

The Emerging Trends of Multi-Label Learning

Weiwei Liu, Haobo Wang, Xiaobo Shen, and Ivor W. Tsang

Abstract—Exabytes of data are generated daily by humans, leading to the growing need for new efforts in dealing with the grand challenges for multi-label learning brought by big data. For example, extreme multi-label classification is an active and rapidly growing research area that deals with classification tasks with an extremely large number of classes or labels; utilizing massive data with limited supervision to build a multi-label classification model becomes valuable for practical applications, etc. Besides these, there are tremendous efforts on how to harvest the strong learning capability of deep learning to better capture the label dependencies in multi-label learning, which is the key for deep learning to address real-world classification tasks. However, it is noted that there has been a lack of systemic studies that focus explicitly on analyzing the emerging trends and new challenges of multi-label learning in the era of big data. It is imperative to call for a comprehensive survey to fulfill this mission and delineate future research directions and new applications.

Index Terms—Extreme Multi-label Learning, Multi-label Learning with Limited Supervision, Deep Learning for Multi-label Learning, Online Multi-label Learning, Statistical Multi-label Learning, New Applications.

1 INTRODUCTION

MULTI-LABEL classification (MLC), which assigns multiple labels for each instance simultaneously, is of paramount importance in a variety of fields ranging from protein function classification and document classification, to automatic image categorization. For example, an image may have Cloud, Tree and Sky tags; the output for a document may cover a range of topics, such as News, Finance and Sport; a gene can belong to the functions of Protein Synthesis, Metabolism and Transcription.

The traditional multi-label classification methods are not coping well with the increasing needs of today's big and complex data structure. As a result, there is a pressing need for new multi-label learning paradigms and new trends are emerging. This paper aims to provide a comprehensive survey on these emerging trends and the state-of-the-art methods, and discuss the possibility of future valuable research directions.

With the advent of the big data era, extreme multi-label classification (XMLC) becomes a rapidly growing new line of research that focuses on multi-label problems with an

extremely large number of labels. Many challenging applications, such as image or video annotation, web page categorization, gene function prediction, language modeling can benefit from being formulated as multi-label classification tasks with millions, or even billions, of labels. The existing MLC techniques can not address the XMLC problem due to the prohibitive computational cost given the large number of labels. One of the most pioneering work in XMLC is SLEEC [1], which learns a small ensemble of local distance preserving embeddings. The authors in SLEEC contribute a popular public Extreme Classification Repository¹, which promote the development of XMLC. The state-of-the-art XMLC techniques are mostly based on one-vs-all classifiers [2], [3], [4], [5], trees [6], [7], [8], [9], [10] and embeddings [1], [11], [12], [13], [14]. Unfortunately, the theoretical results in XMLC under the very high dimensional settings remain relatively under-explored. Moreover, the labels are extremely sparse, which leads to the problem of the long-tail distribution. How to precisely predict all the positive labels to testing examples pose a serious challenge in XMLC.

As the data volume grows quickly these days, it is usually expensive and time-consuming to acquire full supervision. In MLC tasks, the high dimensional output space makes it even harder. To mitigate this problem, a wealth of works have proposed various settings of MLC with limited supervision. For example, multi-label learning with missing labels (MLML) [15] assumes that only a subset of labels is obtained; semi-supervised MLC (SS-MLC) [16] admits a few fully labeled data and a large amount of unlabeled data; partial multi-label learning (PML) [17] studies an ambiguous setting that a superset of labels is given. Many effective models are also proposed based on graph [15], [18], [19], embedding [11], [20], [21], probability models [22], [23] and so on. More interesting improperly-supervised MLC settings are also considered recently, such as MLC with noisy labels [24], multi-label zero-shot learning [25] and multi-label active learning [26]. These settings make MLC

- *Weiwei Liu is with the School of Computer Science, Wuhan University, Wuhan 430079, China. E-mail: liuweimei863@gmail.com.*
- *Haobo Wang is with College of Computer Science and Technology, Zhejiang University. E-mail: wanghaobo@zju.edu.cn.*
- *Xiaobo Shen is with the School of Computer and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China. E-mail: njust.shenxiaobo@gmail.com.*
- *Ivor W. Tsang is with the Centre for Artificial Intelligence, FEIT, University of Technology Sydney, NSW, Australia. E-mail: ivor.tsang@uts.edu.au.*
- *This work is supported by the National Natural Science Foundation of China under Grant No. 61976161, 62176126 and 61906091, the Natural Science Foundation of Jiangsu Province, China (Youth Fund Project) under Grant No. BK20190440, the Fundamental Research Funds for the Central Universities under Grant No. 30921011210, the ARC under Grant No. DP180100106 and DP200101328. (Corresponding author: Weiwei Liu.)*

1. <http://manikvarma.org/downloads/XC/XMLRepository.html>

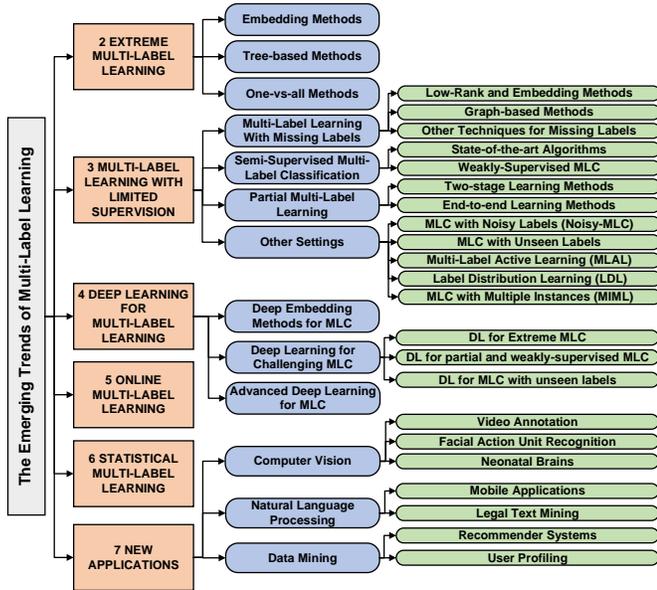


Fig. 1. The structure of this paper.

practical in real-world applications by saving supervision costs, and thus, deserve more attention.

Deep learning has shown excellent potential since 2012 when AlexNet presents surprising performance on the single-label image classification of ILSVRC2012 challenge. As most natural images usually contain multiple objects, it is more practical that each image is associated with multiple tags or labels. Thus developing deep learning techniques that can address MLC problem is more practically demanding in real-world image classification tasks. Some large-scale multi-label image databases, e.g., Open Images [27], newly released Tencent ML-Images [28] promote deep learning for MLC problem. In this area, BP-MLL [29] is the first method to utilize neural network (NN) architecture for MLC problem. Canonical Correlated AutoEncoder (C2AE) [30] is the first Deep NN (DNN) based embedding method for MLC problem. In addition, some deep learning methods are also developed for the Challenging MLC problems, such as Extreme MLC [31], [31], [32], [33], partial and weakly-supervised MLC [19], [23], [23], [34], MLC with unseen labels [35], [36], [36]. Recently advanced deep learning architectures [37], [38], [39], [40] for MLC problems are studied. How to harvest the strong learning capability of deep learning to better capture the label dependencies is key for deep learning to address MLC problems.

The Web continues to generate quintillion bytes of streaming data daily, leading to the key challenges for MLC tasks. Firstly, the existing off-line MLC algorithms are impractical for streaming data sets, since they require to store all data sets in memory; secondly, it is non-trivial to adapt off-line multi-label methods to the sequential data. Therefore, several approaches for online multi-label classification have recently been proposed, including [41], [42], [43]. However, both the experimental and theoretical results obtained so far are still not satisfactory and very limited. There is a real pressing need for credible research into online multi-label learning. Many references [44], [45] have shown

that methods of multi-label learning which explicitly capture label dependency will usually achieve better prediction performance. Therefore, in the past few years, modeling the label dependency is one of the major challenges in multi-label classification problems. A plethora of methods have been motivated to model the dependency. For example, the classifier chain (CC) model [46] captures label dependency by using binary label predictions as extra input attributes for the following classifiers in a chain. CCA [47] uses canonical correlation analysis for learning label dependency. CPLST [48] uses principal component analysis to capture both the label and the feature dependencies. Unfortunately, the statistical properties and asymptotic analysis of all these methods are still not well explored. One of the emerging trends is to develop statistical theory for understanding multi-label dependency modelings.

During the past decade, multi-label classification has been successfully applied in computer vision, natural language processing and data mining. This paper will briefly review these emerging applications, which may inspire the community to explore more interesting applications. The structure of this paper is shown in Figure 1. Some evaluation metrics and important notations used in this paper can be found in the Supplementary Materials.

2 EXTREME MULTI-LABEL LEARNING

Extreme multi-label classification (XMLC) aims to learn a classifier that is able to automatically annotate a data point with the most relevant subset of labels from an extremely large number of labels, which has opened up a new research frontier in data mining and machine learning. For example, there are millions of people who upload their selfies on the Facebook every day, based on these selfies, one might wish to build a classifier that recognizes who appear in the figure. Many XMLC applications have been found in various domains ranging from language modeling, document classification and face recognition to gene function prediction. The main challenging issue of XMLC is that XMLC learns with hundreds of thousands, or even millions, of labels, features and training points. To address this issue, the state-of-the-art XMLC techniques are mostly based on embeddings, trees and one-vs-all classifiers. We will review these advanced techniques in this section. Note that there are also some new deep learning-based XMLC models, but we leave the discussion until §4.

2.1 Embedding Methods

To deal with many labels, [49] assume that label vectors have a little support. In other words, each label vector can be projected into a lower dimensional compressed label space, which can be deemed as encoding. A regression is then learned for each compressed label. Lastly, the compressed sensing technique is used to decode the labels from the regression outputs of each testing instance. Many embedding based works have recently been developed in this learning paradigm. These works mainly differ in compression and decompression methods such as canonical correlation analysis (CCA) [47] and bloom filters [50]. Amongst them, SLEEC [1] is one of the seminal embedding methods in XMLC due to its simplicity and promising experimental results [1].

SLEEC learns low dimensional embeddings which non-linearly capture label correlations by preserving the pairwise distances between only the closest (rather than all) label vectors. Regressors are then trained in the embedding space. SLEEC uses a k -nearest neighbor (k NN) classifier in the embedding space for prediction.

Assume $x_i \in \mathbb{R}^{d \times 1}$ is a real vector representing an input or instance (feature), $y_i \in \{0, 1\}^{L \times 1}$ is the corresponding output or label vector ($i \in \{1, \dots, n\}$). n denotes the number of training data. The input matrix is $X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ and the output matrix is $Y = [y_1, \dots, y_n] \in \{0, 1\}^{L \times n}$. SLEEC maps the label vector y_i to ϖ -dimensional vector $z_i \in \mathbb{R}^{\varpi \times 1}$ ($\varpi < L$ is a small constant) and learns a set of regressors $V \in \mathbb{R}^{\varpi \times d}$ s.t. $z_i \approx Vx_i, \forall i \in \{1, \dots, n\}$. During the prediction, for a testing instance x , SLEEC first computes its embedding Vx and then perform k NN over the set $[Vx_1, \dots, Vx_n]$. We denote the transpose of the vector/matrix by the superscript T and the logarithms to base 2 by \log . Let $\|\cdot\|_F$ and $\|\cdot\|_1$ represent the Frobenius norm and ℓ_1 norm of a matrix.

SLEEC aims to learn a embedding matrix $Z = [z_1, \dots, z_n] \in \mathbb{R}^{\varpi \times n}$ through the following formula:

$$\min_{Z \in \mathbb{R}^{\varpi \times n}} \|P_{\Omega}(Y^T Y) - P_{\Omega}(Z^T Z)\|_F^2 \quad (1)$$

where the index set Ω denotes the set of neighbors: $(i, j) \in \Omega$ iff $j \in \mathcal{N}_i$. \mathcal{N}_i denotes a set of nearest neighbors of i . $P_{\Omega}(\cdot)$ is defined as:

$$(P_{\Omega}(Y^T Y))_{(i,j)} = \begin{cases} y_i^T y_j, & \text{if } (i, j) \in \Omega \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

Based on embedding matrix Z , SLEEC minimizes the following objective with ℓ_1 and ℓ_2 regularization to find regressors V , which is able to reduce the prediction time and the model size, and avoid overfitting.

$$\min_{V \in \mathbb{R}^{\varpi \times d}} \|Z - VX\|_F^2 + \mu \|V\|_F^2 + \lambda \|VX\|_1 \quad (3)$$

where $\mu > 0$ and $\lambda > 0$ are the regularization parameters.

To scale to large-scale data sets, SLEEC clusters the training set into smaller local region just based on features and does not consider label information. Therefore, the instances that have similar labels are not guaranteed to be split into the same region. This partitioning may affect the quality of embeddings learned in SLEEC.

Many methods have been developed to address this issue. For example, AnnexML [12] shows a novel graph embedding method based on the k -nearest neighbor graph (KNNG). AnnexML aims to construct the KNNG of label vectors in the embedding space to improve both the prediction accuracy and speed of the k -nearest neighbor classifier. DEFrag [51] represents each feature $j \in [d]$ as an L -dimensional vector $q^j = \sum_{i=1}^n x_i^j y_i$, which is a weighted aggregate of the label vectors of data points where the feature j is non-zero. After creating these representative vectors, DEFrag performs hierarchical clustering on them to obtain feature clusters, and then performs agglomeration by summing up the coordinates of the feature vectors within a cluster. [51] shows that DEFrag offers faster and better performance.

Word embeddings have been successfully used for learning non-linear representations of text data for natural language processing (NLP) tasks, such as understanding word and document semantics and classifying documents. Recently, [52] first proposes to use word embedding techniques to learn the label embedding of instances. [52] treats each instance as a ‘‘word’’, and define the ‘‘context’’ as k -nearest neighbors of a given instance in the space formed by the training label vectors y_i . Based on Skip Gram Negative Sampling (SGNS) technique, [52] learns embeddings z_1, \dots, z_n through the following formula:

$$\max_{z_1, \dots, z_n} \sum_{i=1}^n \left(\sum_{j \in \mathcal{N}_i} \log(\sigma(\langle z_i, z_j \rangle)) + C \sum_{j'} \log(\sigma(-\langle z_i, z_{j'} \rangle)) \right) \quad (4)$$

where $j' \in \{1, \dots, n\}$, $\sigma(\cdot)$ is a sigmoid function, $\langle \cdot, \cdot \rangle$ denotes the inner product and C is a constant. After learning label embeddings z_1, \dots, z_n , [52] follows the learning algorithm of SLEEC to learn V and make the prediction. [52] shows competitive prediction accuracies compared to state-of-the-art embedding methods, and provides the new insight for XMLC from the popular word2vec in NLP, which may open a new line of research.

The embedding matrix $Z = [z_1, \dots, z_n] \in \mathbb{R}^{\varpi \times n}$ of existing embedding methods is in real space. Hence we need to use regressors for training and may involve solving expensive optimization problems. To break this limitation, many references leverage coding technique for efficiently training the model. For example, based on Bloom filters [53], a well-known space-efficient randomized data structure designed for approximate membership testing, [50] designs a simple scheme to select the k representative bits for labels for training and proposes a robust decoding algorithm for prediction. However, Bloom filters may yield many false positives.

To address this issue, [54] transforms MLC to a popular group testing problem. In the group testing problem, one wish to identify a small number k of defective elements in a population of large size L . The idea is to test the items in groups with the premise that most tests will return negative results. Only few $\varpi < L$ tests are needed to detect the k defective elements. [54] trains ϖ binary classifiers on z_i and learn to test whether the data belongs to a group (of labels) or not, and then uses a simple inexpensive decoding scheme to recover the label vector from the predictions of the classifiers. Recently, [55] develops a novel sparse coding tree framework for XMLC based on Huffman coding and Shannon-Fano coding. [50], [54], [55] introduce the coding theory into MLC which is very novel and worth further research and exploration in this direction.

Remark. Embedding methods are the most popular strategies for addressing XMLC. SLEEC is a seminal work among them and recommended for the beginners to try. The major limitation of existing embedding methods is that the correlations between the input and output are ignored, such that their learned embeddings are not well aligned, which leads to degradation in prediction performance. How to build an embedding space that can preserve the relations between the input and output is an important research topic in the future. For example, [30], [56], [57] explore the

correlations between the input and output. They propose to jointly learn a semantic common subspace and view-specific mappings within one framework. The semantic similarity structure among the embeddings is further preserved, ensuring that close embeddings share similar labels. Another limitation of existing embedding methods is that both the training and testing time complexity are too high (See Table 1). Some techniques, such as random projection, hashing and parallelization, may be able to accelerate the training and testing process. Tree-based methods are able to obtain fast testing speed, which is discussed in the following paragraph.

2.2 Tree-based Methods

For tree-based methods, the original large-scale problem is divided into a sequence of small-scale subproblems by hierarchically partitioning the instance set or the label set. The root node is initialized to contain the entire set. A partitioning formulation is then optimized to partition a set in a node into a fixed number k of subsets which are linked to k child nodes. Nodes are recursively decomposed until a stopping condition is checked on the subsets. Each node involves two optimization problems: optimizing the partition criterion, and defining a condition or building a classifier on the feature space to decide which child node an instance belongs to. In the prediction phase, an instance is passed down the tree until it reaches a leaf (instance tree) or several leaves (label tree). For a label tree, the reached leaves contain the predicted labels. For an instance tree, the prediction is made by a classifier trained on the instances in the leaf node. Thus, the main advantage of tree-based methods is that the prediction costs are sub-linear or even logarithmic if the tree is balanced.

FastXML [6] presents to learn the hierarchy by optimizing the ranking loss function, normalized Discounted Cumulative Gain (nDCG). nDCG brings two main benefits to XMLC. Firstly, nDCG is a measurement which is sensitive to ranking and relevance and therefore ensures that the relevant positive labels are predicted with ranks that are as high as possible. This cannot be guaranteed by rank insensitive measures such as the Gini index or the clustering error. Second, by being rank sensitive, nDCG can be optimized across all L labels at the current node thereby ensuring that the local optimization is not myopic. The experiments show that nDCG is more suitable for extreme multi-label learning.

Based on FastXML, PfastreXML [7] studies how to improve the prediction accuracy of tail labels. The labels in XMLC follow a power law distribution. Infrequently occurring labels usually convey more information, but have little training data and are harder to predict than frequently occurring ones. PfastreXML improves FastXML by replacing the nDCG loss with its unbiased propensity scored variant, and assigns higher rewards for predicting accurate tail labels. Moreover, PfastreXML re-ranks PfastreXML's predictions using tail label classifiers. [7] shows that PfastreXML achieves promising performance in predicting tail labels and successfully applies to tagging, recommendation and ranking problems. SwiftXML [9] maintains all the scaling properties of PfastreXML, but improves the prediction accuracy of PfastreXML by considering more information about

TABLE 1

The training and testing time complexity of XMLC methods ($\text{nnz}(X)$ denotes the number of non-zeros in X , C is a constant, $O(\zeta)$ denotes the computational complexity of ϖ -bit Hamming distance calculation. T is the number of trees. h is the number of levels in the tree. c is the number of top-scoring items being reranked by the base-classifiers. $k \ll L$ is a small constant).

Methods	Training Time	Testing Time
Embedding: SLEEC [1]	$O(n\varpi^2 + n\varpi C)$	$O(nd + kL)$
Embedding: DEFrag [51]	$O(\text{nnz}(X) \log d)$	$O(nd + kL)$
Embedding: CoH [56]	$O(n(d^2 + L^2))$	$O(nc + kL)$
Tree: FastXML [6]	$O(nT \log L + \text{nnz}(X)nT)$	$O(Td \log L)$
Tree: SwiftXML [9]	$O(\text{nnz}(X)Tn \log n)$	$O((T \log n + c)\text{nnz}(X))$
Tree: GBDT-SPARSE [58]	$O(\text{nnz}(X)dTh \log k)$	$O(Tk \log k)$
OVA: PD-Sparse [59]	$O(ndC)$	$O(dL)$
OVA: LF [60]	$O(nd + L \log L + n \log n + nL)$	$O(d + \log(2L))$
OVA: Parabel [4]	$O((nd \log L)/L)$	$O(\text{nnz}(X)Tk \log L)$
OVA: Slice [5]	$O(nd \log L)$	$O(d \log L)$

revealed item preferences and item features. SwiftXML proposes a novel node partitioning function by optimizing two separating hyperplanes in the user and item feature spaces respectively. Experiments on tagging on Wikipedia and item-to-item recommendation on Amazon reveal that SwiftXML is more accurate than leading extreme classifiers by 14%.

FastXML, PfastreXML and SwiftXML have studied the ranking-based measures such as nDCG and its variants. Recently, [8] focuses on F-measure, which is a commonly used performance measure in multi-label classification as well as other fields, such as information retrieval and natural language processing. [8] proposes a novel sparse probability estimates (SPEs) to reduce the complexity of threshold tuning in XMLC. Then, they develop three algorithms for maximizing the F-measure in the Empirical Utility Maximization (EUM) framework by using SPEs. Moreover, Probabilistic label trees (PLTs) and FastXML are discussed for computing SPEs. Recently, the theory in [10] shows that the pick-one-label is inconsistent with respect to the Precision@ k , and PLTs model can get zero regret (i.e., it is consistent) in terms of marginal probability estimation and Precision@ k in the multi-label setting. Inspired by [10], [61] further studies the consistency of other reduction strategies based on a different Recall@ k metric.

Remark. Tree-based methods are the efficient strategy for addressing XMLC with the logarithmic dependence to the number of labels (See Table 1). FastXML is a popular method and recommended for the practitioners. However, one of the major problems for tree-based methods, such as FastXML, PfastreXML and SwiftXML, is that they involve complex non-convex optimization problem at each node. How to obtain cheap and scalable tree structure is an important research topic in the future. For example, GBDT-SPARSE [58] studies the gradient boosted decision trees (GBDT) for XMLC. In each node, the feature is firstly projected into a low-dimensional space and then a simple inexact search strategy is used to find a good split. They significantly reduce the training and prediction time and model size of GBDT to make it suitable for XMLC. CRAFTML [62] tries to use fast partitioning strategies and exploit random forest algorithm. CRAFTML first randomly projects the feature and label into lower dimensional spaces. A k -means algorithm is then used in the projected labels to partition the instances into k temporary subsets. Moreover,

GBDT-SPARSE and CRAFTML also open the way to parallelization, which are able to motivate further research.

2.3 One-vs-all Methods

One-vs-all (OVA) methods are one of the most popular strategies for multi-label classification which independently trains a binary classifier for each label. However, this technique suffers two major limitations for XMLC: 1) Training one-vs-all classifiers for XMLC problems using off-the-shelf solvers such as Liblinear can be infeasible for computation and memory. 2) The model size for XMLC data set can be extremely large, which leads to slow prediction. Recently, many works have been developed to address the above issues of the one-vs-all methods in XMLC.

By exploiting the sparsity of the data, some sub-linear algorithms are proposed to adapt one-vs-all methods in the extreme classification setting. For example, PD-Sparse [59] proposes to minimize the separation ranking loss and ℓ_1 penalty in an Empirical Risk Minimization (ERM) framework for XMLC. The separation ranking loss penalizes the prediction on an instance by the highest response from the set of negative labels minus the lowest response from the set of positive labels. PD-Sparse obtains an extremely sparse solution both in primal and in dual with the sub-linear time cost, while yields higher accuracy than SLEEC, FastXML and some other XMLC methods. By introducing separable loss functions, PPDSparse [3] parallelizes PD-Sparse with sub-linear algorithms to scale out the training. PPDSparse can also reduce the memory cost of PDSparse by orders of magnitude due to the separation of training for each label. DiSMEC [2] also presents a sparse model with a parameter thresholding strategy, and employs a double layer of parallelization to scale one-vs-all methods for problems involving hundreds of thousand labels. ProXML [63] proposes to use ℓ_1 -regularized Hamming loss to address the tail label issues, and reveals that minimizing one-vs-all method based on Hamming loss works well for tail-label prediction in XMLC based on the graph theory.

PD-Sparse, PPDSparse, DiSMEC and ProXML have obtained high prediction accuracies and low model sizes. However, they still train a separate linear classifier for each label and linear scan every single label to decide whether it is relevant or not. Thus the training and testing cost of these methods grow linearly with the number of labels. Some advanced methods are presented to address this issue. For example, to reduce the linear prediction cost of one-vs-all methods, [60] proposes to predict on a small set of labels, which is generated by projecting a test instance on a filtering line, and retaining only the labels that have training instances in the vicinity of this projection. The candidate label set should keep most of the true labels of the testing instances, and be as small as possible. They train the label filters by optimizing these two principles as a mixed integer problem. The label filters can reduce the testing time of existing XMLC classifiers by orders of magnitude, while yields comparable prediction accuracy. [60] shows an interesting technique to find a small number of potentially relevant labels, instead of going through a very long list of labels. How to use label filters to speed up the training time is left as an open problem.

Parabel [4] reduces training time of one-vs-all methods from $O(ndL)$ to $O((nd \log L)/L)$ by learning balanced binary label trees based on an efficient and informative label representation. They also present a probabilistic hierarchical multi-label model for generalizing hierarchical softmax to the multi-label setting. The logarithmic prediction algorithm is also proposed for dealing with XMLC. Experiments show that Parabel could be orders of magnitude faster at training and prediction compared to the state-of-the-art one-vs-all extreme classifiers. However, Parabel is not accurate in low-dimension data set, because Parabel can not guarantee that similar labels are divided into the same group, and the error will be propagated in the deep trees. To reduce the error propagation, Bonsai [64] shows a shallow k -ary label tree structure with generalized label representation. A novel negative sampling technique is also presented in Slice [5] to improve the prediction accuracy for low-dimensional dense feature representations. Slice is able to cut down the training time cost of one-vs-all methods from linear to $O(nd \log L)$ by training classifier on only $O(n/L \log L)$ of the most confusing negative examples rather than on all n training set. Slice employs generative model to estimate $O(n/L \log L)$ negative examples for each label based on approximate nearest neighbor search (ANNS) in time $O((n+L)d \log L)$, and conduct the prediction on $O(\log L)$ of the most probable labels for each testing data. Slice is up to 15% more accurate than Parabel, and able to scale to 100 million labels and 240 million training points. The experiments in [5] show that negative sampling is a powerful tool in XMLC, and the performance gain of some advanced negative sampling technique may be explored for future research.

Remark. One-vs-all methods are the simple strategies for dealing with XMLC, and PD-Sparse is the first choice for the beginners to try. As mentioned before, one-vs-all methods independently train a binary classifier for each label, so computation and memory cost pose an intractable issue for XMLC, and one-vs-all methods do not consider the correlations between labels. Although the reviewed methods in this subsection are able to ease the computation issue, how to use the correlations between labels to boost the performance of one-vs-all methods could pose a serious problem in the future. One possible way is to design some one-vs-all learning models which consider various label correlations.

3 MULTI-LABEL LEARNING WITH LIMITED SUPERVISION

Collecting fully-supervised data is usually hard and expensive and thus a critical bottleneck in real-world classification tasks. In MLC problems, there exist many ground-truth labels and the output space can be very large, which further aggravates the difficulty of precise annotation. To mitigate this problem, plenty of works have studied different settings of MLC with limited supervision. How to model label dependencies and handle incomplete supervision pose two major challenges in these tasks. In this section, we concentrate on several advanced topics. Amongst them, multi-label learning with missing labels (MLML) assumes only a subset of labels are given; semi-supervised MLC (SS-MLC) assumes a large set of unlabeled data as well

labeled data are given; partial multi-label learning (PML) allows the annotators to provide a superset of labels as the candidates. We illustrate the connections between these different supervision types in Figure 2. Note that in these settings, though trained with imperfect supervised signals, the classifier is still evaluated on a perfectly supervised testing data set to quantify the predictive performance.

3.1 Multi-Label Learning With Missing Labels

In real-world scenarios, it is intractable for the annotators to figure out all the ground-truth labels, due to the complicated structure or the high volume of the output space. Instead, a subset of labels can be obtained, which is called multi-label learning with missing labels (MLML). There are two main settings in MLML. The first setting [15] only obtains a subset of relevant labels. It views the MLML problem as a positive-unlabeled learning task such that the remaining labels are all regarded as negative labels. The other setting [65] explicitly indicates which labels are missing. Formally, given a feature vector x_i , we denote the label vector of these two settings by $\hat{y}_i \in \{-1, +1\}^{L \times 1}$ (-1 can be missing or negative labels) and $\tilde{y}_i \in \{-1, 0, +1\}^{L \times 1}$ (0 represents missing labels) respectively. We distinguish these two settings in Figure 2. Moreover, two different learning targets may be considered. One is transductive that only learns to complete the missing entries. The other is inductive where a classifier is trained for unseen data. For simplicity, we do not explicitly distinguish these differences.

Next, we will review state-of-the-art MLML methods which are mainly based on low-rank and graph assumptions.

3.1.1 Low-Rank and Embedding Methods

As discussed in §2.1, the existence of label correlations usually implies the output space is low-rank. Interestingly, this assumption has been widely used to complement the missing entries of a matrix in matrix completion tasks [66]. Since it benefits the two key targets in MLML, i.e. label correlation extraction and missing label completion, many low-rank assumption-based MLML methods have been developed.

In [66], the MLML problem is regarded as a low-rank matrix completion problem with the existence of side information, i.e. the features. To accelerate the learning task, the label matrix is decomposed to be $Y = AWB$ where A and B are side information matrices. W is assumed to be low-rank. In fact, in MLML problems, A is exactly the feature matrix X and B is the identity matrix since there is no side information for the labels. Therefore, W can be viewed as a linear classifier that enables the predicted labels $Y = XW$ to be low-rank. Then, LEML [11] generalizes this paradigm to a flexible empirical risk minimization framework. The formula is as follows,

$$W = \arg \min_W \mathcal{L}(\hat{Y}, XW) + \lambda r(W), \quad \text{s.t. } \text{rank}(W) \leq k \quad (5)$$

where λ and k are constants, $r(W)$ is the regularizer. \mathcal{L} can be any empirical risk that is evaluated on observed entries. To solve this problem, [11] decomposes the classifier to two rank- k ($k \ll L$) matrices V and U such that $W = VU$. Then, an alternative optimization method is used to efficiently

handle large-scale problems. Nevertheless, the presence of tail labels may break the low-rank property. Hence, [20] treats the tail labels as outliers and decompose the label matrix to, $\hat{Y} \approx Y_1 + Y_2$. Here Y_1 is low-rank and Y_2 is sparse. These two components can be obtained by solving the following objective,

$$\begin{aligned} \min_{U, V, H} & \|\hat{Y} - Y_1 - Y_2\|_F^2 + \lambda_1 \|H\|_F^2 \\ & + \lambda_2 (\|U\|_F^2 + \|V\|_F^2) + \lambda_3 \|XH\|_1 \quad (6) \\ \text{s.t. } & Y_1 = XUV, \quad Y_2 = XH \end{aligned}$$

These two learning frameworks are followed by many works. For example, [67] studies the problem that both features and labels are incomplete. The proposed solution, ColEmbed, requires the classifier as well as the recovered feature matrix to be low rank. Moreover, the kernel trick is used to enable the non-linearity of the classifier. Some recent works [68], [69], [70] further integrate the graph-based technique to get more effective models, which we will discuss in §3.1.2.

The low-rank assumption is rather flexible and may be exploited in various ways. For example, COCO [71] considers a more complex setting that the features and labels are missing simultaneously. It imposes the concatenation of recovered feature matrix and label matrix to be low-rank via trace norm. Some works also utilize the assumption through a low-rank label correlation matrix. ML-LRC [72] assumes the label matrix can be reconstructed using a correlation matrix U such that $Y = \hat{Y}^T U$, where $U \in \mathbb{R}^{L \times L}$ is low-rank. Then, the loss is measured using the output and the reconstructed labels $\|XW - YU\|_F^2$. Based on this assumption, ML-LEML [73] further involves an instance-wise label correlation matrix V such that $Y = \hat{Y}V$, where $V \in \mathbb{R}^{n \times n}$ is also low-rank.

Another popular way is to follow the paradigm of embedding methods that projects the label vectors to a low-dimensional space. [30] proposes a deep neural network-based model C2AE. The features and labels are jointly embedded to a latent space using two neural networks F_x and F_e such that their codewords are maximally correlated. Then, the feature codewords are decoded by another neural network F_d , which is also used for prediction. For MLML problems, the decoded labels are evaluated on the observed entries. Though the labels are decoded from a low-dimensional space, the low-rank assumption need not be satisfied with non-linear projection. Thus, REFDHF [74] added a trace norm regularization term on the decoded label matrix. [74] also proposes a novel hypergraph fusion technology that explores and utilizes the complementary between feature space and label space. Compared to low-rank classifier-based methods, the embedding methods are more flexible since the classifier can be non-linear and thus are worthy to be explored.

3.1.2 Graph-based Methods

To handle missing labels, one of the most popular solutions is graph-based model. Denote a weighted graph by $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$, where $\mathcal{V} = \{x_i | 1 \leq i \leq n\}$ denotes the vertex set and $\mathcal{E} = \{(x_i, x_j)\}$ denotes the edge set. $\mathcal{W} = [w_{ij}]_{n \times n}$ is a weight matrix where $w_{ij} = 0$ if $(x_i, x_j) \notin \mathcal{E}$. With the graph being defined, the most typical

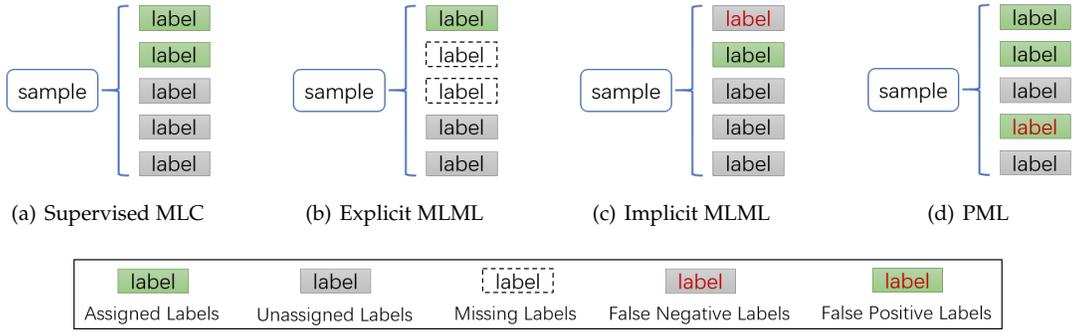


Fig. 2. Illustration of some MLC settings with different types of supervision. (a) instances with full supervision; (b) instances with explicitly missing labels; (c) instances with implicitly missing labels; (d) instances with a set of candidate relevant labels. Here (b) and (c) are two different settings of MLML problems. Semi-supervised MLC is a special case of (b) where some instances miss all the labels.

strategy is adding a manifold regularization term to the empirical risk minimization framework. Note that in this section, we slightly abuse the notation w_{ij} to represent the graph weight entry for the sake of simplicity.

The pioneering work [15] is the first to propose the concept of multi-label learning with weak labels, i.e. the implicit setting of MLML. The proposed method, named WELL, constructs a label-specific graph for each label from a feature-induced similarity graph. Then, the manifold regularization terms are added separately for each label. [65] formalized the other setting of MLML and involves three assumptions into MLML according to [16],

- **Label Consistency:** the predicted labels should be consistent with the initial labels, which is usually achieved by empirical risk minimization principle;
- **Sample-level Smoothness:** if two samples x_i and x_j are close to each other, so are their predicted label vectors;
- **Label-level Smoothness:** if two incomplete label vectors y_i and y_j are semantically similar, so are the predicted label vectors.

Formally, a k -nearest neighbor graph is constructed to satisfy the sample-level smoothness, where weight matrix \mathcal{W}^x is computed by $w_{ij}^x = \exp(-\frac{\|x_i - x_j\|_2^2}{\|x_i - x_h\|_2 \|x_j - x_h\|_2})$, where x_h is the h -th nearest neighbor of x_i (h is a fixed constant). For the label-level smoothness, the authors constructs a L -square weight matrix \mathcal{W}^y where $w_{ij}^y = \exp(-\eta[1 - \frac{\langle \tilde{Y}_{i \cdot}, \tilde{Y}_{j \cdot} \rangle}{\|\tilde{Y}_{i \cdot}\|_2 \|\tilde{Y}_{j \cdot}\|_2}])$. $\tilde{Y}_{i \cdot}$ is the i -th row vector of incomplete matrix \tilde{Y} . Finally, the predicted score matrix \hat{Y} is learned by,

$$\min_{\hat{Y}} \|\hat{Y} - \tilde{Y}\|_F^2 + \frac{\lambda_x}{2} \text{Tr}(\hat{Y} L_x \hat{Y}^T) + \frac{\lambda_y}{2} \text{Tr}(\hat{Y}^T L_y \hat{Y}) \quad (7)$$

where L_x and L_y is the laplacian matrix of \mathcal{W}^x and \mathcal{W}^y . λ_x and λ_y are trade-off parameters. This learning paradigm is followed by some recent works. [75] proposes an inductive version that the trained classifier can also predict on unseen data. [76] chooses the hinge loss as the empirical risk instead of squared loss. To tackle the severe class imbalance problem in MLML, [77] add two class cardinality constraints to Eq. (7) that enforces the number of positive labels is in a predefined range. With hierarchical label information being provided, MLMG-GO [70] involves a semantic hierarchical constraints such that the score of a label y_a is smaller than its parent label y_b . In [19], a new regularization framework

IMCL is proposed that interactively learns the two similarity graphs.

Many graph-based methods only concentrate on the sample-level smoothness principle. That is, the graph information is mainly used for disambiguating the incomplete supervision, and different techniques are involved to utilize the label correlations. [69] treats the problem of one-class matrix factorization with side information as an MLML task. Inspired by [11], the linear classifier is restricted to be low-rank and the predicted label matrix is smoothed by a manifold regularization term. Since the low-rank assumption fails in many applications, MLMG-SL [70] further assumes that the output of graph model can be decomposed to a low-rank matrix and a sparse matrix. There are also several recent works that focus on the label-level smoothness. LSML [78] proposes to learn a label correlation matrix, i.e. a label graph, that can be used to complement the missing labels, smooth the label prediction and guide the learning of label-specific features simultaneously. GLOCAL [79] trains a low-rank model with manifold regularization that exploits global label-level smoothness. In addition, as label correlations may vary from one local region to another, GLOCAL partitions the instances to several groups and learns local label correlations by group-wise manifold regularization. In [34], a fully connected graph is built whose vertices are the labels and then, a graph neural network (GNN) is trained to model the label dependencies. The input of GNN is the L -sized feature vector of the image extracted by a convolutional neural network, and the outputs are the predicted labels. To disambiguate the missing labels, [34] proposes two novel strategies. For the known labels, the authors propose partial binary cross-entropy loss (Partial-BCE) that reduces the normalization factor according to the label proportion. To complete the missing entries, [34] adopts a curriculum learning strategy that learns a self-paced model.

Besides, some studies are also interested in different graph information. APG-Graph [68] proposes a novel semantic descriptor-based approach for visual tasks to construct an instance-instance correlation graph. Specifically, [68] makes use of the posterior probabilities of the classifications on other public large-scale data sets. Then, a k -NN graph is constructed by these predicted tags. [69] regards the user-item interaction in the recommender system as a

bipartite graph.

In the past few years, graph mining techniques have received huge attention. We believe the future graph-based MLML models will involve more expressive graph models, e.g. graph neural networks [80], and various types of graphs, e.g. social networks [81].

3.1.3 Other Techniques for Missing Labels

There are many other techniques can be used for MLML tasks, such as co-regularized learning [82], binary coding embedding [83]. In what follows, we focus on some advanced MLML algorithms.

Due to the capability of exploring the data distribution, probability graphical models (PGMs) have been popular for MLML problems since we can complement the missing labels in a generative manner. SSC-HDP [84] involves a correspondence hierarchical Dirichlet process (Corr-HDP) that enables the dimension of latent factors to be chosen dynamically. Based on Corr-HDP, SSC-HDP iteratively updates the likelihood $P(y^j|x)$ for an instance x whose j -th label is missing, while the likelihood of remaining labels is fixed to 1. CRBM [85] proposes a conditional restricted Boltzmann machine model to capture the high-order label dependence relationships. In specific, a latent layer is added above the labels layer to form a restricted Boltzmann machine, while the features are the conditions. Based on a latent factor model, GenEML [22] proposes a scalable generative model that involves an exposure variable for each missing labels. BMLS [86] jointly learns a low-rank embedding of the label matrix and the label co-occurrence matrix using an Poisson-Dirichlet-gamma non-negative factorization method [87]. Note that [85], [86] are also capable of incorporating auxiliary label relatedness information, such as Wikipedia.

Reweighting empirical risks is also a common strategy. [88] notices that in MLML setting, the traditional multi-label ranking error may overestimate the classification error. Hence, a slack variable is introduced to account for the error of ranking an unassigned class before the assigned class. [7] proposes an unbiased propensity scored variant of nDCG loss and [34] presents Partial-BCE, which we have discussed in previous sections. [89] assigns a weight factor for each term in binary cross-entropy loss. In particular, the weights of the positive labels are fixed to 1. The weights of missing entries are set as $P(\hat{y}_c|y)$, i.e. the probability of having a negative label for the c -th label given the vector of labels y . Specifically, the probability is estimated from the ground-truth label matrix based on label co-occurrences.

Recently, bandit learning-based approaches are also introduced. Specifically, one pioneering work [90] considers the contextual bandits problem in the extreme multi-label learning context. It modifies the inverse gap weighting sampling strategy to select top- k arms, which results in good generalization performance. Besides, this work proposes a tree-based algorithm by grouping similar arms and thus, the model enjoys a poly-logarithm computational cost w.r.t. the number of arms.

Remark To date, graph-based methods and embedding-based methods are still dominant in the MLML context. Though recent works [19] have tried to involve deep models to promote performance, they mainly involve trivial convolutional networks and autoencoders. It would be

promising to design more tailored model architectures for MLML. Other techniques are also worth to be explored. For example, with the success of existing PGM-based MLML methods, we believe that bayesian deep learning (BDL) [91] can further improve the performance due to its superiority on high-dimensional data and complex uncertainty.

3.2 Semi-Supervised Multi-Label Classification

In semi-supervised MLC (SS-MLC) [92], the data set is comprised of two sets: fully labeled data and unlabeled data. Though SS-MLC has a far longer history than MLML, we can regard it as a special case of MLML, i.e. the labels of some instances are totally missing. In fact, similar to MLML, a plenty of SS-MLC algorithms are also based on graph models [93], [94], and low-rank assumptions [79], [95]. In what follows, we first review some state-of-the-art SS-MLC algorithms and then, discuss a novel learning setting called weakly-supervised MLC.

3.2.1 State-of-the-art Algorithms

Graph-based methods are very popular in SS-MLC, which mainly differ in the strategy of utilizing the label-correlation. SLRM [96] enforces the classifier to be low-rank, while a manifold-regularization term is added to ensure the sample-level smoothness. [95] proposes a triple low-rank regularization approach where the graph is dynamically updated using a low-rank feature-recovery matrix. Based on curriculum learning, ML-TLLT [97] forces a teacher pair to generate similar curriculums if the corresponding two labels are highly correlated over the labeled examples. CMLP [94] makes use of collaboration technique [98] to design an scalable multi-label propagation method. Specifically, it breaks the predicted label into two parts: 1) its own prediction part; 2) the prediction of others, i.e. collaborative part.

As mentioned above, other techniques may also be used. COIN [99] adapts the well-known co-training strategy to SS-MLC setting. In each co-training round, a dichotomy over the feature space is learned by maximizing the diversity between the two classifiers induced on either dichotomized feature subset. Then, pairwise ranking predictions on unlabeled data are iteratively communicated for model refinement. Based on COIN, [100] further proposes an ensemble method to accommodate streamed SS-MLC data. DRML [101] designs a dual-classifier domain adaptation network to align the features in a latent space. In order to model label dependencies, DRML generates the final prediction by feeding the outer-product of the dual predicted label vectors to a relation extraction network.

3.2.2 Weakly-Supervised MLC

Due to the large output space, even in the SS-MLC problems, collecting precisely labeled data would take extensive efforts and costs. Hence, a new setting called weakly-supervised multi-label classification (WS-MLC) has attracted enormous attention, i.e. there might be fully labeled data, incompletely-labeled data and unlabeled data in the data set simultaneously. In this survey, we follow the definition of WS-MLC in [23]. However, weakly-supervised MLC may also have other meanings in the literature. In a broad

sense, any noisy supervision can be termed as weakly-supervision. The readers should also be careful about the difference between WS-MLC and multi-label learning with weak labels [15], [102]. The latter sometimes indicates the implicit setting of MLML problems.

Many effective approaches have been developed to deal with WS-MLC problems. For example, WeSed [103] handles the missing labels by a weighted ranking loss and integrates the unlabeled data via a triplet similarity loss. In [104], missing labels are first estimated by a correlation matrix. Then, a linear classifier is trained by minimizing a graph regularized model. SSWL [105] proposes a novel dual similarity regularizer $\|Y - VYU\|$ to characterize both sample-level and label-level smoothness. Here V is the weight matrix of k NN graph over training data and U is a trainable variable that represents the label similarity. Moreover, SSWL also utilizes an ensemble of multiple models to improve the robustness. Though these works have demonstrated promising results, they directly use logical labels, and thus, ignore the relative importance of each label to an instance. To bridge this gap, WSMLLE [106] transforms the original problem to a label distribution learning problem [107]. In specific, a new label enhancement method is proposed that marries the concept of local correlation [79] and dual similarity regularizer [105]. The label enhancement technique is also adopted by fully-supervised MLC [108] and PML [109] models, and we will give a detailed discussion about the latter one in §3.3.

Probabilistic models are also popular in solving WS-MLC tasks, since the distribution of unlabeled data can be seamlessly integrated into a probabilistic framework. DSGM [23] proposes a deep sequential generative model which assumes an instance x is generated from its label y as well as a latent variable z . DSGM leverages information from observed labels in a sequential manner. Then, the model is trained by maximizing the likelihood,

$$\max_{\theta} \sum_{i \in D_l} \log p_{\theta}(x_i, y_i) + \sum_{j \in D_o} \log p_{\theta}(x_j, \tilde{y}_j) + \sum_{k \in D_u} \log p_{\theta}(x_k) \quad (8)$$

where θ is the model parameter. D_l , D_o and D_u are the index sets of fully labeled data, incompletely-labeled data and unlabeled data respectively. [23] also proposes a variational inference method that minimizes the evidence lower bound of the objective. [110] designs an embedding-based probability model called ESMC, which addresses some key issues in WS-MLC tasks. Since the low-rank assumption may be broken by tail labels, ESMC uses the gaussian processes to perform non-linear projection. To handle missing labels, ESMC introduces a set of auxiliary random variable, a.k.a. experts, to model the relationship between the real-valued probability score and the observed logical labels. Finally, the unlabeled data can also be integrated to learn a smooth mapping from the feature space to the label space.

Remark. Compared to MLML problems, the presence of a large amount of unlabeled data in SS-MLC can highly restrict the representation ability of the model. However, few efforts have been made to apply tricks in state-of-the-art deep semi-supervised learning to SS-MLC. We recommend involving techniques such as consistency regulariza-

tion [111] and self-supervised pretraining [112], which have demonstrated exciting ability to utilize the unlabeled data.

3.3 Partial Multi-Label Learning

In practice, the complicated structure of the label space usually makes it hard to decide some *hard* labels are relevant or not. For example, it is usually hard to decide whether a dog is a malamute or a husky. One might naively drop these labels and regard the original problem as an MLML task. However, missing labels provides no information to the user at all. Hence, partial multi-label learning (PML) [17] is proposed to address this issue, which preserves all the potentially correct labels. Formally speaking, each instance x is equipped with a set of candidate labels S , only some of which are the true relevant labels. The remaining labels are called *false positive labels* or *distractor labels*. Technically, PML can be regarded as a dual problem of MLML and solved by existing MLML techniques. However, it is worth noting that this strategy may be less practical owing to the sparsity of the label space. Moreover, PML also provides a safe way to protect data privacy since no label can be determined as the ground-truth, as opposite to MLML data.

3.3.1 Two-stage Learning Methods

In PML, while label correlation still matters, the other key issue becomes identifying the ground-truth from the candidate label set instead of completion. To handle these issues, some PML algorithms adopt a two-stage learning framework. Formally, an enriched label representation $\Lambda = [\lambda_{ij}] \in \mathbb{R}^{L \times n}$ will be learned where λ_{ij} is a real-valued number. The sign of λ_{ij} indicates whether the label is positive or negative, while the magnitude reflects the confidence of the relevance. Then, the PML problem is transformed into a canonical supervised learning problem and the classifier can be easily induced. To obtain Λ , PARTICLE [113] uses the label propagation technique that aggregates the information from the k -nearest neighbors. After that, the confidences are converted back to logical labels by thresholding. To train an MLC classifier, [113] adopts a pairwise label ranking model coupled with virtual label splitting or maximum a posteriori (MAP) reasoning. PARTICLE has two main drawbacks. First, the confidences have richer information than logical labels, but, it is trimmed when thresholding. Second, only the second-order label correlation is considered.

To tackle these problems, DRAMA [18] generates the label confidence matrix under the guidance of feature manifold and the candidate label set. Then, a novel gradient boosting decision tree (GBDT) based multi-output regressor is directly trained on the transformed data set $\tilde{D} = \{(x_i, \lambda_i) | i \in \{1, \dots, n\}\}$ where λ_i is the i -th column vector of Λ . On t -th boosting round, DRAMA augments the feature space using previously learned labels. Therefore, high-order label correlations are automatically exploited to improve performance.

The major limitation of the aforementioned methods is that the disambiguation is achieved purely by features. However, label correlation itself can help to identify the correct labels. Insufficient disambiguation makes the induced MLC classifier error-prone. To this end, PML-LD [109] proposes a novel label enhancement method that transforms

the PML problem to a label distribution learning problem [107]. When learning the label confidence matrix, PML-LD leverages the sample-level smoothness and local label-level smoothness [79] such that the candidate label set can be fully disambiguated. Then, the confidences are normalized by softmax to form an LDL problem and a multi-output support vector machine is induced.

The advantages of two-stage PML methods are two-folds. First, since the label confidences are obtained, we can apply well-studied multi-output learning methods [114]. Second, the real-valued confidences reflect the relative intensity of the relevance or irrelevance, which may give us more information about our data.

3.3.2 End-to-end Learning Methods

As we have mentioned, two-stage learning PML methods usually need be carefully designed, or the induced MLC classifier may be error-prone due to insufficient disambiguation. Hence, many PML algorithms are developed in an end-to-end fashion, which vary from one to another.

[17] proposes a ranking model, which employs the label confidence as a weight for the ranking loss. To estimate the label confidences, [17] provides two practical ways based on label correlation and feature prototypes respectively. Moreover, the classification model along with the ground-truth confidence are optimized in a unified framework such that the two subproblems can benefit from each other. [115] presents a soft sign thresholding method to measure the discrepancy between the real-valued confidences and the candidate labels. Similar to [17], the classifier training and disagreement minimization are performed at the same time. Nevertheless, [115] does not well utilize the label correlations, and thus the performance is limited.

Some methods adopt the low-rank assumption. fPML [116] introduces the matrix factorization technique to obtain a shared latent space for both features and labels. The classifier is then trained by fitting the recovered labels. PML-LRS [21] utilizes the low-rank and sparse decomposition scheme. That is, it assumes the distractor label matrix is sparse while the ground-truth matrix is low-rank. Both fPML and PML-LRS treat the false-positive labels as randomly generated noise. However, in real-world applications, the false-positive labels may be caused by some ambiguous contents of the instance. Therefore, [117] divides the classifier W to two parts $W = U + V$. Here U is the multi-label classifier and V is the distractor label identifier. Meanwhile, U is constrained to be low-rank to utilize label correlations. Since distractor labels usually correlate to only a few ambiguous features, V is regularized to be sparse. MUSER [118] takes redundant labels together with noisy features into account by jointly exploring feature and label subspaces. Furthermore, it uses a manifold regularizer to ensure the consistency between features and latent labels.

Remark. The PML problem is drawing increasing attention in the community. However, the assumption that all labels are equally being candidates can be less practical, since some ground-truth labels can be easily distinguished. Therefore, the key assumptions of PML should be carefully revisited. Here we recommend a more practical setting that besides providing the candidate set, the annotators should also provide partial ranks that which labels are more likely

to be correct. Besides, existing PML data sets are mainly built upon multi-instance multi-label [119] data sets, and thus, there is also an urgent need to establish a benchmark for PML problems.

3.4 Other Settings

The complexity of the label space has expedited various kinds of improperly-supervised MLC settings. In what follows, we briefly review some more state-of-the-art settings in the literature.

MLC with Noisy Labels (Noisy-MLC). While MLML and PML consider single-side noise, Noisy-MLC assumes that noisy labels occur in both relevant and irrelevant labels. Many effective Noisy-MLC algorithms have been proposed to address this problem, including graph based methods [120], probability models [121], teacher-student model [122]. In [106], the WS-MLC framework is extended and noisy labels are assumed to be contained in the data set. Some works [24], [123] maintain a small set of clean data to reduce the noise in the large data set. Since learning from label noise have been a hot topic in the community, Noisy-MLC deserve more attention.

MLC with Unseen Labels. In the aforementioned settings, the label spaces is fixed during training and testing. However, in practice, the label space may be dynamically expanded. For instance, [124] studies an online MLC setting that an arriving data instance may be associated with unknown labels. In [35], knowledge distillation method is used to handle streaming labels. Multi-label zero-shot learning (ML-ZSL) [36] requires the prediction of unknown labels which are not defined during training. To make ML-ZSL feasible, external semantic information is usually involved, such as word vectors [125] and knowledge graphs [36]. In [126], few-shot labels is also considered, which relates to only few instances in the data set, i.e. nearly unseen.

Multi-Label Active Learning (MLAL). Active learning is a notable way to alleviate the difficulty of multi-label tagging. The idea is to carefully select the most informative data instances for labeling such that better models can be trained with less labeling effort. A variety of works have studied MLAL problems. For example, [127] adopts maximum loss reduction with maximal confidence as the sampling criterion for MLAL. [26] solves MLAL problems via a probability model. Moreover, MLAL is also considered in crowdsourcing [128] and novel queries [129] tasks.

Label Distribution Learning (LDL). LDL [107] is a general framework that assigns L normalized real values to label description degree. It aims to tackle inherent ambiguity in data annotation, e.g. a facial expression usually conveys a complex mixture of basic emotions. Since it is difficult to obtain the label distribution directly, many works [108], [130], [131] focus on recovering label distributions from logical labels, which is also known as label enhancement (LE). LE is an effective learning strategy to deal with label ambiguity. LE is also applied in WS-MLC [106] and PML [109] to handle imperfect supervision signals.

MLC with Multiple Instances (MIML). MIML [119] is a popular setting which assumes each example is described by multiple instances as well as associated with multiple binary labels. Recent studies in MIML [132], [133] have

developed many deep learning models such that noisy instances can be effectively figured out. Nevertheless, MIML mainly focuses on the instance-level ambiguity instead of the labels. Hence, we do not further discuss it.

Remark. Intelligent systems are enrolled in increasingly difficult and complicated tasks, and thus new settings like PML and LDL deserve more attention. Moreover, there remain more challenging and complicated settings in real-world applications to be explored. For instance, there might be out-of-distribution detection [134], domain shift [135] and other problems arise in MLC problems.

4 DEEP LEARNING FOR MULTI-LABEL LEARNING

Due to the powerful learning capability, deep learning has achieved state-of-the-art performance in many real-world multi-label applications, e.g., multi-label image classification. In MLC problems, it is key to harvest the advantage of deep learning to better capture the label dependencies. In this section, we first introduce some representative deep embedding methods for MLC, then present deep learning for challenging MLC, and finally review advanced deep learning for MLC.

4.1 Deep Embedding Methods for MLC

Different from conventional multi-label methods, deep neural networks (Deep NNs) often seek a new feature space and employ a multi-label classifier on the top. To our knowledge, BP-MLL [29] is the first method to utilize NN architecture for multi-label learning problem. To explicitly exploit the dependencies among labels, given the neural network F , BP-MLL introduces a pairwise loss function for each instance x_i :

$$E_i = \frac{1}{|y_i^1| |y_i^0|} \sum_{(p,q) \in y_i^1 \times y_i^0} \exp(-(F(x_i)^p - F(x_i)^q)) \quad (9)$$

where y_i^1 and y_i^0 denote the sets of positive and negative labels for the i -th instance x_i respectively, $(F(x_i))^p$ denotes the p -th entry of $F(x_i)$. $F(x_i)^p - F(x_i)^q$ measures the difference between the outputs of the network on the positive and negative labels, and the exponential function is used to severely penalize the difference. Thus the minimization of (9) leads to output larger values for positive labels, and smaller values for the negative labels. [29] further shows that (9) is closely related to *ranking loss*.

Later, [136] finds that BP-MLL does not perform as expected on data sets in textual domain. To address the issue, based on BP-MLL, [136] proposes to use a comparably simple NN approach that can achieve the state-of-the-art performance in large-scale multi-label text classification. They show that the ranking loss in BP-MLL can be efficiently and effectively replaced by the commonly used cross-entropy function, and several NN tricks, i.e., rectified linear units (ReLUs), Dropout, and AdaGrad can be effectively employed in this setting.

Embedding methods have been effective to capture the label dependency and reduce the computation costs. However, existing embedding methods are shallow models, which may not be powerful to discover high order dependency among labels. To fulfill this gap, [30] proposes

Canonical Correlated AutoEncoder (C2AE), which is the first DNN-based embedding method for MLC to our knowledge. The basic idea of C2AE is to seek a deep latent space to jointly embed the instances and labels. C2AE performs feature-aware label embedding and label-correlation aware prediction. The former is realized by joint learning of deep canonical correlation analysis (DCCA) and the encoding stage of autoencoder, while the latter is achieved by the introduced loss function for the decoding outputs.

C2AE consists of two DNN modules, i.e., DCCA and autoencoder, and seeks three mapping functions, i.e., feature mapping F_x , encoding function F_e , and decoding function F_d . For training, C2AE receives instance X and labels Y , associates them in the latent space L , and enforces the recover of Y using autoencoder. The objective function of C2AE is defined as follows:

$$\min_{F_x, F_e, F_d} \Phi(F_x, F_e) + \alpha \Gamma(F_e, F_d) \quad (10)$$

where $\Phi(F_x, F_e)$ and $\Gamma(F_e, F_d)$ denote the losses in the latent and output spaces respectively, α is used to balance the two terms. Inspired by the CCA, C2AE learns the deep latent space by maximizing the correlation between instances and labels. Thus $\Phi(F_x, F_e)$ can be defined as:

$$\begin{aligned} \min_{F_x, F_e} \quad & \|F_x(X) - F_e(Y)\|_F^2 \\ \text{s.t.} \quad & F_x(X)F_x(X)^T = F_e(Y)F_e(Y)^T = I \end{aligned} \quad (11)$$

In addition, C2AE recovers the labels using autoencoder with aim of preserving label dependency. Inspired by [29], $\Gamma(F_e, F_d)$ is defined as follows:

$$\begin{aligned} \Gamma(F_e, F_d) = & \sum_{i=1}^N E_i \\ E_i = & \frac{1}{|y_i^1| |y_i^0|} \sum_{(p,q) \in y_i^1 \times y_i^0} \exp(-(F_d(F_e(x_i))^p - F_d(F_e(x_i))^q)) \end{aligned} \quad (12)$$

where N is the number of the instances, $F_d(F_e(x_i))$ is the recovered label of x_i using the autoencoder. For prediction, given a test instance \hat{x} , C2AE performs prediction as $\hat{y} = F_d(F_x(\hat{x}))$.

Later, inspired by C2AE, [137] presents a two-stage label embedding model based on neural factorization machine model. It first exploits second-order label correlation via a factorization layer and then learns high-order correlation by additional fully-connected layers. [57] proposes another deep embedding method, i.e., Deep Correlation Structure Preserved Label Space Embedding (DCSPE). In addition to DCCA, DCSPE further develops deep multidimensional scaling (DMDS) to preserve the intrinsic structure of the latent space. Finally, DCSPE transforms test instance into the latent space, searches its nearest neighbor, and finally regards label of this neighbor as prediction. However, as the k NN search is time-consuming, the k NN embedding methods are computationally expensive in the large-scale setting. To solve the above issue, [138] proposes a novel deep binary prototype compression (DBPC) for fast multi-label prediction. DBPC compresses the database into a small set of short binary prototypes, and uses the prototypes for prediction.

For multi-label emotion classification, [139] recently proposes latent emotion memory (LEM) to learn latent emotion distribution without external knowledge. LEM includes latent emotion and memory modules to learn emotion distribution and emotional features respectively, and the concatenation of the two is fed into Bi-directional Gated Recurrent Unit (BiGRU) for prediction. For multi-label image classification, [140] proposes a unified deep neural network that exploits both semantic and spatial relations between labels with only image-level supervision. Specifically, the authors propose Spatial Regularization Network (SRN) that generates attention maps for all labels and captures the underlying relations between them via learnable convolutions. [141] finds the consistency of attention regions of CNN classifiers under many transforms are not preserved. To address the issue, the authors propose a two-branch network with original and transformed images as inputs and introduce a new attention consistency loss that measures the attention heatmap consistency between two branches. Later [142] proposes Adjacency-based Similarity Graph Embedding (ASGE) and Cross-modality Attention (CMA) to capture the dependencies between labels and discover locations of discriminative features respectively. Specifically, ASGE learns semantic label embedding that can explicitly exploit label correlations, and CMA generates the meaningful attention maps by leveraging more prior semantic information. Instead of requiring laborious object-level annotations, [143] proposes to distill knowledge from weakly-supervised detection (WSD) task to boost MLC performance. The authors construct an end-to-end MLC framework augmented by a knowledge distillation module that guides the classification model by the WSD model for object RoIs. WSD and MLC are the teacher and student models respectively.

Remark. Deep embedding methods are the most widely-used deep methods for MLC. Among them, C2AE is pioneer deep embedding work for MLC and has been applied in many real-world applications including multi-label emotion classification, which deserves exploration for the beginners to understand basic mechanism. As we know, label correlation is key for MLC, and some objectives, e.g., (12) have been used to model label correlation in deep embedding methods for MLC. However, existing research shows that (12) may not be effective for textual domain. In the future, how to effectively capture label correlation is an important research topic of deep embedding methods for MLC. Some advanced techniques, e.g., graph convolutional network (GCN), recurrent neural network (RNN) open doors to better capture label correlation, and can motivate further research of deep embedding methods for MLC.

4.2 Deep Learning for Challenging MLC

In real world applications, multi-label learning is often challenging due to the complex setting of labels. For instance, the number of labels is very large known as XMLC; the labels are often partially or weakly given; labels emerge continuously or are even unseen before. This section reviews the recent advances of deep learning to address these challenging MLC problems.

DL for Extreme MLC. To our knowledge, [31] is the first attempt at applying deep learning to XMLC. XML-CNN [31] applies convolutional neural network (CNN)

and dynamic pooling to learn the text representation, and a hidden bottleneck layer much smaller than the output layer is used to achieve computational efficiency. However, XML-CNN still suffers from effectiveness of capturing the important subtext for each label. To address this issue, AttentionXML [32] is proposed with two unique features: 1) a multi-label attention mechanism with raw text as input, which allows to capture the most relevant part of text to each label, 2) a shallow and wide probabilistic label tree (PLT), which allows to handle millions of labels, especially for "tail labels". Meanwhile, based on C2AE, a new deep embedding method, i.e., Ranking-based Auto-Encoder (Rank-AE) [33] is proposed for XMLC. Rank-AE first uses an efficient attention mechanism to learn rich representations from any type of input features, learns a latent space for instance and labels, and finally develops a margin-based ranking loss that is more effective for XMLC and noisy labels. [144] empirically demonstrates that overfitting leads to the poor performance of the DNN based embedding methods for XMLC. Based on this finding, [144] further proposes a new regularizer, i.e., GLaS for embedding-based neural network approaches. [145] finetunes a pretrained deep transformer for better feature representation. They propose a novel label clustering model for XMLC and the transformer serves as a neural matcher. With the proposed techniques, the state-of-the-art performance is achieved on several widely-used extreme data sets. Very recently [146] develops DeepXML framework that can generate a family of algorithms by including four sub-tasks, i.e., intermediate representation, negative sampling, transfer learning, and classifier learning. It yields Accelerated Short Text Extreme Classifier (Astec) that is more accurate and faster than state-of-the-art deep-XMLs on public short text data sets. DECAF [147] considers label metadata, e.g., textual descriptions of labels, which is informative but usually ignored by existing methods. DECAF jointly learns model enriched by label metadata and feature representation, and predicts accurately with millions of labels. ECLARE [148] incorporates label text and label correlations, and develops frugal architecture and scalable techniques to train model with label correlation graph with millions of labels. Similarly, GalaXC [149] collaboratively learns over joint document-label graphs that can incorporate various sources, e.g., label metadata. GalaXC further introduces label-wise attention to obtain high-capacity extreme classifiers. [149] shows that GalaXC is up to 18% more accurate than state-of-the-arts while it trains 2-50 times faster and predicts 10 times faster on benchmark data sets.

DL for partial and weakly-supervised MLC. Several efforts [19], [34], [150], [151] have been made towards MLC with partial labels. [34] empirically shows that partially annotating all images is better than fully annotating a small subset. Thus [34] generalizes the standard binary cross-entropy loss by exploiting label proportion information, and develops an approach based on Graph Neural Networks (GNNs) to explicitly model the correlation between categories. Later, [19] regularizes the cross-entropy loss with a cost function that measures the smoothness of labels and features of images on data manifold, and develops an efficient interactive learning framework where similarity learning and CNN training interact and improve each another. [23] is the first deep generative model to tackle

weakly-supervised MLC (WS-MLC). [23] proposes a probabilistic framework that integrates sequential prediction and generation processes to exploit information from unlabeled or partially labeled data.

DL for MLC with unseen labels. In conventional MLC, all the labels are assumed to be fixed and static; however, it is ignored that labels emerge continuously in changing environments. To fulfill this gap, a novel DNN-based method, i.e., Deep Streaming Label Learning (DSLL) [35] is proposed to deal with MLC with newly emerged labels effectively. DSLL uses streaming label mapping, streaming feature distillation, and senior student network to explore the knowledge from past labels and historical models to understand new labels. In addition, [35] further theoretically proves that DSLL admits tight generalization error bounds for new labels in the DNN framework. Different from DSLL, [36] incorporates the additional knowledge graphs for multi-label zero-shot learning (ML-ZSL). [36] advances label propagation mechanism in the semantic space, enabling the reasoning of the learned model for predicting unseen labels.

Remark. MLC problem is challenging due to the high complexity of labels. [149], [34], [23], and [35] are representative deep works for beginners to address extreme MLC, partial MLC, weakly-supervised MLC, and MLC with unseen labels, respectively. The above attempts only focus on challenges of label space in MLC problem. However, in real-world MLC problems, there are some challenges in feature space, e.g., some features may be vanished or augmented, the distribution may change. How to simultaneously address challenges in label and feature spaces is more challenging and can be regarded as future research of challenging MLC problem.

4.3 Advanced Deep Learning for MLC

Recently some advanced deep learning architectures have been developed for MLC problems.

To exploit the underlying rich label structure, [37] proposes Deep In the Output Space (ADIOS) to partition the labels into a Markov Blanket Chain and then apply a novel deep architecture that exploits the partition. In multi-label image classification, CNN-RNN [152] utilizes recurrent neural networks (RNNs) to better exploit the higher-order label dependencies of an image. CNN-RNN learns a joint image-label embedding to characterize the semantic label dependency as well as the image-label relevance, and it can be trained end-to-end from scratch to integrate both information in a unified framework. In addition, instead of using classifier chain, [38] proposes to use RNN to convert MLC into a sequential prediction problem, where the labels are first ordered in an arbitrary fashion. The key advantage is to allow focusing on the prediction of only positive labels, a much smaller set than the full set of possible labels. [153] employs Long-Short Term Memory (LSTM) sub-network to sequentially predict semantic labeling scores on the located regions and capture the global dependencies of these regions, and achieve superior performance in large-scale multi-label image classification. [154] does not require pre-defined label orders. It integrates and learns visual attention and LSTM layers for multi-label image

classification. Instead of a fixed, static label ordering, [39] assumes a dynamic, context-dependent label ordering. [39] consists of a simple EM-like algorithm that bootstraps the learned model, and a more principled approach based on reinforcement learning. The experiments empirically show dynamic label ordering approach based on reinforcement learning outperforms RNN with fixed label ordering. [155] proposes a new framework based on optimal completion distillation and multitask learning that also does not require a predefined label order. Recently [156] proposes predicted label alignment (PLA) and minimal loss alignment (MLA) to dynamically order the ground truth labels with the predicted label sequence. This allows for faster training of more optimal LSTM models, and obtains state-of-the-art results in large-scale image classification.

Graph Convolutional Network (GCN) [80] has been also used to successfully model label dependency in MLC problem. In multi-label image classification problem, [157] first builds a directed graph over the object labels, employs GCN to model the correlations between labels, and maps label representation to inter-dependent object classifiers. Similarly, Semantic-Specific Graph Representation Learning (SSGRL) [158] includes semantic decoupling and interaction modules to learn and correlate semantic-specific representations respectively. The correlation is achieved by GCN on a graph built on label co-occurrence. Later, [159] adds lateral connections between GCN and CNN at shallow, middle and deep layers such that label information can be better injected into backbone CNN for label awareness. For multi-label patent classification, which is regarded as multi-label text classification problem, [160] proposes a new deep learning model based on GCN to capture rich semantic information. The authors design an adaptive non-local second-order attention layer to model long-range semantic dependencies in text content as label attention for patent categories.

As an alternative of DNN, deep forest [161] is a recent deep learning framework based on tree model ensembles, which does not rely on backpropagation. [40] introduces deep forest for MLC due to the advantages of deep forest models. The proposed Multi-Label Deep Forest (MLDF) can handle two challenging problems in MLC: optimizing different performance measures and reducing overfitting. The extensive experiments show that MLDF achieves the best performance over hamming loss, one-error, coverage, ranking loss, average precision and macro-AUC measures.

Remark. Advanced deep architecture has more powerful learning capability, and thus can be more effective for MLC problem. Beginners can try representative deep works of advanced RNN [152], GCN [157] [40] to address MLC problems. However, these advanced deep methods usually contain very large amounts of parameters, and require high complexity in terms of training and prediction costs. To devise lightweight architecture for efficient training and prediction is worthy to be explored for advanced deep methods for MLC.

5 ONLINE MULTI-LABEL LEARNING

Many real-world applications generate a massive volume of streaming data. For example, many web-related applications, such as Twitter and Facebook posts and RSS feeds,

are attached with multiple essential forms of categorization tags. In the search industry, revenue comes from clicks on ads embedded in the result pages. Ad selection and placement can be significantly improved if ads are tagged correctly. This scenario, referred to as online multi-label learning, is a popular tool for addressing large-scale multi-label classification tasks.

The current off-line MLC methods assume that all data are available in advance for learning. However, there are two major limitations of designing MLC methods under such an assumption: firstly, these methods are impractical for large-scale data sets, since they require all data sets to be stored in memory; secondly, it is non-trivial to adapt off-line multi-label methods to the sequential data. In practice, data is collected sequentially, and data that is collected earlier in this process may expire as time passes. Therefore, it is non-trivial to propose new online multi-label learning methods to deal with streaming data. This section presents a review of the latest algorithms on online multi-label classification.

[162] proposes an online universal classifier (OUC) to handle binary, multi-class and multi-label classification problems. To adapt all types of classification, OUC pre-processes the data set that the target label of all three classification types is represented as a vector with dimension equal to the number of output labels. A deep learning model is then employed for online training.

Based on ELM [163], which is a single hidden layer feedforward neural network model, [42] proposes the OSML-ELM approach to handle streaming multi-label data. OSML-ELM uses a sigmoid activation function and outputs weights to predict the labels. In each step, the output weight is learned from the specific equation. OSML-ELM converts the label set from single to multiple representation in order to solve multi-label classification problems.

OLANS GD [41] is proposed based on label ranking, where the ranking functions are learned by minimizing the ranking loss in the large margin framework. However, the memory and computational costs of this process are expensive on large-scale data sets. Stochastic gradient descent (SGD) approaches update the model parameters using only the gradient information calculated from a single label at each iteration. OLANS GD minimizes the primal form using Nesterov's smoothing, which has recently been extended to the stochastic setting.

However, none of these methods analyze the loss function, and do not use the correlations between labels and features. Some works have been developed to address this issue. For example, [164] presents a novel cost-sensitive dynamic principal projection (CS-DPP) method for online MLC. Inspired by matrix stochastic gradient, they develop an efficient online dimension reducer, and provide the theoretical guarantee for their carefully-designed online regression learner. Moreover, the cost information is embedded into label weights to achieve cost-sensitivity along with theoretical guarantees. However, CS-DPP can not capture the joint information between features and labels. To capture such joint information, [43] proposes a novel online metric learning paradigm for MLC. They first project features and labels into the same embedding space, and then the distance metric is learned by enforcing the constraint that the distance between embedded instance and its correct label

must be smaller than the distance between the embedded instance and other labels. Moreover, an efficient optimization algorithm is present for the online MLC. Theoretically, the upper bound of cumulative loss is analyzed in the paper. The experiment results show that their proposed algorithm outperforms the aforementioned baselines.

Recently, some works [165], [166] study online SS-MLC problem, where data examples can be unlabeled. [165] proposes a growing neural gas-based method, which constructs a dynamic graph with incoming data. OnSeML [166] adopts a label embedding fashion that a regression model is learned to fit the latent label vectors. To incorporate the unlabeled data, it extends the regularized moving least-square model [167] with a local smoothness regularizer. It is noteworthy that online learning from weakly-supervised data has long been a difficult issue since the global data structure is not given [168]. Hence, it is valuable to develop online MLC classifiers with limited supervision.

Remark. Online multi-label learning opens a new way to address large-scale MLC issues with limited memory. Unfortunately, the model, algorithm and theoretical results obtained so far are very limited. It is imperative to put more effort to explore this direction.

6 STATISTICAL MULTI-LABEL LEARNING

The generalization error of multi-label learning is analyzed by many papers. For example, [11] formulates MLC as the problem of learning a low-rank linear model in the standard ERM framework that could use a variety of loss functions and regularizations. They analyze the generalization error bounds for low-rank promoting trace norm regularization. There are also some statistical theoretical works which focus on the consistency of multi-label learning: whether the expected loss of a learned classifier converges to the Bayes loss as the training set size increases. For example, [169] studies two well-known multi-label loss functions: ranking loss and hamming loss. They provide a sufficient and necessary condition for the consistency of multi-label learning based on surrogate loss functions. For hamming loss, they propose a surrogate loss function which is consistent for the deterministic case. However, no convex surrogate loss is consistent with the ranking loss. [170] transforms MLC into the bipartite ranking problem, and proposes a simple univariate convex surrogate loss (exponential or logistic) defined on single labels, which is consistent with the ranking loss with explicit regret bounds and convergence rates. Recently, [10] shows that the pick-one-label can not achieve zero regret with respect to the Precision@ k , and PLTs model can get zero regret (i.e., it is consistent) in terms of marginal probability estimation and Precision@ k in the multi-label setting. Inspired by [10], [61] further studies the consistency of one-versus-all, pick-all-labels, normalised one-versus-all and normalised pick-all-labels reduction methods based on a different Recall@ k metric. All these works study the generalization error and consistency of learning approaches which address multi-label learning by decomposing into a set of binary classification problems. However, the existing theory of the generalization error and consistency does not consider label correlations, and desire for more effort to explore.

As mentioned above, several XMLC methods, such as PD-Sparse [59] and SLEEC [1], use ℓ_1 regularization to exploit the sparsity of the data. However, ℓ_1 regularization suffers two major limitations: 1) [171], [172], [173] show that the ℓ_1 regularization produces a bias into the resulting estimator, and harms the estimation accuracy. 2) [174] proves that the oracle property does not hold for ℓ_1 regularization. To address these issues, [175] proposes a unified framework for SLEEC with nonconvex penalty, such as minimax concave penalty (MCP) [173] and smoothly clipped absolute deviation (SCAD) penalty [174], which have recently attracted much attention because they can eliminate the estimation bias and attain attractive statistical properties. Theoretically, they show that their proposed estimator enjoys oracle property, which performs as well as if the underlying model were known beforehand, as well as attains a desirable statistical convergence rate of $\mathcal{O}\left(\frac{\sigma\sqrt{\varpi}+\sqrt{s^*}}{\mu\sqrt{n}}\right)$, where σ , ϖ , μ are positive constants, n is the sample size and s^* denotes the cardinality of the true support of underlying model. Considering the magnitude of the entries in the underlying model, they can achieve a refined convergence rate of $\mathcal{O}\left(\frac{\sqrt{s^*}}{\mu\sqrt{n}}\right)$ under suitable conditions. This paper could inspire the community to bring more powerful statistical penalty method and theory into MLC.

Remark. A key challenging issue in MLC is to model the interdependencies between labels and features. Existing methods, such as classifier chain, CCA and CPLST, attempt to model the correlations between labels and features. However, the statistical properties of these multi-label dependency modelings are less explored, and how to do theoretical analysis for them is an important research topic in the future. Copulas is an influential statistical tool for modeling dependence of multivariate data, and first brought into MLC for modeling label and feature dependencies [176]. In particular, [176] first constructs continuous distribution in the output space via employing the kernel trick, and then develops an unbiased and consistent estimator. Moreover, they also present the asymptotic analysis and mean squared error in the paper. However, the biggest problem for this paper is that it can not handle high dimension issues. The use of copula for modeling label and feature dependencies reveals new statistical insights in multi-label learning, and could orient more high dimension driven works in this direction.

7 NEW APPLICATIONS

During the past decade, multi-label classification has been successfully applied in various applications, such as protein function classification, music categorization and semantic scene classification. Recently, some new applications in computer vision (CV), natural language processing (NLP) and data mining (DM) are emerging, which are summarized in the Supplementary Materials. This section will briefly review some of them.

7.1 Computer Vision

7.1.1 Video Annotation

With the development of considerable videos on the Internet (e.g., Youtube, Flickr and Facebook), efficient and effective indexing and searching these video corpus becomes

more and more important for the research and industry community. In many real-world video corpus, the videos are multi-labeled. For instance, most of the videos in the popular TRECVID data set [177] are annotated by more than one label from a set of 39 different concepts. Currently, semantic-level video annotation (i.e., the semantic video concept detection) has been an important research topic in the multimedia research community, which aims to tag videos with a set of concepts of interest, including scenes (e.g., garden, sky, tree), objects (e.g., animals, people, airplane, car), events (e.g., election, ceremony) and certain named entities (e.g., university, person, home). [178] attempts to capture the correlations between different labels to improve the annotation performance on video concepts. [179] proposes a novel online multi-label learning method for large-scale video annotation.

7.1.2 Facial Action Unit Recognition

Thoughts and feelings are revealed in the face. The facial muscle movements tell a person's social behavior, psychopathology and internal states. Facial Action Unit (AU) Recognition plays an important role in describing comprehensive facial expressions, and has been successfully applied in mental state analysis. Some works [180] have provided the evidence that the occurrence of AUs are strongly correlated, and the sample distribution of AUs is unbalanced. Based on these properties, multi-label learning methods are well-matched to this learning scenario. For example, [181] introduces joint-patch and multi-label learning (JPML) to leverage group sparsity by selecting a sparse subset of facial patches while learning a multi-label classifier. [182] presents deep region and multi-label learning (DRML) for AU detection. Recently, [183] proposes a semi-supervised multi-label approach for AU recognition utilizing a large number of web face images without AU labels.

7.1.3 Neonatal Brains

Effective and consistent segmentation of brain white matter bundles at neonatal stage plays a vital role in detecting white matter abnormalities and understanding brain development for the prediction of psychiatric disorders. Because of the complexity of white matter anatomy and the spatial resolution of diffusion-weighted MR imaging, multiple fiber bundles can pass through one voxel. [184] aims to assign one or multiple anatomical labels of white matter bundles to each voxel to reflect complex white matter anatomy of the neonatal brain. To achieve this goal, [184] explores the supervised multi-label learning algorithm in Riemannian diffusion tensor spaces, which considers diffusion tensors lying on the Log-Euclidean Riemannian manifold of symmetric positive definite (SPD) matrices and their corresponding vector space as feature space. [184] demonstrates that they are able to automatically learn the number of white matter bundles at a location and provide anatomical annotation of the neonatal white matter. Recently, [185] and [186] present some weakly-supervised multi-label learning methods for neonatal brain extraction.

7.2 Natural Language Processing

7.2.1 Mobile Applications

Recently, the development of mobile applications has become one of the most important topics in communications [187]. Under this field, advanced high performance algorithms for mobile applications have attracted the attention of researchers. Recommendation systems are widely used to predict the “rating” or “preference” that a user would give to an item. A good recommendation system with high performance is able to attract users to the service for 5G applications. [188] focuses on high performance multi-label classification methods and their applications for medical recommendations in the domain of 5G communication. [189] develops a deep convolutional neural network for iris segmentation of noisy images acquired by mobile devices. A novel multi-label active learning method is proposed by [190] for mobile reviews classification tasks. Mobile applications involve language understandings, we group it to NLP.

7.2.2 Legal Text Mining

MLC has been widely used in the legal domain, especially for legal text mining tasks. In 2008, Mencia and Fürnkranz [191] collects a data set EUR-Lex, which comprises of documents about European Union law, including treaties, legislation, case-law and legislative proposals. The documents are categorized into several orthogonal concepts according to the European Vocabulary (EUROVOC), to allow for multiple search facilities. Recently, there arises new interest. In [192], a new legal MLC data set, dubbed EURLEX-57K is released. This is a large-scale version of EUR-Lex data set (19.6k documents, 4k EUROVOC labels) that contains 57k EU legislative documents from the EUR-Lex portal, each of which is labeled by 4.3k concepts from EUROVOC. In the Chinese AI and Law challenge [193], MLC is also applied to the legal judgment prediction (LJP) task, which aims to empower the machine to predict the judgment results of legal cases after reading fact descriptions. Since each criminal case can be relevant to multiple law articles, charges and prison terms, the LJP task can be regarded as a multi-label text classification problem. XMLC [192] and DL MLC [193] are proposed to address this task. Based on syntactic and grammatical features, legal text mining is categorized as NLP.

7.3 Data Mining

7.3.1 Recommender Systems

The recommender system can be naturally regarded as an MLC tasks since we usually recommend multiple items simultaneously to the users. For example, [194] develops an MLC model to automatically recommend bid phrases to an advertiser from a given ad landing page; [195] approaches the item-to-item recommendation task on Amazon, which aims at predicting the subset of items (labels) that a user might buy along with a given item. A recent work [145] regarded the keyword recommendation as an XMLC task, that provides keyword suggestions for advertisers to create campaigns. The MLC model receives the product-query customer purchase records and then suggests queries that are relevant to any given product by utilizing product

information, like title, description, brand, and so on. The applications of XMLC in recommendation have been widely studied in the literature [145], [194], [195].

7.3.2 User Profiling

In many applications, such as social media and e-commerce, it is essential to provide adaptive and personalized services to users. Therefore, user profiling, which infers user characteristics and personal interests from user-generated data, has been widely adopted by many online platforms. Some works regard this problem as a single-label learning task. However, obviously, more user characteristics lead to better personalization and the correlations between different user profiles can help improve the quality of user profiling. Hence, some works try to infer multiple attributes simultaneously. For example, Farnadi [196] proposes a hybrid deep learning framework to infer multiple types of user-profiles from multiple modalities of user data. Their experiments on 5K Facebook users also validates the superiority of the multi-label learning fashion to single-label learning. [197] explores the user profiles on Weibo, a famous social network platform in China, by using graph information in social networks. Another example is fraud detection in e-commerce platforms [94], since fraud users usually have different spam behaviors simultaneously. [94] presents a collaboration based multi-label propagation method to utilized the correlations among different fraud behaviors.

8 CONCLUSION

Multi-label classification has attracted significant attention from the community over the last decade. This paper provides a comprehensive review of the emerging topics of multi-Label learning, which include extreme multi-label classification, multi-label learning with limited supervision, deep learning for multi-label learning, online multi-label learning, statistical multi-label learning and new applications. We provide an overview of the representative works referenced throughout. In addition, we emphasize the challenges of these emerging topics and some future research directions and the promising extensions that are worthy of further study.

APPENDIX A

EVALUATION METRICS AND NOTATIONS AND NEW APPLICATIONS

Assume $x_i \in \mathbb{R}^{d \times 1}$ is a real vector representing an input or instance (feature), $y_i = (y_{i,1}, \dots, y_{i,L}) \in \{0, 1\}^{L \times 1}$ is the corresponding output or label vector ($i \in \{1, \dots, n\}$). n , d and L denote the number of training data, feature dimensions and the number of labels, respectively. The input matrix is $X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ and the output matrix is $Y = [y_1, \dots, y_n] \in \{0, 1\}^{L \times n}$. MLC aims to learn a classifier which predicts the testing instance as accurate as possible with the set of proper labels. Let $\check{Y} = [\check{y}_1, \dots, \check{y}_n] \in \{0, 1\}^{L \times n}$ be the predicted label. We first introduce some evaluation metrics for MLC.

Hamming loss. Hamming loss is defined as follows:

$$1/n \sum_{i=1}^n |\Upsilon(y_i) \Delta \Upsilon(\check{y}_i)|/L$$

TABLE 2
Important notations used in the main paper.

Notations	Explanations
x_i, y_i	Input and output vectors
X, Y	Input and output matrices
\tilde{D}	Transformed data set
\hat{Y}, \tilde{Y}	Label matrices of implicit and explicit missing labels
$Z = [z_1, \dots, z_n]$	Embedding matrix and vectors
\hat{y}	Predicted score vector
\tilde{Y}, \tilde{y}	Predicted logical label matrix and vector
$\Upsilon(y_i)$	The indices of the positive labels of y_i
$\Lambda = [\lambda_{ij}]$	Enriched real-value label representation
S	The candidate label set in PML
Ω	The index set of neighbors
N_i	The index set of neighbors of the i -th instance
D_l, D_o, D_u	The index sets of labeled, incompletely-labeled and unlabeled data
Δ	The symmetric difference between two sets
$ \cdot $	The set cardinality
$\langle \cdot, \cdot \rangle$	Inner product
$O(\cdot)$	Computational complexity
T	Matrix Transpose
$\sigma(\cdot)$	Sigmoid function
$\text{nnz}(\cdot)$	The number of non-zero entries
$\ \cdot\ _F, \ \cdot\ _2, \ \cdot\ _1$	Frobenius norm, ℓ_2 and ℓ_1 norm of a matrix (vector)
$\text{Tr}(\cdot)$	Trace operator
$r(\cdot)$	The regularizer function
$\mathcal{L}(\cdot)$	Empirical risk function
n	The number of training data
d, L	Feature dimensions and the number of labels
ϖ	Dimension of embedding vectors
\mathbb{R}	Set of real numbers
s^*	The cardinality of the true support of the underlying model
α, μ, λ, C	Trade-off hyperparameters
I	Identity Matrix
A, B	Side information matrices w.r.t input and output
W, U, V, H	Projection or similarity Matrix
F_e, F_x, F_d	Label encoding, feature encoding and decoding network of C2AE
$\mathcal{W} = [w_{ij}]_{n \times n}$	Graph weight matrix
L_x, L_y	The laplacian matrix of \mathcal{W}^x and \mathcal{W}^y
$\Phi(F_x, F_e), \Gamma(F_e, F_d)$	The losses of C2AE in the latent and output space

where $\Upsilon(y_i)$ denotes the indices of the positive labels of y_i , Δ stands for the symmetric difference between two sets, $|\cdot|$ means the cardinality. The hamming loss evaluates the fraction of misclassified instance-label pairs.

Ranking loss. Let f be the real-valued function. Ranking loss is defined as follows:

$$1/n \sum_{i=1}^n \frac{|\{(a, b) : f(x_i, a) \leq f(x_i, b), (a, b) \in \Upsilon(y_i) \times \bar{y}_i\}|}{|\Upsilon(y_i)| |\bar{y}_i|}$$

where \bar{y}_i is the complementary set of $\Upsilon(y_i)$ in the label space. The ranking loss evaluates the fraction of reversely ordered label pairs.

F-measure.

$$\text{F-measure} = 1/L \sum_{j=1}^L \frac{2 \sum_{i=1}^n y_{i,j} \check{y}_{i,j}}{\sum_{i=1}^n y_{i,j} + \sum_{i=1}^n \check{y}_{i,j}}$$

F-measure computes true positives, true negatives, false positives and false negatives over labels, and then calculates an overall F-1 score.

Precision@k.

$$\text{Precision@}k = 1/k \sum_{k \in \text{rank}_k(\hat{y})} y_k$$

where $\hat{y} \in \mathbb{R}^{L \times 1}$ is a predicted score vector, y is a ground truth label vector and $\text{rank}_k(\hat{y})$ returns the k largest indices of \hat{y} ranked in descending order.

Recall@k.

$$\text{Recall@}k = 1/|\Upsilon(y)| \sum_{k \in \text{rank}_k(\hat{y})} y_k$$

Precision@ k and Recall@ k evaluate top- k precision and recall over labels respectively, and both of them are the standard measures for XMLC. F-measure and Ranking loss are usually used in recommender system. Some CV and NLP applications, such as facial action unit recognition and web page categorization, usually use the Hamming loss

and Ranking loss as the performance metric. The important notations and new applications in the main paper are summarized in Tables 2 and 3, respectively.

REFERENCES

- [1] K. Bhatia, H. Jain, P. Kar, M. Varma, and P. Jain, "Sparse local embeddings for extreme multi-label classification," in *NeurIPS*, 2015, pp. 730–738.
- [2] R. Babbar and B. Schölkopf, "Dismec: Distributed sparse machines for extreme multi-label classification," in *WSDM*, 2017, pp. 721–729.
- [3] I. E. Yen, X. Huang, W. Dai, P. Ravikumar, I. S. Dhillon, and E. P. Xing, "Pdpdpsparse: A parallel primal-dual sparse method for extreme classification," in *KDD*, 2017, pp. 545–553.
- [4] Y. Prabhu, A. Kag, S. Harsola, R. Agrawal, and M. Varma, "Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising," in *WWW*, 2018, pp. 993–1002.
- [5] H. Jain, V. Balasubramanian, B. Chunduri, and M. Varma, "Slice: Scalable linear extreme classifiers trained on 100 million labels for related searches," in *WSDM*, 2019, pp. 528–536.
- [6] Y. Prabhu and M. Varma, "Fastxml: a fast, accurate and stable tree-classifier for extreme multi-label learning," in *KDD*, 2014, pp. 263–272.
- [7] H. Jain, Y. Prabhu, and M. Varma, "Extreme multi-label loss functions for recommendation, tagging, ranking & other missing label applications," in *KDD*, 2016, pp. 935–944.
- [8] K. Jasinska, K. Dembczynski, R. Busa-Fekete, K. Pfannschmidt, T. Klerx, and E. Hüllermeier, "Extreme f-measure maximization using sparse probability estimates," in *ICML*, 2016, pp. 1435–1444.
- [9] Y. Prabhu, A. Kag, S. Gopinath, K. Dahiya, S. Harsola, R. Agrawal, and M. Varma, "Extreme multi-label learning with label features for warm-start tagging, ranking & recommendation," in *WSDM*, 2018, pp. 441–449.
- [10] M. Wydmuch, K. Jasinska, M. Kuznetsov, R. Busa-Fekete, and K. Dembczynski, "A no-regret generalization of hierarchical softmax to extreme multi-label classification," in *NeurIPS*, 2018, pp. 6358–6368.
- [11] H. Yu, P. Jain, P. Kar, and I. S. Dhillon, "Large-scale multi-label learning with missing labels," in *ICML*, 2014, pp. 593–601.
- [12] Y. Tagami, "Annexml: Approximate nearest neighbor search for extreme multi-label classification," in *KDD*, 2017, pp. 455–464.
- [13] W. Liu, D. Xu, I. W. Tsang, and W. Zhang, "Metric learning for multi-output tasks," *TPAMI*, vol. 41, no. 2, pp. 408–422, 2019.
- [14] X. Gong, D. Yuan, and W. Bao, "Fast multi-label learning," in *IJCAI*, 2021, pp. 2432–2438.
- [15] Y. Sun, Y. Zhang, and Z. Zhou, "Multi-label learning with weak label," in *AAAI*, 2010.
- [16] G. Chen, Y. Song, F. Wang, and C. Zhang, "Semi-supervised multi-label learning by solving a sylvester equation," in *SDM*, 2008, pp. 410–419.
- [17] M. Xie and S. Huang, "Partial multi-label learning," in *AAAI*, 2018, pp. 4302–4309.
- [18] H. Wang, W. Liu, Y. Zhao, C. Zhang, T. Hu, and G. Chen, "Discriminative and correlative partial multi-label learning," in *IJCAI*, 2019, pp. 3691–3697.
- [19] D. Huynh and E. Elhamifar, "Interactive multi-label cnn learning with partial labels," in *CVPR*, 2020, pp. 9423–9432.
- [20] C. Xu, D. Tao, and C. Xu, "Robust extreme multi-label learning," in *KDD*, 2016, pp. 1275–1284.
- [21] L. Sun, S. Feng, T. Wang, C. Lang, and Y. Jin, "Partial multi-label learning by low-rank and sparse decomposition," in *AAAI*, 2019, pp. 5016–5023.
- [22] V. Jain, N. Modhe, and P. Rai, "Scalable generative models for multi-label learning with missing labels," in *ICML*, 2017, pp. 1636–1644.
- [23] H. Chu, C. Yeh, and Y. F. Wang, "Deep generative models for weakly-supervised multi-label classification," in *ECCV*, 2018, pp. 409–425.
- [24] M. Hu, H. Han, S. Shan, and X. Chen, "Weakly supervised image classification through noise regularization," in *CVPR*, 2019, pp. 11 517–11 525.
- [25] Z. Ji, B. Cui, H. Li, Y.-G. Jiang, T. Xiang, T. Hospedales, and Y. Fu, "Deep ranking for image zero-shot multi-label classification," *TIP*, 2020.

TABLE 3
The new applications of multi-label learning.

Reference	New Applications	Approaches	Evaluation Metrics
[178]	CV: automatic video annotation	XMLC [198], online MLC [179]	Precision@k, Recall@k and Hamming loss
[133]	CV: action recognition and localization in videos	multi-instance MLC [133]	Hamming loss
[183]	CV: facial action unit recognition	DL MLC [182], semi-supervised MLC [183]	Hamming loss and Ranking loss
[199]	CV: visual object recognition	online MLC [199], [200]	Hamming loss
[119]	CV: visual mobile robot navigation	multi-instance MLC [119]	Hamming loss
[201]	CV: biomedical image segmentation	semi-supervised MLC [201]	Hamming loss
[184]	CV: neonatal brains	semi-supervised MLC [185], [186]	Hamming loss
[188]	NLP: 5G mobile medical recommendations	DL MLC [189], MLAL [190]	Hamming loss and Ranking loss
[202]	NLP: social network analysis	DL MLC [81]	Hamming loss and Ranking loss
[203]	NLP: high-speed streaming data	online MLC [203]	Hamming loss
[204]	NLP: web page categorization	DL MLC [204]	Hamming loss and Ranking loss
[205]	NLP: protein subcellular localization	XMLC [205]	Precision@k and Recall@k
[205]	NLP: legal text mining	XMLC [192], DL MLC [193]	Precision@k, Recall@k and Hamming loss
[196]	DM: recommender system	XMLC [145], [194], [195]	Precision@k, Recall@k, F-measure and Ranking loss
[196]	DM: user profiling in social media	DL MLC [196], semi-supervised MLC [206]	Hamming loss and Ranking loss
[94]	DM: e-commercial fraud user detection	semi-supervised MLC [94]	Hamming loss and Ranking loss

- [26] W. Shi and Q. Yu, "Fast direct search in an optimally compressed continuous target space for efficient multi-label active learning," in *ICML*, 2019, pp. 5769–5778.
- [27] A. Kuznetsova, H. Rom, N. Alldrin, J. R. R. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, T. Duerig, and V. Ferrari, "The open images dataset V4: unified image classification, object detection, and visual relationship detection at scale," *CoRR*, vol. abs/1811.00982, 2018.
- [28] B. Wu, W. Chen, Y. Fan, Y. Zhang, J. Hou, J. Liu, and T. Zhang, "Tencent ml-images: A large-scale multi-label image database for visual representation learning," *IEEE Access*, vol. 7, pp. 172 683–172 693, 2019.
- [29] M.-L. Zhang and Z.-H. Zhou, "Multilabel neural networks with applications to functional genomics and text categorization," *TKDE*, vol. 18, no. 10, pp. 1338–1351, 2006.
- [30] C. Yeh, W. Wu, W. Ko, and Y. F. Wang, "Learning deep latent space for multi-label classification," in *AAAI*, 2017, pp. 2838–2844.
- [31] J. Liu, W. Chang, Y. Wu, and Y. Yang, "Deep learning for extreme multi-label text classification," in *SIGIR*, 2017, pp. 115–124.
- [32] R. You, Z. Zhang, Z. Wang, S. Dai, H. Mamitsuka, and S. Zhu, "Attentionxml: Label tree-based attention-aware deep model for high-performance extreme multi-label text classification," in *NeurIPS*, 2019, pp. 5812–5822.
- [33] B. Wang, L. Chen, W. Sun, K. Qin, K. Li, and H. Zhou, "Ranking-based autoencoder for extreme multi-label classification," in *NAACL-HLT*, 2019, pp. 2820–2830.
- [34] T. Durand, N. Mehrasa, and G. Mori, "Learning a deep convnet for multi-label classification with partial labels," in *CVPR*, 2019, pp. 647–657.
- [35] Z. Wang, L. Liu, and D. Tao, "Deep streaming label learning," in *ICML*, 2020.
- [36] C. Lee, W. Fang, C. Yeh, and Y. F. Wang, "Multi-label zero-shot learning with structured knowledge graphs," in *CVPR*, 2018, pp. 1576–1585.
- [37] M. Cissé, M. Al-Shedivat, and S. Bengio, "ADIOS: architectures deep in output space," in *ICML*, 2016, pp. 2770–2779.
- [38] J. Nam, E. L. Mencia, H. J. Kim, and J. Fürnkranz, "Maximizing subset accuracy with recurrent neural networks in multi-label classification," in *NeurIPS*, 2017, pp. 5413–5423.
- [39] J. Nam, Y. Kim, E. L. Mencia, S. Park, R. Sarikaya, and J. Fürnkranz, "Learning context-dependent label permutations for multi-label classification," in *ICML*, 2019, pp. 4733–4742.
- [40] L. Yang, X. Wu, Y. Jiang, and Z. Zhou, "Multi-label learning with deep forest," *CoRR*, vol. abs/1911.06557, 2019.
- [41] S. Park and S. Choi, "Online multi-label learning with accelerated nonsmooth stochastic gradient descent," in *ICASSP*, 2013, pp. 3322–3326.
- [42] R. Venkatesan, M. J. Er, M. Dave, M. Pratama, and S. Wu, "A novel online multi-label classifier for high-speed streaming data applications," *Evolving Systems*, vol. 8, no. 4, pp. 303–315, 2017.
- [43] X. Gong, D. Yuan, and W. Bao, "Online metric learning for multi-label classification," in *AAAI*, 2020, pp. 4012–4019.
- [44] M.-L. Zhang and Z.-H. Zhou, "ML-KNN: A lazy learning approach to multi-label learning," *PR*, vol. 40, no. 7, pp. 2038–2048, 2007.
- [45] W. Liu, I. W. Tsang, and K.-R. Müller, "An easy-to-hard learning paradigm for multiple classes and multiple labels," *JMLR*, vol. 18, no. 94, pp. 1–38, 2017.
- [46] J. Read, B. Pfahringer, G. Holmes, and E. Frank, "Classifier chains for multi-label classification," in *ECML/PKDD*, 2009, pp. 254–269.
- [47] Y. Zhang and J. G. Schneider, "Multi-label output codes using canonical correlation analysis," in *AISTATS*, 2011, pp. 873–882.
- [48] Y.-N. Chen and H.-T. Lin, "Feature-aware label space dimension reduction for multi-label classification," in *NeurIPS*, 2012, pp. 1538–1546.
- [49] D. Hsu, S. Kakade, J. Langford, and T. Zhang, "Multi-label prediction via compressed sensing," in *NeurIPS*, 2009, pp. 772–780.
- [50] M. Cissé, N. Usunier, T. Artières, and P. Gallinari, "Robust bloom filters for large multilabel classification tasks," in *NeurIPS*, 2013, pp. 1851–1859.
- [51] A. Jalan and P. Kar, "Accelerating extreme classification via adaptive feature agglomeration," in *IJCAI*, 2019, pp. 2600–2606.
- [52] V. Gupta, R. Wadbude, N. Natarajan, H. Karnick, P. Jain, and P. Rai, "Distributional semantics meets multi-label learning," in *AAAI*, 2019, pp. 3747–3754.
- [53] B. H. Bloom, "Space/time trade-offs in hash coding with allowable errors," *Communications of the ACM*, vol. 13, no. 7, pp. 422–426, 1970.
- [54] S. Ubaru and A. Mazumdar, "Multilabel classification with group testing and codes," in *ICML*, 2017, pp. 3492–3501.
- [55] W. Liu and I. W. Tsang, "Making decision trees feasible in ultrahigh feature and label dimensions," *JMLR*, vol. 18, no. 81, pp. 1–36, 2017.
- [56] X. Shen, W. Liu, I. W. Tsang, Q. Sun, and Y. Ong, "Multilabel prediction via cross-view search," *TNNLS*, vol. 29, no. 9, pp. 4324–4338, 2018.
- [57] K. Wang, M. Yang, W. Yang, and Y. Yin, "Deep correlation structure preserved label space embedding for multi-label classification," in *ACML*, 2018, pp. 1–16.
- [58] S. Si, H. Zhang, S. S. Keerthi, D. Mahajan, I. S. Dhillon, and C. Hsieh, "Gradient boosted decision trees for high dimensional sparse output," in *ICML*, 2017, pp. 3182–3190.
- [59] I. E. Yen, X. Huang, P. Ravikumar, K. Zhong, and I. S. Dhillon, "Pd-sparse: A primal and dual sparse approach to extreme multiclass and multilabel classification," in *ICML*, 2016, pp. 3069–3077.
- [60] A. Niculescu-Mizil and E. Abbasnejad, "Label filters for large scale multilabel classification," in *AISTATS*, 2017, pp. 1448–1457.
- [61] A. K. Menon, A. S. Rawat, S. J. Reddi, and S. Kumar, "Multilabel reductions: what is my loss optimising?" in *NeurIPS*, 2019, pp. 10 599–10 610.
- [62] W. Sibli, F. Meyer, and P. Kuntz, "Craftml, an efficient clustering-based random forest for extreme multi-label learning," in *ICML*, 2018, pp. 4671–4680.
- [63] R. Babbar and B. Schölkopf, "Data scarcity, robustness and extreme multi-label classification," *Machine Learning*, vol. 108, no. 8–9, pp. 1329–1351, 2019.
- [64] S. Khandagale, H. Xiao, and R. Babbar, "Bonsai - diverse and shallow trees for extreme multi-label classification," *CoRR*, vol. abs/1904.08249, 2019.
- [65] B. Wu, Z. Liu, S. Wang, B. Hu, and Q. Ji, "Multi-label learning with missing labels," in *ICPR*, 2014, pp. 1964–1968.
- [66] M. Xu, R. Jin, and Z. Zhou, "Speedup matrix completion with side information: Application to multi-label learning," in *NeurIPS*, 2013, pp. 2301–2309.

- [67] Y. Han, G. Sun, Y. Shen, and X. Zhang, "Multi-label learning with highly incomplete data via collaborative embedding," in *KDD*, 2018, pp. 1494–1503.
- [68] H. Yang, J. T. Zhou, and J. Cai, "Improving multi-label learning with missing labels by structured semantic correlations," in *ECCV*, 2016, pp. 835–851.
- [69] H. Yu, H. Huang, I. S. Dhillon, and C. Lin, "A unified algorithm for one-class structured matrix factorization with side information," in *AAAI*, 2017, pp. 2845–2851.
- [70] B. Wu, F. Jia, W. Liu, B. Ghanem, and S. Lyu, "Multi-label learning with missing labels using mixed dependency graphs," *IJCV*, vol. 126, no. 8, pp. 875–896, 2018.
- [71] M. Xu, G. Niu, B. Han, I. W. Tsang, Z. Zhou, and M. Sugiyama, "Matrix co-completion for multi-label classification with missing features and labels," *CoRR*, vol. abs/1805.09156, 2018.
- [72] L. Xu, Z. Wang, Z. Shen, Y. Wang, and E. Chen, "Learning low-rank label correlations for multi-label classification with missing labels," in *ICDM*, 2014, pp. 1067–1072.
- [73] J. Ma, Z. Tian, H. Zhang, and T. W. S. Chow, "Multi-label low-dimensional embedding with missing labels," *KBS*, vol. 137, pp. 65–82, 2017.
- [74] K. Wang, "Robust embedding framework with dynamic hypergraph fusion for multi-label classification," in *ICME*, 2019, pp. 982–987.
- [75] B. Wu, S. Lyu, B. Hu, and Q. Ji, "Multi-label learning with missing labels for image annotation and facial action unit recognition," *PR*, vol. 48, no. 7, pp. 2279–2289, 2015.
- [76] Y. Liu, K. Wen, Q. Gao, X. Gao, and F. Nie, "SVM based multi-label learning with missing labels for image annotation," *PR*, vol. 78, pp. 307–317, 2018.
- [77] B. Wu, S. Lyu, and B. Ghanem, "Constrained submodular minimization for missing labels and class imbalance in multi-label learning," in *AAAI*, 2016, pp. 2229–2236.
- [78] J. Huang, F. Qin, X. Zheng, Z. Cheng, Z. Yuan, and W. Zhang, "Learning label-specific features for multi-label classification with missing labels," in *Fourth IEEE International Conference on Multimedia Big Data*, 2018, pp. 1–5.
- [79] Y. Zhu, J. T. Kwok, and Z. Zhou, "Multi-label learning with global and local label correlation," *TKDE*, vol. 30, no. 6, pp. 1081–1094, 2018.
- [80] J. Zhou, G. Cui, Z. Zhang, C. Yang, Z. Liu, and M. Sun, "Graph neural networks: A review of methods and applications," *CoRR*, vol. abs/1812.08434, 2018.
- [81] H. Dong, W. Wang, K. Huang, and F. Coenen, "Joint multi-label attention networks for social text annotation," in *NAACL-HLT*, 2019, pp. 1348–1354.
- [82] M. Chen, A. X. Zheng, and K. Q. Weinberger, "Fast image tagging," in *ICML*, 2013, pp. 1274–1282.
- [83] Q. Wang, B. Shen, S. Wang, L. Li, and L. Si, "Binary codes embedding for fast image tagging with incomplete labels," in *ECCV*, 2014, pp. 425–439.
- [84] Z. Qi, M. Yang, Z. M. Zhang, and Z. Zhang, "Mining partially annotated images," in *KDD*, 2011, pp. 1199–1207.
- [85] X. Li, F. Zhao, and Y. Guo, "Conditional restricted boltzmann machines for multi-label learning with incomplete labels," in *AISTATS*, 2015.
- [86] H. Zhao, P. Rai, L. Du, and W. L. Buntine, "Bayesian multi-label learning with sparse features and labels, and label co-occurrences," in *AISTATS*, 2018, pp. 1943–1951.
- [87] M. Zhou, L. Hannah, D. B. Dunson, and L. Carin, "Beta-negative binomial process and poisson factor analysis," in *AISTATS*, 2012, pp. 1462–1471.
- [88] S. S. Bucak, R. Jin, and A. K. Jain, "Multi-label learning with incomplete class assignments," in *CVPR*, 2011, pp. 2801–2808.
- [89] K. M. Ibrahim, E. V. Epure, G. Peeters, and G. Richard, "Confidence-based weighted loss for multi-label classification with missing labels," in *ICMR*, 2020, pp. 291–295.
- [90] R. Sen, A. Rakhlin, L. Ying, R. Kidambi, D. P. Foster, D. N. Hill, and I. S. Dhillon, "Top-k extreme contextual bandits with arm hierarchy," *CoRR*, vol. abs/2102.07800, 2021.
- [91] M. E. Khan, D. Nielsen, V. Tangkaratt, W. Lin, Y. Gal, and A. Srivastava, "Fast and scalable bayesian deep learning by weight-perturbation in adam," in *ICML*, 2018, pp. 2616–2625.
- [92] Y. Liu, R. Jin, and L. Yang, "Semi-supervised multi-label learning by constrained non-negative matrix factorization," in *AAAI*, 2006, pp. 421–426.
- [93] X. Kong, M. K. Ng, and Z. Zhou, "Transductive multilabel learning via label set propagation," *TKDE*, vol. 25, no. 3, pp. 704–719, 2013.
- [94] H. Wang, Z. Li, J. Huang, P. Hui, W. Liu, T. Hu, and G. Chen, "Collaboration based multi-label propagation for fraud detection," in *IJCAI*, 2020.
- [95] L. Sun, S. Feng, G. Lyu, and C. Lang, "Robust semi-supervised multi-label learning by triple low-rank regularization," in *PAKDD*, 2019, pp. 269–280.
- [96] L. Jing, L. Yang, J. Yu, and M. K. Ng, "Semi-supervised low-rank mapping learning for multi-label classification," in *CVPR*, 2015, pp. 1483–1491.
- [97] C. Gong, D. Tao, J. Yang, and W. Liu, "Teaching-to-learn and learning-to-teach for multi-label propagation," in *AAAI*, 2016, pp. 1610–1616.
- [98] L. Feng, B. An, and S. He, "Collaboration based multi-label learning," in *AAAI*, 2019, pp. 3550–3557.
- [99] W. Zhan and M. Zhang, "Inductive semi-supervised multi-label learning with co-training," in *KDD*, 2017, pp. 1305–1314.
- [100] Z. Chu, P. Li, and X. Hu, "Co-training based on semi-supervised ensemble classification approach for multi-label data stream," in *ICBK*, 2019, pp. 58–65.
- [101] L. Wang, Y. Liu, C. Qin, G. Sun, and Y. Fu, "Dual relation semi-supervised multi-label learning," in *AAAI*, 2020, pp. 6227–6234.
- [102] Q. Wang, L. Yang, and Y. Li, "Learning from weak-label data: A deep forest expedition," in *AAAI*, 2020, pp. 6251–6258.
- [103] F. Wu, Z. Wang, Z. Zhang, Y. Yang, J. Luo, W. Zhu, and Y. Zhuang, "Weakly semi-supervised deep learning for multi-label image annotation," *IEEE Transactions on Big Data*, vol. 1, no. 3, pp. 109–122, 2015.
- [104] Q. Tan, Y. Yu, G. Yu, and J. Wang, "Semi-supervised multi-label classification using incomplete label information," *Neurocomputing*, vol. 260, pp. 192–202, 2017.
- [105] H. Dong, Y. Li, and Z. Zhou, "Learning from semi-supervised weak-label data," in *AAAI*, 2018, pp. 2926–2933.
- [106] J. Lv, N. Xu, R. Zheng, and X. Geng, "Weakly supervised multi-label learning via label enhancement," in *IJCAI*, 2019, pp. 3101–3107.
- [107] X. Geng, "Label distribution learning," *TKDE*, vol. 28, no. 7, pp. 1734–1748, 2016.
- [108] R. Shao, N. Xu, and X. Geng, "Multi-label learning with label enhancement," in *ICDM*, 2018, pp. 437–446.
- [109] N. Xu, Y. Liu, and X. Geng, "Partial multi-label learning with label distribution," in *AAAI*, 2020, pp. 6510–6517.
- [110] A. Akbarnejad and M. S. Baghshah, "An efficient semi-supervised multi-label classifier capable of handling missing labels," *TKDE*, vol. 31, no. 2, pp. 229–242, 2019.
- [111] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *NeurIPS*, 2017, pp. 1195–1204.
- [112] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, vol. 119, 2020, pp. 1597–1607.
- [113] M. Zhang and J. Fang, "Partial multi-label learning via credible label elicitation," *TPAMI*, 2020.
- [114] D. Xu, Y. Shi, I. W. Tsang, Y. Ong, C. Gong, and X. Shen, "A survey on multi-output learning," *CoRR*, vol. abs/1901.00248, 2019.
- [115] S. He, K. Deng, L. Li, S. Shu, and L. Liu, "Discriminatively relabel for partial multi-label learning," in *ICDM*, 2019, pp. 280–288.
- [116] G. Yu, X. Chen, C. Domeniconi, J. Wang, Z. Li, Z. Zhang, and X. Wu, "Feature-induced partial multi-label learning," in *ICDM*, 2018, pp. 1398–1403.
- [117] M. Xie and S. Huang, "Partial multi-label learning with noisy label identification," in *AAAI*, 2020, pp. 6454–6461.
- [118] Z. Li, G. Lyu, and S. Feng, "Partial multi-label learning via multi-subspace representation," in *IJCAI*, 2020, pp. 2612–2618.
- [119] J. He, H. Gu, and Z. Wang, "Multi-instance multi-label learning based on gaussian process with application to visual mobile robot navigation," *Information Sciences*, vol. 190, pp. 162–177, 2012.
- [120] C. Zhang, Z. Yu, H. Fu, P. Zhu, L. Chen, and Q. Hu, "Hybrid noise-oriented multilabel learning," *TCYB*, vol. 50, no. 6, pp. 2837–2850, 2020.
- [121] Z. Cui, Y. Zhang, and Q. Ji, "Label error correction and generation through label relationships," in *AAAI*, 2020, pp. 3693–3700.

- [122] M. Hu, H. Han, S. Shan, and X. Chen, "Multi-label learning from noisy labels with non-linear feature transformation," in *ACCV*, 2018, pp. 404–419.
- [123] A. Veit, N. Alldrin, G. Chechik, I. Krasin, A. Gupta, and S. J. Belongie, "Learning from noisy large-scale datasets with minimal supervision," in *CVPR*, 2017, pp. 6575–6583.
- [124] Y. Zhu, K. M. Ting, and Z. Zhou, "Multi-label learning with emerging new labels," *TKDE*, vol. 30, no. 10, pp. 1901–1914, 2018.
- [125] Y. Zhang, B. Gong, and M. Shah, "Fast zero-shot image tagging," in *CVPR*, 2016, pp. 5985–5994.
- [126] A. Alfassy, L. Karlinsky, A. Aides, J. Shtok, S. Harary, R. S. Feris, R. Giryes, and A. M. Bronstein, "Laso: Label-set operations networks for multi-label few-shot learning," in *CVPR*, 2019, pp. 6548–6557.
- [127] B. Yang, J. Sun, T. Wang, and Z. Chen, "Effective multi-label active learning for text classification," in *KDD*, 2009, pp. 917–926.
- [128] S. Li, Y. Jiang, N. V. Chawla, and Z. Zhou, "Multi-label learning from crowds," *TKDE*, vol. 31, no. 7, pp. 1369–1382, 2019.
- [129] S. Huang, S. Chen, and Z. Zhou, "Multi-label active learning: Query type matters," in *IJCAI*, 2015, pp. 946–952.
- [130] N. Xu, J. Shu, Y. Liu, and X. Geng, "Variational label enhancement," in *ICML*, vol. 119, 2020, pp. 10597–10606.
- [131] N. Xu, Y. Liu, and X. Geng, "Label enhancement for label distribution learning," *TKDE*, vol. 33, no. 4, pp. 1632–1643, 2021.
- [132] H. Yang, J. T. Zhou, J. Cai, and Y. Ong, "MIML-FCN+: multi-instance multi-label learning via fully convolutional networks with privileged information," in *CVPR*, 2017, pp. 5996–6004.
- [133] X. Zhang, H. Shi, C. Li, and P. Li, "Multi-instance multi-label action recognition and localization based on spatio-temporal pre-trimming for untrimmed videos," in *AAAI*, 2020, pp. 12886–12893.
- [134] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in *ICLR*, 2017.
- [135] M. Wang and W. Deng, "Deep visual domain adaptation: A survey," *Neurocomputing*, vol. 312, pp. 135–153, 2018.
- [136] J. Nam, J. Kim, E. L. Mencia, I. Gurevych, and J. Fürnkranz, "Large-scale multi-label text classification - revisiting neural networks," in *ECML-PKDD*, 2014, pp. 437–452.
- [137] C. Chen, H. Wang, W. Liu, X. Zhao, T. Hu, and G. Chen, "Two-stage label embedding via neural factorization machine for multi-label classification," in *AAAI*, 2019, pp. 3304–3311.
- [138] X. Shen, W. Liu, Y. Luo, Y. Ong, and I. W. Tsang, "Deep discrete prototype multilabel learning," in *IJCAI*, 2018, pp. 2675–2681.
- [139] H. Fei, Y. Zhang, Y. Ren, and D. Ji, "Latent emotion memory for multi-label emotion classification," in *AAA*, 2020, pp. 7692–7699.
- [140] F. Zhu, H. Li, W. Ouyang, N. Yu, and X. Wang, "Learning spatial regularization with image-level supervisions for multi-label image classification," in *CVPR*, 2017, pp. 2027–2036.
- [141] H. Guo, K. Zheng, X. Fan, H. Yu, and S. Wang, "Visual attention consistency under image transforms for multi-label image classification," in *CVPR*, 2019, pp. 729–739.
- [142] R. You, Z. Guo, L. Cui, X. Long, Y. Bao, and S. Wen, "Cross-modality attention with semantic graph embedding for multi-label classification," in *AAAI*, 2020, pp. 12709–12716.
- [143] Y. Liu, L. Sheng, J. Shao, J. Yan, S. Xiang, and C. Pan, "Multi-label image classification via knowledge distillation from weakly-supervised detection," in *ACM MM*, 2018, pp. 700–708.
- [144] C. Guo, A. Mousavi, X. Wu, D. N. Holtmann-Rice, S. Kale, S. J. Reddi, and S. Kumar, "Breaking the glass ceiling for embedding-based classifiers for large output spaces," in *NeurIPS*, 2019, pp. 4944–4954.
- [145] W. Chang, H. Yu, K. Zhong, Y. Yang, and I. S. Dhillon, "Taming pretrained transformers for extreme multi-label text classification," in *KDD*, 2020.
- [146] K. Dahiya, D. Saini, A. Mittal, A. Shaw, K. Dave, A. Soni, H. Jain, S. Agarwal, and M. Varma, "Deepxml: A deep extreme multi-label learning framework applied to short text documents," in *WSDM*, 2021, pp. 31–39.
- [147] A. Mittal, K. Dahiya, S. Agrawal, D. Saini, S. Agarwal, P. Kar, and M. Varma, "DECAF: deep extreme classification with label features," in *WSDM*, 2021, pp. 49–57.
- [148] A. Mittal, N. Sachdeva, S. Agrawal, S. Agarwal, P. Kar, and M. Varma, "ECLARE: extreme classification with label graph correlations," in *WWW*, 2021, pp. 3721–3732.
- [149] D. Saini, A. K. Jain, K. Dave, J. Jiao, A. Singh, R. Zhang, and M. Varma, "Galax: Graph neural networks with labelwise attention for extreme classification," in *WWW*, 2021, pp. 3733–3744.
- [150] X. Gong, J. Yang, D. Yuan, and W. Bao, "Generalized large margin knn for partial label learning," *TMM*, 2021.
- [151] X. Gong, D. Yuan, and W. Bao, "Top-k partial label machine," *TNNLS*, 2021.
- [152] J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, and W. Xu, "CNN-RNN: A unified framework for multi-label image classification," in *CVPR*, 2016, pp. 2285–2294.
- [153] Z. Wang, T. Chen, G. Li, R. Xu, and L. Lin, "Multi-label image recognition by recurrently discovering attentional regions," in *ICCV*, 2017, pp. 464–472.
- [154] S. Chen, Y. Chen, C. Yeh, and Y. F. Wang, "Order-free RNN with visual attention for multi-label classification," in *AAAI*, 2018, pp. 6714–6721.
- [155] C. Tsai and H. Lee, "Order-free learning alleviating exposure bias in multi-label classification," in *AAAI*, 2020, pp. 6038–6045.
- [156] V. O. Yazici, A. Gonzalez-Garcia, A. Ramisa, B. Twardowski, and J. van de Weijer, "Orderless recurrent models for multi-label classification," in *CVPR*, 2020, pp. 13437–13446.
- [157] Z. Chen, X. Wei, P. Wang, and Y. Guo, "Multi-label image recognition with graph convolutional networks," in *CVPR*, 2019, pp. 5177–5186.
- [158] T. Chen, M. Xu, X. Hui, H. Wu, and L. Lin, "Learning semantic-specific graph representation for multi-label image recognition," in *ICCV*, 2019, pp. 522–531.
- [159] Y. Wang, D. He, F. Li, X. Long, Z. Zhou, J. Ma, and S. Wen, "Multi-label classification with label graph superimposing," in *AAAI*, 2020, pp. 12265–12272.
- [160] P. Tang, M. Jiang, B. N. Xia, J. W. Pitera, J. Welser, and N. V. Chawla, "Multi-label patent categorization with non-local attention-based graph convolutional network," in *AAAI*, 2020.
- [161] Z. Zhou and J. Feng, "Deep forest: Towards an alternative to deep neural networks," in *IJCAI*, 2017, pp. 3553–3559.
- [162] M. J. Er, R. Venkatesan, and N. Wang, "An online universal classifier for binary, multi-class and multi-label classification," in *IEEE International Conference on Systems, Man, and Cybernetics*, 2016, pp. 3701–3706.
- [163] S. Ding, H. Zhao, Y. Zhang, X. Xu, and R. Nie, "Extreme learning machine: algorithm, theory and applications," *Artificial Intelligence Review*, vol. 44, no. 1, pp. 103–115, 2015.
- [164] H. Chu, K. Huang, and H. Lin, "Dynamic principal projection for cost-sensitive online multi-label classification," *Machine Learning*, vol. 108, no. 8-9, pp. 1193–1230, 2019.
- [165] S. Boulbazine, G. Cabanes, B. Matei, and Y. Bennani, "Online semi-supervised growing neural gas for multi-label data classification," in *IJCNN*, 2018, pp. 1–8.
- [166] P. Li, H. Wang, C. Böhm, and J. Shao, "Online semi-supervised multi-label classification with label compression and local smooth regression," in *IJCAI*, 2020, pp. 1359–1365.
- [167] D. Yeung and H. Chang, "Locally smooth metric learning with application to image retrieval," in *ICCV*, 2007, pp. 1–7.
- [168] H. Wang, Y. Qiang, C. Chen, W. Liu, T. Hu, Z. Li, and G. Chen, "Online partial label learning," in *ECML/PKDD*, 2020.
- [169] W. Gao and Z. Zhou, "On the consistency of multi-label learning," *Artificial Intelligence*, vol. 199-200, pp. 22–44, 2013.
- [170] K. Dembczynski, W. Kotłowski, and E. Hüllermeier, "Consistent multilabel ranking through univariate losses," in *ICML*, 2012.
- [171] H. Zou, "The adaptive lasso and its oracle properties," *Journal of the American Statistical Association*, vol. 101, no. 476, pp. 1418–1429, 2006.
- [172] C.-H. Zhang and J. Huang, "The sparsity and bias of the lasso selection in high-dimensional linear regression," *Annals of Statistics*, vol. 36, no. 4, pp. 1567–1594, 2008.
- [173] C.-H. Zhang, "Nearly unbiased variable selection under minimax concave penalty," *Annals of Statistics*, vol. 38, no. 2, pp. 894–942, 2010.
- [174] J. Fan and R. Li, "Variable selection via nonconcave penalized likelihood and its oracle properties," *Journal of the American Statistical Association*, vol. 96, no. 456, pp. 1348–1360, 2001.
- [175] W. Liu and X. Shen, "Sparse extreme multi-label learning with oracle property," in *ICML*, 2019, pp. 4032–4041.
- [176] W. Liu, "Copula multi-label learning," in *NeurIPS*, 2019, pp. 6334–6343.
- [177] C. Snoek, M. Worring, J. Geusebroek, D. Koelma, F. J. Seinstra, and A. W. M. Smeulders, "The semantic pathfinder: Using an

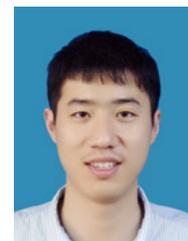
- authoring metaphor for generic multimedia indexing," *TPAMI*, vol. 28, no. 10, pp. 1678–1689, 2006.
- [178] G.-J. Qi, X.-S. Hua, Y. Rui, J. Tang, T. Mei, M. Wang, and H.-J. Zhang, "Correlative multi-label video annotation with temporal kernels," *ACM Transactions on Multimedia Computing, Communications, and Applications*, vol. 5, no. 1, 2008.
- [179] X.-S. Hua and G.-J. Qi, "Online multi-label active learning for large-scale multimedia annotation," Microsoft, Tech. Rep., 2008.
- [180] Z. Wang, Y. Li, S. Wang, and Q. Ji, "Capturing global semantic relationships for facial action unit recognition," in *ICCV*, 2013, pp. 3304–3311.
- [181] K. Zhao, W. Chu, F. D. la Torre, J. F. Cohn, and H. Zhang, "Joint patch and multi-label learning for facial action unit detection," in *CVPR*, 2015, pp. 2207–2216.
- [182] K. Zhao, W. Chu, and H. Zhang, "Deep region and multi-label learning for facial action unit detection," in *CVPR*, 2016, pp. 3391–3399.
- [183] X. Niu, H. Han, S. Shan, and X. Chen, "Multi-label co-regularization for semi-supervised facial action unit recognition," in *NeurIPS*, 2019, pp. 907–917.
- [184] N. Ratnarajah and A. Qiu, "Multi-label segmentation of white matter structures: Application to neonatal brains," *NeuroImage*, vol. 102, pp. 913–922, 2014.
- [185] N. Noorizadeh, K. Kazemi, H. Danyali, and A. Aarabi, "Multi-atlas based neonatal brain extraction using a two-level patch-based label fusion strategy," *Biomedical Signal Processing and Control*, vol. 54, 2019.
- [186] N. Noorizadeh, K. Kazemi, H. Danyali, A. Babajani-Feremi, and A. Aarabi, "Multi-atlas based neonatal brain extraction using atlas library clustering and local label fusion," *Multimedia Tools and Applications*, vol. 79, no. 27–28, pp. 19411–19433, 2020.
- [187] G. Wunder, P. Jung, M. Kasparick, T. Wild, F. Schaich, Y. Chen, S. ten Brink, I. Gaspar, N. Michailow, A. Festag, L. L. Mendes, N. Cassiau, D. Ktenas, M. Dryjanski, S. Pietrzyk, B. Eged, P. Vago, and F. Wiedmann, "5gnow: non-orthogonal, asynchronous waveforms for future mobile applications," *IEEE Communications Magazine*, vol. 52, no. 2, pp. 97–105, 2014.
- [188] L. Guo, B. Jin, R. Yu, C. Yao, C. Sun, and D. Huang, "Multi-label classification methods for green computing and application for mobile medical recommendations," *IEEE Access*, vol. 4, pp. 3201–3209, 2016.
- [189] C. Wang, Y. Wang, B. Xu, Y. He, Z. Dong, and Z. Sun, "A lightweight multi-label segmentation network for mobile iris biometrics," in *ICASSP*, 2020, pp. 1006–1010.
- [190] M. B. Messaoud, I. Jenhani, N. B. Jemaa, and M. W. Mkaouer, "A multi-label active learning approach for mobile app user review classification," in *KSEM*, 2019, pp. 805–816.
- [191] E. L. Mencia and J. Fürnkranz, "Efficient multilabel classification algorithms for large-scale problems in the legal domain," in *Semantic Processing of Legal Texts: Where the Language of Law Meets the Law of Language*, ser. Lecture Notes in Computer Science, vol. 6036, 2010, pp. 192–215.
- [192] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, and I. Androutopoulos, "Large-scale multi-label text classification on EU legislation," in *ACL*, 2019, pp. 6314–6322.
- [193] C. Xiao, H. Zhong, Z. Guo, C. Tu, Z. Liu, M. Sun, Y. Feng, X. Han, Z. Hu, H. Wang, and J. Xu, "CAIL2018: A large-scale legal dataset for judgment prediction," *CoRR*, vol. abs/1807.02478, 2018.
- [194] R. Agrawal, A. Gupta, Y. Prabhu, and M. Varma, "Multi-label learning with millions of labels: recommending advertiser bid phrases for web pages," in *WWW*, 2013, pp. 13–24.
- [195] J. J. McAuley, R. Pandey, and J. Leskovec, "Inferring networks of substitutable and complementary products," in *KDD*, 2015, pp. 785–794.
- [196] G. Farnadi, J. Tang, M. D. Cock, and M. Moens, "User profiling through deep multimodal fusion," in *WSDM*, 2018, pp. 171–179.
- [197] J. Wen, L. Wei, W. Zhou, J. Han, and T. Guo, "GCN-IA: user profile based on graph convolutional network with implicit association labels," in *ICCS*, vol. 12139, 2020, pp. 355–364.
- [198] M. R. Naphade, J. R. Smith, J. Tesic, S. Chang, W. H. Hsu, L. S. Kennedy, A. G. Hauptmann, and J. Curtis, "Large-scale concept ontology for multimedia," *TMM*, vol. 13, no. 3, pp. 86–91, 2006.
- [199] S. S. Bucak, R. Jin, and A. K. Jain, "Multi-label multiple kernel learning by stochastic approximation: Application to visual object recognition," in *NeurIPS*, 2010, pp. 325–333.
- [200] D. Y. Kim, B. Vo, and B. Vo, "Online visual multi-object tracking via labeled random finite set filtering," *CoRR*, vol. abs/1611.06011, 2016.
- [201] L. Grady and G. Funka-Lea, "Multi-label image segmentation for medical applications based on graph-theoretic electrical potentials," in *Computer Vision and Mathematical Methods in Medical and Biomedical Image Analysis*, 2004, pp. 230–245.
- [202] A. Schulz, L. M. Eneldo, and B. Schmidt, "A rapid-prototyping framework for extracting small-scale incident-related information in microblogs: Application of multi-label classification on tweets," *Information Systems*, vol. 57, pp. 88–110, 2016.
- [203] R. Venkatesan, M. J. Er, M. Dave, M. Pratama, and S. Wu, "A novel online multi-label classifier for high-speed streaming data applications," *CoRR*, vol. abs/1609.00086, 2016.
- [204] P. M. Ciarelli, E. Oliveira, and E. O. T. Salles, "Multi-label incremental learning applied to web page categorization," *Neural Computing and Applications*, vol. 24, no. 6, pp. 1403–1419, 2014.
- [205] S. Wan, M. Mak, B. Zhang, Y. Wang, and S. Kung, "Ensemble random projection for multi-label classification with application to protein subcellular localization," in *ICASSP*, 2014, pp. 5999–6003.
- [206] L. Wei, W. Zhou, J. Wen, M. Lin, J. Han, and S. Hu, "MLP-IA: multi-label user profile based on implicit association labels," in *ICCS*, vol. 11536, 2019, pp. 548–561.



Weiwei Liu received his PhD degree under the supervision of Prof. Ivor W. Tsang in computer science from University of Technology Sydney, Australia in 2017. He is currently a full professor with the School of Computer Science, Wuhan University, China. His current research interest is machine learning. His research results have been published at prestigious journals and leading conferences such as JMLR, IEEE TPAMI, IEEE TNNLS, IEEE TIP, IEEE TCYB, NeurIPS, ICML, ACL, AACL, IJCAI and so on.



Haobo Wang received his B.S. degree in Computer Science and Technology from Zhejiang University, China, in 2018. He is currently working toward the Ph.D. degree in the College of Computer Science and Technology, Zhejiang University. His research interests include machine learning and data mining, especially on multi-label learning and weakly-supervised learning.



Xiaobo Shen received his BSc and PhD from School of Computer Science and Engineering, Nanjing University of Science and Technology in 2011 and 2017 respectively. He is currently a full professor with the School of Computer Science and Engineering, Nanjing University of Science and Technology, China. He has authored over 30 technical papers in prominent journals and conferences, such as IEEE TNNLS, IEEE TIP, IEEE TCYB, NeurIPS, ACM MM, AACL, and IJCAI. His primary research interests are Multi-view Learning, Multi-label Learning, Network Embedding and Hashing.



Ivor W. Tsang is an ARC Future Fellow and Professor of Artificial Intelligence with the University of Technology Sydney, Australia. He is also the Research Director of the Australian Artificial Intelligence Institute. His research interests include transfer learning, generative models, and big data analytics for data with extremely high dimensions. In 2013, Prof Tsang received his prestigious ARC Future Fellowship for his research regarding Machine Learning on Big Data. In 2019, his JMLR paper titled "Towards ultra-

high dimensional feature selection for big data" received the International Consortium of Chinese Mathematicians Best Paper Award. In 2020, Prof Tsang was recognized as the AI 2000 AAIL/JCAI Most Influential Scholar in Australia for his outstanding contributions to the field of AAIL/JCAI between 2009 and 2019. His research on transfer learning granted him the Best Student Paper Award at CVPR 2010 and the 2014 IEEE TMM Prize Paper Award. In addition, he received the IEEE TNN Outstanding 2004 Paper Award in 2007. He serves as a Senior Area Chair/Area Chair for NeurIPS, ICML, AISTATS, AAIL and JCAI, and the Editorial Board for JMLR, MLJ, and IEEE TPAMI.