

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/267154292>

Multilabel Learning: A Review of the State of The Art and Ongoing Research

Article in *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* · November 2014

DOI: 10.1002/widm.1139

CITATIONS

87

READS

3,548

2 authors:



Eva Gibaja

University of Cordoba (Spain)

57 PUBLICATIONS 489 CITATIONS

SEE PROFILE



Sebastian Ventura

University of Cordoba (Spain)

317 PUBLICATIONS 9,993 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



artificial intelligence [View project](#)



Emerging Trends in Data Analysis (EMERALD) [View project](#)

Multi-Label Learning. A Review of the State of the Art and Ongoing Research

Eva Gibaja · Sebastián Ventura

Received: date / Accepted: date

Abstract Multi-label learning is quite a recent supervised learning paradigm. Owing to its capabilities to improve performance in problems where a pattern may have more than one associated class, it has attracted the attention of researchers, producing an increasing number of publications. The paper presents an up-to-date overview about multi-label learning with the aim of sorting and describing the main approaches developed till now. The formal definition of the paradigm, the analysis of its impact on the literature, its main applications, works developed and ongoing research are presented.

Keywords multi-label learning · review

1 Introduction

Multi-label learning (MLL) is a supervised learning paradigm which has attracted a great deal of attention in recent years due to its capabilities of improving performance in many current applications such as the classification of multimedia, the prediction of gene and protein functions, the direct marketing or the social network mining. All of these applications have in common that not one, but multiple outputs are required. Therefore the *only-one-label-per-pattern* restriction of classical supervised learning (a.k.a. single-label learning) is not satisfied. As well as dealing with multiple outputs, MLL has to deal with trending challenges such as relationships between labels, the computational costs of generating the models, presence of imbalanced labels or high dimensionality of data. The key challenge has been recently identified as dealing with the high dimensionality of the output space, especially in domains with a large number of labels [261]. This challenge involves exploring label correlations efficiently. Besides, specialised workshops [2], [3], [5], special issues [4], repositories of benchmark data sets and software [1], [39], [167], [213], have contributed to the progress in this field. All of these factors

University of Córdoba, Spain · University of Córdoba, Spain – King Abdulaziz University, Saudi Arabia
E-mail: {egibaja,sventura}@uco.es

have become MLL into a relevant supervised learning paradigm with an increasing number of papers published on it per year.

First tutorials about MLL were published in its earlier stages [33], [208]. They introduced the setting of multi-label learning and compiled the contributions and methods developed up to 2007. Next, in [210], a review that has become a reference for the MLL community was published. It included the main proposals developed up to 2008 and also the description of the main evaluation metrics. It is also worth citing two other recent papers. First, in [135] an experimental comparison of 12 well-known MLL methods was carried out using 16 evaluation measures on 11 benchmark data sets. Its aim was to provide a better understanding of the performance of these methods. Second, Zhang and Zhou [261] have recently published a paper, whose main aim is describing the setting, evaluation metrics and 8 representative MLL algorithms in an elaborated and formal way. Due to the high number of proposals developed the latest years (around 700 new works only from 2009 to 2012) the aim of our paper is filling this gap with an update for the topic.

The rest of the paper is organised as follows: first, the formal definition of MLL and the main fields of application are described (section 2) followed by the analysis of its impact on the literature (section 3). Next the state of art (section 4) and ongoing research (section 5) are analysed. Section 6 summarises a series of pitfalls and guidelines mainly focused on the selection of a proper multi-label learner. The paper finishes with the set of more relevant conclusions.

2 Multi-Label Learning

This section presents the main fields of application of MLL, a formal definition of the paradigm, a summary of evaluation metrics and the description of other learning settings that may share some features with MLL.

2.1 Key applications of MLL

- **TEXT CATEGORISATION** consists of assigning a set of predefined categories to documents. Since a document can belong simultaneously to more than one category it can be tackled with MLL. It has been applied to many kinds of documents such as legal texts [129], [130], web documents [175], [215], news [180], research papers [145], narrative clinical text [190], patents [54] or aeronautics reports [149]. Other related applications are document indexing [119], tag suggestion [110], [189], e-mail filtering [243], medical coding [238], query categorisation [201], or the classification of news sentences into multiple emotion categories [16], [17].
- **MULTIMEDIA**. MLL techniques have been applied to many types of resources such as images, videos and sound. Examples of applications are: automatic image annotation [141], [226], face verification [117], video annotation [224], object recognition [187], detection of emotions into music [133], [206], music metadata extraction [151] and speech emotion classification [188].
- **BIOLOGY**. It is worth noting gene function prediction [15], [36], [37], [45], [72], [217], [257] and protein function prediction [11], [12], [38], [148] applications.

Other recent applications are the prediction of proteins' 3D structures [71] and, finally, the problem of protein sub-cellular multi-location [48] (proteins may simultaneously exist at, or move between, two or more different subcellular locations).

- CHEMICAL DATA ANALYSIS. MLL has also been applied to predict adverse drug reactions [136], to identify the drugs that have two or more different biological actions (drug discovery) [111] and to detect contaminants in machine lubricants by using spectral images (vision-based metal spectral analysis) [216].
- SOCIAL NETWORK MINING has become a new area of interest. Collective behaviour learning consists of inferring behaviour or preferences of individuals [199], [200]. Social networking advertising [116] or the automatic annotation of the nodes of a partially labelled multi-relational graph [155] are other fields of application.
- E-LEARNING. MLL has been also applied to classify learning styles based on learners' profiles [150] and to tag learning objects [128].
- OTHER APPLICATIONS. Other fields of application worth citing are direct marketing, where potential buyers of certain products are identified [264], and medical diagnosis [183] (many symptoms may be associated with more than one syndrome). Finally, in [6] MLL was applied to classify dermoscopy images of skin lesions which could contain several pattern lesions.

2.2 Formal definition

According to [180] and [256], given $\mathcal{X} = X_1 \times \dots \times X_d$ a d -dimensional input space of numerical or categorical features, and an output space of q labels, $\mathcal{Y} = \{\lambda_1, \lambda_2, \dots, \lambda_q\}$, a multi-label pattern can be defined as a pair (\mathbf{x}, Y) , where $\mathbf{x} = (x_1, \dots, x_d) \in \mathcal{X}$ and $Y \subseteq \mathcal{Y}$ is called *label set*. Label associations can be also represented as a q dimensional binary vector $\mathbf{y} = (y_1, y_2, \dots, y_q) = \{0, 1\}^q$ where each element is 1 if the label is relevant and 0 otherwise. Three different tasks can be included in MLL [210]:

- *Label Ranking* (LR) consists of producing a function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ that induces an ordering of all the possible labels which express the relevance of labels to a given instance \mathbf{x} . Thus label λ_1 is considered to be ranked higher than λ_2 if $f(\mathbf{x}, \lambda_1) > f(\mathbf{x}, \lambda_2)$. For each instance $\mathbf{x} \in \mathcal{X}$, a rank function, $\tau_{\mathbf{x}} : \mathcal{Y} \rightarrow \{1, 2, \dots, q\}$, can be defined using the output real value of the classifier f , such that if $f(\mathbf{x}, \lambda_1) > f(\mathbf{x}, \lambda_2)$ then $\tau_{\mathbf{x}}(\lambda_1) < \tau_{\mathbf{x}}(\lambda_2)$. The lower the value, the better the position in the ranking is.
- *Multi-Label Classification* (MLC) consists of defining a function $h_{\text{MLC}} : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$ that returns the set of relevant labels. So for each $\mathbf{x} \in \mathcal{X}$, we have a bipartition (Y, \overline{Y}) of the label set \mathcal{Y} , where $Y = h_{\text{MLC}}(\mathbf{x})$ is the set of relevant labels and \overline{Y} is the set of irrelevant ones. Here $\overline{Y} = \mathcal{Y} \setminus Y$ denotes the set theoretic complement of Y in \mathcal{Y} . Multi-class (MCC) and binary classification (BC) can be seen as a particular case of MLC where $h_{\text{MCC}} : \mathcal{X} \rightarrow \mathcal{Y}$ and $h_{\text{BC}} : \mathcal{X} \rightarrow \{0, 1\}$. A multi-label classifier can be derived from a multi-label ranking model by using a threshold function. Strategies for thresholding can be found in [72], [76], [100], [139], [201], [241].
- *Multi-Label Ranking* (MLR) is a generalisation of MLC and LR consisting of producing, at the same time, both a bipartition and a consistent ranking. In

other words, if Y is the set of labels associated with an instance, then, in a consistent ranking, labels in Y will have higher rank than labels in \overline{Y} .

According to [261], multi-label learning can be considered a particular setting into a wider framework, called *multi-target learning*, where a pattern is associated with multiple outputs. Depending on the kind of outputs different instantiations are considered, calling the settings: i) multi-label learning when output variables are binary, ii) *multi-dimensional learning* [19] if output variables are multi-class, or iii) *multi-output regression* [125] for numerical outputs. Besides, combination of different types of output variables can be considered.

2.3 Metrics to evaluate the models

Tsoumakas et al. [210] distinguish two kinds of metrics to evaluate MLL methods: *label-based* metrics and *example-based* metrics. The idea of the label-based approach is computing a single-label metric for each label based on the number of true positives (tp), true negatives (tn), false positives (fp) and false negatives (fn) and then obtaining an average value. Given B any binary evaluation measure, two different averaging approaches can be used: the *macro approach* computes one metric for each label and then the values are averaged over all the categories, while the *micro approach* considers predictions of all instances together (aggregating the tp , tn , fp , fn values of all classes) and then calculates the measure across all labels.

On the other hand, example-based metrics are calculated for each test example and then averaged across the test set. They are categorised commonly into two groups: metrics to evaluate rankings and metrics to evaluate bipartitions. Table 1 summarises the main metrics described in literature according to the categorization described. Besides, metrics to evaluate confidence scores and hierarchies of labels have been included.

Let $T = \{(\mathbf{x}_i, Y_i) | 1 \leq i \leq t\}$ be a multi-label test set with t instances. Given an instance, \mathbf{x} , let Y and Z be the set of true and predicted labels, and let $\mathbf{w} = (w_{\lambda_1}, w_{\lambda_2}, \dots, w_{\lambda_q})$ be a vector with normalised output confidence scores in $[0, 1]$. For any predicate, π , $\llbracket \pi \rrbracket$ returns 1 if the predicate is true and 0 otherwise. Let τ^* be the true ranking, Δ stands for the symmetric difference of two sets, \arg function returns a label, and given a hierarchy of labels, $anc(i)$ returns the set of ancestors of a node i .

2.4 Other related learning settings

This section discusses other learning settings, which may share some features with MLL and are worth to be briefly discussed.

- MULTI-TASK LEARNING OR LEARNING PARALLEL TASKS (MTL) [138] tries learning in parallel several tasks which share a common representation. Thus what is learned for one task can help the others to be learnt better. According to Zhang and Zhou [261], there are three main differences between multi-label and multi-task learning. First, in MLL all the examples have the same feature space while in multi-task learning it can be the same or different. The

LABEL-BASED	[210]	Macro approach	$B_{macro} = \frac{1}{q} \sum_{i=1}^q B(tp_i, fp_i, tn_i, fn_i)$ e.g. $recall_{mac} = \frac{1}{q} \sum_{i=1}^q \frac{tp_i}{tp_i + fn_i}$
	[210]	Micro approach	$B_{micro} = B(\sum_{i=1}^q tp_i, \sum_{i=1}^q fp_i, \sum_{i=1}^q tn_i, \sum_{i=1}^q fn_i)$ e.g. $recall_{mic} = \frac{\sum_{i=1}^q tp_i}{\sum_{i=1}^q tp_i + \sum_{i=1}^q fn_i}$
METRICS TO EVALUATE BIPARTITIONS			
EXAMPLE-BASED	[86]	↑ Subset accuracy	$= \frac{1}{t} \sum_{i=1}^t \mathbb{I}[Z_i = Y_i]$
	[62]	↓ Subset 0/1 loss	$= \frac{1}{t} \sum_{i=1}^t \mathbb{I}[Z_i \neq Y_i]$
	[179]	↓ Hamming loss	$= \frac{1}{t} \sum_{i=1}^t \frac{1}{q} Z_i \Delta Y_i $
	[87]	↑ Recall	$= \frac{1}{t} \sum_{i=1}^t \frac{ Z_i \cap Y_i }{ Y_i }$
	[87]	↑ Precision	$= \frac{1}{t} \sum_{i=1}^t \frac{ Z_i \cap Y_i }{ Z_i }$
	[87]	↑ Accuracy	$= \frac{1}{t} \sum_{i=1}^t \frac{ Z_i \cap Y_i }{ Z_i \cup Y_i }$
	[87]	↑ F1-Score	$= \frac{1}{t} \sum_{i=1}^t \frac{2 Z_i \cap Y_i }{ Z_i + Y_i }$
METRICS TO EVALUATE RANKINGS			
	[180]	↓ One-error	$= \frac{1}{t} \sum_{i=1}^t \mathbb{I}[\arg \min_{\lambda \in \mathcal{Y}} \tau_i(\lambda) \notin Y_i]$
	[180]	↓ Coverage	$= \frac{1}{t} \sum_{i=1}^t \max_{\lambda \in Y_i} \tau_i(\lambda) - 1$
	[179]	↓ Ranking loss	$= \frac{1}{t} \sum_{i=1}^t \frac{1}{ Y_i \overline{Y_i} } E $ where $E = \{(\lambda, \lambda') \tau_i(\lambda) > \tau_i(\lambda'), (\lambda, \lambda') \in Y_i \times \overline{Y_i}\}$
	[56][130]	↓ IsError	$= \frac{1}{t} \sum_{i=1}^t \mathbb{I}[\sum_{\lambda \in \mathcal{Y}} \tau_i^*(\lambda) - \tau_i(\lambda) \neq 0]$
	[180]	↑ Average precision	$= \frac{1}{t} \sum_{i=1}^t \frac{1}{ Y_i } \sum_{\lambda \in Y_i} \frac{ \{\lambda' \in Y_i \tau_i(\lambda') \leq \tau_i(\lambda)\} }{\tau_i(\lambda)}$
	[130]	↓ Margin loss	$= \frac{1}{t} \sum_{i=1}^t \max(0, \max_{\lambda \in Y_i} \{\tau(\lambda) \lambda \in Y_i\} - \min_{\lambda' \notin Y_i} \{\tau(\lambda') \lambda' \notin Y_i\})$
	[153]	↓ Ranking error	$= \frac{1}{t} \sum_{i=1}^t \sum_{\lambda \in \mathcal{Y}} \tau_i^*(\lambda) - \tau_i(\lambda) ^2$
METRICS TO EVALUATE CONFIDENCE SCORES			
	[165]	↓ Log loss	$= \frac{1}{tq} \sum_{i=1}^t \sum_{\lambda \in \mathcal{Y}} \min(-\text{LOGLOSS}(\lambda, \mathbf{w}_i), \ln(t))$ where $\text{LOGLOSS}(\lambda, \mathbf{w}) = \ln(w_\lambda)$ if $\lambda \in Y$ $\ln(1 - w_\lambda)$ if $\lambda \in \overline{Y}$
METRICS FOR HMC			
	[18]	↓ 0/1 loss	$= \frac{1}{t} \sum_{i=1}^t \mathbb{I}[Z_i \neq Y_i]$
	[18]	↓ Symmetric diff.	$= \frac{1}{t} \sum_{i=1}^t Z_i \Delta Y_i $
	[18]	↓ Hierarchical loss	$= \frac{1}{t} \sum_{i=1}^t \sum_{j=1}^q \mathbb{I}[Z_{ij} \neq Y_{ij} \wedge Z_{ik} = Y_{ik}, k \in \text{anc}(j)]$

Table 1 Taxonomy of metrics to evaluate MLL algorithms. ↑ means the metric has to be maximised and ↓ means the metric has to be minimised

purpose is also different, thus in MLL one task is learned (i.e. predicting the label set associated with an object) while in multi-task learning, multiple tasks are learned simultaneously. Finally, the label space in MLL is large while in multi-task learning is not reasonable to consider a large number of tasks.

- MULTIPLE-LABELS LEARNING OR PARTIAL LABELLING [146]. In this kind of problem each pattern has multiple candidate labels, and only one of them is the correct one. Real problems such as disagreement between assessors can be viewed from this point of view.
- MULTI-INSTANCE LEARNING (MIL) [248] consists of learning a concept where training labels are associated with sets (*bags*) of patterns (*instances*) rather than with individual patterns. A bag is positive if, at least, one of its patterns is positive and negative otherwise. Numerous real-world tasks (e.g. drug activity prediction or web index page recommendation) can be naturally represented as multiple instance problems. Section 5.3 describes the MIML setting, in which a bag may have associated not one, but a set of labels.
- REVERSE MULTI-LABEL LEARNING (RMLL) [157] consists of carrying out a reverse prediction, that is, to predict sets of relevant instances given a set of labels.

- PREFERENTIAL TEXT CLASSIFICATION [10] is a problem where primary (central topics of the document) and secondary categories associated with a document can be distinguished. Misclassifications related to primary categories should be penalised more severely than those related to secondary ones.
- WEAK LABEL PROBLEM OR LEARNING WITH INCOMPLETE CLASS ASSIGNMENT [31], [196]. In this kind of problem only a partial labelling associated with each training example is provided. In other words, if a label has been assigned to an instance, it is a proper label of the instance but, if a label has not been assigned it cannot be concluded that it is not a proper label.
- MULTI-VALUED MULTI-LABEL (M^2) [121]. In this kind of problem samples are not only associated with a set of labels, they may also have some attributes of the pattern presenting several values.
- GRADED MULTI-LABEL CLASSIFICATION (GMLC) [45]. The aim is to obtain not only a set of labels, but also a membership value for each label in the sense of fuzzy set theory.
- MULTI-VIEW LEARNING [24]. One object has different representations in the form of several disjoint subsets of features (each sub-set is a view), each of which is sufficient for learning the target concept. Thus several classifiers are trained on subsets of features or views. This setting has been combined with active learning in order to reduce the annotation effort in multi-label tasks (see section 5.4).

3 MLL in the literature

In this section, the visibility of research in MLL is going to be studied by analysing the number of publications and citations at the two scientific citation databases more broadly used: Thompson ISI and Elsevier Scopus. Table 2 shows the queries made and the overall results. In the case of ISI, the search has been carried out by title and topic while in Scopus the search has been carried out by title, abstract and keywords. The search was restricted to articles and conference papers.

DATABASE	ISI	SCOPUS
QUERY	TOPIC (("multi-label" OR "multilabel") AND ("classification" OR "learning")) OR TITLE: (("multi-label" OR "multilabel") AND ("classification" OR "learning")) Refined by: DOCUMENT TYPES: (ARTICLE OR MEETING) Timespan: all years	TITLE-ABS-KEY((((multi-label) OR (multilabel)) AND ((classification) OR (learning)))) AND (LIMIT-TO(DOCTYPE, "cp") OR LIMIT-TO(DOCTYPE, "ar"))
DATE	6th July 2014	2nd July 2014
#DOCUMENTS	638	1116
H-INDEX	28	39

Table 2 Queries made to ISI and Scopus

The *h-index* reflects both the number of publications and the number of citations per publication. The h-index for the 1116 papers considered in Scopus is 39 which means that 39 of these papers have been cited at least 39 times. In ISI the h-index of the 638 retrieved references is 28. Two graphics summarising the total of articles and citations both in ISI and Scopus databases are shown in Figures 1

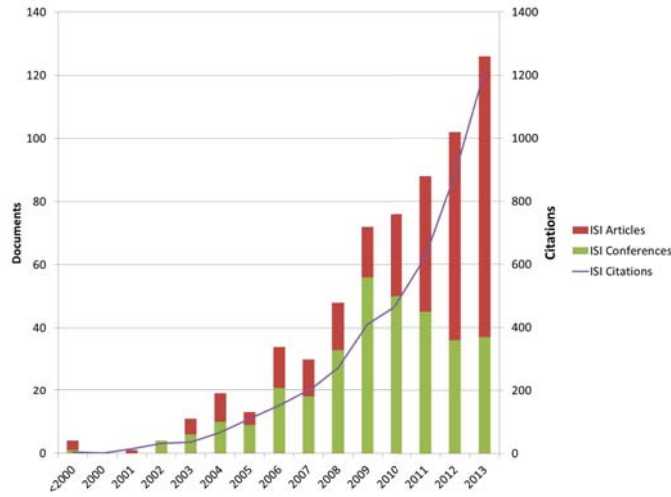


Fig. 1 Multi-label learning in ISI

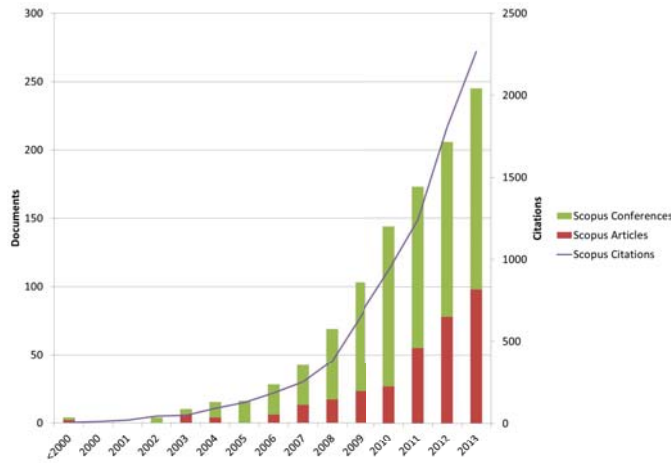


Fig. 2 Multi-label learning in Scopus

and 2. It is observed that MLL is a quite recent topic and also the exponential increase of papers and citations in the latest years. Finally, the ten most cited papers in ISI and Scopus databases with their respective number of citations are listed in Tables 3 and 4.

4 MLL methods

This section is going to follow the taxonomy of MLL methods presented in [210] that distinguishes between: *problem transformation methods* and *algorithm adaptation methods*. The former are algorithm independent and transform the multi-label

#	CIT.	TITLE	AUTHORS	PUBLICATION	YEAR
1	911	Improved boosting algorithms using confidence-rated predictions [179]	Shapire and Singer	Machine Learning Theory	1999
2	378	RCV1: A new benchmark collection for text categorization research [120]	Lewis et al.	Journal of Machine Learning Research	2004
3	259	Learning multi-label scene classification [26]	Boutell et al.	Pattern Recognition	2004
4	228	ML-KNN: A lazy learning approach to multi-label learning [259]	Zhang and Zhou	Pattern Recognition	2007
5	151	iLoc-Euk: A multi-label classifier for predicting the subcellular localization of singleplex and multiplex eukaryotic proteins [48]	Chou et al.	PLoS ONE	2011
6	143	Multilabel neural networks with applications to functional genomics and text categorization [257]	Zhang and Zhou	IEEE Transactions on Knowledge and Data Engineering	2006
7	120	Optimization method based extreme learning machine for classification [96]	Huang et al.	Neurocomputing	2010
8	97	A kernel method for multi-labelled classification [72]	Elisseeff and Weston	NIPS	2001
9	82	iLoc-Virus: a multi-label learning classifier for identifying the subcellular localization of virus proteins with both single and multiple sites [230]	Xiao et al.	Journal of Theoretical Biology	2011
10	81	Decision trees for hierarchical multi-label classification [220]	Vens et al.	Machine Learning	2008

Table 3 ISI top ten most cited papers

#	CIT.	TITLE	AUTHORS	PUBLICATION	YEAR
1	1360	Improved boosting algorithms using confidence-rated predictions [179]	Shapire and Singer	Machine Learning	1999
2	415	Learning multi-label scene classification [26]	Boutell et al.	Pattern Recognition	2004
3	381	Multi-label classification: An overview [208]	Tsoumakas and Katakis	International Journal of Data Warehousing and Mining	2007
4	359	ML-KNN: A lazy learning approach to multi-label learning [259]	Zhang and Zhou	Pattern Recognition	2007
5	193	Multilabel neural networks with applications to functional genomics and text categorization [257]	Zhang and Zhou	IEEE Transactions on Knowledge and Data Engineering	2006
6	155	iLoc-Euk: A multi-label classifier for predicting the subcellular localization of singleplex and multiplex eukaryotic proteins [48]	Chou et al.	PLoS ONE	2011
7	149	Improved boosting algorithms using confidence-rated predictions [178]	Shapire and Singer	ACM Conference on Computational Learning Theory	1998
8	142	Optimization method based extreme learning machine for classification [96]	Huang et al.	Neurocomputing	2010
9	128	Semantic annotation and retrieval of music and sound effects [214]	Turnbull et al.	IEEE Trans. Audio, Speech and Language Processing	2008
10	127	A kernel method for multi-labelled classification [72]	Elisseeff and Weston	NIPS	2001

Table 4 Scopus top ten most cited papers

problem into one or more single-label ones in which is then applied a single-label classification algorithm, whilst the latter extend a single-label algorithm in order to directly deal with multi-label data. A summary of the main problem transformation methods is found in Table 5. Table 6 summarises the main algorithm adaptation methods developed according to the adapted paradigm.

RANKING VIA SINGLE LABEL LEARNING	
Ignore	
Select (<i>min</i> , <i>max</i> , <i>random</i>)	[210]
Copy, copy-weight	
BINARY	
BR	[210]
CC	[171]
PCC	[60]
2BR (MBR)	[13] [151] [207]
BR+	[47]
PAIRWISE	
RPC	[99]
CLR	[28] [83]
QCLR	[131]
DLVM	[134]
LABEL COMBINATION	
LP	[26] [210]
PS	[164]
LPBR(ChiDep)	[203]
ENSEMBLES OF MLL METHODS	
ECC	[171]
EBR	[171]
EPS	[164]
RAkEL	[211]
RAkEL-CSML	[127]
CDE	[203]
RAKEL++	[174]
RF-PCT	[114]
RFML-C4.5	[135]
TREMLC	[143]
DST-fusion	[177]
OTHER TRANSFORMATIONS	
InsDif	[259]
OBO	[233]

Table 5 Taxonomy of the MLL transformation methods

4.1 Problem transformation methods

4.1.1 Ranking via single-label learning

This approach consists of transforming the instances in a multi-label data set in order to obtain a single-label one. Then a single-label classifier, which is able to produce a score for each label (e.g. probability), will allow to obtain a ranking [210]. Three strategies can be followed: ignoring all multi-label instances (*ignore* method), transforming every multi-label instance into several ones, one per label (*copy* and *copy-weight* methods) or selecting one of the labels of the multiple-labelled patterns (*select-max*, *select-min* or *select-random* methods). Despite being simple, these methods produce problems of information loss in terms of labels or label relationships and so they may not be very useful.

4.1.2 Binary methods

The *Binary Relevance* method (BR) [210] follows the *one-versus-all* (OVA) philosophy and builds one binary data set for each label. Patterns predicting the label are considered positive patterns and the rest are considered to be negative. Once an unknown pattern is presented to the model, the output will be the set of positive classes predicted. The main problem of BR is the assumption of label

independence that ignores the relationships between labels [267] and may lead to failure to predicting label combinations or rankings of labels [207]. Nevertheless it is computationally simple, scales linearly with the number of labels and can be parallelised [171].

Some approaches have been developed in order to overcome the label independence assumption of BR while maintaining a reasonable complexity. They are described below.

The *Classifier Chains* (CC) model [171] generates q binary classifiers, but they are linked in such a way that each classifier incorporates the labels predicted by the previous classifiers in the chain as additional features. Therefore, label correlations are considered in a random manner while the complexity is linear with the number of labels. CC can be parallelised in the training stage. Experiments proved CC overall improved over BR with a similar complexity. As the order of the chain itself can influence the performance, an *Ensemble of Classifier Chains* (ECC), which trained a set of CC classifiers with a random chain ordering and a random subset of training patterns sampled with replacement, was also proposed. These authors also proposed an *Ensemble of Binary Relevance classifiers* (EBR) developed identically to ECC but without chaining. ECC and EBR obtained promising predictive results in a wide range of metrics and data sets and demonstrated to be efficient on large data sets without significant losses in predictive performance. In [60] *Probabilistic Classifier Chains* (PCC), a Bayes optimal way of forming classifier chains that outperformed CC, was described. Its drawback was the computational complexity at prediction time (it must look at each of 2^q possible combinations) being recommendable only with a small to moderate number of labels ($q < 15$).

A different approach is 2BR, also called *Meta-BR* (MBR) in [171], which basically consists of applying BR twice [13] [151] [207]. It follows the philosophy of *Stacking* [228] and maintains the linear time complexity with respect to the number of labels in the data set. During the first step (the base-level), a BR classifier is learnt and the second BR step (the meta-level) implements a meta-learning stage. To do this, the input feature set of the meta-level is augmented with q extra features consisting of the predictions of each binary classifier in the base-level. After that, q new binary classifiers are trained with this extended data set considering the desired outputs as targets. For classification, outputs of the first level of classifiers are used to extend the examples and these new examples are classified by the second level. In [207] only the predictions of the base-level models regarding correlated labels (according to the ϕ coefficient) were considered when constructing the meta-level training examples for a certain label. Decision trees and Support Vector Machines (SVM) were used in both base and meta-learning levels and linear regression at meta-learning level. This approach was able to improve the efficiency substantially, without significant loss in predictive performance.

Finally, *BRplus* (BR+) [47] also incremented the feature space of the BR classifiers with labels, but in this case with $q - 1$ features corresponding to the other labels in the data set. During the classification stage, features of unlabelled examples were augmented by using a BR predictor trained with the original training data. The reported experiments showed it effectively improved BR and maintained results comparable to CC.

4.1.3 Pairwise methods

The *Ranking by Pairwise Comparison* (RPC) approach [99] follows a *one-versus-one* (OVO) philosophy and transforms a data set with q labels into $q(q-1)/2$ binary data sets, one per each pair of labels. Each data set uses the examples belonging to one of the two classes as positive or negative examples respectively (patterns belonging to both labels are not considered). Then, a binary classifier is built for each data set. Given an unknown pattern, the prediction is obtained by invoking all models and obtaining a rank from the counting votes for each label. It is also worth mentioning *Calibrated Label Ranking* (CLR) [28], [83] which adds to the RPC transformation q data sets, λ_i vs. λ_0 , which are identical to a BR transformation. The key idea of this approach is the virtual label, λ_0 , that acts as a split point separating the relevant for the irrelevant labels obtaining a consistent ranking and bipartition. Empirical results in the area of text categorization, image classification and gene analysis showed it outperformed BR, MMP and MLPP (these last two methods will be described in section 4.3.4). Authors also found that CLR is typically a bit more conservative in predicting the number of relevant labels. This translates into lower recall, but precision and F1-score values tend to be high. The main drawbacks of pairwise methods described are the space complexity and the need to query all the generated (q^2) binary models at runtime, which may become impractical for large number of labels. *QWeighted Calibrated Label Ranking* (QCLR) [131] combined CLR with the *QWeighted* [129] voting schema. Multi-class QWeighted efficiently computes the class with the highest accumulated voting mass (i.e. the top ranked class) without evaluating all pairwise classifiers. In the multi-label case, the process is iteratively applied until the virtual label is returned, which means that all remaining labels are irrelevant. It speeded up the voting process reducing the evaluations needed from $q(q-1)/2$ to $q \log(q)$ in practice, which is near the q evaluations processed by BR. The remaining bottleneck is the need to store a quadratic number of base classifiers. Finally, *Dual layer Voting Method* (DLVM) [134] is another voting strategy consisting of a two-stage architecture. In the first layer BR models output a probability about the relevance of every label. Then, if this probability is above a certain threshold, the pairwise models of the second layer are consulted. DLVM outperformed the CLR's majority voting in terms of speed whilst maintaining the prediction performance.

4.1.4 Label Combination methods

The *Label Powerset* (LP) approach [26], [210] builds a single-label data set where each possible combination of labels is considered as a class itself. Then a multi-class algorithm is applied. Given a new instance, LP outputs a class, which actually is a label set in the original data set. Despite being able to model label correlations, the drawback is its complexity, in the worst case, exponential with the number of labels. *Pruned Problem Transformation* (PPT) or *Pruned Sets* (PS) [164] tries to reduce this complexity focusing on the most important combinations of labels by pruning examples with less frequent label sets; to compensate for such information loss it reintroduces the pruned examples along with subsets of their label sets. After that, LP is applied. Like LP, PS is not able to output label sets that are not in the training set. A method to tackle this issue is to combine the results of several classifiers in an ensemble. For each classifier in the ensemble a subset of

the training set (i.e. 63%) is sampled without replacement and a PS classifier is trained. During the prediction stage, outputs are combined by a voting scheme and a threshold separates relevant and irrelevant labels. EPS proved to be competitive with LP and RAKEL in terms of efficiency and predictive accuracy.

Random k-label sets (RAkEL) [211] is based on random projections of the label space and builds an ensemble of LP classifiers, each one trained with a random subset of k labels. Thus it is able to deal with label correlations whilst avoiding the computational complexity of LP. During classification, the output of the classifiers is averaged per label and thresholding is used to assign the label set. Disjoint and overlapping subsets of labels were studied and experiments showed that both improved over LP, even in large data sets. RAKEL with overlapping label sets and C4.5 as base classifier improved RAKEL with disjoint subsets. In order to minimise the noise or errors in collaborative filtering cost-sensitive approaches of RAKEL and stacking [228] have been applied to audio tagging by considering the tag count as a misclassification cost [127]. In general the cost sensitive version outperformed the cost-insensitive counterpart.

LPBR [203] is an hybrid method which carried out combinations of LP and BR rounds. After a first round with BR, the most dependent labels were clustered into a new label. LP was applied within the groups of dependent labels (a group with a limited number of labels) while BR was applied to the independent groups of labels. The process was repeated until the accuracy did not improve. The approach where the dependence between labels is computed by using the χ^2 score was called *ChiDep*. In order to improve the performance, *ChiDep Ensemble* (CDE) was also proposed. A high number of random label sets partitions were generated (e.g. 10000), and for each partition a score was computed based on sum of the χ^2 of the pairs of labels in the partition. The top high-scored partitions were selected as members of the ensemble. It obtained competitive prediction results. Train time of ChiDep and CDE was relatively long and approximate to that of RAKEL and 2BR. Test time of ChiDep was comparable to BR while test time of CDE was longer.

Finally, in [172], [173] an improvement of RAKEL, where the subsets were not selected randomly but in advance was described. The subset selection problem was formulated as a *Set Covering Problem* (SCP) (i.e. cover the set of labels by a set of label subsets of a given size k) and an approximation greedy algorithm was employed to derive the subsets. The work has been recently generalised [174] by proposing a general framework that allows the application of a wide variety of optimization criteria (i.e. balanced representation of each label, coverage of inter-label correlations or both). Authors have also proposed in this paper RAKEL++, an improved version of RAKEL in which, instead of voting, the confidence values are taken into account by thresholding the average of the probabilities provided by the base-classifier for each label. Besides, it employs built-in cross validation for deriving the threshold value. The proposed strategies perform in a efficient and stable manner overall obtaining better results than RAKEL.

4.1.5 Other transformations

Other transformations that do not fit into the previously defined categories have been described. An example is INSTANCE DIFFERENTIATION (InsDif) [259] whose

aim is exploiting the relationship between the input ambiguity and output ambiguity by supposing that an object belongs to several labels due to the diverse information embodied in the object (in the input space). For each label, a prototype vector was calculated by averaging all instances belonging to such label. Then each example was transformed into a bag of q instances each one being equal to the difference between the original example and one of the prototype vectors. A two-level classification strategy learnt from the transformed data set. Finally, in [233] a *one-by-one* (OBO) transformation generated a new data set for each label where instances were only the positive ones and each sub-problem was solved by means of a *one-class* classifier (i.e. Support Vector Data Description [202]). To compensate missing correlations between labels, linear ridge regression integrated predictions of all sub-classifiers.

4.2 Ensembles of MLL methods

Ensemble methods whose base classifiers are multi-label learners are considered by Madjarov et al. [135] a special group of methods because they are developed on top of problem transformation and algorithm adaptation approaches. The RAKEL, EPS, ECC, EBR, PCC and CDE approaches previously described are examples of ensembles of MLL methods. Binary methods are occasionally referred to as ensemble methods because they employ multiple binary models. As none of these models is multi-label capable, the term ensemble is preferable in the sense of an ensemble of MLL methods [171].

Other than the above cited methods, *Random Forest of Predictive Clustering Trees* (RF-PCT) [114] is an ensemble that uses PCT (described in section 4.3.1) as base classifier. The diversity among the base classifiers is obtained by using bagging and selecting, at each node, the best feature from a random subset of the input attributes. The outputs are combined using a voting scheme. Finally, *Random forest of ML-C4.5* (RFML-C4.5) [135] follows the same philosophy but uses ML-C4.5 trees (described in section 4.3.1) as base classifiers. In [135] an intensive experimental evaluation involving a wide variety of algorithms, metrics and statistical tests was carried out in which RF-PCT obtained very competitive prediction results. The main finding of this experimental evaluation are analysed in section 6.

Triple Random ensemble multi-label classification (TREMLC) [143] integrated the ideas of random subspace method, bagging and RAKEL by randomly selecting feature subsets, instance subsets and label subsets respectively to build an ensemble of LP multi-label classifiers. Its performance was competitive but with a high computational cost maybe due to this triple randomization. In [177] the feature space was divided into a number of subsets and a baseline induction algorithm (i.e. AdaBoost.MH described in section 4.3.8) was run over each one. Confidence outputs of AdaBoost.MH were combined by an approach based on the Dempster-Shafer theory [182] outperforming traditional voting schemes.

DECISION TREES	
ML-C4.5 [52]	M2 [147]
PCT [22]	IS-MLT [133]
SUPPORT VECTOR MACHINES	
SVM-HF [87]	BandSVM, ConfMat [87]
ML-PC [156]	SVM-ML [234]
PSVM [225]	SSVM [225]
OVO3C-SVM [221]	OVODL-SVM [122]
Rank-SVM [72]	Calibrated-RankSVM [102]
INSTANCE-BASED	
ML-kNN [256]	IBLR [46]
BRkNN [192]	LPkNN [192]
DML-kNN [244]	KNNMLC [226]
kNN-MLR* [29]	FkNN [17]
FV-kNN [246]	EML-kNN [245]
FSKNN [103]	Mr.kNN [124]
NEURAL NETWORKS	
MMP [56]	BP-MLL [257]
DMLPP [129]	CMLPP [131]
QCMLPP [131]	ML-RBF [250]
PNN [50]	PNN-centroid [49]
ML-FAM, ML-ARAM [176]	MLPP [131]
GENERATIVE AND PROBABILISTIC MODELS	
Multi-label Mixture Model* [137]	PMM1 [215]
PMM2 [215]	EPMM [109]
CoLModel [223]	MadGen [193]
Flat-LDA [175]	Prior-LDA [175]
Dependency-LDA [175]	CML, CMLF [86]
ASSOCIATIVE CLASSIFICATION	
MMAC [204]	RMR [205]
BF-TP, DF-TP [162]	CLAC [219]
BIO-INSPIRED APPROACHES	
MuLAM [38]	G3P-ML [32]
GC [104]	GEP-MLC [14]
ML-2OKM [232]	MoML [185]
EnML [186]	GACC [88]
ENSEMBLES	
AdaBoost.MH [180] [179]	AdaBoost.MR [180] [179]
AdaBoost.MH ^{KR} [181]	AdaBoost.MH with discr. [140]
AdaBoost.MH ^{KR} with discr. [140]	ADTboost.MH [57]
AdaBoost.SZ [66]	AdaBoost.SP [66]
MLBoost [105]	MP-Boost [73]
MSSBoost [237]	ML-RDT [263]
CCA-OC [265]	ML-BCHRF, ML-CRF [115]
MCSP-ECOC [108]	FDT [218]

Table 6 Taxonomy of the MLL adaptation methods. Names with a final * have been given by the authors of the present paper

4.3 Algorithm adaptation methods

4.3.1 Decision trees

Clare and King [52] proposed ML-C4.5, an adaptation of the C4.5 algorithm to the MLL setting. It allowed multiple labels in the leaves and modified the definition of entropy in order to consider not only membership, but also non-membership of each class. This method has become a reference and has been mainly used as base classifier in ensembles of MLL methods (e.g. RFML-C4.5 described in section 4.2).

On the other hand, *Predictive Clustering Tree* (PCT) [22] is framework for prediction that can be instantiated to a particular task by defining a distance metric and a prototype. Thus PCTs have been used for predicting tuples of variables,

time series and even classes organised into a hierarchy (this issue is described in section 5.6). Particularly, in MLL each label is a component of the target tuple. A PCT is top-down generated. At each node, data is partitioned into clusters in such a way that the intra-cluster variation is minimised. The result of the induction process is a decision tree in which each leaf contains the prototype of the instances belonging to that leaf. PCTs have yielded very good classification results combined with random forest (RF-PCTs are described in section 4.2).

The two described approaches are the most widely used, besides other tree-based algorithms have been proposed. First, in [147] the fact that the greedy search of predictors tends to select those with many splits was analysed and the *M2 tree* was introduced. It was a two-stage method that separated the splitting-variable selection (using the statistic test of Nettleton and Banerjee [144]) and the splitting-point selection (that generates binary partitions of data) steps. M2 reduced the bias in predictor selection, but accuracy values were lower when the target vector did not have any significant correlation structure. Finally, it is worth citing *Iterative Split Multi-Label decision Tree* (IS-MLT) [133] that was built following a top-down strategy. At each node, the set of labels was split into two groups by clustering. These two sets become the target objective over which to find the best split with a SVM. The leaf node frequencies were considered scores and labels were assigned by threshold.

4.3.2 Support vector machines

Many approaches have used single-label SVMs with OVA approach [26], [89], [90], [123]. In [87] 2BR (described in section 4.1.2) was called *SVMs with heterogeneous feature kernels* (SVM-HF) as it considered SVM in both base and meta-learning levels. Besides, two mechanisms for improving the margin quality of SVMs in an OVA setting were presented. The first one, *Band-removal method* (BandSVM), operated at the instance level. Once an OVA SVM had been learnt, it removed similar negative examples that were within a threshold distance (band) from the learned decision hyperplane. Then, the model was re-trained. On the other hand, *confusion-matrix based pruning method* (ConfMat) worked at the label level. A confusion matrix, M , over the original learning problem was obtained with any fast, moderately accurate classifier. Each M_{ij} represented the percentage of class i misclassified as class j . Then all training examples of confusing classes were removed and an OVA SVM was learnt. ConfMat was faster to train than BandSVM and required only one OVA SVM step. Experiments on text classification showed that SVM-HF and BandSVM were comparable in their results, being better than ConfMat and OVA SVM.

Other authors have developed OVO decomposition with strategies to treat the special case occurring when samples have double labels at the same time. With this aim, *Multi-Label Paired Comparisons* (ML-PC) [156] separated each pair of overlapping classes by using two probabilistic binary classifiers (e.g. SVM). The individual probabilities of the binary classifiers were combined with the extended Bradley-Terry model with ties [163]. It yielded competitive predictive results but the method to combine predictions was computationally expensive, so it may not be proper for large data sets. Later, in [225], two algorithms known as *Parallel Support Vector Machines* (PSVMs) and *Sequential Support Vector Machines*

(SSVMs) were devised. Double labelled instances were considered as a new independent class, and two parallel hyperplanes were used to separate the three possible classes. The drawback of this strategy was the number of binary classifiers needed for data sets with large number of labels. In [221] the *triple class SVM* (OVO3C-SVM), that was able to deal with positive, negative and mixed classes simultaneously, was proposed. It performed well and the algorithmic complexity was less than PSVMs. Finally, the so-called *Double Label SVM* (OVODL-SVM) [122] considered double labelled instances to be located at marginal region between positive and negative instances by building a double label SVM in which a bias term was needed. The predictive performance was comparable to OVO3C-SVM but faster at train and test procedures.

The algorithm adaptation approach has also been used. Thus Elisseeff and Weston [72] proposed *Rank-SVM*, a ranking method that has become a reference in MLL. It defines a set of q linear classifiers which are optimised to minimise a measure that evaluates the average fraction of label pairs that are reversely ordered for the instance (i.e. the empirical ranking loss defined in Table 1). In addition, it handles non-linear cases with kernel trick. The optimization is carried out under quadratic programming framework in its dual form. This quadratic programming problem has high computational complexity due to a huge number of variables to be solved. In order to obtain a bipartition from the ranking generated, a threshold selection stage, named *set size predictor*, is carried out. In this stage, a proper threshold is determined on the basis of the obtained rank using linear least squares. This stage has no interactive effects with the rank stage, so the threshold selected may be not the optimal. Due to this reason *Calibrated-RankSVM* [102] added a virtual label to determine the relevant labels whose optimal coefficients were determined embedded during the ranking learning process. Nevertheless training procedure is still time consuming. Finally *SVM-ML* [234] was proposed to overcome these drawbacks. It consisted of adding a zero label to detect relevant labels and simplified the original form of Rank-SVM. This led to a novel quadratic programming problem in which each class had an independent equality constraint. It reduced the computational cost in comparison to Rank-SVM and obtained competitive performance.

4.3.3 Instance-based algorithms

Multi-label K-nearest neighbour (ML-kNN) [256] was one of the first lazy MML proposals. After determining the k nearest neighbours, a membership counting vector with the number of neighbours belonging to each possible class was computed. Based on this statistical information, gained from the label sets of the neighbours, the set of labels for the unseen instance was identified by using the *maximum a posteriori* (MAP) principle. The experiments showed that it performed well on several real-world data sets. Nevertheless it is often criticized because it does not take into account label correlations. Therefore, *dependent multi-label k-nearest neighbour* (DML-kNN) [244] was proposed as a version of ML-kNN where dependencies between classes were considered by defining a *global* MAP rule. Unlike ML-kNN, that for each label to be predicted considered only the number of neighbours containing such label, DML-kNN took into account the numbers of all labels in the neighbourhood. The experiments conducted proved the effectiveness of the method as compared to ML-kNN.

Instance Based Learning by Logistic Regression (IBLR) [46] combined instance-based learning (IBL) and logistic regression. It took label correlations into account by using the labels of neighbour examples as extra attributes in a logistic regression scheme. The experimental results proved its performance on several real-world problems. Together with ML-kNN it can be considered the state-of-art in instance-based MLL.

BRkNN [192] is also worth citing. It was equivalent to using BR with kNN as the base classifier, but much faster because, instead of computing q times the k nearest neighbours, it only searched once. Two extensions, dubbed BRkNN-a and BRkNN-b, to tackle the case where BR outputs an empty set for any test instance, were proposed. Defining the confidence of a label as the percentage of the k nearest neighbours that included it, the first method output the label with the highest confidence, and the second one calculated the average size of the label sets of the k nearest neighbours, s , and then output the $[s]$ (nearest integer of s) labels with the highest confidence. The experiments carried out supported that both extensions were beneficial. Besides, authors proposed LPkNN, which consisted of an LP transformation with kNN as the base classifier and also yielded good performance results. In [29] multi-label ranking was described as a special case of rank aggregation (with ties) supplemented with an additional virtual label within a case-based framework.

In the field of image annotation, *kNN multi-label classification* (KNNMLC) [226] was developed with the aim of reducing the bias between semantic similarity of concepts and visual similarity of images annotated with these concepts. It combined a weighted version of kNN and multiple SVM classifiers, used for jointly finding optimal margins in both spaces.

Approximate reasoning techniques have also been incorporated in many proposals. For instance, FkNN [17] adapted the well-known *Fuzzy kNN algorithm* [112]. First, the k nearest neighbours were found and then, a membership degree for each label was computed based on the memberships of the neighbours to each class and on the distances between the test instance and the neighbours. It output a membership degree for each label and the bipartition was obtained applying an α -cut. Other example was *Mr.kNN* [124], a two-step approach. First, a modified fuzzy c-means clustering algorithm, that treated each label as a cluster, assigned each training example a soft relevance value for each label that indicates the strength of an instance related to a label. Then, a kNN algorithm implemented a voting factor based on the soft relevance of each neighbour and on the distances between a test instance and its neighbours. Reported experiments showed it outperformed ML-kNN, nevertheless its computational complexity was higher.

Evidential Multi-Label k-Nearest Neighbour (EML-kNN) [245] used an evidence-theoretic rule that extended the Dempster-Shafer framework [182] to the MLL setting with only a moderate increase in complexity as compared to the classical case. *Fuzzy Veristic k-Nearest Neighbor* (FV-kNN) [246] used a fuzzy kNN rule for MLL based on the theory of *Veristic Variables* [236]. It was able to generate fuzzy label sets for instances that have been originally labelled by crisp ones and obtained competitive results. Finally, FSKNN [103] was based on a *fuzzy similarity measure* (FSM) and kNN. The main aim was reducing the computational power required for finding the neighbours in kNN-like algorithms. First, a FSM grouped the training patterns into clusters. Given a test instance, only those clusters whose similarities to the test instance exceeded a predefined threshold were used to calcu-

late the nearest neighbours. An unseen document was labelled based on its nearest neighbours using the MAP estimate. Experiments in text categorization showed the reduction of the computational costs while the performance was maintained.

4.3.4 Neural networks

Crammer and Singer [56] proposed *Multi-label Multi-class Perceptron* (MMP), a family of on-line algorithms for topic-ranking on text documents. One perceptron was used for each label but, unlike BR, the performance of the whole ensemble was considered to update each perceptron. They showed that MMP outperformed BR on text classification tasks.

The pairwise approach is often regarded as superior to BR because the former profits from simpler decision boundaries in the subproblems. Thus, while in MMP one perceptron was trained for each class, in [131], *multi-label pairwise perceptron* (MLPP) was described as the instantiation of the RPC transformation with perceptrons as base classifiers. It was less efficient than MMP, but it resulted in a gain of accuracy and was able to tackle large text corpora (i.e. RCV1) in a pairwise approach. In order to alleviate the complexity of the RPC approach, quadratic with the number of labels, *dual multi-label pairwise perceptrons algorithm* (DMLPP) [129] formulated the perceptrons in dual form. Thus the prediction time depended linearly on the number of labels. Authors recommend this approach when the number of classes is high. Nevertheless, it is still less efficient than MMP and, as it keeps the whole training set in memory, it has problems to handle training sets with many instances. On the contrary, MLPP is advisable if the number of classes is low and the number of examples high. Instead the RPC approach, CMLPP [131] used a CLR decomposition strategy, and QCMLPP [131] used iteratively the QWeighted voting schema (described in section 4.1.3) until the calibrated label was found. Reported experiments showed that QCMLPP achieved a good trade-off between predictive performance and time complexity.

The algorithm adaptation approach has also been used. Zhang and Zhou [257] developed *Backpropagation for Multilabel Learning* (BP-MLL), an adaptation of the traditional multilayer feed-forward neural network. The key idea is the definition of an error function, closely related to the ranking loss (defined in Table 1). The error function is minimised with gradient descent combined with the error backpropagation. The net has one input unit per input feature, one output unit per label, and the hidden layer is fully connected with weights to the input and output layers. Its computational complexity in the training phase is high, but the time cost of predictions is quite trivial.

Multi Label Radial Basis Function (ML-RBF) [250] was inspired in the well-known RBF method. The first layer was obtained through k-means clustering on instances of each possible class, the centre of each cluster being the prototype vector of a basis function. Weights of the second layer were learnt by minimising a sum-of-squares function. Experimental results showed it outperformed methods as Rank-SVM and BP-MLL in a wide range of metrics and data sets. On the other hand, while the prediction time was similar to BP-MLL, training time was much less in ML-RBF.

It is worth citing works adapting other models of artificial neural networks. *Probabilistic Neural Network* (PNN) [191] has been adapted to the MLL setting in the field of text categorisation [50]. Later, *PNN-centroid* [49] used a technique

of centroids with good performance results that was faster and needed less memory than PNN. Finally ML-FAM and ML-ARAM [176] were two extensions of the neuro-fuzzy models based on *Adaptive Resonance Theory* (ART) that outperformed the single-label proposal.

4.3.5 Generative and probabilistic models

Generative models have been developed, mainly related to text categorisation, under the assumption that a document is generated by a mixture of single-label document models (one per category). In [137] McCallum presented a probabilistic generative model which was based on naive Bayes with *expectation-maximisation* (EM) [64] to learn the mixture weights and the word distributions in each mixture component. The proposed model tried to capture the relationship between the classes and word occurrences, but it did not consider the correlation within the classes.

On the other hand, *Parametric Mixture Models* PMM1 and PMM2 [215] were two probabilistic generative approaches in which documents were modelled by a single multinomial distribution for each class. Experiments showed that PMM obtained better classification performance than the binary decomposition in several text categorization problems. The PMM was a very efficient model but tended to underfit, which led to the development of the *Extended PMM* (EPMM) [109]. It incorporated *latent categories* into PMM so that the model complexity could be adaptively controlled according to the given data. The experiments conducted showed classification performance higher than PMM. *Correlated Labelling Model* (COLMODEL) [223] was proposed to formulate the correlation between different classes. It captured the underlying structures via the latent random variables in a supervised manner. In *Additive-Generative Multi-Label Model* (MAdGen) [193] a deconvolution approach estimated the individual contribution of each label to a given data item and, in [175], a set of three models based on the *latent Dirichlet allocation* (LDA) [21] framework was presented.

Finally, *Conditional Random Fields* (CRFs) [118] have been used in *Collective multi-label* (CML) and *Collective multi-label with features* (CMLF) [86], two multi-label graphical models for text classification that parametrised label co-occurrences. In the field of image categorization CRFs have been used to capture associations between labels in [187] and [227].

4.3.6 Associative classification

Multi-class multi-label associative classification (MMAC) [204] was one of the first algorithms for MLL based on associative classification. The antecedent of a rule was a set of attribute-value pairs in conjunctive form (i.e. an item), and the consequent was a list of ranked class labels. MMAC first applied association rule mining over the training data to discover and generate an initial set of classification rules in which each rule was associated with the most obvious class label. After that, the process was repeated and new rules sets were generated from the remaining unclassified instances, until no more frequent items could be discovered. The rules sets derived at each iteration were merged to form a multi-label classifier. As in different iterations the same antecedent, but with different class labels, might have

been discovered, the algorithm merged such antecedents and ranked the labels according to their frequency in the training patterns satisfying the antecedent. The same authors proposed *Ranked Multi-label Rule* (RMR) [205], similar to MMAC but including a pruning step that removed training objects shared by rules.

Although MMAC tackled with multiple labels, it generated rules by iteratively repeating the adopted method for generating single-label rules. In [162] an algorithm that generated multi-label rules in a single run was proposed. It was based on the *tree-projection-based frequent pattern mining* algorithm [7] and had two possible versions based on *breadth-first* (BF-TP) and *depth-first* (DF-TP) search algorithms.

Finally, *Correlated Lazy Associative classification* (CLAC) [219] was a lazy approach which delayed the inductive process until a test instance, that was used as a filter to remove irrelevant examples from the training data, was given for classification. After that, a greedy heuristic built a specific model for the test instance. This model was composed of class association rules that satisfied thresholds of support and confidence and included labels into the antecedent. By including labels in the antecedent, interactions among labels were explored.

4.3.7 Bio-inspired approaches

Several bio-inspired approaches have been described that built a MLL classifier. One example is *Multi-Label Ant-Miner* (MuLAM) [38], that was based on the *Ant Colony Optimisation* (ACO) Ant-Miner algorithm [154]. Unlike the original Ant-Miner, a pheromone matrix was created for each class, each ant discovered a set of rules (at most one rule for each label) and more than one class in the rule consequent was allowed. Experimental results did not showed significant differences with the single-label proposal. A subsequent version for hierarchical MLL is described in section 5.6. GEP-MLC [14] was another proposal based on *Gene Expression Programming* (GEP) [79]. Each individual codified a discriminant function that was applied to the input features of the pattern to produce a numerical value in such a way that a threshold determined membership to a class. A population of discriminant functions evolved and a niching algorithm was used to guarantee diversity in the solutions and to determine the functions that finally composed the classifier. Later authors proposed GC [104], also based in GEP, but in this case each individual codified a rule, a model more interpretable than a discriminant function. The final classifier consisted of a set of rules and several labels were allowed as consequent. Reported experiments in both works showed competitive results in data sets on different domains. Finally, it is worth citing G3P-ML [32], a Grammar-Guided Genetic Programming algorithm in which a population of classification rules that evolved following a classical generational and elitist evolutionary algorithm.

On the other hand, other approaches have been focused on the optimisation of a population of MML algorithms over several objectives simultaneously. Thus, ML-2OKM [232] was a multi-label kernel algorithm, derived from Rank-SVM (described in section 4.3.2). Two objectives, the model regularization term and the ranking loss, were optimised by using NSGA-II [58], an elitist multi-objective genetic algorithm (MOGA) that provided a Pareto optimal set of solutions to implement ML-2OKM. Unlike ML-2OKM, which optimised two particular objectives,

Multi-Objective Multi-label (MOML) [185] allowed the optimisation of any evaluation metric. Its base model was ML-RBF (described in section 4.3.4) with an additional regularization term added to reduce overfitting risks. The best models were selected and the final prediction was obtained with a majority vote. The reported experiments showed an improvement in performance accuracy, not only limited to the optimisation objectives. *EnML* [186] was another multi-objective evolutionary approach whose aim was to obtain an optimal ensemble consisting of a group of accurate and diverse multi-label base learners (i.e. ML-RBF). Finally, in [88] a genetic algorithm, dubbed GACC, has been proposed to optimise the chain ordering in Classifier Chains classifiers (described in section 4.1.2). Experiments on diverse benchmark data sets indicated the improvement of the output classifiers. Bio-inspired approaches have been also developed in the context of hierarchical multi-label learning. A description of these proposals can be found in section 5.6.

4.3.8 Ensembles

Schapire and Singer proposed *AdaBoost.MH* and *AdaBoost.MR* [179], [180] two boosting algorithms for text categorisation inspired in the popular AdaBoost [82]. The purpose was to find a highly accurate classifier (*final hypothesis*) by combining many classifiers (*weak hypotheses*), each of which might be only moderately accurate. *AdaBoost.MH* tried to minimise the number of misclassified labels (i.e. the Hamming loss described in Table 1), for which it maintained a set of weights not only over the training examples (as AdaBoost does), but also over the labels. Thus, each round, training examples and their corresponding labels that were harder to predict got incrementally higher weights while examples and labels that were easy to predict got lower weights. In practice, it mapped the original multi-label learning problem into a binary learning problem, which was solved by the traditional AdaBoost algorithm with one-level decision trees as base learners. On the other hand, the purpose of *AdaBoost.MR* was minimising the number of labels misorderings (i.e. the ranking loss described in Table 1), so relevant labels would be ranked above the irrelevant ones, for which a weight distribution was maintained over examples and pairs of labels.

Many variants of AdaBoost.MH have been subsequently proposed. For instance, *AdaBoost.MH with K-fold real-valued predictions* (AdaBoost.MH^{KR}) [181]. Its main idea was producing, at each iteration, not a single weak hypothesis, but a complex weak hypothesis consisting of a sub-committee of simple weak hypotheses (committees of decision stumps), which, at that iteration, looked the most promising. AdaBoost.MH and AdaBoost.MH^{KR} required documents to be represented by binary vectors corresponding to the presence or absence of the terms in the document. In order to overcome this drawback, in [140] the potential of weighted (*tf * idf*) representations was studied with AdaBoost-like algorithms by means of two different entropy-based discretization methods. Experiments showed that binary representations obtained after discretization outperformed the original binary representation. Other example is ADTboost.MH [57], that can be viewed as an extension of AdaBoost.MH, that allows a better readability of the classification rules, as well as an extension of ADTBoost [81], that extends the formalism of ADTrees to deal with multi-label problems. Finally, other boosting-type algorithms have been proposed in [66], [73] and [105].

Multi-Label RDT (ML-RDT) [263] was based on *Random Decision Tree* (RDT) [77] and built multiple decision trees by selecting, at each step, a random feature until the number of examples of a node was under a threshold or the depth exceeded a limit. Apart from the label probability distribution counting on each leaf node it did not use any information of labels being independent from the number of labels. This fact reduced the computation time and made it effectively handle large number of labels. On the other hand, *Fast Decision Tree induction* (FDT) [218], employed feature pre-selection and data partitioning to induce a set of C4.5 decision trees from different subsets of the training data. These subsets were non-overlapping, equally-sized and maintained the same ratio between positive and negative classes. Multi-label data was handled by inducing a binary classifier for each class separately. The final output was computed by a strategy called *one-vote* that biased the output toward the minority class by considering the label as positive if at least one of the trees said so. *Model-shared subspace boosting* (MSSBoost) [237] was a boosting-type algorithm that worked at feature and data level. Thus, each model was learnt from random feature subspace and bootstrap data samples. It exploited the label space redundancy by sharing base models across multiple labels.

Error Correcting Output Codes (ECOCs or ECC¹) have been used in multi-class learning due to its capabilities to reduce a multi-class problem to a set of binary-classification problems and also to the improvement in accuracy that their error-correcting properties may provide [67]. The immediate translation to MLL is to consider each combination of labels in the data set as a new one and to apply the ECOC framework in order to build a set of binary-classification problems (e.g. with OVO or OVA coding schemes). This approach was developed in *Multiple Classifier Method for Structured Output Prediction based on Error Correcting Output Codes* (MCSP-ECOC) [108]. Other coding schemes have been studied in the field of MLL, for instance, BCH [115], convolution code [115], canonical correlation analysis [265], repetition code [78], Hamming code [78] and low-density parity-check code [78]. The framework for MLL ECOCs was formalized in [78] and [107]. The main idea was to transform the original multi-label problem into another (larger) multi-label classification task. This was carried out by an ECOC encoder that expanded the original label sets to codewords with redundant information in such a way that provided correcting capabilities. Then a multi-label classifier was learnt considering this augmented label space. During prediction, an ECOC decoder transformed the outputs of the classifier to the original set of labels taking advantage of its error-correcting properties. Experimental results showed the validity of this framework when coupled with classic ECOC coding schemes. Experiments carried out in [107] showed that LP and RAKEL did not benefit from ECOC while BR did. The redundancy is just the main drawback of this approach. The coding of original label information always provides a higher dimensionality and this implies a computational cost. Finally, in [84] multi-class ECOC was interpreted as a way to map a conventional multi-class problem into a multi-label one. This obvious correspondence, and the solution of multi-class problems via MLL, had not yet been noted in the literature.

¹ We will use ECOC instead ECC in order to avoid confusion with Ensemble of Classifier Chains

5 Ongoing research

5.1 Label dependence

Exploring correlations between labels may reduce complexity when data have a moderate or high number of labels. From a probabilistic point of view, two types of dependence in multi-label data have been identified [60], [61]. *Conditional label dependence* captures the dependences between labels given a specific instance, reflecting how likely or unlikely labels are to occur together given the attribute values of a specific instance. This kind of dependence has been explored in [60], [210] and [211]. On the other hand, *unconditional (marginal) label dependence* is independent of a certain instance and refers to the idea that certain labels are likely or unlikely to occur together. This kind of dependence has been explored in [46], [94] and [168].

In [255] a categorisation of multi-label techniques based on the order of correlations was considered. Thus, *first order* approaches ignore label dependences decomposing the problem into a set of independent binary problems (e.g. [26], [53], [57], [258]), *second order* approaches consider the pairwise relations between labels such as their ranking constraint (e.g. [83], [180], [257]) or their co-occurrence patterns (e.g. [86]) and finally, *high-order* approaches consider high-order relations between labels (e.g. [46], [169], [171], [211]).

It is also worth citing that recent studies have pointed to the presence of *asymmetric* [97] and *local* [98] correlations, more consistent with realistic situations. Asymmetric correlations are present when the influence of one label to the other is not necessary the same in the inverse direction (e.g. the label *bear* implies *mammal* but the inverse may not be true). On the other hand, local correlations are shared by subsets of instances rather than all the instances, and exploiting such correlations globally could affect the performance by predicting some irrelevant labels.

Many approaches have already been described in this paper that tackle with label dependence. Nevertheless it is still worth noting other ones. For instance, in the field of automatic video annotation, *Correlative Multi-Label* (CML) [158] simultaneously classified concepts and modelled correlations between them in a single step. Instances were transformed into high-dimensional vectors by encoding correlation information between inputs and outputs, and a maximum-margin type algorithm learned from these transformed vectors. On the other hand, in [270] a method for modelling the correlations between categories based in the principle of *maximum entropy model* (MEM) was presented. According to [159] it could be effective on a set of samples that vary linearly, but it might fail to capture the structure of the feature space if the variations among the samples were non-linear (e.g. image classification). Therefore, it was extended to a non-linear case by incorporating a kernel function into the model [159]. Later, in [101], a general framework to encode class dependences was described. A subspace was assumed to be shared among multiple labels that was computed by solving a generalized eigenvalue problem. *Multi-label Learning by Exploiting lAbel Dependency* (LEAD) [255] was an hybrid generative-discriminative approach, that first built a Bayesian network to encode the conditional dependencies. Then it decomposed the problem into a set of single-label ones considering as additional features its parental labels in the Bayesian network. Its complexity was linear with the number of labels. The

reported experiments over a wide range of data sets showed it was comparable to the state-of-the-art especially on large-scale problems with large number of labels and examples.

Other authors have opted for explicitly represent the dependence structure between the classes via *multi-dimensional Bayesian network classifiers* (MBC). A MBC is a Bayesian network of restricted topology consisting of three subgraphs: one for the class variables, other for the feature variables, and a bridge subgraph that interconnects the class and feature subgraphs allowing only arcs from classes to features. The parameter set defines the conditional probability distribution of each variable given its parents. Several approaches have been recently proposed to learn MBCs [19], [25], [85]. The main drawback of this approach is the high computational cost of determining the optimal network structure and computing the most probable explanation for any instance with unknown values for the classes. Due to this reason in [194], *Bayesian Chain Classifiers* (BCC), combined a first stage, in which Bayesian network with a simple tree-based structure captured the dependency relationships between labels, with a second stage in which a CC (described in section 4.1.2) was built based on the dependence structure. The network constrained the possible chaining orders and reduced the number of classes. This approach yielded competitive results both in complexity and predictive performance.

5.2 Dimensionality reduction

5.2.1 Reduction of the input space (features)

This section tackles *feature selection* and *feature extraction* methods. The former obtains a new feature set which is a subset of the original one and the latter obtains new features by combinations and transformations of the original ones.

FEATURE SELECTION

The *wrapper* approach, in which the subset of features to be used is determined by the learning method, can be directly applicable to MLL by searching for a subset of features optimising a multi-label loss function on an evaluation data set [210]. Another strategy is transforming the data set into one or more single-label data sets and using any existing *filter* method, in which features are scored depending on a measure and the best ones are selected [242]. Thus, in [209] the χ^2 method was separately applied to each label in order to obtain a ranking of all features for each label. Then the top 500 features based on the maximum rank over all labels were selected. In [206] an LP transformation was applied and then a common attribute evaluation statistic (i.e. χ^2) was used, while in [70] a PS transformation with a greedy feature selection algorithm, based on a nearest neighbour estimator of multidimensional mutual information, was applied. Once the features had been ranked, the original multi-label problem was considered again with all the samples. A two-stage filter-wrapper feature selection strategy was described in [253]. It first removed irrelevant and redundant features by using *Principal Component Analysis* (PCA) [106] and then selected the more appropriate subset of features by means of a genetic algorithm whose fitness function addressed the label correlations. Then, multi-label classification was solved using OVA decomposition and binary naive Bayes classifiers with Gaussian density estimation. Finally, *Hybrid Optimisation*

based Multi-Label feature selection (HOML) [183] combined simulated annealing, hill climbing and a genetic algorithm. It reduced the data dimension, improved the classification performance and yielded competitive results compared to benchmark feature selection/feature reduction techniques.

FEATURE EXTRACTION

PCA and *Linear Discriminant Analysis* (LDA) [80] are two common single-label unsupervised techniques for data classification and dimensionality reduction. In [222] a novel method, dubbed *Multi-label Linear Discriminant Analysis* (MLDA), that generalised the single-label framework and surpassed label correlations, was proposed. Another generalised LDA for multi-label problems was studied in [152] concluding the effectiveness of dimension reduction for high dimensional data. In [132] *Self Organising Feature Map* (SOM) and *Latent Semantic Indexing* (LSI) were studied as unsupervised representations of the input space for text classification. The experiments carried out showed that SOM approach obtained better results with multi-label data while LSI performed better with single-label data.

A criticism of unsupervised methods is that they ignore information in labels, thus *Multi-label informed Latent Semantic Indexing* (MLSI) [247] was presented as a multi-label extension of LSI that carried out a mapping from the input features into a new feature space by taking into account not only the information of the inputs, but also capturing the correlations in the label space. Then a set of q linear SVMs was trained on this projected space. *Multi-label Dimensionality reduction via Dependence Maximisation* (MDDM) [266] was another supervised method. Based on the Hilbert-Schmidt independence criterion (HSIC) [91], it performed dimensionality reduction by maximising the dependence between the feature description and the associated class labels. The reported experiments showed it was slightly superior to PCA and significantly superior to MLSI. Finally, *multi-label learning with Label specific FeaTures* (LIFT) [252], considered a different feature set for each label. In a first stage, cluster analysis was conducted for each label considering its positive and negative instances in the training data and new specific features for each label were constructed based on the clustering results. After that, q classifiers were trained taking into account only these specific features. Experiments validated its effectiveness as compared to other well established multi-label learning algorithms.

5.2.2 Reduction of the output space (labels)

Hierarchy Of Multilabel classifiERS (HOMER) [209] is an algorithm whose complexity depends on q . The main idea is the transformation of a MLL problem into a tree hierarchy of simpler MLL tasks, each one dealing with a small number of classes. At each node, labels are split with a balanced clustering algorithm that groups similar labels into a *meta-label*. A multi label classifier is then built in order to predict one or more of these meta-labels. Each child node filters the data of his parent, keeping only the examples that are annotated with at least one of his own labels. Leaves represent one label of the data set. Experimental results showed that HOMER with BR as classifier at the nodes outperformed BR in time and accuracy. Further experimental research has identified HOMER as a computationally efficient MLL method specifically designed for large multi-label data sets [135]. Later, in [212], the use of HOMER with a QCLR classifier (described in section 4.1.3) at each node improved the predictive performance results of HOMER with

BR at a small expense in training time and classification time. It also improved simple BR and QCLR in prediction performance, train and classification time and QCLR also in terms of memory usage.

Label Reduction with Association Rules (LRwAR) [41] applied a different approach. It first run an association rule algorithm over the label space (i.e. FP-Growth [92]) in order to obtain a set of association rules representing the presence of certain labels when others were present. Then, the data set was pre-processed by hiding the inferred labels. After the classification stage, these rules were used to retrieve the relevant labels in a post-processing phase. Results maintained or even improved the evaluation measures of BR, LP, BP-MLL, IBLR and ML-kNN and were obtained in a shorter run time.

Other methods map labels into a reduced label space to reduce computational and space complexities. Examples are *Canonical Correlation Analysis* (CCA) [195], *Compressed Sensing* (ML-CS) [94], *Principal Label Space Transform* (PLST) [198] or *Compressed Labelling* (CL) [267].

Finally, it is worth highlighting that recent research has shifted its focus to problems on large-scale problems where the number of labels is assumed to be extremely large [8], [59]. The key challenge being the design of scalable algorithms that offer real-time predictions, have a small memory footprint and even are able to accommodate missing labels (human annotators tag only with categories they know about).

5.3 MIML

In *Multi-Instance Multi-Label learning* (MIML) a pattern (*bag*) is described by multiple instances and each pattern is associated with a set of labels. For example, an image (bag) can contain multiple regions (instances) and the image can belong to several classes simultaneously. With a MIML representation the relation between the input patterns and the semantic meanings may become more easily discoverable. For example, it can be discovered that one object (bag) has a certain label because it contains a certain instance, and even that the occurrence of several instances triggers other labels [269]. Besides, in many MLL problems, different labels are often tied to the different parts of the object, so, developing classifiers based on the whole object would incur too much noise and harm the performance [93].

In [268] and [269] two approaches based on a simple degeneration strategy for MIML problems were proposed. The first one used MIL as the bridge and obtained a MIL task by applying a category-wise decomposition. After that, any MIL algorithm could be used or the problem could be transformed again to obtain a traditional supervised learning task. Particularly, Zhou and Zhang proposed using MiBOOSTING [235] calling this approach MIMLBOOST. Some authors have noticed that this approach is time-consuming [184]. The second approach used MLL as the bridge. First, the problem was transformed into a MLL task and then any MLL technique could be used. In this case constructive clustering combined all instances in a bag into a single instance. Then a BR approach was used with SVM as base classifier. This approach was called MIMLSVM. Besides authors proved that MIML setting is even useful in tasks where data are represented as single-instance multi-label or as multi-instance single-label examples by using two algorithms,

INstance DIFFerentiation (InsDif) and *SUB-CONcept Discovery* (SubCod), which transformed, respectively, examples in single-instance multi-label or multi-instance single-label format to MIML.

Although these reduction processes are feasible, the performance of the resultant algorithms may suffer from the information loss incurred during the reduction process [251]. Thus, several methods have been proposed to learn from MIML examples directly in order to avoid this loss of information and to model the connections between instances and labels. Zhou et al. proposed D-MIMLSVM [269], a regularisation method. It achieved better performance than degeneration methods, but it could only deal with moderate size of training set due to the associated demanding optimization problem [254]. On the other hand, MIML- k NN [251] not only considered the neighbours of an example, but also those training examples (bags) that counted the example as a neighbour (i.e. citers). In this way, the correlations between instances and labels of an example were exploited. After that, a label counting vector was constructed from its neighbours and citers, and then fed to q trained linear classifiers for prediction. Experimental results showed that MIML- k NN outperformed MIMLBOOST and MIMLSVM in predictive performance. Its training and testing efficiency was slightly worse than MIMLSVM and far superior than MIMLBOOST.

Maximum Margin Method for Multi-Instance Multi-Label (M^3 MIML) [260] was based on a maximum margin method which directly exploited connections between instances and labels. The task was formulated as a quadratic programming problem and implemented in its dual form. Experimental results showed that this algorithm achieved superior performance than MIMLSVM and MIMLBOOST, but the cost of this learning algorithm was quite high [184]. *Multi-instance multi-label radial basis function* (MIMLRBF) [254] used RBF neural networks to learn from MIML patterns and to exploit connections between instances and labels. The input of the first layer was a bag of d -dimensional instances. The first layer consisted of medoids (i.e. bags of instances) formed by performing k -medoids clustering on MIML examples for each possible class. Second layer weights of MIMLRBF network were optimised by minimizing a sum-of-squares error function. Each output unit was related to a possible class label. The reported experiments on two real-world data sets showed it was competitive with MIMLBOOST and MIMLSVM. Finally, it is worth citing other approaches for MIML such as the probabilistic generative model called *Dirichlet-Bernoulli alignment* (DBA) [240] or the *Multi-Instance Multi-Label Gaussian Process* (MIMLGP) [93].

5.4 Semi-supervised and active learning

In many applications data is unlabelled or labelling is expensive or impractical. This fact is even more challenging in MLL. Thus efforts have been also focused in *semi-supervised* (using large amounts of unlabelled data to augment limited labelled data) and *Active Learning* (AL) (the algorithm iteratively asks for labelling examples carefully chosen with the goal of minimizing the labelling effort).

5.4.1 Semi-supervised learning

Semi-supervised methods get benefit of the information provided by unlabelled instances outperforming supervised learning when the number of training data is relatively small and the number of classes is large. In the context of MLL, it is worth citing CNMF [126], in which the key assumption is that two examples tend to have large overlap in their assigned class memberships (class-based similarity) if they share high similarity in their input patterns (input-based similarity). Based in this assumption, CNMF computed two similarity matrices for input patterns and labels. By minimizing the difference of these two matrices, CNMF determined the labels of unlabelled data. The optimization problem was formulated as a *Constrained Non-negative Matrix Factorization*. Experiments on text categorization showed that CNMF had more stable performance than the competing MLSI approach (cited in section 5.2).

While traditional graph-based semi-supervised methods only construct a graph on instance level (in which each node represents one instance and each edge the similarity between corresponding pairwise instances), in SMSE [44] two graphs, on instance (based on both labelled and unlabelled instances) and category level respectively, were constructed. By combining the regularization terms for the two graphs, a regularization framework for multi-label learning was suggested and the labels of unlabelled instances were obtained by solving a *Sylvester Equation* [95]. Experiments on text categorization showed competitive results against binary SVMs, CNMF and MLSI (cited in section 5.2). Zha et al. [249] proposed other graph-based framework for video annotation which employed one loss function and two types of regulariser. One was used to tackle the label consistency on the graph and the other was adopted to tackle the correlations of multiple labels. Based on the proposed framework, two novel graph-based algorithms were developed and experiments showed this framework outperformed key existing graph-based methods and semi-supervised MLL approaches. It is worth citing a different approach, based on *Semi-supervised Impurity based Subspace Clustering Multi-Label* (SISC-ML) [9]. During subspace clustering, labelled and unlabelled examples were used and prototypes maintained the summary about the percentage of each label within each cluster. To obtain predictions, a kNN approach was used considering k nearest neighbour clusters. It was applied to text categorization and performed well even when a very limited amount of labelled training data was available.

Most of the works on semi-supervised multi-label learning work under the transductive setting. Nonetheless, *inductive Multi-Label Classification with Unlabelled data* (iMLCU) [229] is one recent work which tackles semi-supervised multi-label learning under the inductive setting. The inductive semi-supervised MLL is formulated as an optimization problem of learning q linear models, which fits labeled data by exploiting pairwise label correlations and uses unlabeled data via appropriate regularizations. After that the resulting optimization, which is non-convex, is solved via the *ConCave Convex Procedure* (CCCP) [40].

5.4.2 Active-learning

The key in active learning is the sample selection strategy whose aim is choosing the most informative instance to obtain the best classification performance. Bin-Min [27] selected unlabelled examples with respect to the most uncertain label

and OVA was used for multi-label classification with SVM as the base classifier. This method did not take advantages of the multi-label information. In the field of image classification, it was proposed the *Mean Max Loss* (MML) or 1DAL strategy [123], that selected the unlabelled instance which had the maximum mean loss value over the predicted classes. One SVM was trained for each label and a threshold cutting method decided the relevant labels. The overall loss value was averaged over the labels. As this strategy selected only along the sample dimension, it did not take advantages of the label correlations to reduce human labelling cost. Besides, when one sample was selected, all its labels had to be labelled. 1DAL was improved in [160], where *2 Dimensional Active Learning* (2DAL) was described. It considered both relationships between samples and between labels. Sample-label pairs were selected in order to minimise the multi-label Bayesian error bound. This allowed the annotation of a subset of labels and the inference of the rest of labels was performed from labels correlations with expectation-maximisation. It improved 1DAL and random selection in image classification tasks and was later extended in [262] to the multi-view learning framework (described in section 2.4). By taking advantage of both active learning and multi-view learning, the annotation effort was reduced in comparison with random selection, 1DAL and 2DAL. As multi-view learning and active learning can be effectively integrated, this may be a line to be explored in the MLL setting.

An approach for text classification, so-called MMC, was proposed in [239]. One SVM was trained for each label and the overall loss of the classifier was measured by gathering the loss of all binary classifiers. Instead of estimating the labels for each instance, the number of labels was estimated by applying *Logistic Regression* (LR). The training features of the LR model were the probabilities obtained by the binary classifiers and the number of labels was the categorical target to be predicted. It outperformed random selection, 1DAL and BinMin in the domain of text classification and reduced significantly the labelling cost. Finally in [75] several strategies to carry out a global labelling in text classification, in which a unique ranking of unlabelled patterns combined the outputs of q individual binary classifiers, have been proposed.

5.5 On-line MLL and data streams

Many challenging real world problems involve multi-label data streams and learning in such scenarios has special requirements: i) one example is processed at a time and only once, ii) there are limitations of memory and time, iii) data distribution evolves producing concept drift, and also iv) the model must be ready to predict at any time. Very few authors have explored this task in a MLL setting. In [161] the incoming data was partitioned into chunks and, to take advantage of label correlations, a SVM-HF classifier (described in section 4.1.2) was built for each chunk. Concept drift was managed by keeping only the latest trained classifiers and discarding the oldest ones. The empirical results showed that it outperformed partitioning data into chunks with BR, that does not consider label correlations, as classifier.

Later, Read et al. [168] discussed the use of BR, LP, CC and PS transformation methods in evolving scenarios by instantiating incremental base classifiers. Besides, an on-line *Multi Label Hoeffding Tree* that used the Hoeffding bound [69]

to determine the number of instances needed to split a node was described. The tree used the multi-label definition of entropy proposed by Care & King [52] and applied a PS classifier at the leaves. It achieved faster and more accurate performance than the cited transformation methods. Besides, authors used *ADWIN Bagging* [20] to deal with concept drift. In [170] a framework to generate synthetic multi-label data streams can be found. Its aim is facilitating the study and evaluation of MLL algorithms in a data stream scenario.

Finally, it is worth citing *Multiple Windows* (MW) [231], which used a double window mechanism for each label (one for positive examples and one for negative examples) to deal with the fact that labels do not drift simultaneously and with the same rate (multiple concept drift). Besides, this mechanism was able to deal with the imbalance problem by oversampling the positive examples and undersampling the negative ones according to a user defined parameter. It outperformed SVM-HF.

5.6 Hierarchical multi-label classification

In *hierarchical multi-label classification*, in contrast to *flat classification*, examples can be associated with multiple labels and labels are organised in a hierarchical structure such as a tree or a *directed acyclic graph* (DAG), which allows a child category to have more than one parent category. This entails several challenges. First, classes in the bottom of the hierarchy tend to be more difficult to identify because the number of samples is usually less than in upper classes. Secondly, the closer a category is to the root, the more a wrong decision affects lower levels. And finally, predictions must respect the class hierarchy. Thus, according to the *true path rule* (TPR), borrowed from the Gene Ontology and FunCat taxonomies [55], an example that belongs to some class automatically belongs to all its ancestors, and negative predictions for a given node are propagated to the descendants to preserve the consistency of the hierarchy. Typical examples of hierarchical domains are protein function prediction and text categorization.

In [34] and [113] two approaches for hierarchical multi-label classification are described: the *local* one (top-down) consists of training a hierarchy of classifiers (e.g. SVM or decision trees), which are used in a top-down fashion for the classification of new examples; and the *global* one (one-shot, big-bang), that induces a unique classifier using all classes of the hierarchy at once and is able to predict just in one step. On the other hand, three baseline approaches were identified in [220], the two first corresponded to local methods and the last one corresponded to global methods: *Single-label Classification approach* (SC), *Hierarchical Single-label Classification* (HSC) and *Hierarchical Multi-label Classification* (HMC).

SC trains a binary classifier for each label in the hierarchy considering as positive examples those labelled with such class and the rest are considered to be negative. This approach has several drawbacks. First, it needs to train one classifier per class (which can be hundreds or thousands). Besides, as the hierarchy is not taken into account, it is also possible having inconsistencies. Finally, it is very likely having a problem of imbalanced data at lower levels of the hierarchy with only a few positive patterns and too many negative ones.

HSC, also called *hierarchical binary relevance* (HBR) [210], consists of adapting SC by considering the hierarchy during the prediction in such a way that a classifier only predicts positive if the classifier for the parent class also makes a

positive prediction. The hierarchy can also be considered during the training step by restricting the training set of a classifier to those instances belonging to the parent class. HSC has been followed in [15] and [18].

Finally, HMC consists of generating a global model which is able to predict the class associated with a pattern at any level of the hierarchy. This approach has been followed in [23], [51] and [113].

In the group of local methods, it can be cited *Bayes-optimal classifier* (HBAYES) [35]. It followed the HBR approach but the output for a given pattern was computed by a bottom-up process. Later, HBAYES-CS [37] tackled the problem of sparsity of annotations by a cost sensitive parameter to control the trade-off between precision and recall. According to the authors, HBAYES-CS resulted more suitable for dealing with skewed data sets. Esuli et al. presented *TreeBoost.MH* [74], a recursive algorithm that generated an AdaBoost.MH classifier for each non-root category. As multi-label (instead of binary) classifiers were built, it is considered a generalisation of HBR [210]. In [145] three local approaches for text classification with DAG categories were proposed: *flat*, *tree-based* and *DAG-based* that, respectively, transformed the DAG into a flat structure (in which categories were treated in isolation), an equivalent tree or generated one classifier per parent when the node had several parents. SVM were used as base classifiers. The results showed that the flat approach had a comparable performance to the hierarchical approaches when the number of categories involved was small and tree-based and DAG-based approaches had nearly the same classification accuracy, but the former tended to produce larger trees. To finish with local methods, in [36], the convenience of combining techniques of hierarchical classification (to take into account relationships between classes), data fusion techniques (to integrate multiple sources of data) and cost-sensitive methods (to address the imbalance between positive and negative examples) has been studied. It was found that the combined action of some of these techniques achieved better performance than the average of the performances of the strategies used separately.

Regarding to global methods, Clare [53] proposed one based on the ML-C4.5 described in section 4.3.1. It modified the definition of entropy to take into account both the multi-label aspect and the hierarchical relationship between labels. In [113] a global approach applied to a DAG hierarchy was presented. The main idea was transforming an initial (possibly single-label) task into a multi-label task by expanding the label set of each training example with the corresponding ancestor labels. After that AdaBoost.MH algorithm was applied and, finally, inconsistently classified instances were re-labelled. Experiments carried out in annotation of gene functions demonstrated that the approach improved the flat AdaBoost.MH (i.e. AdaBoost without considering any hierarchical information) and was comparable to local AdaBoost.MH. Clus-HMC [23], based on the PCT algorithm described in section 4.3.1, is other global approach, which trains only one decision-tree to cope with the entire classification problem. Experiments concluded that it was fast, identified features relevant for all the labels together and performed similar to Clus-SC (that learnt a separate PCT for each class). Clus-SC, Clus-HSC and Clus-HMC were later compared in [220] finding that Clus-HMC performed better for tree and DAG class hierarchies.

HMC has also been tackled with bio-inspired approaches. Thus, *Hierarchical Classification Ant-Miner* (*hmAnt-Miner*) [148] was based on ACO and discovered a global classification model in the form of an ordered list of IF-THEN rules. The

construction of a rule was divided into two ant colonies which worked in a cooperative manner, one for the rule antecedents and the other for the rule consequents. At each iteration, a rule was built by the pairing of an antecedent ant with a consequent ant. It operated both with tree or DAG hierarchies and the results obtained were competitive in terms of accuracy and complexity of the model. *Artificial immune systems* (AIS) were used in [11] and [12] defining a local and a global version of an algorithm called *Multi-label Hierarchical classification with AIS* (HMC-AIS) to discover classification rules for protein function prediction. A sequential covering procedure iteratively called a rule evolution procedure, that evolved classification rules (antibodies) and added the best one to the set of discovered rules, until all (or almost all) training examples (antigens) were covered by the discovered rules. *Hierarchical Multi-Label Classification with Genetic Algorithm* (HMC-GA) [34] was another bio-inspired approach. It evolved the antecedents of decision rules with a sequential covering strategy, removing from the training set examples already covered by the generated rules. The fitness function was biased towards rules with high example coverage. Experiments showed that *hmAnt-Miner*, *Clus-HMC* and *HMC-GA* produced considerably less rules than the local methods (i.e. *Clus-SC* and *Clus-HSC*), resulting in much simpler and interpretable final models.

Finally an innovative approach which was able to obtain not only the classifier, but also the hierarchy from the multi-label prediction was described in [30]. It is more complex than standard HMC where the class hierarchy is known a priori.

5.7 Dealing with class imbalance

It is commonly accepted that multi-label data may suffer from imbalance. Thus, *label skew* [165] is defined as a relatively high number of examples associated with the most common label sets, while a relatively high number of examples are associated with infrequent label sets. When each label is considered separately, label skew becomes *class imbalance*. In this case, not only may some labels be more frequent than others (inter class), but a strong imbalance between positive and negative examples for each label (inner or intra class) may also occur [231]. Measures regarding imbalance in MLL are found in [42], [65]. The proposals developed till now are sparse and have been focused on algorithmic adaptations of MLL algorithms, the use of ensembles of classifiers or preprocessing methods [42], [65], [197].

6 Pitfalls and guidelines

Given the high number of available algorithms, selecting the most suitable set of learners to be applied to a given data set is an important issue. This section intends to give some guidelines in the light of the results reported in literature.

Despite that some empirical evaluations have been carried out in concrete domains as image annotation [142], video annotation [68] or within a concrete family of algorithms [192], some authors recently suggested that the development of empirical comparisons needs to be more explored [261]. In our opinion, these comparisons should take into account a wide range of data sets, algorithms and

performance metrics, conduct statistical tests to determine significant differences between proposals and consider execution time in training and test stages.

In [135] an extensive experimental comparison with statistical proofs was conducted on 12 state-of-the-art MLL algorithms, 16 evaluation metrics and 11 benchmark data sets. The overall conclusion was that the best performing methods were RF-PCT, HOMER, BR, and CC and thus, they were recommended as good methods for MLL and as benchmarks for testing new algorithms. Particularly, for each measure, the best algorithm was RF-PCT followed by HOMER. This last method was specially recommended for large data sets [135], [209]. ML-kNN performed poor across all evaluation measures, but in spite of its limitations, it resulted more efficient than ECC for large training sets. More specifically, in example-based and label-based measures RF-PCT was best according to precision (the prediction was more exact) and HOMER was better in recall and poor in precision (the prediction was more complete). With ranking-based measures RF-PCT was the best followed by CC and BR while HOMER performance was poor. Regarding base classifiers, SVMs and trees were used. Experiments lead to think that SVMs work better in domains with large number of features as text classification, typically $d > 500$, and small number of patterns, while trees perform better in domains with large number of examples. This may be due to the fact that SMVs are able to exploit the information of all features while trees only exploit a subset.

In [171] the experiments carried out yielded interesting conclusions about complexity and limits of the algorithms that can help researchers to choose the proper method. First, CC and BR time complexity is similar (up to $q > 128$) and overall, were considered the best candidates for very large problems. RAKEL run out of memory when $q > 256$. CLR scaled well with respect to the number of patterns but not with respect to the number of labels (it became intractable for $q > 64$) and vice versa for IBLR. CC, BR and CLR were able to complete the task for 819200 patterns.

Chekina et al. [43] postulated that additional information (besides the target evaluation measure) was needed to select the appropriate MML method and contemplated a meta-learning approach based on 49 descriptive characteristics of the data sets (e.g. number of instances, number of labels, number of distinct labels in the data set, etc.). A set of 10 algorithms, 12 data sets and 18 evaluation measures were used. Almost in two-thirds of the cases the meta-learner was able to predict which multi-label classification method outperformed other methods. Besides, in the light of experiments, authors suggested a set of meta-features of data sets that may determine the performance of the MLL methods (e.g. number of labels, average examples per class etc.).

In [61] and [63] conditional and unconditional dependences were studied (see definitions in section 5.1) obtaining interesting findings that are listed below:

- Most MLL algorithms learn by explicitly or implicitly optimise a specific metric (the main metrics are summarized in Table 1). The same MLL method rarely will be optimal for different types of losses. Depending on the metric being minimised conditional dependence treatment can improve or not the performance.
- Minimizing the subset 0/1 loss requires modelling dependences between labels. A classifier that minimises the subset 0/1 loss may cause low values for the Hamming loss and vice versa.

- The conditional dependence is more related to non-decomposable losses (e.g. subset 0/1 loss) than to decomposable ones (e.g. Hamming loss).

A general conclusion is that algorithms perform different according to the data sets and evaluation metrics, and the algorithm to use will depend on the needs of the problem. A guideline could be decision trees for efficiency, ensembles for predictive performance and transformation methods for the flexibility of using any single-label classifier [166]. Besides, other features as scalability (in patterns, features or labels) and interpretability of the model could be taken into account.

7 Concluding remarks

This paper has carried out a review of the state-of-the-art in MLL and ongoing research. The descriptions of the multi-label framework and the main areas of application have provided us with the background needed to understand the works reviewed. The review shows that MLL has been successfully applied to fields such as text, image and video annotation, detection of emotions in music, medical diagnosis, gene and protein function prediction and even new areas of application are arising (e.g. speech emotion recognition or social network mining). A bibliometric study has revealed the increasing interest and number of papers that have been published in this field. The more accepted taxonomy of MLL methods has been used to categorise, sort and describe the relevant literature showing that many of the main traditional single-label classification models have been adapted to the MLL framework. Guidelines and pitfalls to choose the proper method have also been provided. The key challenge and open issue has been recently identified as dealing with the high dimensionality of the output space, especially in domains with a large number of labels. This challenge involves exploring label correlations efficiently. The study of recent work also reveals that there exist many challenges that could merit further attention in the future, such as dimensionality reduction, MIML learning, semi-supervised and active learning, structured prediction or managing imbalanced data.

Acknowledgements This work is supported by the Ministry of Science and Technology project TIN-2011-22408.

References

1. Lamda. learning and mining from data. data & code. URL <http://lamda.nju.edu.cn/Data.ashx>
2. 1st International Workshop on Learning from Multi-Label Data (MLD'09). <http://lps.csd.auth.gr/workshops/mld09/mld09.pdf> (2009)
3. 2nd International Workshop on Learning from Multi-Label Data (MLD'10). <http://cse.seu.edu.cn/conf/MLD10/files/MLD'10.pdf> (2010)
4. Machine Learning. Special Issue on Learning from Multi-Label Data 88(1-2) (2012)
5. Extreme Classification: Multi-Class & Multi-Label Learning with Millions of Categories. <http://nips.cc/Conferences/2013/Program/event.php?ID=3707> (2013)
6. Abbas, Q., Celebi, M., Serrano, C., García, I.F., Ma, G.: Pattern classification of dermoscopy images: A perceptually uniform model. *Pattern Recognition* **46**(1), 86 – 97 (2013)
7. Agarwal, R., Aggarwal, C., Prasad, V.: A tree projection algorithm for generation of frequent item sets. *Journal of Parallel and Distributed Computing* **61**(3), 350–371 (2001)

8. Agrawal, R., Gupta, A., Prabhu, Y., Varma, M.: Multi-label Learning with Millions of Labels: Recommending Advertiser Bid Phrases for Web Pages. In: Proceedings of the 22nd International Conference on World Wide Web (WWW13), pp. 13–24 (2013)
9. Ahmed, M.S., Khan, L., Oza, N.C., Rajeswari, M.: Multi-label ASRS Dataset Classification Using Semi Supervised Subspace Clustering. In: Proceedings of the 2010 Conference on Intelligent Data Understanding, (CIDU), pp. 285–299 (2010)
10. Aiolli, F., Cardin, R., Sebastiani, F., Sperduti, A.: Preferential text classification: learning algorithms and evaluation measures. *Information Retrieval* **12**(5), 559–580 (2009)
11. Alves, R.T., Delgado, M.R., Freitas, A.A.: Multi-label hierarchical classification of protein functions with artificial immune systems. In: Proc. Brazilian Symposium in Bioinformatics (BSB-2008), *Lecture Notes in Bioinformatics*, vol. 5167, pp. 1–12 (2008)
12. Alves, R.T., Delgado, M.R., Freitas, A.A.: Knowledge discovery with Artificial Immune Systems for hierarchical multi-label classification of protein functions. In: IEEE International Conference on Fuzzy Systems (FUZZ-IEEE), pp. 1–8. Barcelona (2010)
13. Antenreiter, M., Ortner, R., Auer, P.: Combining Classifiers for Improved Multilabel Image Classification. In: Proceedings of the 1st workshop on learning from multilabel data (MLD) held in conjunction with ECML/PKDD, pp. 16–27. Bled, Slovenia (2009)
14. Ávila, J.L., Gibaja, E.L., Zafra, A., Ventura, S.: A Gene Expression Programming Algorithm for Multi-Label Classification. *Journal of Multiple-Valued Logic and Soft Computing* **17**, 183–206 (2011)
15. Barutcuoglu, Z., Schapire, R.E., Troyanskaya, O.G.: Hierarchical multi-label prediction of gene function. *Bioinformatics* **22**(7), 830–836 (2006)
16. Bhowmick, P.K., Basu, A., Mitra, P.: Reader Perspective Emotion Analysis in Text through Ensemble based Multi-Label Classification Framework. *Computer and Information Science* **2**(4), 64–74 (2009)
17. Bhowmick, P.K., Basu, A., Mitra, P., Prasad, A.: Sentence Level News Emotion Analysis in Fuzzy Multi-label Classification Framework. *Research in Computer Science*, special issue: Natural Language Processing and its Applications **46**, 143–154 (2010)
18. Bianchi, N.C., Gentile, C., Zaniboni, L.: Incremental Algorithms for Hierarchical Classification. *J. Mach. Learn. Res.* **7**, 31–54 (2006)
19. Bielza, C., Li, G., Larrañaga, P.: Multi-dimensional classification with Bayesian networks. *Int. J. Approx. Reasoning* **52**(6), 705–727 (2011)
20. Bifet, A., Holmes, G., Pfahringer, B., Kirkby, R., Gavaldà, R.: New ensemble methods for evolving data streams. In: 15th ACM SIGKDD international conference on knowledge discovery and data mining, pp. 139–148 (2009)
21. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent Dirichlet Allocation. *Journal of Machine Learning Research* **3** (2003)
22. Blockeel, H., Raedt, L.D., Ramon, J.: Top-Down Induction of Clustering Trees. In: Proceedings of the Fifteenth International Conference on Machine Learning (ICML '98), pp. 55–63. San Francisco, CA, USA (1998)
23. Blockeel, H., Schietgat, L., Struyf, J., Dzēroski, S., Clare, A.: Decision trees for hierarchical multilabel classification: A case study in functional genomics. In: 10th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD), *Lecture Notes in Computer Science*, vol. 4213, pp. 18–29 (2006)
24. Blum, A., Mitchell, T.: Combining labeled and unlabeled data with co-training. In: 11th Annual Conference on Computational Learning Theory, pp. 92–100. Madison, WI. (1998)
25. Borchani, H., Bielza, C., Toro, C., Larrañaga, P.: Predicting human immunodeficiency virus inhibitors using multi-dimensional bayesian network classifiers. *Artificial Intelligence in Medicine* **57**(3), 219–229 (2013)
26. Boutell, M., Luo, J., Shen, X., Brown, C.: Learning multi-label scene classification. *Pattern Recognition* **37**(9), 1757–1771 (2004)
27. Brinker, K.: On Active Learning in Multi-label Classification. In: From Data and Information Analysis to Knowledge Engineering, Studies in Classification, Data Analysis, and Knowledge Organization, pp. 206–213. Springer Berlin Heidelberg (2006)
28. Brinker, K., Fürnkranz, J., Hüllermeier, E.: A Unified Model for Multilabel Classification and Ranking. In: Proceeding of the ECAI 2006: 17th European Conference on Artificial Intelligence, pp. 489–493 (2006)
29. Brinker, K., Hüllermeier, E.: Case-based multilabel ranking. In: Proceedings of the 20th international joint conference on Artificial intelligence (IJCAI'07), pp. 702–707. San Francisco, CA, USA (2007)

30. Brucker, F., Benites, F., Sapozhnikova, E.: Multi-label classification and extracting predicted class hierarchies. *Pattern Recognition* **44**(3), 724–738 (2010)
31. Bucak, S., Jin, R., Jain, A.: Multi-label learning with incomplete class assignments. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2801–2808 (2011)
32. Cano, A., Zafra, A., Galindo, E.L.G., Ventura, S.: A grammar-guided genetic programming algorithm for multi-label classification. In: *16th European Conference, EuroGP, Lecture Notes in Computer Science*, vol. 7831, pp. 217–228 (2013)
33. de Carvalho, A., Freitas, A.: A Tutorial on Multi-label Classification Techniques. In: *Foundations of Computational Intelligence*, vol. 5, pp. 177–195. Springer Berlin / Heidelberg (2009)
34. Cerri, R., Barros, R.C., de Carvalho, A.C.P.L.F.: A genetic algorithm for Hierarchical Multi-Label Classification. In: *Proceedings of the 27th Annual ACM Symposium on Applied Computing (SAC '12)*, pp. 250–255. New York, NY, USA (2012)
35. Cesa-Bianchi, N., Gentile, C., Zaniboni, L.: Hierarchical classification: combining Bayes with SVM. In: *Proceedings of the Twenty-Third International Conference on Machine Learning (ICML 2006)*, pp. 177–184 (2006)
36. Cesa-Bianchi, N., Re, M., Valentini, G.: Synergy of multi-label hierarchical ensembles, data fusion, and cost-sensitive methods for gene functional inference. *Machine Learning* **88**, 209–241 (2012)
37. Cesa-Bianchi, N., Valentini, G.: Hierarchical Cost-Sensitive Algorithms for Genome-Wide Gene Function Prediction. *Journal of Machine Learning Research - Proceedings Track* **8**, 14–29 (2010)
38. Chan, A., Freitas, A.A.: A new ant colony algorithm for multi-label classification with applications in bioinformatics. In: *GECCO '06: Proceedings of the 8th annual conference on Genetic and evolutionary computation*, pp. 27–34. New York, USA (2006)
39. Chang, C.C., Lin, C.J.: LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* **2**, 27:1–27:27 (2011). Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
40. Chapelle, O., Sindhwani, V., Keerthi, S.S.: Optimization techniques for semisupervised support vector machines. *Journal of Machine Learning Research* **9**, 203–233 (2008)
41. Charte, F., Rivera, A., del Jesus, M., Herrera, F.: Improving Multi-label Classifiers via Label Reduction with Association Rules. In: *Hybrid Artificial Intelligent Systems, Lecture Notes in Computer Science*, vol. 7209, pp. 188–199. Springer Berlin / Heidelberg (2012)
42. Charte, F., Rivera, A., del Jesus, M., Herrera, F.: A First Approach to Deal with Imbalance in Multi-label Datasets. In: *HAIS 2013 - LNAI 8073*, pp. 150–160 (2013)
43. Chekina, L., Rokach, L., Shapira, B.: Meta-learning for Selecting a Multi-label Classification Algorithm. In: *IEEE 11th International Conference on Data Mining Workshops (ICDMW)*, pp. 220–227 (2011)
44. Chen, G., Song, Y., Wang, F., Zhang, C.: Semi-supervised Multi-label Learning by Solving a Sylvester Equation. In: *Proceedings of the SIAM International Conference on Data Mining (SDM)*, pp. 410–419 (2008)
45. Cheng, W., Dembczynski, K., Hüllermeier, E.: Graded Multilabel Classification: The Ordinal Case. In: *Proceedings of the 27 th International Conference on Machine Learning (ICML)*, pp. 223–230 (2010)
46. Cheng, W., Hüllermeier, E.: Combining instance-based learning and logistic regression for multilabel classification. *Mach. Learn.* **76**, 211–225 (2009)
47. Cherman, E.A., Metz, J., Monard, M.C.: Incorporating label dependency into the binary relevance framework for multi-label classification. *Expert Syst. Appl.* **39**(2), 1647–1655 (2012)
48. Chou, K.C., Wu, Z.C., Xiao, X.: iLoc-Euk: A Multi-Label Classifier for Predicting the Subcellular Localization of Singleplex and Multiplex Eukaryotic Proteins. *PLoS ONE* **6**(3) (2011)
49. Ciarelli, P., Oliveira, E.: An Enhanced Probabilistic Neural Network Approach Applied to Text Classification. In: *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications, Lecture Notes in Computer Science*, vol. 5856, chap. 78, pp. 661–668. Springer Berlin / Heidelberg (2009)
50. Ciarelli, P.M., Oliveira, E., Badue, C., Souza, A.F.: Multi-Label Text Categorization Using a Probabilistic Neural Network. *International Journal of Computer Information Systems and Industrial Management Applications (IJCISIM)* **1**, 133–144 (2009)

51. Clare, A.: Machine learning and data mining for yeast functional genomics. Ph.D. thesis, University of Wales (2003)
52. Clare, A., King, R.D.: Knowledge Discovery in Multi-label Phenotype Data. In: PKDD '01 Proceedings of the 5th European Conference on Principles of Data Mining and Knowledge Discovery, *Lecture Notes in Computer Science*, vol. 2168, pp. 42 – 53 (2001)
53. Clare, A., King, R.D.: Predicting gene function in *Saccharomyces cerevisiae*. *Bioinformatics* **2**, 42–49 (2003)
54. Cong, H., Tong, L.H.: Grouping of TRIZ Inventive Principles to facilitate automatic patent classification. *Expert Systems with Applications* **34**(1), 788–795 (2008)
55. Consortium, T.G.O.: Gene ontology: tool for the unification of biology. *Nature Genet* **25**, 25–29 (2000)
56. Crammer, K., Singer, Y.: A family of additive online algorithms for category ranking. *J. Mach. Learn. Res.* **3**, 1025–1058 (2003)
57. De Comit  , F., Gilleron, R., Tommasi, M.: Learning multi-label alternating decision trees from texts and data. In: Proceedings of the 3rd international conference on Machine learning and data mining in pattern recognition (MLDM'03), pp. 35–49. Berlin, Heidelberg (2003)
58. Deb, K., Pratap, A., Agarwal, S., Meyarivan, T.: A fast and elitist multiobjective genetic algorithm: NSGA-II. *Evolutionary Computation, IEEE Transactions on* **6**(2), 182 –197 (2002)
59. Dekel, O., Shamir, O.: Multiclass-multilabel classification with more labels than examples. In: Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS), pp. 137–144 (2010)
60. Dembczyński, K., Cheng, W., H  llermeier, E.: Bayes Optimal Multilabel Classification via Probabilistic Classifier Chains. In: Proceedings of the 27 th International Conference on Machine Learning (ICML), pp. 279–286 (2010)
61. Dembczyński, K., Waegeman, W., Cheng, W., H  llermeier, E.: On Label Dependence in Multi-Label Classification. In: Proceedings of the 2nd International Workshop on Learning from Multi-Label Data (MLD'10), pp. 5–12 (2010)
62. Dembczyński, K., Waegeman, W., Cheng, W., H  llermeier, E.: Regret Analysis for Performance Metrics in Multi-Label Classification: The Case of Hamming and Subset Zero-One Loss. In: Machine Learning and Knowledge Discovery in Databases, *Lecture Notes in Computer Science*, vol. 6321, pp. 280–295. Springer Berlin / Heidelberg (2010)
63. Dembczyński, K., Waegeman, W., Cheng, W., H  llermeier, E.: On label dependence and loss minimization in multi-label classification. *Machine Learning* **88**, 5–45 (2012)
64. Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistics Society -B* **39**(1), 1–38 (1977)
65. Dendamrongvit, S., Kubat, M.: Undersampling Approach for Imbalanced Training Sets and Induction from Multi-label Text-Categorization Domains. In: New Frontiers in Applied Data Mining, *LNCs*, vol. 5669, pp. 40–52. Springer Berlin / Heidelberg (2010)
66. Diao, L., Hu, K., Lu, Y., Shi, C.: Boosting simple decision trees with Bayesian learning for text categorization. In: Proceedings of the 4th World Congress on Intelligent and Automation, pp. 321–325. Shanghai, P.R. China (2002)
67. Dietterich, T.G., Bakiri, G.: Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research* **2**, 263–286 (1995)
68. Dimou, A., Tsoumakas, G., Mezaris, V., Kompatsiaris, I., Vlahavas, I.: An Empirical Study of Multi-label Learning Methods for Video Annotation. In: International Workshop on Content-Based Multimedia Indexing, pp. 19–24. IEEE Computer Society, Los Alamitos, CA, USA (2009)
69. Domingos, P., Hulten, G.: Mining high-speed data streams. In: Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '00), pp. 71–80. New York, NY, USA (2000)
70. Doquire, G., Verleysen, M.: Feature selection for multi-label classification problems. In: 11th International Work-Conference on on Artificial Neural Networks (IWANN), *Lecture Notes in Computer Science*, vol. 6691, pp. 9–16 (2011)
71. Duwairi, R., Kassawneh, A.: A framework for predicting proteins 3D structures. In: IEEE/ACS International Conference on Computer Systems and Applications (AICCSA '08), pp. 37 –44. Washington, DC, USA (2008)
72. Elisseeff, A., Weston, J.: A kernel method for multi-labelled classification. In: Advances in Neural Information Processing Systems (NIPS), vol. 14, pp. 681–687 (2001)

73. Esuli, A., Fagni, T., Sebastiani, F.: MP-Boost: A Multiple-Pivot Boosting Algorithm and Its Application to Text Categorization. In: String Processing and Information Retrieval (SPIRE), *Lecture Notes in Computer Science*, vol. 4209, pp. 1–12. Springer Berlin / Heidelberg (2006)
74. Esuli, A., Fagni, T., Sebastiani, F.: TreeBoost.MH: A Boosting Algorithm for Multi-label Hierarchical Text Categorization. In: String Processing and Information Retrieval (SPIRE), *Lecture Notes in Computer Science*, vol. 4209, pp. 13–24. Springer Berlin / Heidelberg (2006)
75. Esuli, A., Sebastiani, F.: Active Learning Strategies for Multi-Label Text Classification. In: Advances in Information Retrieval, *Lecture Notes in Computer Science*, vol. 5478, pp. 102–113. Springer Berlin / Heidelberg (2009)
76. Fan, R.E., Lin, C.J.: A Study on Threshold Selection for Multi-label Classification. Tech. rep., National Taiwan University (2007)
77. Fan, W., Wang, H., Yu, P.S., Ma, S.: Is random model better? On its accuracy and efficiency. In: Proceedings of the Third IEEE International Conference on Data Mining (ICDM '03) (2003)
78. Ferng, C.S., Lin, H.T.: Multi-label Classification with Error-correcting Codes. *Journal of Machine Learning Research - Proceedings Track* **20**, 281–295 (2011)
79. Ferreira, C.: Gene expression programming: A new adaptive algorithm for solving problems. *Complex Systems* **13**(2), 87–129 (2001)
80. Fisher, R.A.: The Use of Multiple Measurements in Taxonomic Problems. *Annals of Eugenics* **7**, 179–188 (1936)
81. Freund, Y., Mason, L.: The Alternating Decision Tree Learning Algorithm. In: Proceedings of the Sixteenth International Conference on Machine Learning (ICML '99), pp. 124–133. San Francisco, CA, USA (1999)
82. Freund, Y., Schapire, R.E.: A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.* **55**, 119–139 (1997)
83. Fürnkranz, J., Hüllermeier, E., Loza mencia, E., Brinker, K.: Multilabel classification via calibrated label ranking. *Machine Learning* **73**(2), 133 – 153 (2008)
84. Fürnkranz, J., Park, S.H.: Error-Correcting Output Codes as a Transformation from Multi-Class to Multi-Label Prediction. In: Discovery Science, *Lecture Notes in Computer Science*, vol. 7569, pp. 254–267. Springer Berlin Heidelberg (2012)
85. van der Gaag, L.C., de Waal, P.R.: Multi-dimensional bayesian network classifiers. In: Third European Workshop on Probabilistic Graphical Models, pp. 107–114 (2006)
86. Ghamrawi, N., McCallum, A.: Collective multi-label classification. In: CIKM '05: Proceedings of the 14th ACM international conference on Information and knowledge management, pp. 195–200 (2005)
87. Godbole, S., Sarawagi, S.: Discriminative Methods for Multi-Labeled Classification. In: Proceedings of the 8th Pacific-Asia Conference on Knowledge Discovery and Data Mining, pp. 22–30 (2004)
88. Gonçalves, E.C., Plastino, A., Freitas, A.A.: A genetic algorithm for optimizing the label ordering in multi-label classifier chains. In: IEEE 25th International Conference on Tools with Artificial Intelligence (ICTAI), pp. 469–476 (2013)
89. Gonçalves, T., Quaresma, P.: A Preliminary Approach to the Multilabel Classification Problem of Portuguese Juridical Documents. In: Progress in Artificial Intelligence, *Lecture Notes in Computer Science*, vol. 2902, pp. 435–444 (2003)
90. Gonçalves, T., Quaresma, P.: The impact of NLP techniques in the multilabel text classification problem. In: Proceedings of Intelligent Information Processing and Web Mining (IIPWM'04), *Advances in Soft Computing*, pp. 424–428 (2004)
91. Gretton, A., Bousquet, O., Smola, A., Schölkopf, B.: Measuring statistical dependence with hilbert-schmidt norms. In: Proceedings of the 16th international conference on Algorithmic Learning Theory (ALT'05), pp. 63–77 (2005)
92. Han, J., Pei, J., Yin, Y., Mao, R.: Mining Frequent Patterns without Candidate Generation: A Frequent-Pattern Tree Approach. *Data Min. Knowl. Discov.* **8**(1), 53–87 (2004)
93. He, J., Gu, H., Wang, Z.: Multi-instance multi-label learning based on Gaussian process with application to visual mobile robot navigation. *Information Sciences* **190**, 162 – 177 (2012)
94. Hsu, D., Kakade, S., Langford, J., Zhang, T.: Multi-Label Prediction via Compressed Sensing. In: Advances in Neural Information Processing Systems (NIPS), pp. 772–780 (2009)

95. Hu, D.Y., Reichel, L.: Krylov-subspace methods for the sylvester equation. *Linear Algebra and Its Applications* **172**, 283 – 313 (1992)
96. Huang, G.B., Ding, X., Zhou, H.: Optimization method based extreme learning machine for classification. *Neurocomputing* **74**(1-3), 155–163 (2010)
97. Huang, S.J., Yu, Y., Zhou, Z.H.: Multi-label Hypothesis Reuse. In: Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD12), pp. 525–533. Beijing, China (2012)
98. Huang, S.J., Zhou, Z.H.: Multi-Label Learning by Exploiting Label Correlations Locally. In: Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence (2012)
99. Hüllermeier, E., Fürnkranz, J., Cheng, W., Brinker, K.: Label Ranking by Learning Pairwise Preferences. *Artificial Intelligence* **172**, 1897–1916 (2008)
100. Ioannou, M., Sakkas, G., Tsoumakas, G., Vlahavas, I.P.: Obtaining Bipartitions from Score Vectors for Multi-Label Classification. In: 22nd IEEE International Conference on Tools with Artificial Intelligence (ICTAI), pp. 409–416 (2010)
101. Ji, S., Tang, L., Yu, S., Ye, J.: A shared-subspace learning framework for multi-label classification. *ACM Trans. Knowl. Discov. Data* **4**(2), 1–29 (2010)
102. Jiang, A., Wang, C., Zhu, Y.: Calibrated Rank-SVM for multi-label image categorization. In: IEEE World Congress on Computational Intelligence (IJCNN), pp. 1450–1455 (2008)
103. Jiang, J.Y., Tsai, S.C., Lee, S.J.: FSKNN: Multi-label text categorization based on fuzzy similarity and k nearest neighbors. *Expert Syst. Appl.* **39**(3), 2813–2821 (2012)
104. Ávila Jiménez, J., Gibaja, E., Ventura, S.: Evolving Multi-label Classification Rules with Gene Expression Programming: A Preliminary Study. In: Hybrid Artificial Intelligence Systems (HAIS), *Lecture Notes in Computer Science*, vol. 6077, pp. 9–16 (2010)
105. Johnson, M., Cipolla, R.: Improved Image Annotation and Labelling through Multi-Label Boosting. In: British Machine Vision Association (BMVC) (2005)
106. Jolliffe, I.: *Principal Component Analysis*. Springer (1986)
107. Kajdanowicz, T., Kazienko, P.: Multi-label classification using error correcting output codes. *Int. J. Appl. Math. Comput. Sci* **22**(4), 829–840 (2012)
108. Kajdanowicz, T., Wozniak, M., Kazienko, P.: Multiple Classifier Method for Structured Output Prediction Based on Error Correcting Output Codes. In: Intelligent Information and Database Systems, *Lecture Notes in Computer Science*, vol. 6592, pp. 333–342 (2011)
109. Kaneda, Y., Ueda, N., Saito, K.: Extended Parametric Mixture Model for Robust Multi-labeled Text Categorization. In: Knowledge-Based Intelligent Information and Engineering Systems, *Lecture Notes in Computer Science*, vol. 3214, pp. 616–623 (2004)
110. Katakis, I., Tsoumakas, G., Vlahavas, I.: Multilabel Text Classification for Automated Tag Suggestion. In: Proceedings of the ECML/PKDD 2008 Discovery Challenge (2008)
111. Kawai, K., Takahashi, Y.: Identification of the Dual Action Antihypertensive Drugs Using TFS-Based Support Vector Machines. *Chem-Bio Informatics Journal* **4**, 44–51 (2009)
112. Keller, J.M., Gray, M.R., Givens, J.A.: Fuzzy K-Nearest neighbor algorithm. *IEEE Transactions on Systems, Man and Cybernetics* **15**(14), 580–585 (1985)
113. Kiritchenko, S., Matwin, S., Famili, A.F.: Functional annotation of genes using hierarchical text categorization. In: Proc. of the BioLINK SIG: Linking Literature, Information and Knowledge for Biology (held at ISMB-05) (2005)
114. Kocev, D., Vens, C., Struyf, J., Džeroski, S.: Ensembles of Multi-Objective Decision Trees. In: Proceedings of the 18th European conference on Machine Learning (ECML '07), pp. 624–631. Berlin, Heidelberg (2007)
115. Kouzani, A.: Multilabel Classification Using Error Correction Codes. In: Advances in Computation and Intelligence, *Lecture Notes in Computer Science*, vol. 6382, pp. 444–454 (2010)
116. Krohn-Grimberghe, A., Drumond, L., Freudenthaler, C., Schmidt-Thieme, L.: Multi-relational matrix factorization using Bayesian personalized ranking for social network data. In: Proceedings of the fifth ACM international conference on Web search and data mining (WSDM '12), pp. 173–182. New York, NY, USA (2012)
117. Kumar, N., Berg, A.C., Belhumeur, P.N., Nayar, S.K.: Attribute and Simile Classifiers for Face Verification. In: IEEE International Conference on Computer Vision (ICCV) (2009)
118. Lafferty, J., McCallum, A., Pereira, F.: Conditional random fields: Probabilistic models for segmenting and labelling sequence data. In: Proceedings of the Eighteenth International Conference on Machine Learning (ICML), pp. 282 – 289 (2001)
119. Lauser, B., Hotho, A.: Automatic Multi-label Subject Indexing in a Multilingual Environment. In: European Conference on Digital Libraries (ECDL), *Lecture Notes in Computer Science*, vol. 2769, pp. 140–151 (2003)

120. Lewis, D.D., Yang, Y., Rose, T.G., Li, F.: RCV1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research* **5**, 361–397 (2005)
121. Li, H., Guo, Y.J., Wu, M., Li, P., Xiang, Y.: Combine multi-valued attribute decomposition with multi-label learning. *Expert Syst. Appl.* **37(12)**, 8721–8728 (2010)
122. Li, J., Xu, J.: A Fast Multi-label Classification Algorithm Based on Double Label Support Vector Machine. In: *Proceedings of the 2009 International Conference on Computational Intelligence and Security (CIS '09)*, pp. 30–35. Washington, DC, USA (2009)
123. Li, X., Wang, L., Sung, E.: Multilabel SVM active learning for image classification. In: *International Conference on Image Processing (ICIP '04)*, pp. 2207–2210 (2004)
124. Lin, X., Chen, X.W.: Mr.KNN: soft relevance for multi-label classification. In: *Proceedings of the 19th ACM international conference on Information and knowledge management (CIKM '10)*, pp. 349–358. New York, NY, USA (2010)
125. Liu, G., Lin, Z., Yu, Y.: Multi-output regression on the output manifold. *Pattern Recogn.* **42(11)**, 2737–2743 (2009)
126. Liu, Yi and Jin, Rong and Yang, Liu: Semi-supervised Multi-label Learning by Constrained Non-negative Matrix Factorization. In: *Proceedings of the 21st National Conference on Artificial Intelligence (AAAI06) - Volume 1*, pp. 421–426 (2006)
127. Lo, H., Wang, J., Wang, H., Lin, S.: Cost-sensitive multi-label learning for audio tag annotation and retrieval. *IEEE Transactions on Multimedia* **13(3)**, 518–529 (2011)
128. López, V.F., de la Prieta, F., Ogihara, M., Wong, D.D.: A model for multi-label classification and ranking of learning objects. *Expert Systems with Applications* **39(10)**, 8878 – 8884 (2012)
129. Loza, E., Fürnkranz, J.: Efficient Pairwise Multilabel Classification for Large-Scale Problems in the Legal Domain. In: *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD-2008)*, pp. 50–65. Springer-Verlag (2008)
130. Loza, E., Fürnkranz, J.: Efficient Multilabel Classification Algorithms for Large-Scale Problems in the Legal Domain. In: *Semantic Processing of Legal Texts, Lecture Notes in Computer Science*, vol. 6036, pp. 192–215. Springer Berlin / Heidelberg (2010)
131. Loza, E., Park, S.H., Fürnkranz, J.: Efficient voting prediction for pairwise multilabel classification. *Neurocomput.* **73**, 1164–1176 (2010)
132. Luo, X., Heywood, Z.A.N.: Evaluation of Two Systems on Multi-class Multi-label Document Classification. In: *Proceedings of the 15th International Symposium on Methodologies for Intelligent Systems*, pp. 161–169 (2005)
133. Ma, A., Sethi, I., Patel, N.: Multimedia Content Tagging Using Multilabel Decision Tree. In: *11th IEEE International Symposium on Multimedia (ISM '09)*, pp. 606–611 (2009)
134. Madjarov, G., Gjorgjevikj, D., Džeroski, S.: Dual layer voting method for efficient multi-label classification. In: *Proceedings of the 5th Iberian conference on Pattern recognition and image analysis, Lecture Notes in Computer Science*, vol. 6669, pp. 232–239 (2011)
135. Madjarov, G., Kocev, D., Gjorgjevikj, D., Džeroski, S.: An extensive experimental comparison of methods for multi-label learning. *Pattern Recognition* **45(9)**, 3084 – 3104 (2012)
136. Mammadov, M.A., Rubinov, A.M., Yearwood, J.: The study of drug-reaction relationships using global optimization techniques. *Optimization Methods Software* **22**, 99–126 (2007)
137. McCallum, A.K.: Multi-label text classification with a mixture model trained by EM. In: *AAAI 99 Workshop on Text Learning* (1999)
138. Mencía, E.L.: Multi-Label Classification in Parallel Tasks. In: *2nd International Workshop on Learning from Multi-Label Data (MLD'10)*, pp. 29–36 (2010)
139. Montejo-Ráez, A., Ureña López, L.: Selection Strategies for Multi-label Text Categorization. In: *Advances in Natural Language Processing, Lecture Notes in Computer Science*, vol. 4139, pp. 585–592 (2006)
140. Nardiello, P., Sebastiani, F., Sperduti, A.: Discretizing Continuous Attributes in AdaBoost for Text Categorization. In: *Advances in Information Retrieval, Lecture Notes in Computer Science*, vol. 2633, pp. 320–334. Springer Berlin / Heidelberg (2003)
141. Nasierding, G., Kouzani, A.: Image to Text Translation by Multi-Label Classification. In: *Advanced Intelligent Computing Theories and Applications. With Aspects of Artificial Intelligence, Lecture Notes in Computer Science*, vol. 6216, pp. 247–254. Springer Berlin / Heidelberg (2010)
142. Nasierding, G., Kouzani, A.Z.: Empirical Study of Multi-label Classification Methods for Image Annotation and Retrieval. In: *Digital Image Computing: Techniques and Applications*, pp. 617–622 (2010)

143. Nasierding, G., Kouzani, A.Z., Tsoumakas, G.: A Triple-Random Ensemble Classification Method for Mining Multi-label Data. In: Proceedings of the 2010 IEEE International Conference on Data Mining Workshops (ICDMW '10), pp. 49–56. Washington, DC, USA (2010)
144. Nettleton, D., Banerjee, T.: Testing the equality of distributions of random vectors with categorical components. *Comput. Statist. Data Anal.* **37**, 195–208 (2001)
145. Nguyen, C.D., Dung, T.A., Cao, T.H.: Text Classification for DAG-Structured Categories. In: *Advances in Knowledge Discovery and Data Mining, Lecture Notes in Computer Science*, vol. 3518, chap. 36, pp. 1–18. Springer Berlin / Heidelberg (2005)
146. Nguyen, Nam and Caruana, Rich: Classification with partial labels. In: Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '08), pp. 551–559. New York, NY, USA (2008)
147. Noh, H.G., Song, M.S., Park, S.H.: An unbiased method for constructing multilabel classification trees. *Computational Statistics & Data Analysis* **47**(1), 149–164 (2004)
148. Otero, F., Freitas, A., Johnson, C.: A hierarchical multi-label classification ant colony algorithm for protein function prediction. *Memetic Computing* **2**(3), 165–181 (2010)
149. Oza, N., Castle, J.P., Stutz, J.: Classification of Aeronautics System Health and Safety Documents. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* **39**(6), 670–680 (2009)
150. Özpolat, E., Akar, G.B.: Automatic detection of learning styles for an e-learning system. *Comput. Educ.* **53**, 355–367 (2009)
151. Pachet, F., Roy, P.: Improving Multilabel Analysis of Music Titles: A Large-Scale Validation of the Correction Approach. *IEEE Transactions on Audio, Speech, and Language Processing* **17**(2), 335–343 (2009)
152. Park, C.H., Lee, M.: On applying linear discriminant analysis for multi-labeled problems. *Pattern Recogn. Lett.* **29**(7), 878–887 (2008)
153. Park, S.H., Fürnkranz, J.: Multi-Label Classification with Label Constraints. Tech. Rep. TUD-KE-2008-04, Knowledge Engineering Group, TU Darmstadt (2008). URL <http://www.ke.tu-darmstadt.de/publications/reports/tud-ke-2008-04.pdf>
154. Parpinelli, R., Lopes, H., Freitas, A.: Data Mining with an Ant Colony Optimization Algorithm. *EEE Trans. On Evolutionary Computation* **6**(4), 321–332 (2002)
155. Peters, S., Denoyer, L., Gallinari, P.: Iterative Annotation of Multi-relational Social Networks. In: International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pp. 96–103 (2010)
156. Petrovskiy, M.: Paired Comparisons Method for Solving Multi-Label Learning Problem. In: Sixth International Conference on Hybrid Intelligent Systems (HIS '06), p. 42 (2006)
157. Petterson, J., Caetano, T.S.: Reverse Multi-Label Learning. In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 1912–1920 (2010)
158. Qi, G.J., Hua, X.S., Rui, Y., Tang, J., Mei, T., Zhang, H.J.: Correlative multi-label video annotation. In: Proceedings of the 15th international conference on Multimedia, pp. 17–26. New York, NY, USA (2007)
159. Qi, G.J., Hua, X.S., Rui, Y., Tang, J., Zhang, H.J.: Two-Dimensional Active Learning for image classification. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2008)*, pp. 1–8 (2008)
160. Qi, G.J., Hua, X.S., Rui, Y., Tang, J., Zhang, H.J.: Two-Dimensional Multilabel Active Learning with an Efficient Online Adaptation Model for Image Classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **31**(10), 1880–1897 (2009)
161. Qu, W., Zhang, Y., Zhu, J., Qiu, Q.: Mining Multi-label Concept-Drifting Data Streams Using Dynamic Classifier Ensemble. In: *Advances in Machine Learning, Lecture Notes in Computer Science*, vol. 5828, pp. 308–321. Springer Berlin / Heidelberg (2009)
162. Rak, R., Kurgan, L., Reformat, M.: A tree-projection-based algorithm for multi-label recurrent-item associative-classification rule generation. *Data & Knowledge Engineering* **64**(1), 171–197 (2008)
163. Rao, P., Kupper, L.: Ties in paired-comparison experiments: A generalization of the bradley-terry model. *Amer. Statist. Assoc.* **62**, 194–204 (1967)
164. Read, J.: A Pruned Problem Transformation Method for Multi-label Classification. In: *Proc. of the NZ Computer Science Research Student Conference*, pp. 143 – 150 (2008)
165. Read, J.: Scalable Multi-label Classification. Ph.D. thesis, University of Waikato (2010)
166. Read, J.: Advances in multi-label classification. <http://users.ics.aalto.fi/jesse/talks/Charla-Malaga.pdf> (2011)

167. Read, J.: MEKA: A Multi-label Extension to WEKA. <http://meka.sourceforge.net/> (2012)
168. Read, J., Bifet, A., Holmes, G., Pfahringer, B.: Scalable and efficient multi-label classification for evolving data streams. *Machine Learning* **88**, 243–272 (2012)
169. Read, J., Pfahringer, B., Holmes, G.: Multi-label Classification Using Ensembles of Pruned Sets. In: *ICDM '08: Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, pp. 995–1000. Washington, DC, USA (2008)
170. Read, J., Pfahringer, B., Holmes, G.: Generating Synthetic Multi-label Data Streams. In: *ECML/PKDD 2009 Workshop on Learning from Multi-label Data (MLD'09)* (2009)
171. Read, J., Pfahringer, B., Holmes, G., Frank, E.: Classifier chains for multi-label classification. *Machine Learning* **85**(3), 1–27 (2011)
172. Rokach, L., Itach, E.: An Ensemble Method for Multi-Label Classification using a Transportation Model. In: *Proceedings of the 1st workshop on learning from multilabel data (MLD) held in conjunction with ECML/PKDD*, pp. 49–60. Bled, Slovenia (2009)
173. Rokach, L., Itach, E.: An Ensemble Method for Multi-Label Classification using an Approximation Algorithm for the Set Covering problem. In: *Proceedings of the 2n international workshop on learning from multilabel data (MLD)*, pp. 37–44. Haifa, Israel (2010)
174. Rokach, L., Schclar, A., Itach, E.: Ensemble methods for multi-label classification. *Expert Systems with Applications* **41**, 7507–7523 (2014)
175. Rubin, T., Chambers, A., Smyth, P., Steyvers, M.: Statistical topic models for multi-label document classification. *Machine Learning* **88**, 157–208 (2012)
176. Sapozhnikova, E.: ART-Based Neural Networks for Multi-label Classification. In: *Advances in Intelligent Data Analysis VIII, Lecture Notes in Computer Science*, vol. 5772, pp. 167–177. Springer Berlin / Heidelberg (2009)
177. Sarinnapakorn, K., Kubat, M.: Induction from multi-label examples in information retrieval systems: a case study. *Applied Artificial Intelligence* **22**(5), 407–432 (2008)
178. Schapire, R.E., Singer, Y.: Improved boosting algorithms using confidence-rated predictions. In: *Proceedings of the eleventh annual conference on Computational learning theory (COLT' 98)*, pp. 80–91. New York, NY, USA (1998)
179. Schapire, R.E., Singer, Y.: Improved Boosting Algorithms Using Confidence-rated Predictions. *Machine Learning* **37**(3), 297 – 336 (1999)
180. Schapire, R.E., Singer, Y.: BoosTexter: A Boosting-based System for Text Categorization. *Machine Learning* **39**, 135–168 (2000)
181. Sebastiani, F., Sperduti, A., Valdambrini, N.: An improved boosting algorithm and its application to text categorization. In: *Proceedings of the ninth international conference on Information and knowledge management (CIKM '00)*, pp. 78–85. New York, NY, USA (2000)
182. Shafer, G.: *A Mathematical Theory of Evidence*. Princeton University Press (1976)
183. Shao, H., Li, G., Liu, G., Wang, Y.: Symptom selection for multi-label data of inquiry diagnosis in traditional Chinese medicine. *Science China Information Sciences* **1**, 1–13 (2010)
184. Shen, C., Jing, L., Ng, M.: Sparse-MIML: A Sparsity-Based Multi-Instance Multi-Learning Algorithm. In: *Energy Minimization Methods in Computer Vision and Pattern Recognition, Lecture Notes in Computer Science*, vol. 8081, pp. 294–306 (2013)
185. Shi, C., Kong, X., Yu, P., Wang, B.: Multi-Objective Multi-Label Classification. In: *Proceedings of the SIAM International Conference on Data Mining*, pp. 355–366. Anaheim, California (2012)
186. Shi, C., Kong, X., Yu, P.S., Wang, B.: Multi-label Ensemble Learning. In: *European conference on Machine learning and knowledge discovery in databases (ECML/PKDD)*, pp. 223–239 (2011)
187. Shotton, J., Winn, J., Rother, C., Criminisi, A.: TextonBoost for Image Understanding: Multi-Class Object Recognition and Segmentation by Jointly Modeling Texture, Layout, and Context. *Int. J. Comput. Vision* **81**, 2–23 (2009)
188. Sobol-Shikler, T., Robinson, P.: Classification of Complex Information: Inference of Co-Occurring Affective States from Their Expressions in Speech. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **32**(7), 1284–1297 (2010)
189. Song, Y., Zhang, L., Giles, C.L.: Automatic tag recommendation algorithms for social recommender systems. *ACM Trans. Web* **5**(1), 4:1–4:31 (2011)
190. Spat, S., Cadonna, B., Rakovac, I., Gütl, C., Leitner, H., Stark, G., Beck, P.: Enhanced Information Retrieval from Narrative German-language Clinical Text Documents using Automated Document Classification. In: *eHealth Beyond the Horizon - Get IT There*,

- Proceedings of MIE2008, The XXist International Congress of the European Federation for Medical Informatics, pp. 473–478. Göteborg, Sweden (2008)
191. Specht, D.F.: Probabilistic neural networks. *Neural Netw.* **3**(1), 109–118 (1990)
 192. Spyromitros, E., Tsoumakas, G., Vlahavas, I.: An Empirical Study of Lazy Multilabel Classification Algorithms. In: SETN '08: Proceedings of the 5th Hellenic conference on Artificial Intelligence, pp. 401–406. Berlin, Heidelberg (2008)
 193. Streich, A., Buhmann, J.: Classification of Multi-labeled Data: A Generative Approach. In: Proceedings of the European conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD '08), pp. 390–405. Springer-Verlag (2008)
 194. Sucar, L.E., Bielza, C., Morales, E.F., Hernandez-Leal, P., Zaragoza, J.H., Larrañaga, P.: Multi-label classification with bayesian network-based chain classifiers. *Pattern Recognition Letters* **41**, 14–22 (2014)
 195. Sun, L., Ji, S., Ye, J.: Canonical Correlation Analysis for Multilabel Classification: A Least-Squares Formulation, Extensions, and Analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **33**(1), 194–200 (2011)
 196. Sun, Y.Y., Zhang, Y., Zhou, Z.H.: Multi-label learning with weak label. In: Proceedings of the National Conference on Artificial Intelligence, pp. 593–598 (2010)
 197. Tahir, M.A., Kittler, J., Bouridane, A.: Multilabel classification using heterogeneous ensemble of multi-label classifiers. *Pattern Recognition Letters* **33**(5), 513 – 523 (2012)
 198. Tai, F., Lin, H.T.: Multi-Label Classification with Principal Label Space Transformation. In: 2nd International Workshop on Learning from Multi-Label Data (MLD'10), pp. 45–52 (2010)
 199. Tang, L., Liu, H.: Relational learning via latent social dimensions. In: Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '09), pp. 817–826. New York, NY, USA (2009)
 200. Tang, L., Liu, H.: Scalable learning of collective behavior based on sparse social dimensions. In: Proceedings of the 18th ACM conference on Information and knowledge management (CIKM '09), pp. 1107–1116. New York, NY, USA (2009)
 201. Tang, L., Rajan, S., Narayanan, V.K.: Large scale multi-label classification via metabeler. In: Proceedings of the 18th international conference on World wide web (WWW '09), pp. 211–220. New York, NY, USA (2009)
 202. Tax, D., Duan, R.P.W.: Support vector data description. *Machine Learning* **54**(1), 45–66 (2004)
 203. Tenenboim, L., Rokach, L., Shapira, B.: Identification of Label Dependencies for Multi-Label Classification. In: 2nd International Workshop on Learning from Multi-Label Data (MLD'10), pp. 53–60 (2010)
 204. Thabtah, F.A., Cowling, P., Peng, Y., Rastogi, R., Morik, K., Bramer, M., Wu, X.: MMAC: A new multi-class, multi-label associative classification approach. In: Proceedings - Fourth IEEE International Conference on Data Mining, ICDM 2004, pp. 217–224 (2004)
 205. Thabtah, F.A., Cowling, P.I.: A greedy classification algorithm based on association rule. *Appl. Soft Comput.* **7**(3), 1102–1111 (2007)
 206. Trohidis, K., Tsoumakas, G., Kalliris, G., Vlahavas, I.: Multi-label Classification of Music into Emotions. In: Proc. 9th International Conference on Music Information Retrieval (ISMIR 2008), pp. 325–330 (2008)
 207. Tsoumakas, G., Dimou, A., Spyromitros, E., Mezaris, V., Kompatsiaris, I., Vlahavas, I.: Correlation-Based Pruning of Stacked Binary Relevance Models for Multi-Label Learning. In: Proceedings of the 1st International Workshop on Learning from Multi-Label Data (MLD'09), pp. 101–116. Bled, Slovenia (2009)
 208. Tsoumakas, G., Katakis, I.: Multi Label Classification: An Overview. *International Journal of Data Warehousing and Mining* **3**, 1–13 (2007)
 209. Tsoumakas, G., Katakis, I., Vlahavas, I.: Effective and Efficient Multilabel Classification in Domains with Large Number of Labels. In: Proc. ECML/PKDD 2008 Workshop on Mining Multidimensional Data (MMD'08) (2008)
 210. Tsoumakas, G., Katakis, I., Vlahavas, I.: *Data Mining and Knowledge Discovery Handbook*, Part 6, chap. Mining Multi-label Data, pp. 667–685. Springer (2010)
 211. Tsoumakas, G., Katakis, I., Vlahavas, I.: Random k-Labelsets for Multi-Label Classification. *IEEE Transactions on Knowledge and Data Engineering* **23**(7), 1079–1089 (2010)
 212. Tsoumakas, G., Mencia, E.L., Katakis, I., Park, S., Fürnkranz, J.: On the combination of two decompositive multi-label classification methods. In: Workshop on Preference Learning, ECML PKDD 09, pp. 114–133 (2009)

213. Tsoumakas, G., Spyromitros-Xioufis, E., Vilcek, J., Vlahavas, I.: Mulan: A Java Library for Multi-Label Learning. *Journal of Machine Learning Research* **12**, 2411–2414 (2011)
214. Turnbull, D., Barrington, L., Torres, D., Lanckriet, G.: Semantic Annotation and Retrieval of Music and Sound Effects. *Audio, Speech, and Language Processing, IEEE Transactions on* **16**(2), 467–476 (2008)
215. Ueda, N., Saito, K.: Parametric Mixture Models for Multi-Labeled Text. In: *Neural Information Processing Systems (NIPS)*, pp. 721–728 (2002)
216. Ukwatta, E., Samarabandu, J.: Vision Based Metal Spectral Analysis Using Multi-label Classification. In: *Canadian Conference on Computer and Robot Vision (CRV '09)*, pp. 132–139 (2009)
217. Valentini, G., Re, M.: Weighted True Path Rule: a multilabel hierarchical algorithm for gene function prediction. In: *1st International Workshop on learning from Multi-Label Data (MLD-ECML 2009)*, pp. 132–145. Bled, Slovenia (2009)
218. Vateekul, P., Kubat, M.: Fast Induction of Multiple Decision Trees in Text Categorization from Large Scale, Imbalanced, and Multi-label Data. In: *IEEE International Conference on Data Mining Workshops (ICDMW)*, pp. 320–325 (2009)
219. Veloso, A., Jr., W.M., Gonçalves, M.A., Zaki, M.J.: Multi-label Lazy Associative Classification. In: *Proceedings of the 11th European Conference on Principles and Practice of Knowledge Discovery in Databases, PKDD*, pp. 605–612. Warsaw, Poland (2007)
220. Vens, C., Struyf, J., Schietgat, L., Džeroski, S., Blockeel, H.: Decision trees for hierarchical multi-label classification. *Machine Learning* **73**(2), 185–214 (2008)
221. Wan, S.P., Xu, J.H.: A multi-label classification algorithm based on triple class support vector machine. In: *International Conference on Wavelet Analysis and Pattern Recognition (ICWAPR '07)*, vol. 4, pp. 1447–1452 (2007)
222. Wang, H., Ding, C., Huang, H.: Multi-label Linear Discriminant Analysis. In: *Computer Vision – ECCV 2010, Lecture Notes in Computer Science*, vol. 6316, pp. 126–139. Springer Berlin / Heidelberg (2010)
223. Wang, H., Huang, M., Zhu, X.: A Generative Probabilistic Model for Multi-label Classification. In: *ICDM '08: Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, pp. 628–637. Washington, DC, USA (2008)
224. Wang, J., Zhao, Y., Wu, X., Hua, X.S.: A transductive multi-label learning approach for video concept detection. *Pattern Recognition* **44**, 2274–2286 (2010)
225. Wang, L., Chang, M., Feng, J.: Parallel and sequential support vectormachines for multi-label classification. *International Journal of Information Technology* **11**(9), 11–18 (2005)
226. Wang, M., Zhou, X., Chua, T.S.: Automatic image annotation via local multi-label classification. In: *Proceedings of the 2008 international conference on Content-based image and video retrieval (CIVR '08)*, pp. 17–26. New York, NY, USA (2008)
227. Wang, X., Liu, X., Shi, Z., Shi, Z., Sui, H.: Voting conditional random fields for multi-label image classification. In: *3rd International Congress on Image and Signal Processing (CISP)*, pp. 1984–1988 (2010)
228. Wolpert, D.H.: Stacked Generalization. *Neural Networks* **5**, 241–259 (1992)
229. Wu, L., Zhang, M.L.: Multi-label classification with unlabeled data: An inductive approach. In: *Proceedings of the 5th Asian Conference on Machine Learning (ACML'13)*, pp. 197–212. Canberra, Australia (2013)
230. Xiao, X., Wu, Z.C., Chou, K.C.: iLoc-Virus: a multi-label learning classifier for identifying the subcellular localization of virus proteins with both single and multiple sites. *J Theor Biol* **284**(1), 42–51 (2011)
231. Xioufis, E.S., Spiliopoulou, M., Tsoumakas, G., Vlahavas, I.P.: Dealing with Concept Drift and Class Imbalance in Multi-Label Stream Classification. In: *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1583–1588. Barcelona, Spain (2011)
232. Xu, H., Xu, J.: Designing a multi-label kernel machine with two-objective optimization. In: *Proceedings of the 2010 international conference on Artificial intelligence and computational intelligence (AICI'10): Part I*, pp. 282–291. Berlin, Heidelberg (2010)
233. Xu, J.: Constructing a Fast Algorithm for Multi-label Classification with Support Vector Data Description. In: *Proceedings of the IEEE International Conference on Granular Computing (GrC)*, pp. 817–821 (2010)
234. Xu, J.: An efficient multi-label support vector machine with a zero label. *Expert Systems with Applications* **39**(5), 4796–4804 (2012)
235. Xu, X., Frank, E.: Logistic Regression and Boosting for Labeled Bags of Instances. In: *Advances in Knowledge Discovery and Data Mining, Lecture Notes in Computer Science*, vol. 3056, pp. 272–281. Springer Berlin / Heidelberg (2004)

236. Yager, R.R.: Veristic variables. *IEEE Transactions on Systems, Man, and Cybernetics* **30**(1), 71–84 (2000)
237. Yan, R., Tesic, J., Smith, J.R.: Model-shared subspace boosting for multi-label classification. In: *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '07)*, pp. 834–843. New York, NY, USA (2007)
238. Yan, Y., Fung, G., Dy, J.G., Rosales, R.: Medical coding classification by leveraging inter-code relationships. In: *Proceedings of the 16th international conference on Knowledge discovery and data mining (KDD '10)*, pp. 193–202. New York, NY, USA (2010)
239. Yang, B., Sun, J.T., Wang, T., Chen, Z.: Effective multi-label active learning for text classification. In: *KDD '09: Proceedings of the 15th international conference on Knowledge discovery and data mining*, pp. 917–926. New York, NY, USA (2009)
240. Yang, S.H., Zha, H., Hu, B.G.: Dirichlet-Bernoulli Alignment: A Generative Model for Multi-Class Multi-Label Multi-Instance Corpora. In: *Annual Conference on Neural Information Processing Systems*, pp. 2143–2150. Vancouver, British Columbia, Canada (2009)
241. Yang, Y.: A study of thresholding strategies for text categorization. In: *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR '01)*, pp. 137–145. New York, NY, USA (2001)
242. Yang, Y., Pedersen, J.O.: A Comparative Study on Feature Selection in Text Categorization. In: *Proceedings of the Fourteenth International Conference on Machine Learning (ICML '97)*, pp. 412–420. San Francisco, CA, USA (1997)
243. Yearwood, J., Mammadov, M., Banerjee, A.: Profiling Phishing Emails Based on Hyperlink Information. In: *International Conference on Advances in Social Networks Analysis and Mining*, pp. 120–127 (2010)
244. Younes, Z., Aballah, F., Denoeux, T.: Multi-label classification algorithm derived from k-nearest neighbor rule with label dependencies. In: *Proceedings of the 16th European Signal Processing Conference* (2008)
245. Younes, Z., Abdallah, F., Denoeux, T.: Evidential Multi-Label Classification Approach to Learning from Data with Imprecise Labels. In: *Computational Intelligence for Knowledge-Based Systems Design, Lecture Notes in Computer Science*, vol. 6178, pp. 119–128. Springer Berlin / Heidelberg (2010)
246. Younes, Z., Abdallah, F., Denoux, T.: Fuzzy multi-label learning under veristic variables. In: *IEEE International Conference on Fuzzy Systems, FUZZ-IEEE 2010*, pp. 1–8 (2010)
247. Yu, K., Yu, S., Tresp, V.: Multi-label informed latent semantic indexing. In: *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 258–265. New York, NY, USA (2005)
248. Zafra, A., Gibaja, E., Ventura, S.: Multiple Instance Learning with Multiple Objective Genetic Programming for Web Mining. *Applied Soft Computing* **11**(1), 93 – 102 (2011)
249. Zha, Z.J., Mei, T., Wang, J., Wang, Z., Hua, X.S.: Graph-based semi-supervised learning with multiple labels. *Journal of Visual Communication and Image Representation* **20**(2), 97 – 103 (2009). Special issue on Emerging Techniques for Multimedia Content Sharing, Search and Understanding
250. Zhang, M.L.: ML-rbf: RBF Neural Networks for Multi-Label Learning. *Neural Processing Letters* **29**(2), 61–74 (2009)
251. Zhang, M.L.: A k-Nearest Neighbor Based Multi-Instance Multi-Label Learning Algorithm. In: *Proceedings of the 22nd IEEE International Conference on Tools with Artificial Intelligence, ICTAI (2)*, pp. 207–212. Arras, France (2010)
252. Zhang, M.L.: LIFT: Multi-Label Learning with Label-Specific Features. In: *Proceedings of the 22nd International Joint Conference on Artificial Intelligence IJCAI*, pp. 1609–1614 (2011)
253. Zhang, M.L., Peña, J.M., Robles, V.: Feature selection for multi-label naive Bayes classification. *Information Sciences* **179**(19), 3218–3229 (2009)
254. Zhang, M.L., Wang, Z.J.: MIMLRBF: RBF neural networks for multi-instance multi-label learning. *Neurocomputing* **72**(16–18), 3951 – 3956 (2009)
255. Zhang, M.L., Zhang, K.: Multi-label learning by exploiting label dependency. In: *Proceedings of the 16th international conference on Knowledge discovery and data mining (KDD '10)*, pp. 999–1008. New York, NY, USA (2010)
256. Zhang, M.L., Zhou, Z.H.: A k-Nearest Neighbor Based Algorithm for Multi-label Classification. In: *Proceedings of the IEEE International Conference on Granular Computing (GrC)*, pp. 718–721. Beijing, China (2005)
257. Zhang, M.L., Zhou, Z.H.: Multilabel Neural Networks with Applications to Functional Genomics and Text Categorization. *IEEE Transactions on Knowledge and Data Engineering* **18**(10), 1338–1351 (2006)

258. Zhang, M.L., Zhou, Z.H.: ML-KNN: A lazy learning approach to multi-label learning. *Pattern Recognition* **40**(7), 2038–2048 (2007)
259. Zhang, M.L., Zhou, Z.H.: Multi-Label Learning by Instance Differentiation. In: *AAAI Conference on Artificial Intelligence*, pp. 669–674 (2007)
260. Zhang, M.L., Zhou, Z.H.: M3MIML: A Maximum Margin Method for Multi-instance Multi-label Learning. In: *ICDM '08: Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, pp. 688–697. Washington, DC, USA (2008)
261. Zhang, M.L., Zhou, Z.H.: A Review On Multi-Label Learning Algorithms. *IEEE Transactions on Knowledge and Data Engineering* **26**(8), 1819 – 1837 (2014)
262. Zhang, X., Cheng, J., Xu, C., Lu, H., Ma, S.: Multi-view multi-label active learning for image classification. In: *IEEE International Conference on Multimedia and Expo*, pp. 258–261 (2009)
263. Zhang, X., Yuan, Q., Zhao, S., Fan, W., Zheng, W., Wang, Z.: Multi-label classification without the multi-label cost. In: *Proceedings of the 10th SIAM International Conference on Data Mining* (2010)
264. Zhang, Y., Burer, S., Street, W.N., Bennett, K., Parrado-hern, E.: Ensemble Pruning Via Semi-definite Programming. *Journal of Machine Learning Research* **7**, 1315–1338 (2006)
265. Zhang, Y., Schneider, J.: Multi-Label Output Codes using Canonical Correlation Analysis. In: *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 873–882 (2011)
266. Zhang, Y., Zhou, Z.H.: Multilabel dimensionality reduction via dependence maximization. *ACM Transactions on Knowledge Discovery from Data (TKDD)* **4**, paper 14 (2010)
267. Zhou, T., Tao, D., Wu, X.: Compressed labeling on distilled labelsets for multi-label learning. *Machine Learning* **88**, 69–126 (2012)
268. Zhou, Z.H., Zhang, M.L.: Multi-Instance Multi-Label Learning with Application to Scene Classification. In: *NIPS*, pp. 1609–1616 (2006)
269. Zhou, Z.H., Zhang, M.L., Huang, S.J., Li, Y.F.: Multi-instance multi-label learning. *Artificial Intelligence* **176**(1), 2291 – 2320 (2012)
270. Zhu, S., Ji, X., Xu, W., Gong, Y.: Multi-labelled classification using maximum entropy method. In: *SIGIR '05: Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 274–281. New York, NY, USA (2005)