbf{w}_i, \mathbf{w}_j \rangle + \gamma \mathbf{V} . \end{aligned}$$

So, the new optimization task becomes

$$\begin{aligned}
\min_{\mathbf{A}, \boldsymbol{\xi}, \boldsymbol{\delta}, \mathbf{b}} \quad & \frac{1}{8T^2} \sum_{i=1}^T \sum_{j=1}^T \left( \frac{16\mu_{ij}^2 + 10\mu_{ij} + 1}{(2\mu_{ij} + 1)^2} \boldsymbol{\alpha}'_i \mathbf{K} \boldsymbol{\alpha}_i + \frac{2\mu_{ij} + 1}{(2\mu_{ij} + 1)^2} \boldsymbol{\alpha}'_j \mathbf{K} \boldsymbol{\alpha}_j \right) \\
& - \frac{1}{4T^2} \sum_{i=1}^T \sum_{j=1}^T \frac{2\mu_{ij} + 1}{(2\mu_{ij} + 1)^2} \boldsymbol{\alpha}'_i \mathbf{K} \boldsymbol{\alpha}_j + \frac{\gamma}{mT} \boldsymbol{\xi}' \mathbf{1} + \frac{\gamma\lambda}{mT} \boldsymbol{\delta}' \mathbf{1} \\
\text{s.t.} \quad & y_{it}(\mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t + b_t) \geq 1 - \xi_{it}, \\
& \boldsymbol{\xi} \geq \mathbf{0}, \\
& \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t - \delta_{it} \leq \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t, \\
& \mathbf{k}'_{\mathcal{I}(X_i)} \boldsymbol{\alpha}_t - \max_{j=1, \dots, n_i} \mathbf{k}'_{\mathcal{I}(\mathbf{x}_{ij})} \boldsymbol{\alpha}_t \leq \delta_{it}.
\end{aligned} \tag{28}$$

By solving Eq. 28 we can get not only an MIML learner, but also some understanding on the relatedness between pairs of labels from  $\mathbf{W}_{ij}$ , and some understanding on the different importance of the  $\mathbf{W}_{ij}$ 's in determining the concerned class label from  $\mu_{ij}$ 's; this may be very helpful for understanding the complicated concepts underlying the task. Eq. 28, however, is difficult to solve since it involves too many variables. Thus, how to exploit/understand the pairwise relatedness between different pairs of labels remains an open problem.

## 6 Solving Single-Instance Multi-Label Problems through MIML Transformation

The previous sections show that when we have access to the real objects and are able to represent complicated objects as MIML examples, using the MIML framework is beneficial. However, in many practical tasks we are given observational data where each object has already been represented by a single instance, and we do not have access to the real objects. In such case, we cannot capture more information from the real objects using the MIML representation. Even in this situation, however, MIML is still useful. Here we propose the INSDIF (i.e., INSTance DIFFerentiation) algorithm which transforms single-instance multi-label examples into MIML examples to exploit the power of MIML.

### 6.1 INSDIF

For an object associated with multiple class labels, if it is described by only a single instance, the information corresponding to these labels are mixed and thus difficult to learn. The basic assumption of INSDIF is that the spatial distribution of

instances with different labels encodes information helpful for discriminating these labels, and such information will become more explicit by breaking the single-instances into a number of instances each corresponds to one label.

INSDIF is a two-stage algorithm, which is based on *instance differentiation*. In the first stage, INSDIF transforms each example into a bag of instances, by deriving one instance for each class label, in order to explicitly express the ambiguity of the example in the input space; in the second stage, an MIML learner is utilized to learn from the transformed data set. For the consistency with our previous description of the algorithm [81], in the current version of INSDIF we use a two-level classification strategy, but note that other MIML algorithms such as D-MIMLSVM can also be applied.

Using the same denotation as that in Sections 3 and 4, that is, given data set  $S = \{(\mathbf{x}_1, Y_1), (\mathbf{x}_2, Y_2), \dots, (\mathbf{x}_m, Y_m)\}$ , where  $\mathbf{x}_i \in \mathcal{X}$  is an instance and  $Y_i \subseteq \mathcal{Y}$  a set of labels  $\{y_{i1}, y_{i2}, \dots, y_{i, l_i}\}$ ,  $y_{ik} \in \mathcal{Y}$  ( $k = 1, 2, \dots, l_i$ ). Here  $l_i$  denotes the number of labels in  $Y_i$ .

In the first stage, INSDIF derives a prototype vector  $\mathbf{v}_l$  for each class label  $l \in \mathcal{Y}$  by averaging all the training instances belonging to  $l$ , i.e.,

$$\mathbf{v}_l = \frac{1}{|S_l|} \left( \sum_{\mathbf{x}_i \in S_l} \mathbf{x}_i \right), \quad (29)$$

where

$$S_l = \{\mathbf{x}_i | \{\mathbf{x}_i, Y_i\} \in S, l \in Y_i\}, l \in \mathcal{Y}.$$

Here  $\mathbf{v}_l$  can be approximately regarded as a profile-style vector describing common characteristics of the class  $l$ . Actually, this kind of prototype vectors have already shown their usefulness in solving text categorization problems. For example, the ROCCHIO method [34, 59] forms a prototype vector for each class by averaging all the documents (represented by weight vectors) of this class, and then classifies the test document by calculating the dot-products between the weight vector representing the document and each of the prototype vectors. Here we use such prototype vectors to facilitate bag generation. After obtaining the prototype vectors, each example  $\mathbf{x}_i$  is re-represented by a bag of instances  $B_i$ , where each instance in  $B_i$  expresses the difference between  $\mathbf{x}_i$  and a prototype vector according to Eq. 30. In

this way, each example is transformed into a bag whose size equals to the number of class labels.

$$B_i = \{\mathbf{x}_i - \mathbf{v}_l | l \in \mathcal{Y}\} \quad (30)$$

In fact, such a process attempts to exploit the spatial distribution since  $\mathbf{x}_i - \mathbf{v}_l$  in Eq. 30 is a kind of distance between  $\mathbf{x}_i$  and  $\mathbf{v}_l$ . The transformation can also be realized in other ways. For example, other than referring to the prototype vector of each class, one could also consider the following approach. For each possible class  $l$ , identify the  $k$ -nearest neighbors of  $\mathbf{x}_i$  among training instances that have  $l$  as a proper label. Then, the mean vector of these neighbors can be regarded as an instance in the bag. Note that the transformation of a single instance into a bag of instances can be realized as a general pre-processing method which can be plugged into many learning systems.

In the second stage, INSDIF learns from the transformed training set  $S^* = \{(B_1, Y_1), (B_2, Y_2), \dots, (B_m, Y_m)\}$ . This task can be realized by any MIML learning algorithm. By default we use the MIMLNN algorithm introduced in Section 4.3.2. The use of other MIML algorithms for this stage will also be studied in the next section.

The pseudo-code of INSDIF is summarized in Appendix A (Table A.4). In the first stage (Steps 1 to 2), INSDIF transforms each example into a bag of instances by querying the class prototype vectors. In the second stage (Step 3), an MIML algorithm is used to learn from the transformed data set. A test example  $\mathbf{x}^*$  is then transformed into the corresponding bag representation  $B^*$  and then fed to the learned MIML model.

## 6.2 Experiments

We compare INSDIF with several state-of-the-art multi-label learning algorithms, including ADTBOOST.MH [22], RANKSVM [27], MLSVM [11], ML- $k$ NN [80] and CNMF [43]; these algorithms have been introduced briefly in Section 2. In addition, by using MIMLBOOST, MIMLSVM and MIMLSVM $_{mi}$  respectively to replace MIMLNN for realizing the second stage of INSDIF, we get three variants of INSDIF, i.e., INSDIF<sub>MIMLBOOST</sub>, INSDIF<sub>MIMLSVM</sub> and INSDIF<sub>MIMLSVM $_{mi}$</sub> . These variants are also evaluated for comparison.

Note that the experiments here are very different from that in Sections 4.3 and 5.6. In Sections 4.3 and 5.6, it is assumed that the data are MIML examples; while in this section, it is assumed that we are given observational data where each real object has already been represented as a single instance. In other words, in this section we are trying to learn from single-instance multi-label examples, and therefore the experimental data sets are different from those used in Sections 4.3 and 5.6.

### 6.2.1 *Yeast Gene Functional Analysis*

The task here is to predict the gene functional classes of the Yeast *Saccharomyces cerevisiae*, which is one of the best studied organisms. Specifically, the Yeast data set investigated in [27,80] is studied. Each gene is represented by a 103-dimensional feature vector generated by concatenating a gene expression vector and the corresponding phylogenetic profile. Each 79-element gene expression vector reflects the expression levels of a particular gene under two different experimental conditions, while the phylogenetic profile is a Boolean string, each bit indicating whether the concerned gene has a close homolog in the corresponding genome. Each gene is associated with a set of functional classes whose maximum size can be potentially more than 190. Elisseeff and Weston [27] have pre-processed the data set where only the known structure of the functional classes are used. In fact, the whole set of functional classes is structured into hierarchies up to 4 levels deep.<sup>7</sup> Illustrations on the first level of the hierarchy used to generate the Yeast data can be found in [27, 79, 80]. The resulting multi-label data set contains 2,417 genes, fourteen possible class labels and the average number of labels for each gene is  $4.24 \pm 1.57$ .

For INSDIF, the parameter  $M$  is set to be 20% of the size of training set; The performance of this algorithm with different  $M$  settings is shown in Appendix B (Fig. B.3), it can be found that its performance is not sensitive to the setting of  $M$ . The boosting rounds of ADTBOOST.MH are set to 25; The performance of this algorithm at different boosting rounds is shown in Appendix B (Fig. B.4), it can be observed that after this round its performance has become stable. (Similar observations are also found in Section 6.2.2.) For RANKSVM, polynomial kernel

---

<sup>7</sup> See <http://mips.gsf.de/proj/yeast/catalogues/funecat/> for more details.

Table 3

Results (mean $\pm$ std.) on Yeast gene data set (‘ $\downarrow$ ’ indicates ‘the smaller the better’; ‘ $\uparrow$ ’ indicates ‘the larger the better’).

Compared Algorithms	Evaluation Criteria						
	<i>hloss</i> $\downarrow$	<i>one-error</i> $\downarrow$	<i>coverage</i> $\downarrow$	<i>rloss</i> $\downarrow$	<i>aveprec</i> $\uparrow$	<i>avgrec</i> $\uparrow$	<i>avgF1</i> $\uparrow$
INSDIF	<b>.189<math>\pm</math>.010</b>	<b>.214<math>\pm</math>.030</b>	6.288 $\pm$ 0.240	<b>.163<math>\pm</math>.017</b>	<b>.774<math>\pm</math>.019</b>	.602 $\pm$ .026	.677 $\pm$ .023
INSDIF <sub>MIMLSVM</sub>	<b>.189<math>\pm</math>.009</b>	.232 $\pm$ .040	6.625 $\pm$ 0.261	.179 $\pm$ .015	.763 $\pm$ .021	.591 $\pm$ .023	.666 $\pm$ .022
INSDIF <sub>MIMLSVM<sub>mi</sub></sub>	.196 $\pm$ .011	.238 $\pm$ .043	6.396 $\pm$ 0.206	.172 $\pm$ .012	.765 $\pm$ .019	<b>.655<math>\pm</math>.024</b>	<b>.706<math>\pm</math>.017</b>
ADTBOOST.MH	.212 $\pm$ .008	.247 $\pm$ .029	6.385 $\pm$ 0.151	N/A	.739 $\pm$ .022	N/A	N/A
RANKSVM	.207 $\pm$ .013	.243 $\pm$ .039	7.090 $\pm$ 0.502	.195 $\pm$ .021	.750 $\pm$ .026	.500 $\pm$ .047	.600 $\pm$ .041
MLSVM	.199 $\pm$ .009	.227 $\pm$ .032	7.220 $\pm$ 0.338	.201 $\pm$ .019	.749 $\pm$ .021	.572 $\pm$ .023	.649 $\pm$ .022
ML- <i>k</i> NN	.194 $\pm$ .010	.230 $\pm$ .030	<b>6.275<math>\pm</math>0.240</b>	.167 $\pm$ .016	.765 $\pm$ .021	.574 $\pm$ .022	.656 $\pm$ .021
CNMF	N/A	.354 $\pm$ .184	7.930 $\pm$ 1.089	.268 $\pm$ .062	.668 $\pm$ .093	N/A	N/A

with degree 8 is used as suggested in [27]. For MLSVM, a Gaussian kernel is used with default LIBSVM setting for kernel width (i.e.  $\frac{1}{\# \text{ features}}$ ). For CNMF, a normalized Gaussian kernel as recommended in [43] is used to compute the pairwise class similarity. For ML-*k*NN, the number of nearest neighbors considered is set to 10. The criteria introduced in Section 4.3.1 are used to evaluate the learning performance. Ten-fold cross-validation is conducted on this data set and the results are summarized in Table 3,<sup>8</sup> where the best performance on each criterion has been highlighted in boldface.

Table 3 shows that INSDIF and its variants achieve good performance on the Yeast gene functional data set. Pairwise *t*-tests with 95% significance level disclose that: a) INSDIF is significantly better than all the compared multi-label learning algorithms (i.e., the second part of Table 3) on all criteria, except that on *coverage* it is worse than ML-*k*NN but the difference is not statistically significant;<sup>9</sup> b)

<sup>8</sup> *Hamming loss*, *average recall* and *average F1* are not available for CNMF; *ranking loss*, *average recall* and *average F1* are not available for ADTBOOST.MH. The performance of INSDIF<sub>MIMLBOOST</sub> is not reported since this algorithm did not terminate within reasonable time on this data.

<sup>9</sup> Note that our implementation of RANKSVM was obtained with the help of the authors of [27], yet our results are somewhat worse than the best results reported in [27]. We think that the performance gap may be caused by the minor implementation differences and the different experimental data partitions. Nevertheless, it is worth mentioning that the results of INSDIF are better than the best results of RANKSVM in [27] in terms of *hamming loss*, *one-error* and *average precision*, and as same as the best results of

INSDIF<sub>MIMLSVM</sub> is significantly better than the compared multi-label learning algorithms for more than 68% cases, and is significantly inferior to them for less than 11% cases; c) INSDIF<sub>MIMLSVM<sub>mi</sub></sub> is significantly better than the compared multi-label learning algorithms for more than 65% cases, and is never significantly inferior to them. Specifically, INSDIF<sub>MIMLSVM<sub>mi</sub></sub> outperforms all the compared algorithms in terms of *average recall* and *average F1*. It is noteworthy that CNMF performs quite poorly compared to other algorithms although it has used test set information. The reason may be that the key assumption of CNMF, i.e., two examples with high similarity in the input space tend to have large overlap in the output space, does not hold on this gene data since there are some genes whose functions are quite different but the physical appearances are similar.

Overall, results on the Yeast gene functional analysis task suggest that MIML can be useful when we are given observational data where each complicated object has already been represented by a single instance.

### 6.2.2 Web Page Categorization

The web page categorization task has been studied in [39, 65, 80]. The web pages were collected from the “yahoo.com” domain and then divided into 11 data sets based on Yahoo’s top-level categories.<sup>10</sup> After that, each page is classified into a number of Yahoo’s second-level subcategories. Each data set contains 2,000 training documents and 3,000 test documents. The simple term selection method based on *document frequency* (the number of documents containing a specific term) was applied to each data set to reduce the dimensionality. Actually, only 2% words with the highest document frequency were retained in the final vocabulary. Other term selection methods such as *information gain* and *mutual information* can also be adopted. After term selection, each document in the data set is described as a feature vector using the “*Bag-of-Words*” representation, i.e., each feature expresses the number of times a vocabulary word appearing in the document.

Characteristics of the web page data sets are summarized in Appendix C (Table C.1). Compared to the Yeast data in Section 6.2.1, here the instances are rep-

---

RANKSVM in [27] in terms of *ranking loss*.

<sup>10</sup> Data set available at <http://www.kecl.ntt.co.jp/as/members/ueda/yahoo.tar.gz>.

Table 4

Results (mean $\pm$ std.) on eleven web page categorization data sets ( $\downarrow$  indicates ‘the smaller the better’;  $\uparrow$  indicates ‘the larger the better’).

Compared Algorithms	Evaluation Criteria						
	<i>hloss</i> $\downarrow$	<i>one-error</i> $\downarrow$	<i>coverage</i> $\downarrow$	<i>rloss</i> $\downarrow$	<i>aveprec</i> $\uparrow$	<i>avgrecl</i> $\uparrow$	<i>aveF1</i> $\uparrow$
INSDIF	<b>.039<math>\pm</math>.013</b>	.381 $\pm$ .118	4.545 $\pm$ 1.285	<b>.102<math>\pm</math>.037</b>	<b>.686<math>\pm</math>.091</b>	.377 $\pm$ .163	.479 $\pm$ .154
INSDIF <sub>MIMLSVM</sub>	.043 $\pm$ .015	.395 $\pm$ .119	6.823 $\pm$ 1.623	.166 $\pm$ .045	.653 $\pm$ .093	<b>.501<math>\pm</math>.105</b>	<b>.566<math>\pm</math>.102</b>
ADTBOOST.MH	.044 $\pm$ .014	.477 $\pm$ .144	4.177 $\pm$ 1.261	N/A	.621 $\pm$ .108	N/A	N/A
RANKSVM	.043 $\pm$ .014	.424 $\pm$ .135	7.228 $\pm$ 2.442	.182 $\pm$ .057	.621 $\pm$ .108	.252 $\pm$ .172	.345 $\pm$ .177
MLSVM	.042 $\pm$ .015	<b>.375<math>\pm</math>.119</b>	6.919 $\pm$ 1.767	.168 $\pm$ .047	.660 $\pm$ .093	.378 $\pm$ .167	.472 $\pm$ .156
ML- <i>k</i> NN	.043 $\pm$ .015	.471 $\pm$ .157	<b>4.097<math>\pm</math>1.236</b>	<b>.102<math>\pm</math>.045</b>	.625 $\pm$ .116	.292 $\pm$ .189	.381 $\pm$ .196
CNMF	N/A	.509 $\pm$ .142	6.717 $\pm$ 1.588	.171 $\pm$ .058	.561 $\pm$ .114	N/A	N/A

resented by much higher-dimensional feature vectors and a large portion of them (about 20-45%) are multi-labeled. Moreover, here the number of categories (21-40) are much larger and many of them are *rare* categories (about 20-55%). So, the web page data sets are more difficult than the Yeast data to learn.

The parameter settings are similar as those in Section 6.2.1. That is, for INSDIF, the parameter  $M$  is set to be 20% of the size of training set; the boosting rounds of ADTBOOST.MH are set to 25; for RANKSVM, polynomial kernel is used where polynomial degrees of 2 to 9 are considered as in [27] and chosen by hold-out tests on training sets; for MLSVM and CNMF, linear and Gaussian kernel are used respectively; for ML- $k$ NN, the number of nearest neighbors considered is set to 10.

Results of the eleven data sets are shown in Appendix C (Fig. C.1), and the average results are summarized in Table 4 where the best performance on each criterion has been highlighted in boldface.<sup>11</sup>

Table 4 shows that INSDIF and INSDIF<sub>MIMLSVM</sub> perform well on the Yahoo data. Pairwise  $t$ -tests with 95% significance level disclose that: a) INSDIF is only inferior to ADTBOOST.MH and ML- $k$ NN in terms of *coverage*, inferior to MLSVM

<sup>11</sup> The performance of INSDIF<sub>MIMLBOOST</sub> and INSDIF<sub>MIMLSVM <sub>$m_i$</sub></sub>  are not reported since these algorithms did not terminate within reasonable time on this data. Note that though the significant differences between some numbers in the table might be subtle at the first glance (e.g., INSDIF vs. RANKSVM in terms of *one-error*), statistical tests based on detailed information (in online supplementary file) justify the significance.

in terms of *one-error*, comparable to ML-*k*NN in terms of *ranking loss*, comparable to MLSVM in terms of *average recall* and *average F1*. Under all the other circumstances (more than 79% cases), the performance of INSDIF is significantly better than the compared multi-label learning algorithms (i.e., the second part of Table 4); b) INSDIF<sub>MIMLSVM</sub> is significantly better than the compared multi-label learning algorithms for more than 44% cases, and is significantly inferior to them for less than 18% cases. Specifically, INSDIF<sub>MIMLSVM</sub> achieves the best performance in terms of *average recall* and *average F1*; on *one-error*, it is only inferior to MLSVM but significantly superior the other compared multi-label learning algorithms.

Overall, results on the web page categorization task suggest that MIML can be useful when we are given observational data where each complicated object has already been represented by a single instance.

## 7 Solving Multi-Instance Single-Label Problems through MIML Transformation

In many tasks we are given observational data where each object has already been represented as a multi-instance single-label example, and we do not have access to the real objects. In such case, we cannot capture more information from the real objects using the MIML representation. Even in this situation, however, MIML is still useful. Here we propose the SUBCOD (i.e., SUB-CONcept Discovery) algorithm which transforms multi-instance single-label examples into MIML examples to exploit the power of MIML.

### 7.1 SUBCOD

For an object that has been described by multi-instances, if it is associated with a label corresponding to a high-level complicated concept such as *Africa* in Fig. 2(a), it may be quite difficult to learn this concept directly. The basic assumption of SUBCOD is that high-level complicated concepts can be derived by a number of lower-level sub-concepts which are relatively clearer and easier for learning, so that we can transform the single-label into a set of labels each corresponds to one sub-concept. Therefore, we can learn these labels at first and then derive the high-level

complicated label based on them, as illustrated in Fig. 2(b).

SUBCOD is a two-stage algorithm, which is based on *sub-concept discovery*. In the first stage, SUBCOD transforms each single-label example into a multi-label example by discovering and exploiting sub-concepts involved by the original label; this is realized by constructing multiple labels through unsupervised clustering all instances and then treating each cluster as a set of instances of a separate sub-concept. In the second stage, the outputs learned from the transformed data set are used to derive the original labels that are to be predicted; this is realized by using a supervised learning algorithm to predict the original labels from the sub-concepts predicted by an MIML learner.

Using the same denotation as that in Sections 3 and 4, that is, given data set  $\{(X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)\}$ , where  $X_i \subseteq \mathcal{X}$  is a set of instances  $\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{i, n_i}\}$ ,  $\mathbf{x}_{ij} \in \mathcal{X}$  ( $j = 1, 2, \dots, n_i$ ), and  $y_i \in \mathcal{Y}$  is the label of  $X_i$ . Here  $n_i$  denotes the number of instances in  $X_i$ .

In the first stage, SUBCOD collects all instances from all the bags to compose a data set  $D = \{\mathbf{x}_{11}, \dots, \mathbf{x}_{1, n_1}, \mathbf{x}_{21}, \dots, \mathbf{x}_{2, n_2}, \dots, \mathbf{x}_{m1}, \dots, \mathbf{x}_{m, n_m}\}$ . For the ease of discussion, let  $N = \sum_{i=1}^m n_i$  and re-index the instances in  $D$  as  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ . A Gaussian mixture model with  $M$  mixture components is to be learned from  $D$  by the EM algorithm, and the mixture components are regarded as sub-concepts. The parameters of the mixture components, i.e., the means  $\boldsymbol{\mu}_k$ , covariances  $\Sigma_k$  and mixing coefficients  $\pi_k$  ( $k = 1, 2, \dots, M$ ), are randomly initialized and the initial value of the log-likelihood is evaluated. In the E-step, the responsibilities are measured according to

$$\gamma_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^M \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \Sigma_j)} \quad (i = 1, 2, \dots, N) . \quad (31)$$

In the M-step, the parameters are re-estimated according to

$$\boldsymbol{\mu}_k^{new} = \frac{\sum_{i=1}^N \gamma_{ik} \mathbf{x}_i}{\sum_{i=1}^N \gamma_{ik}} , \quad (32)$$

$$\Sigma_k^{new} = \frac{\sum_{i=1}^N \gamma_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k^{new})(\mathbf{x}_i - \boldsymbol{\mu}_k^{new})^T}{\sum_{i=1}^N \gamma_{ik}}, \quad (33)$$

$$\pi_k^{new} = \frac{\sum_{i=1}^N \gamma_{ik}}{N}, \quad (34)$$

and the log-likelihood is evaluated according to

$$\ln p(D|\boldsymbol{\mu}, \Sigma, \pi) = \sum_{i=1}^N \ln \left( \sum_{k=1}^M \pi_k^{new} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k^{new}, \Sigma_k^{new}) \right). \quad (35)$$

After the convergence of the EM process (or after a pre-specified number of iterations), we can estimate the associated sub-concept for every instance  $\mathbf{x}_i \in D$  ( $i = 1, 2, \dots, N$ ) by

$$sc(\mathbf{x}_i) = \arg \max_k \gamma_{ik} \quad (k = 1, 2, \dots, M). \quad (36)$$

Then, we can derive the multi-label for each  $X_i$  ( $i = 1, 2, \dots, m$ ) by considering the sub-concept belongingness. Let  $\mathbf{c}_i$  denote an  $M$ -dimensional binary vector where each element is either  $+1$  or  $-1$ . For  $j = 1, 2, \dots, M$ ,  $c_{ij} = +1$  means that the sub-concept corresponding to the  $j$ -th Gaussian mixture component appears in  $X_i$ , while  $c_{ij} = -1$  means that this sub-concept does not appear in  $X_i$ . Here the value of  $c_{ij}$  can be determined according to a simple rule that  $c_{ij} = +1$  if  $X_i$  has at least one instance which takes the  $j$ -th sub-concept (i.e., satisfying Eq. 36); otherwise  $c_{ij} = -1$ . Note that for examples with identical single-label, the derived multi-labels for them may be different.

The above process works in an unsupervised way which does not consider the original labels of the bags  $X_i$ 's. Thus, the derived multi-labels  $\mathbf{c}_i$  need to be polished by incorporating the relation between the sub-concepts and the original label of  $X_i$ . Here the maximum margin criterion is used. In detail, consider a vector  $\mathbf{z}_i$  with elements  $z_{ij} \in [-1.0, +1.0]$  ( $j = 1, 2, \dots, M$ );  $z_{ij} = +1$  means that the label  $c_{ij}$  should not be modified while  $z_{ij} = -1$  means that the label  $c_{ij}$  should be inverted. Denote  $\mathbf{q}_i = \mathbf{c}_i \odot \mathbf{z}_i$  as that for  $j = 1, 2, \dots, M$ ,  $q_{ij} = c_{ij} z_{ij}$ . Let  $\theta$  denote the smallest number of labels that cannot be inverted. SUBCOD attempts to optimize the objective

$$\begin{aligned}
& \min_{\mathbf{w}, b, \boldsymbol{\xi}, \mathbf{Z}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi_i & (37) \\
& \text{s.t. } y_i(\mathbf{w}^\top(\mathbf{c}_i \odot \mathbf{z}_i) + b) \geq 1 - \xi_i, \\
& \quad \boldsymbol{\xi} \geq \mathbf{0}, \quad -1 \leq z_{ij} \leq 1 \\
& \quad \sum_{i,j} z_{ij} \geq 2\theta - mM,
\end{aligned}$$

where  $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m]$ .

By solving Eq. 37 we will get the vector  $\mathbf{z}_i$  which maximizes the margin of the prediction of the proper labels of  $X_i$ . Here we solve Eq. 37 iteratively. We initialize  $\mathbf{Z}$  with all 1's. First, we fix  $\mathbf{Z}$  to get the optimal  $\mathbf{w}$  and  $b$ ; this is a standard QP problem. Then, we fix  $\mathbf{w}$  and  $b$  to get the optimal  $\mathbf{Z}$ ; this is a standard LP problem. These two steps are iterated till convergence. Finally, we set the multi-label vector's elements which correspond to positive  $c_{ij}z_{ij}$ 's ( $i = 1, 2, \dots, m; j = 1, 2, \dots, M$ ) to +1, and set the remaining ones to -1. Thus, we get all the polished multi-label vectors  $\tilde{\mathbf{c}}_i$  for the bags  $X_i$ . Thus, the original data set  $\{(X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)\}$  is transformed to an MIML data set  $\{(X_1, \tilde{\mathbf{c}}_1), (X_2, \tilde{\mathbf{c}}_2), \dots, (X_m, \tilde{\mathbf{c}}_m)\}$ , and any MIML algorithms can be applied.

To map the multi-labels predicted by the MIML classifier for a test example to the original single-labels  $y \in \mathcal{Y}$ , in the second stage of SUBCOD, a traditional classifier  $f : \{+1, -1\}^M \rightarrow \mathcal{Y}$  is generated from the data set  $\{(\tilde{\mathbf{c}}_1, y_1), (\tilde{\mathbf{c}}_2, y_2), \dots, (\tilde{\mathbf{c}}_m, y_m)\}$ . This is relatively simple and traditional supervised learning algorithms can be applied.

The pseudo-code of SUBCOD is summarized in Appendix A (Table A.5). In the first stage (Steps 1 to 3), SUBCOD derives multi-labels via sub-concept discovery and transforms single-label examples into MIML examples, from which an MIML learner is generated. In the second stage (Step 4), a traditional classifier is trained to map the derived multi-labels to the original single-labels. Test example  $X^*$  is fed to the MIML learner to get its multi-labels, and the multi-labels are then fed to the supervised classifier to get the label  $y^*$  predicted for  $X^*$ .

## 7.2 Experiments

We compare SUBCOD with several state-of-the-art multi-instance learning algorithms, including DIVERSE DENSITY [45], EM-DD [83], MI-SVM and MI-SVM [3], and CH-FD [31]; these algorithms have been introduced briefly in Section 2. For SUBCOD, the MIML learner in Step 3 is realized by MIMLSVM and the classifier in Step 4 is realized by SMO with default parameters. In addition, by using MIMLNN and MIMLSVM<sub>mi</sub> respectively to replace MIMLSVM for realizing Step 3 of SUBCOD, we get two variants of SUBCOD, i.e., SUBCOD<sub>MIMLNN</sub> and SUBCOD<sub>MIMLSVM<sub>mi</sub></sub>. They are also evaluated for comparison.<sup>12</sup>

Note that the experiments here are very different from that in Sections 4.3, 5.6 and 6.2. Both Sections 4.3 and 5.6 deal with learning from MIML examples, Section 6.2 deals with learning from single-instance multi-label examples, while this section deals with learning from multi-instance single-label examples, and therefore the experimental data sets in this section are different from those used in Sections 4.3, 5.6 and 6.2.

Five benchmark multi-instance learning data sets are used, including *Musk1*, *Musk2*, *Elephant*, *Tiger* and *Fox*. Both *Musk1* and *Musk2* are drug activity prediction data sets, publicly available at the UCI machine learning repository [8]. Here every bag corresponds to a molecule, while every instance corresponds to a low-energy shape of the molecule [24]. *Musk1* contains 47 positive bags and 45 negative bags, and the number of instances contained in each bag ranges from 2 to 40. *Musk2* contains 39 positive bags and 63 negative bags, and the number of instances contained in each bag ranges from 1 to 1,044. Each instance is a 166-dimensional feature vector. *Elephant*, *Tiger* and *Fox* are three image annotation data sets generated by [3] for multi-instance learning. Here every bag is an image, while every instance corresponds to a segmented region in the image [3]. Each data set contains 100 positive and 100 negative bags, and each instance is a 230-dimensional feature vector. These data sets are popularly used in evaluating the performance of multi-instance learning algorithms.

---

<sup>12</sup> We have also evaluated the variant SUBCOD<sub>MIMLBOOST</sub> which is obtained by employing MIMLBOOST to replace MIMLSVM, however, it did not terminate within reasonable time and so its performance is not reported in this section.

Table 5  
 Predictive accuracy on *Musk1*, *Musk2*, *Elephant*, *Tiger* and *Fox* data sets

Compared Algorithms	Data sets				
	<i>Musk1</i>	<i>Musk2</i>	<i>Elephant</i>	<i>Tiger</i>	<i>Fox</i>
SUBCOD	0.850±0.035	<b>0.921±0.014</b>	<b>0.836±0.010</b>	0.808±0.013	<b>0.616±0.020</b>
SUBCOD <sub>MIMLNN</sub>	0.859±0.025	0.888±0.022	0.815±0.023	0.795±0.018	0.599±0.032
SUBCOD <sub>MIMLSVM<sub>mi</sub></sub>	0.870±0.023	0.869±0.020	0.805±0.017	0.787±0.016	0.590±0.015
DIVERSE DENSITY	0.880	0.840	N/A	N/A	N/A
EM-DD	0.848	0.849	0.783	0.721	0.561
MI-SVM	0.874	0.836	0.820	0.789	0.582
MI-SVM	0.779	0.843	0.814	<b>0.840</b>	0.594
CH-FD	<b>0.888</b>	0.857	0.824	0.822	0.604

Parameters of SUBCOD are determined by hold-out tests on training sets. Specifically, candidate values of  $M$  (the number of Gaussian mixture components) range between  $[10, 70]$ , while candidate values of  $\theta$  (the smallest number of labels that cannot be inverted) range between  $[mM \times 10\%, mM \times 70\%]$ . Ten runs of ten-fold cross validation are performed and the results are summarized in Table 5, where the best performance on each data set has been highlighted in boldface. Note that the results of the compared algorithms (second part of Table 5) are the best performance reported in literatures [3, 31].<sup>13</sup>

Table 5 shows that SUBCOD and its variants are very competitive to state-of-the-art multi-instance learning algorithms. In particular, on *Musk2* their performance are much better than other algorithms. This is expectable because *Musk2* is a complicated data set which has the largest number of instances, while on such data set the sub-concept discovery process of SUBCOD may be more effective.

Overall, the experimental results suggest that MIML could be useful when we are given observational data where each object has already been represented as a multi-instance single-label example.

<sup>13</sup> The tradition of the multi-instance learning community is to compare with the best performance reported in literature. Since the detailed results are not available [3, 17, 18, 31, 32, 45, 67, 83], we do not perform statistical significance tests at here.

## 8 Conclusion

This paper extends our preliminary work [81, 92] to formalize the MIML *Multi-Instance Multi-Label learning* framework, where an example is described by multiple instances and associated with multiple class labels. It was inspired by the recognition that when solving real-world problems, having a good representation is often more important than having a strong learning algorithm because a good representation may capture more meaningful information and make the learning task easier to tackle. Since many real objects are inherited with input ambiguity as well as output ambiguity, MIML is more natural and convenient for tasks involving such objects.

To exploit the advantages of the MIML representation, we propose the MIML-BOOST algorithm and the MIMLSVM algorithm based on a simple degeneration strategy. Experiments on scene classification and text categorization show that solving problems involving complicated objects with multiple semantic meanings under the MIML framework can lead to good performance. Considering that the degeneration process may lose information, we also propose the D-MIMLSVM algorithm which tackles MIML problems directly in a regularization framework. Experiments show that this “direct” SVM algorithm outperforms the “indirect” MIMLSVM algorithm.

In some practical tasks we are given observational data where each complicated object has already been represented by a single instance, and we do not have access to the real objects such that we cannot capture more information from the real objects using the MIML representation. For such scenario, we propose the INSDIF algorithm which transforms single-instances into MIML examples to learn. Experiments on Yeast gene functional analysis and web page categorization show that such algorithm is able to achieve a better performance than learning the single-instances directly. This is not difficult to understand. Actually, by representing the multi-label object using multi-instances, the structure information collapsed in traditional single-instance representation may become easier to exploit, and for each label the number of training instances can be significantly increased. So, transforming multi-label examples to MIML examples for learning may be beneficial in some tasks.

MIML can also be helpful for learning single-label examples involving complicated high-level concepts. Usually it may be quite difficult to learn such concepts directly since many different lower-level concepts are mixed together. If we can transform the single-label into a set of labels corresponding to some sub-concepts, which are relatively clearer and easier to learn, we can learn these labels at first and then derive the high-level complicated label based on them. Inspired by this recognition, we propose the SUBCOD algorithm which works by discovering sub-concepts of the target concept at first and then transforming the data into MIML examples to learn. Experiments show that this algorithm is able to achieve better performance than learning the single-label examples directly in some tasks.

We believe that semantics exist in the connections between atomic input patterns and atomic output patterns; while a prominent usefulness of MIML, which has not been realized in this paper, is the possibility of identifying such connection. As stated in Section 3, in the MIML framework it is possible to understand why a concerned object has a certain class label; this may be more important than simply making an accurate prediction, because the results could be helpful for understanding the source of ambiguous semantics.

## **Acknowledgements**

The authors want to thank the anonymous reviewers for helpful comments and suggestions. The authors also want to thank De-Chuan Zhan and James Kwok for help on D-MIMLSVM, Yang Yu for help on SUBCOD, and André Elisseeff and Jason Weston for providing the Yeast data and the implementation details of RANKSVM. A preliminary Chinese version has been presented at the Chinese Workshop on Machine Learning and Applications 2009.

## Appendix

### A Pseudo-codes of the Learning Algorithms

Table A.1  
The MIMLBOOST algorithm

- 
- 1 Transform each MIML example  $(X_u, Y_u)$  ( $u = 1, 2, \dots, m$ ) into  $|\mathcal{Y}|$  number of multi-instance bags  $\{(X_u, y_1), \Psi(X_u, y_1), \dots, (X_u, y_{|\mathcal{Y}|}), \Psi(X_u, y_{|\mathcal{Y}|})\}$ . Thus, the original data set is transformed into a multi-instance data set containing  $m \times |\mathcal{Y}|$  number of multi-instance bags, denoted by  $\{(X^{(i)}, y^{(i)}), \Psi(X^{(i)}, y^{(i)})\}$  ( $i = 1, 2, \dots, m \times |\mathcal{Y}|$ ).
  - 2 Initialize weight of each bag to  $W^{(i)} = \frac{1}{m \times |\mathcal{Y}|}$  ( $i = 1, 2, \dots, m \times |\mathcal{Y}|$ ).
  - 3 Repeat for  $t = 1, 2, \dots, T$  iterations:
    - 3a Assign the bag's label  $\Psi(X^{(i)}, y^{(i)})$  to each of its instances  $(\mathbf{x}_j^{(i)}, y^{(i)})$  ( $i = 1, 2, \dots, m \times |\mathcal{Y}|$ ;  $j = 1, 2, \dots, n_i$ ), set the weight of the  $j$ -th instance of the  $i$ -th bag  $W_j^{(i)} = W^{(i)}/n_i$ , and build an instance-level predictor  $h_t[(\mathbf{x}_j^{(i)}, y^{(i)})] \in \{-1, +1\}$ .
    - 3b For the  $i$ -th bag, compute the error rate  $e^{(i)} \in [0, 1]$  by counting the number of misclassified instances within the bag, i.e.  $e^{(i)} = \frac{\sum_{j=1}^{n_i} [h_t[(\mathbf{x}_j^{(i)}, y^{(i)})] \neq \Psi(X^{(i)}, y^{(i)})]}{n_i}$ .
    - 3c If  $e^{(i)} < 0.5$  for all  $i \in \{1, 2, \dots, m \times |\mathcal{Y}|\}$ , go to Step 4.
    - 3d Compute  $c_t = \arg \min_{c_t} \sum_{i=1}^{m \times |\mathcal{Y}|} W^{(i)} \exp[(2e^{(i)} - 1)c_t]$ .
    - 3e If  $c_t \leq 0$ , go to Step 4.
    - 3f Set  $W^{(i)} = W^{(i)} \exp[(2e^{(i)} - 1)c_t]$  ( $i = 1, 2, \dots, m \times |\mathcal{Y}|$ ) and re-normalize such that  $0 \leq W^{(i)} \leq 1$  and  $\sum_{i=1}^{m \times |\mathcal{Y}|} W^{(i)} = 1$ .
  - 4 Return  $Y^* = \{y | \text{sign}(\sum_j \sum_t c_t h_t[(\mathbf{x}_j^*, y)]) = +1\}$  ( $\mathbf{x}_j^*$  is  $X^*$ 's  $j$ -th instance).
-

Table A.2  
The MIMLSVM algorithm

- 
- 1 For MIML examples  $(X_u, Y_u)$  ( $u = 1, 2, \dots, m$ ),  $\Gamma = \{X_u | u = 1, 2, \dots, m\}$ .
  - 2 Randomly select  $k$  elements from  $\Gamma$  to initialize the medoids  $M_t$  ( $t = 1, 2, \dots, k$ ), repeat until all  $M_t$  do not change:
    - 2a  $\Gamma_t = \{M_t\}$  ( $t = 1, 2, \dots, k$ ).
    - 2b Repeat for each  $X_u \in (\Gamma - \{M_t | t = 1, 2, \dots, k\})$ :
$$index = \arg \min_{t \in \{1, \dots, k\}} d_H(X_u, M_t), \Gamma_{index} = \Gamma_{index} \cup \{X_u\}.$$
    - 2c  $M_t = \arg \min_{A \in \Gamma_t} \sum_{B \in \Gamma_t} d_H(A, B)$  ( $t = 1, 2, \dots, k$ ).
  - 3 Transform  $(X_u, Y_u)$  into a multi-label example  $(z_u, Y_u)$  ( $u = 1, 2, \dots, m$ ), where  $z_u = (z_{u1}, z_{u2}, \dots, z_{uk}) = (d_H(X_u, M_1), d_H(X_u, M_2), \dots, d_H(X_u, M_k))$ .
  - 4 For each  $y \in \mathcal{Y}$ , derive a data set  $\mathcal{D}_y = \{(z_u, \Phi(z_u, y)) | u = 1, 2, \dots, m\}$ , and then train an SVM  $h_y = SVMTrain(\mathcal{D}_y)$ .
  - 5 Return  $Y^* = \{\arg \max_{y \in \mathcal{Y}} h_y(z^*)\} \cup \{y | h_y(z^*) \geq 0, y \in \mathcal{Y}\}$ , where  $z^* = (d_H(X^*, M_1), d_H(X^*, M_2), \dots, d_H(X^*, M_k))$ .
- 

Table A.3  
Efficient Algorithm for Eq. 24

- 
- Input:**  $K, \lambda, \mu, \gamma, \varepsilon, \{X_i, Y_i\}_{i=1}^m$
- 1  $\forall t, S_t = \emptyset, \mathbf{v}_t = (\boldsymbol{\alpha}_t^T, \boldsymbol{\xi}_{t1}, \dots, \boldsymbol{\xi}_{tm}, \boldsymbol{\delta}_{t1}, \dots, \boldsymbol{\delta}_{tm}, b_t) = \mathbf{0}$
  - 2 **Repeat**
  - 3   **For**  $t = 1, \dots, T$
  - 4     Pick  $p$  indexes of constraints that are not in  $S_t$  randomly, denoted by  $I$ ;
  - 5     Compute  $Loss_i$  for every constraint in  $I$ ;
  - 6     % find out the cutting plane
  - 7      $q = \arg \max_{i \in I} Loss_i$
  - 8     **If**  $Loss_q > \varepsilon$
  - 9        $S_t = S_t \cup \{q\}$ ;
  - 10      $\mathbf{v}_t \leftarrow$  optimized over  $S_t$ ;
  - 11     **End If**
  - 12   **End For**
  - 13 **Until** no  $S_t$  changes
-

Table A.4  
The INSDIF algorithm

- 
- 1 For single-instance multi-label examples  $(x_u, Y_u)$  ( $u = 1, 2, \dots, m$ ), compute the prototype vectors  $\mathbf{v}_l$  ( $l \in \mathcal{Y}$ ) using Eq. 29.
  - 2 Derive the new training set  $S^*$  by transforming each  $\mathbf{x}_i$  into a bag of instances  $B_i$  using Eq. 30.
  - 3 Learning from  $S^* = \{(B_1, Y_1), (B_2, Y_2), \dots, (B_m, Y_m)\}$  by using an MIML algorithm.
- 

Table A.5  
The SUBCOD algorithm

- 
- 1 For multi-instance single-label examples  $(X_u, y_u)$  ( $u = 1, 2, \dots, m$ ), collect all the instances  $\mathbf{x} \in X_u$  together and identify the Gaussian mixture components through the EM process detailed in Eqs. 31 to 35.
  - 2 Determine the sub-concept for every instance  $\mathbf{x} \in X_u$  according to Eq. 36, and then derive the label vector  $\mathbf{c}_u$  for  $X_u$ .
  - 3 Make corrections to  $\mathbf{c}_u$  by optimizing Eq. 37, which results in  $\tilde{\mathbf{c}}_u$  for  $X_u$ , and then train an MIML learner  $h_t(X)$  on  $\{(X_u, \tilde{\mathbf{c}}_u)\}$  ( $u = 1, 2, \dots, m$ ).
  - 4 Train a classifier  $h_y(\tilde{\mathbf{c}})$  on  $\{(\tilde{\mathbf{c}}_u, y_u)\}$  ( $u = 1, 2, \dots, m$ ), which maps the derived multi-labels to the original single-labels.
  - 5 Return  $y^* = h_y(h_t(X^*))$ .
-

## B Parameter Settings of the Learning Algorithms

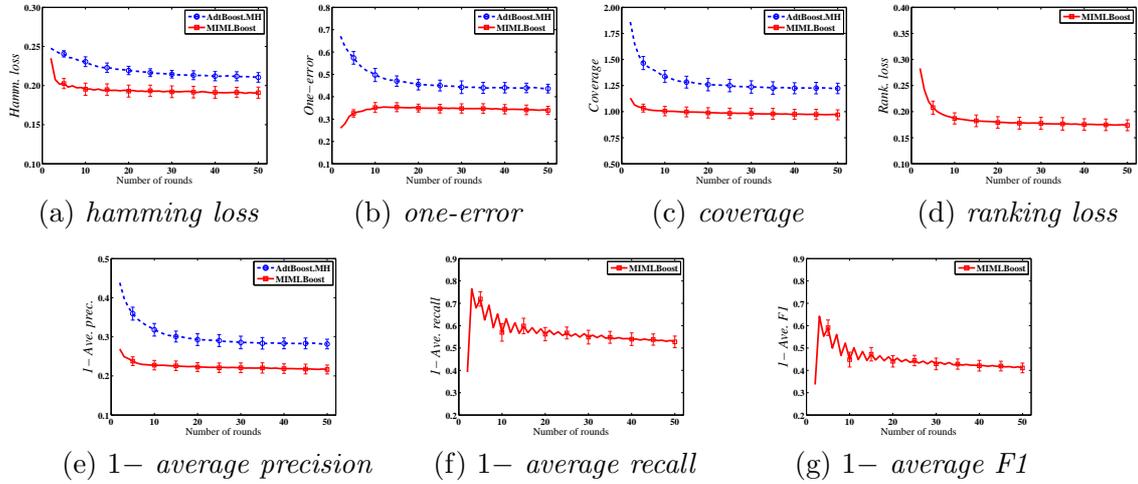


Fig. B.1. Performance of MIMLBOOST and ADTBOOST.MH at different rounds on scene classification data set.

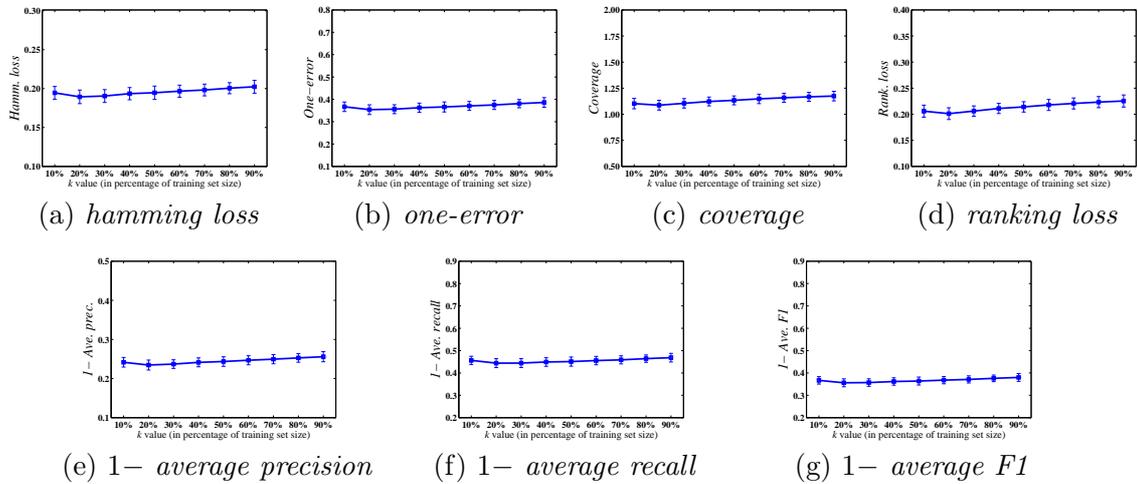


Fig. B.2. Performance of MIMLSVM with different  $k$  values on scene classification data set.

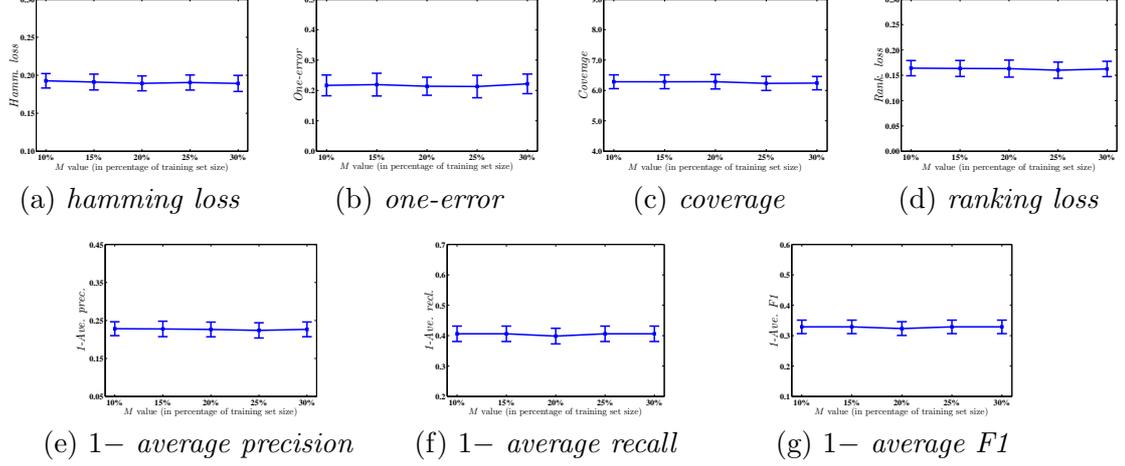


Fig. B.3. Performance of INSDIF with different  $M$  settings on Yeast gene data set.

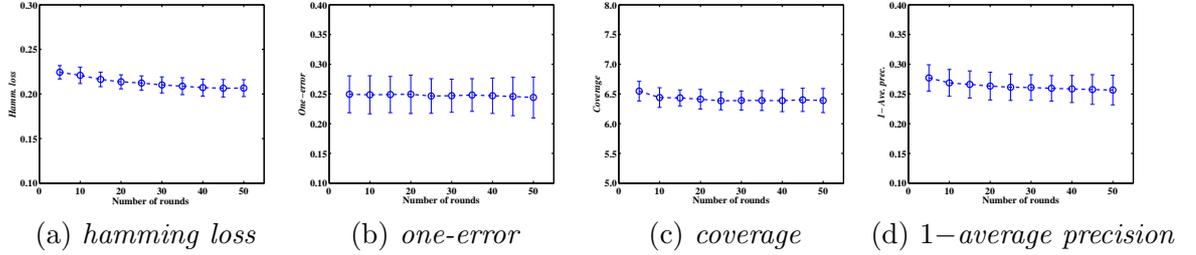


Fig. B.4. Performance of ADTBOOST.MH at different rounds on Yeast gene data set.

## C Web Page Data Sets

Table C.1

Characteristics of the web page data sets (after term selection). *PMC* denotes the percentage of documents belonging to more than one category; *ANL* denotes the average number of labels for each document; *PRC* denotes the percentage of *rare* categories, i.e., the kind of category where only less than 1% instances in the data set belong to it.

Data Set	Number of Categories	Vocabulary Size	Training Set			Test Set		
			<i>PMC</i>	<i>ANL</i>	<i>PRC</i>	<i>PMC</i>	<i>ANL</i>	<i>PRC</i>
Arts&Humanities	26	462	44.50%	1.627	19.23%	43.63%	1.642	19.23%
Business&Economy	30	438	42.20%	1.590	50.00%	41.93%	1.586	43.33%
Computers&Internet	33	681	29.60%	1.487	39.39%	31.27%	1.522	36.36%
Education	33	550	33.50%	1.465	57.58%	33.73%	1.458	57.58%
Entertainment	21	640	29.30%	1.426	28.57%	28.20%	1.417	33.33%
Health	32	612	48.05%	1.667	53.13%	47.20%	1.659	53.13%
Recreation&Sports	22	606	30.20%	1.414	18.18%	31.20%	1.429	18.18%
Reference	33	793	13.75%	1.159	51.52%	14.60%	1.177	54.55%
Science	40	743	34.85%	1.489	35.00%	30.57%	1.425	40.00%
Social&Science	39	1 047	20.95%	1.274	56.41%	22.83%	1.290	58.97%
Society&Culture	27	636	41.90%	1.705	25.93%	39.97%	1.684	22.22%

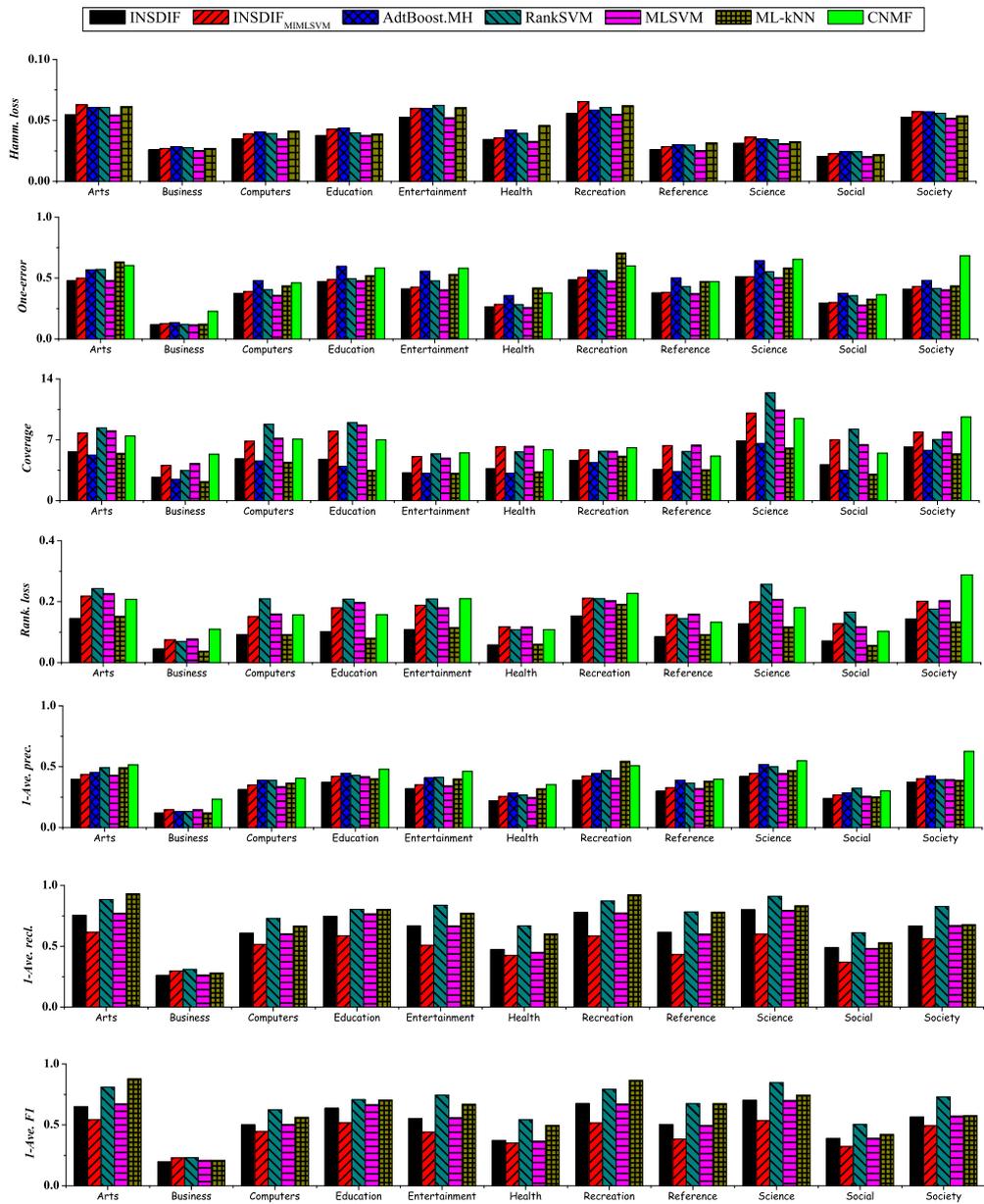


Fig. C.1. Results on the eleven Yahoo data sets.

## References

- [1] É. Alphonse and S. Matwin. Filtering multi-instance problems to reduce dimensionality in relational learning. *Journal of Intelligent Information Systems*, 22(1):23–40, 2004.
- [2] R. A. Amar, D. R. Dooly, S. A. Goldman, and Q. Zhang. Multiple-instance learning of real-valued data. In *Proceedings of the 18th International Conference on Machine Learning*, pages 3–10, Williamstown, MA, 2001.
- [3] S. Andrews, I. Tsochantaridis, and T. Hofmann. Support vector machines for multiple-instance learning. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 561–568. MIT Press, Cambridge, MA, 2003.
- [4] P. Auer. On learning from multi-instance examples: Empirical evaluation of a theoretical approach. In *Proceedings of the 14th International Conference on Machine Learning*, pages 21–29, Nashville, TN, 1997.
- [5] P. Auer, P. M. Long, and A. Srinivasan. Approximating hyper-rectangles: Learning and pseudo-random sets. *Journal of Computer and System Sciences*, 57(3):376–388, 1998.
- [6] P. Auer and R. Ortner. A boosting approach to multiple instance learning. In *Proceedings of the 15th European Conference on Machine Learning*, pages 63–74, Pisa, Italy, 2004.
- [7] Z. Barutcuoglu, R. E. Schapire, and O. G. Troyanskaya. Hierarchical multi-label prediction of gene function. *Bioinformatics*, 22(7):830–836, 2006.
- [8] C. Blake, E. Keogh, and C. J. Merz. UCI repository of machine learning databases. [<http://www.ics.uci.edu/~mllearn/MLRepository.html>], Department of Information and Computer Science, University of California, Irvine, CA, 1998.
- [9] H. Blockeel, D. Page, and A. Srinivasan. Multi-instance tree learning. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 57–64, Bonn, Germany, 2005.
- [10] A. Blum and A. Kalai. A note on learning from multiple-instance examples. *Machine Learning*, 30(1):23–29, 1998.

- [11] M. R. Boutell, J. Luo, X. Shen, and C. M. Brown. Learning multi-label scene classification. *Pattern Recognition*, 37(9):1757–1771, 2004.
- [12] K. Brinker, J. Fürnkranz, and E. Hüllermeier. A unified model for multilabel classification and ranking. In *Proceedings of the 17th European Conference on Artificial Intelligence*, pages 489–493, Riva del Garda, Italy, 2006.
- [13] K. Brinker and E. Hüllermeier. Case-based multilabel ranking. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, pages 702–707, Hyderabad, India, 2007.
- [14] L. Cai and T. Hofmann. Hierarchical document categorization with support vector machines. In *Proceedings of the 13th ACM International Conference on Information and Knowledge Management*, pages 78–87, Washington, DC, 2004.
- [15] N. Cesa-Bianchi, C. Gentile, and L. Zaniboni. Hierarchical classification: Combining Bayes with SVM. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 177–184, Pittsburgh, PA, 2006.
- [16] C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. Technical report, Department of Computer Science and Information Engineering, National Taiwan University, Taipei, 2001.
- [17] Y. Chen, J. Bi, and J. Z. Wang. MILES: Multiple-instance learning via embedded instance selection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(12):1931–1947, 2006.
- [18] Y. Chen and J. Z. Wang. Image categorization by learning and reasoning with regions. *Journal of Machine Learning Research*, 5:913–939, 2004.
- [19] P.-M. Cheung and J. T. Kwok. A regularization framework for multiple-instance learning. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 193–200, Pittsburgh, PA, 2006.
- [20] Y. Chevaleyre and J.-D. Zucker. A framework for learning rules from multiple instance data. In *Proceedings of the 12th European Conference on Machine Learning*, pages 49–60, Freiburg, Germany, 2001.
- [21] A. Clare and R. D. King. Knowledge discovery in multi-label phenotype data. In *Proceedings of the 5th European Conference on Principles of Data Mining and Knowledge Discovery*, pages 42–53, Freiburg, Germany, 2001.

- [22] F. De Comité, R. Gilleron, and M. Tommasi. Learning multi-label alternating decision tree from texts and data. In *Proceedings of the 3rd International Conference on Machine Learning and Data Mining in Pattern Recognition*, pages 35–49, Leipzig, Germany, 2003.
- [23] L. De Raedt. Attribute-value learning versus inductive logic programming: The missing links. In *Proceedings of the 8th International Workshop on Inductive Logic Programming*, pages 1–8, Madison, WI, 1998.
- [24] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez. Solving the multiple-instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89(1-2):31–71, 1997.
- [25] T. Pham Dinh and H. A. Le Thi. A D. C. optimization algorithm for solving the trust-region subproblem. *SIAM Journal on Optimization*, 8(2):476–505, 1998.
- [26] G. A. Edgar. *Measure, Topology, and Fractal Geometry*. Springer, Berlin, 1990.
- [27] A. Elisseeff and J. Weston. A kernel method for multi-labelled classification. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, pages 681–687. MIT Press, Cambridge, MA, 2002.
- [28] T. Evgeniou, C. A. Micchelli, and M. Pontil. Learning multiple tasks with kernel methods. *Journal of Machine Learning Research*, 6:615–637, 2005.
- [29] J. Foulds and E. Frank. A review of multi-instance learning assumptions. *Knowledge Engineering Review*, 25(1):1–25, 2010.
- [30] Y. Freund and L. Mason. The alternating decision tree learning algorithm. In *Proceedings of the 16th International Conference on Machine Learning*, pages 124–133, Bled, Slovenia, 1999.
- [31] G. Fung, M. Dundar, B. Krishnappuram, and R. B. Rao. Multiple instance learning for computer aided diagnosis. In B. Schölkopf, J. Platt, and T. Hofmann, editors, *Advances in Neural Information Processing Systems 19*, pages 425–432. MIT Press, Cambridge, MA, 2007.
- [32] T. Gärtner, P. A. Flach, A. Kowalczyk, and A. J. Smola. Multi-instance kernels. In *Proceedings of the 19th International Conference on Machine Learning*, pages 179–186, Sydney, Australia, 2002.
- [33] S. Godbole and S. Sarawagi. Discriminative methods for multi-labeled classification. In *Proceedings of the 8th Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 22–30, Sydney, Australia, 2004.

- [34] D. J. Ittner, D. D. Lewis, and D. D. Ahn. Text categorization of low quality images. In *Proceedings of the 4th Annual Symposium on Document Analysis and Information Retrieval*, pages 301–315, Las Vegas, NV, 1995.
- [35] R. Jin and Z. Ghahramani. Learning with multiple labels. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 897–904. MIT Press, Cambridge, MA, 2003.
- [36] T. Joachims. Text categorization with support vector machines: Learning with many relevant features. In *Proceedings of the 10th European Conference on Machine Learning*, pages 137–142, Chemnitz, Germany, 1998.
- [37] Z. Jorgensen, Y. Zhou, and M. Inge. A multiple instance learning strategy for combating good word attacks on spam filters. *Journal of Machine Learning Research*, 8:993–1019, 2008.
- [38] F. Kang, R. Jin, and R. Sukthankar. Correlated label propagation with application to multi-label learning. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1719–1726, New York, NY, 2006.
- [39] H. Kazawa, T. Izumitani, H. Taira, and E. Maeda. Maximal margin labeling for multi-topic text categorization. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 649–656. MIT Press, Cambridge, MA, 2005.
- [40] J. E. Kelley. The cutting-plane method for solving convex programs. *Journal of the Society for Industrial and Applied Mathematics*, 8(4):703–712, 1960.
- [41] H. Kück and N. de Freitas. Learning about individuals from group statistics. In *Proceedings of the 21st Conference on Uncertainty in Artificial Intelligence*, pages 332–339, Edinburgh, Scotland, 2005.
- [42] J. T. Kwok and P.-M. Cheung. Marginalized multi-instance kernels. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, pages 901–906, Hyderabad, India, 2007.
- [43] Y. Liu, R. Jin, and L. Yang. Semi-supervised multi-label learning by constrained non-negative matrix factorization. In *Proceedings of the 21st National Conference on Artificial Intelligence*, pages 421–426, Boston, MA, 2006.
- [44] P. M. Long and L. Tan. PAC learning axis-aligned rectangles with respect to product distributions from multiple-instance examples. *Machine Learning*, 30(1):7–21, 1998.

- [45] O. Maron and T. Lozano-Pérez. A framework for multiple-instance learning. In M. I. Jordan, M. J. Kearns, and S. A. Solla, editors, *Advances in Neural Information Processing Systems 10*, pages 570–576. MIT Press, Cambridge, MA, 1998.
- [46] O. Maron and A. L. Ratan. Multiple-instance learning for natural scene classification. In *Proceedings of the 15th International Conference on Machine Learning*, pages 341–349, Madison, MI, 1998.
- [47] A. McCallum. Multi-label text classification with a mixture model trained by EM. In *Working Notes of the AAAI'99 Workshop on Text Learning*, Orlando, FL, 1999.
- [48] G.-J. Qi, X.-S. Hua, Y. Rui, J. Tang, T. Mei, and H.-J. Zhang. Correlative multi-label video annotation. In *Proceedings of the 15th ACM International Conference on Multimedia*, pages 17–26, Augsburg, Germany, 2007.
- [49] R. Rahmani and S. A. Goldman. MISSL: Multiple-instance semi-supervised learning. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 705–712, Pittsburgh, PA, 2006.
- [50] R. Rak, L. Kurgan, and M. Reformat. Multi-label associative classification of medical documents from medline. In *Proceedings of the 4th International Conference on Machine Learning and Applications*, pages 177–186, Los Angeles, CA, 2005.
- [51] S. Ray and M. Craven. Supervised versus multiple instance learning: An empirical comparison. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 697–704, Bonn, Germany, 2005.
- [52] S. Ray and D. Page. Multiple instance regression. In *Proceedings of the 18th International Conference on Machine Learning*, pages 425–432, Williamstown, MA, 2001.
- [53] J. Rousu, C. Saunders, S. Szedmak, and J. Shawe-Taylor. Learning hierarchical multi-category text classification models. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 774–751, Bonn, Germany, 2005.
- [54] G. Ruffo. *Learning single and multiple instance decision trees for computer security applications*. PhD thesis, Department of Computer Science, University of Turin, Torino, Italy, 2000.
- [55] G. Salton. *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer*. Addison-Wesley, Reading, MA, 1989.

- [56] R. E. Schapire and Y. Singer. BoosTexter: A boosting-based system for text categorization. *Machine Learning*, 39(2-3):135–168, 2000.
- [57] B. Schölkopf and A. J. Smola. *Learning With Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, 2002.
- [58] S. D. Scott, J. Zhang, and J. Brown. On generalized multiple-instance learning. Technical Report UNL-CSE-2003-5, Department of Computer Science, University of Nebraska, Lincoln, NE, 2003.
- [59] F. Sebastiani. Machine learning in automated text categorization. *ACM Computing Surveys*, 34(1):1–47, 2002.
- [60] B. Settles, M. Craven, and S. Ray. Multiple-instance active learning. In J. C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 1289–1296. MIT Press, Cambridge, MA, 2008.
- [61] A. J. Smola and B. Schölkopf. Sparse greedy matrix approximation for machine learning. In *Proceedings of the 17th International Conference on Machine Learning*, pages 911–918, San Francisco, CA, 2000.
- [62] A. J. Smola, S. V. N. Vishwanathan, and T. Hofmann. Kernel methods for missing variables. In *Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics*, pages 325–332, Savannah Hotel, Barbados, 2005.
- [63] F. A. Thabtah, P. I. Cowling, and Y. Peng. MMAC: A new multi-class, multi-label associative classification approach. In *Proceedings of the 4th IEEE International Conference on Data Mining*, pages 217–224, Brighton, UK, 2004.
- [64] I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun. Large margin methods for structured and interdependent output variables. *Journal of Machine Learning Research*, 6:1453–1484, 2005.
- [65] N. Ueda and K. Saito. Parametric mixture models for multi-labeled text. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 721–728. MIT Press, Cambridge, MA, 2003.
- [66] P. Viola, J. Platt, and C. Zhang. Multiple instance boosting for object detection. In Y. Weiss, B. Schölkopf, and J. Platt, editors, *Advances in Neural Information Processing Systems 18*, pages 1419–1426. MIT Press, Cambridge, MA, 2006.
- [67] J. Wang and J.-D. Zucker. Solving the multi-instance problem: A lazy learning approach. In *Proceedings of the 17th International Conference on Machine Learning*, pages 1119–1125, San Francisco, CA, 2000.

- [68] N. Weidmann, E. Frank, and B. Pfahringer. A two-level learning method for generalized multi-instance problem. In *Proceedings of the 14th European Conference on Machine Learning*, pages 468–479, Cavtat-Dubrovnik, Croatia, 2003.
- [69] G. M. Weiss. Mining with rarity - problems and solutions: A unifying framework. *SIGKDD Explorations*, 6(1):7–19, 2004.
- [70] X. Xu and E. Frank. Logistic regression and boosting for labeled bags of instances. In *Proceedings of the 8th Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 272–281, Sydney, Australia, 2004.
- [71] C. Yang and T. Lozano-Pérez. Image database retrieval with multiple-instance learning techniques. In *Proceedings of the 16th International Conference on Data Engineering*, pages 233–243, San Diego, CA, 2000.
- [72] Y. Yang. An evaluation of statistical approaches to text categorization. *Information Retrieval*, 1(1-2):67–88, 1999.
- [73] Y. Yang and J. O. Pedersen. A comparative study on feature selection in text categorization. In *Proceedings of the 14th International Conference on Machine Learning*, pages 412–420, Nashville, TN, 1997.
- [74] K. Yu, S. Yu, and V. Tresp. Multi-label informed latent semantic indexing. In *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 258–265, Salvador, Brazil, 2005.
- [75] A. L. Yuille and A. Rangarajan. The concave-convex procedure. *Neural Computation*, 15(4):915–936, 2003.
- [76] C. Zhang and P. Viola. Multiple-instance pruning for learning efficient cascade detectors. In J. C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 1681–1688. MIT Press, Cambridge, MA, 2008.
- [77] M.-L. Zhang and Z.-H. Zhou. Improve multi-instance neural networks through feature selection. *Neural Processing Letters*, 19(1):1–10, 2004.
- [78] M.-L. Zhang and Z.-H. Zhou. Adapting RBF neural networks to multi-instance learning. *Neural Processing Letters*, 23(1):1–26, 2006.
- [79] M.-L. Zhang and Z.-H. Zhou. Multilabel neural networks with applications to functional genomics and text categorization. *IEEE Transactions on Knowledge and Data Engineering*, 18(10):1338–1351, 2006.

- [80] M.-L. Zhang and Z.-H. Zhou. ML-kNN: A lazy learning approach to multi-label learning. *Pattern Recognition*, 40(7):2038–2048, 2007.
- [81] M.-L. Zhang and Z.-H. Zhou. Multi-label learning by instance differentiation. In *Proceedings of the 22nd AAAI Conference on Artificial Intelligence*, pages 669–674, Vancouver, Canada, 2007.
- [82] M.-L. Zhang and Z.-H. Zhou. Multi-instance clustering with applications to multi-instance prediction. *Applied Intelligence*, 31(1):47–68, 2009.
- [83] Q. Zhang and S. A. Goldman. EM-DD: An improved multi-instance learning technique. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, pages 1073–1080. MIT Press, Cambridge, MA, 2002.
- [84] Q. Zhang, W. Yu, S. A. Goldman, and J. E. Fritts. Content-based image retrieval using multiple-instance learning. In *Proceedings of the 19th International Conference on Machine Learning*, pages 682–689, Sydney, Australia, 2002.
- [85] Y. Zhang and Z.-H. Zhou. Multi-label dimensionality reduction via dependency maximization. *ACM Transactions on Knowledge Discovery from Data*, 4(3):Article 14, 2010.
- [86] Z.-H. Zhou, K. Jiang, and M. Li. Multi-instance learning based web mining. *Applied Intelligence*, 22(2):135–147, 2005.
- [87] Z.-H. Zhou and X.-Y. Liu. On multi-class cost-sensitive learning. In *Proceeding of the 21st National Conference on Artificial Intelligence*, pages 567–572, Boston, WA, 2006.
- [88] Z.-H. Zhou and J.-M. Xu. On the relation between multi-instance learning and semi-supervised learning. In *Proceeding of the 24th International Conference on Machine Learning*, pages 1167–1174, Corvallis, OR, 2007.
- [89] Z.-H. Zhou and Y. Yu. AdaBoost. In X. Wu and V. Kumar, editors, *The Top Ten Algorithms in Data Mining*, pages 127–149. Chapman & Hall, Boca Raton, FL, 2009.
- [90] Z.-H. Zhou and M.-L. Zhang. Neural networks for multi-instance learning. Technical report, AI Lab, Department of Computer Science and Technology, Nanjing University, Nanjing, China, August 2002.
- [91] Z.-H. Zhou and M.-L. Zhang. Ensembles of multi-instance learners. In *Proceeding of the 14th European Conference on Machine Learning*, pages 492–502, Cavtat-Dubrovnik, Croatia, 2003.

- [92] Z.-H. Zhou and M.-L. Zhang. Multi-instance multi-label learning with application to scene classification. In B. Schölkopf, J. Platt, and T. Hofmann, editors, *Advances in Neural Information Processing Systems 19*, pages 1609–1616. MIT Press, Cambridge, MA, 2007.
- [93] Z.-H. Zhou and M.-L. Zhang. Solving multi-instance problems with classifier ensemble based on constructive clustering. *Knowledge and Information Systems*, 11(2):155–170, 2007.