

# On Aggregation in Ensembles of Multilabel Classifiers

Vu-Linh Nguyen<sup>1</sup>, Eyke Hüllermeier<sup>1</sup>, Michael Rapp<sup>2</sup>,  
Eneldo Loza Mencía<sup>2</sup>, and Johannes Fürnkranz<sup>3</sup>

<sup>1</sup> Heinz Nixdorf Institute and Department of Computer Science Paderborn University, Germany  
vu.linh.nguyen@uni-paderborn.de, eyke@upb.de

<sup>2</sup> Knowledge Engineering Group, TU Darmstadt, Germany mrapp@ke.tu-darmstadt.de,  
eneldo@ke.tu-darmstadt.de

<sup>3</sup> Computational Data Analytics Group, JKU Linz, Austria juffi@faw.jku.at

**Abstract.** While a variety of ensemble methods for multilabel classification have been proposed in the literature, the question of how to aggregate the predictions of the individual members of the ensemble has received little attention so far. In this paper, we introduce a formal framework of ensemble multilabel classification, in which we distinguish two principal approaches: “predict then combine” (PTC), where the ensemble members first make loss minimizing predictions which are subsequently combined, and “combine then predict” (CTP), which first aggregates information such as marginal label probabilities from the individual ensemble members, and then derives a prediction from this aggregation. While both approaches generalize voting techniques commonly used for multilabel ensembles, they allow to explicitly take the target performance measure into account. Therefore, concrete instantiations of CTP and PTC can be tailored to concrete loss functions. Experimentally, we show that standard voting techniques are indeed outperformed by suitable instantiations of CTP and PTC, and provide some evidence that CTP performs well for decomposable loss functions, whereas PTC is the better choice for non-decomposable losses.

**Keywords:** Ensembles of Multilabel Classifiers, Predict then Combine, Combine then Predict, Hamming loss, F-measure, Subset 0/1 loss.

## 1 Introduction

The setting of *multilabel classification* (MLC), which generalizes standard multi-class classification by relaxing the assumption of mutual exclusiveness of classes, has received a lot of attention in the recent machine learning literature—we refer to [21] and [23] for comprehensive survey articles on this topic.

Like for other types of classification problems, the idea of *ensemble learning* [5] has also been applied to MLC (cf. Section 3). However, somewhat surprisingly, the question of how to aggregate the predictions of the individual members of an ensemble has so far received little attention in MLC. Instead, most approaches are based on simple voting techniques, which are typically applied in a label-wise manner: For each label, the predictions—either binary predictions of relevance or, more generally, label probabilities—of all ensemble members are collected, averaged, and thresholded to obtain a final prediction for this label.

An obvious disadvantage of this simple approach is that the aggregation is independent of the underlying performance measure, i.e., the aggregation procedure is not tailored to a specific loss function. This, however, would supposedly be important: In contrast to standard classification,

where a loss function compares a predicted class label with a ground truth, an MLC loss compares a *subset* of labels predicted to be relevant with a ground-truth subset. As there are various ways in which subsets can be compared with each other, a wide spectrum of loss functions is commonly used in MLC, and it is well known that different losses may call for different (Bayes-optimal) predictions [3,4]. Naturally, the idea of customizing an MLC predictor to a specific loss function should not only be considered at the level of individual predictors, but also at the level of the ensemble as a whole, and hence also concern the way in which the predictions are combined.

In this paper, we study the problem of aggregation in ensembles of multilabel classifiers (EMLC) in a systematic way. To this end, we introduce a formal framework, in which we distinguish two principal approaches: “predict then combine” (PTC), where the ensemble members first make loss minimizing predictions which are then combined, and “combine then predict” (CTP), which first aggregates information such as marginal label probabilities from the individual ensemble members, and then derives a prediction from this aggregation. While both approaches generalize common voting techniques as mentioned above, they also include more general variants and, moreover, allow one to explicitly take the target loss into account. In other words, concrete instantiations of CTP and PTC can be tailored to concrete loss functions. In an extensive experimental study, we demonstrate that such loss-based aggregation functions do indeed outperform simple voting techniques, and also investigate the question which type of aggregation is more suitable for which loss functions.

## 2 Multilabel Classification

Let  $\mathcal{X}$  denote an instance space, and let  $\mathcal{L} = \{\lambda_1, \dots, \lambda_K\}$  be a finite set of class labels. We assume that an instance  $\mathbf{x} \in \mathcal{X}$  is (probabilistically) associated with a subset of labels  $\Lambda = \Lambda(\mathbf{x}) \in 2^{\mathcal{L}}$ ; this subset is often called the set of relevant labels, while the complement  $\mathcal{L} \setminus \Lambda$  is considered as irrelevant for  $\mathbf{x}$ . We identify a set  $\Lambda$  of relevant labels with a binary vector  $\mathbf{y} = (y_1, \dots, y_K)$ , where  $y_k = \llbracket \lambda_k \in \Lambda \rrbracket$ .<sup>4</sup> By  $\mathcal{Y} = \{0, 1\}^K$  we denote the set of possible labelings.

We assume observations to be realizations of random variables generated independently and identically (i.i.d.) according to a probability distribution  $\mathbf{p}$  on  $\mathcal{X} \times \mathcal{Y}$ , i.e., an observation  $\mathbf{y} = (y_1, \dots, y_K)$  is the realization of a corresponding random vector  $\mathbf{Y} = (Y_1, \dots, Y_K)$ . We denote by  $\mathbf{p}(\mathbf{Y} | \mathbf{x})$  the conditional distribution of  $\mathbf{Y}$  given  $\mathbf{X} = \mathbf{x}$ , and by  $\mathbf{p}_k(Y_k | \mathbf{x})$  the corresponding marginal distribution of  $Y_k$ :

$$\mathbf{p}_k(b | \mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}: y_k = b} \mathbf{p}(\mathbf{y} | \mathbf{x}). \quad (1)$$

Moreover, we denote by  $p_k = p_k(\mathbf{x}) = \mathbf{p}_k(1 | \mathbf{x})$  the probability of relevance of the label  $\lambda_k$ .

Given training data in the form of a finite set of observations

$$\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N \subset \mathcal{X} \times \mathcal{Y}, \quad (2)$$

drawn independently from  $\mathbf{p}(\mathbf{X}, \mathbf{Y})$ , the goal in MLC is to learn a predictive model in the form of a multilabel classifier  $\mathbf{h}$ , which is a mapping  $\mathcal{X} \rightarrow \mathcal{Y}$  that assigns a (predicted) label subset to each instance  $\mathbf{x} \in \mathcal{X}$ . Thus, the output of a classifier  $\mathbf{h}$  is a vector of predictions

$$\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_K(\mathbf{x})) \in \{0, 1\}^K, \quad (3)$$

<sup>4</sup>  $\llbracket \cdot \rrbracket$  is the indicator function, i.e.,  $\llbracket A \rrbracket = 1$  if the predicate  $A$  is true and  $= 0$  otherwise.

also denoted as  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_K)$ .

## 2.1 MLC Loss Functions

The main goal in MLC is to induce predictions (3) that generalize well beyond the training data (2), i.e., predictions

$$\hat{\mathbf{y}} = \operatorname{argmin}_{\mathbf{y}} \sum_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}, \bar{\mathbf{y}}) p(\mathbf{y} | \mathbf{x}), \quad (4)$$

that minimize the expected loss with respect to a specific MLC loss function  $\ell : \mathcal{Y}^2 \rightarrow \mathbb{R}$ . Two important loss functions, both generalizing the standard 0/1 loss commonly used in classification, are the *Hamming loss* and the *subset 0/1 loss*:

$$\ell_H(\mathbf{y}, \hat{\mathbf{y}}) := \frac{1}{K} \sum_{k=1}^K \llbracket y_k \neq \hat{y}_k \rrbracket, \quad (5)$$

$$\ell_S(\mathbf{y}, \hat{\mathbf{y}}) := \llbracket \mathbf{y} \neq \hat{\mathbf{y}} \rrbracket. \quad (6)$$

The (*instance-wise*) *F-measure* compares a set of predicted labels to a corresponding set of ground-truth labels via the harmonic mean of precision and recall:

$$F(\mathbf{y}, \hat{\mathbf{y}}) = \frac{2 \sum_{k=1}^K \hat{y}_k y_k}{\sum_{k=1}^K \hat{y}_k + \sum_{k=1}^K y_k}. \quad (7)$$

The goal of classification algorithms in general is to capture dependencies between input features and the target variable. In MLC, dependencies may not only exist between the features and each target, but also between the targets  $Y_1, \dots, Y_K$  themselves. The idea to improve predictive accuracy by capturing such dependencies is a driving force in research on multilabel classification.

Not all loss functions capture label dependencies to the same extent: A *decomposable loss* can be reduced to loss functions for the individual labels, i.e., it can be expressed in the form

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{k=1}^K \ell_k(y_k, \hat{y}_k), \quad (8)$$

with suitable binary loss functions  $\ell_k : \{0, 1\}^2 \rightarrow \mathbb{R}$ . A *non-decomposable loss* does not permit such a representation. It can be shown that, for making optimal predictions  $\hat{\mathbf{y}} = \mathbf{h}(\mathbf{x})$  which minimize the expected loss, knowledge about the marginals (1) is sufficient in the case of a decomposable loss (such as Hamming), but not in the case of a non-decomposable loss [3]. Instead, if a loss is non-decomposable, higher-order probabilities are needed, and in the extreme case even the entire distribution  $p(\mathbf{Y} | \mathbf{x})$  (like in the case of the subset 0/1 loss).

On an algorithmic level, this means that MLC with a decomposable loss can be tackled by what is commonly called binary relevance (BR) learning, i.e., by learning one binary classifier for each individual label, whereas non-decomposable losses call for more sophisticated learning methods that are able to take label dependencies into account.

## 2.2 Risk Minimization

In the most general case, the problem of finding a risk-minimizing (Bayes-optimal) prediction is tackled by producing a prediction  $p(\cdot | \mathbf{x})$  of the conditional joint distribution of labelings, and explicitly solving (4) as a combinatorial optimization problem. Obviously, this approach is infeasible unless the number of class labels is very low. Fortunately, the problem can be solved more efficiently for specific loss functions, including those considered in this paper.

In the case of the *Hamming loss*, the Bayes-optimal prediction can be obtained by thresholding the marginal probabilities, regardless of whether the labels are independent or not:

$$\hat{y}_k = \begin{cases} 0 & \text{if } p_k(\mathbf{x}) \leq 1/2 \\ 1 & \text{if } p_k(\mathbf{x}) > 1/2 \end{cases} \quad (9)$$

Thus, it is sufficient to have good estimates for the marginal probabilities, which can be accomplished by simple techniques such as binary relevance [3].

For *subset 0/1 loss*, the Bayes-optimal prediction is not the marginal but the *joint* mode of the distribution  $p(\cdot | \mathbf{x})$ :

$$\hat{\mathbf{y}} \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y} | \mathbf{x}).$$

In this case, label dependence needs to be taken into account for optimal performance.

The *F-measure* is in a sense in-between these two extremes. It can be shown that, while the entire distribution  $p(\cdot | \mathbf{x})$  is not needed to find a Bayes-optimal prediction for this measure, marginal probabilities (1) do not suffice either. Instead, probabilities on pairwise label combinations are required in the general case, whereas under the assumption of conditional label independence, marginal probabilities again provide sufficient information [22].

## 3 Ensembles of MLC

In general, an ensemble approach to multilabel classification (EMLC) learns a set of  $M$  multilabel classifiers, each of which predicts a binary label vector  $\hat{\mathbf{y}}_j$ . Given a query instance  $\mathbf{x} \in \mathcal{X}$ , these are then aggregated into a final prediction  $\hat{\mathbf{y}} = \operatorname{agg}(\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_M)$ . For this aggregation, variants of label-wise *majority voting* (MV) are typically used:

- **Binary majority voting (BMV)** assigns to each label  $\lambda_k \in \mathcal{L}$  the prediction given by the majority of the classifiers:

$$\hat{y}_k := \operatorname{argmax}_{y_k \in \{0,1\}} \sum_{j=1}^M \mathbb{1}[y_k = \hat{y}_{j,k}]. \quad (10)$$

- **Graded majority voting (GMV)**, also known as *weighted voting*, adds up confidence scores  $p_j = (p_{j,1}, p_{j,2}, \dots, p_{j,K})$  for each label  $\lambda_k \in \mathcal{L}$ :

$$\hat{y}_k := \operatorname{argmax}_{y_k \in \{0,1\}} \sum_{j=1}^M p_{j,k}^{y_k} (1 - p_{j,k})^{1-y_k}. \quad (11)$$

Several ensemble-based multi-label classifiers have been tried in the literature, which typically use the above-mentioned voting techniques for combining the predictions of the ensemble members [6,7,10,18,19]. While we aim at optimizing the predictions for a particular loss function, a different line of work—orthogonal to our approach—aims at simultaneously optimizing for multiple loss functions [16,17]. In the following, we briefly recall some commonly used EMLC methods, which will serve as baselines in our experimental evaluation. We refer to [11] for an extensive discussion on ensembles of MLC classifiers.

- **Ensembles of Binary Relevance Classifiers (EBR)** use bagging [1] to construct  $K$  independent ensembles of binary classifiers, one for each label  $\lambda_k \in \mathcal{L}$  [20]. At prediction time, the predictions of these classifiers are combined for each label using majority voting, as is commonly used in bagged ensembles. Obviously, like all BR methods, EBR ignores any relationships between the labels and implicitly assumes them to be independent. Moreover, EBR is computationally expensive, since  $K \cdot M$  classifiers are required in order to have an “actual ensemble” of cardinality  $M$ .
- **Ensembles of Classifier Chains (ECC).** The classifier chains (CC) method [14] also trains  $K$  binary classifiers  $h_k$ ,  $k \in [K] := \{1, \dots, K\}$ , one for each label. Yet, to capture label dependencies,  $h_k$  is trained on an augmented input space  $\mathcal{X} \times \{0, 1\}^{k-1}$ , taking the (binary) values of the  $k - 1$  previous labels as additional attributes. More specifically,  $h_k$  predicts  $\hat{y}_{\sigma(k)} \in \{0, 1\}$  using

$$(\mathbf{x}, \hat{y}_{\sigma(1)}, \hat{y}_{\sigma(2)}, \dots, \hat{y}_{\sigma(k-1)}) \in \mathcal{X} \times \{0, 1\}^{k-1}$$

as input, where  $\sigma$  is some permutation of  $[K]$ .

Practically, it turns out that the order of labels on the chain, defined by  $\sigma$ , has an impact on predictive performance [2,15]. As finding an optimal order appears to be difficult, [15] suggest to use an ensemble of CCs over a (randomly chosen) set of permutations and combine their predictions. In the original CC, the final prediction is derived in a label-wise manner using BMV. In a probabilistic variant of CC, we allow each classifier  $h_k$ ,  $k \in [K]$ , to produce a score in  $[0, 1]$ , namely an estimation of the conditional probability

$$\mathbf{p}(y_{\sigma(k)} = 1 \mid \mathbf{x}, y_{\sigma(1)}, \dots, y_{\sigma(k-1)}) . \quad (12)$$

The score (12) can be seen as a “dependent” marginal probability, i.e., a marginal probability which to some extent takes label dependence into account.

- **Ensembles of Multi-Objective Decision Trees (EMODT)** are a computationally efficient EMLC method [8]. Similar to conventional decision trees (DT) [12,13], a multi-objective decision tree (MODT) partitions the instance space  $\mathcal{X}$  into (axis-parallel) regions  $R_1, \dots, R_L$  (i.e.,  $\bigcup_{i=1}^L R_i = \mathcal{X}$  and  $R_i \cap R_j = \emptyset$  for  $i \neq j$ ), corresponding to individual leaves of the tree. In a probabilistic setting, each leaf of the MODT is associated with a complete marginal probability vector, where the marginal probability corresponding to a particular label is simply estimated as the proportion of the training instances in the leaf for which the label is relevant. The binary label vector predicted by an EMODT can be derived with GMV on the probability vectors provided by the individual MODTs in a label-wise manner. Due to the label-wise voting, EMODT is also tailored to decomposable performance measures.

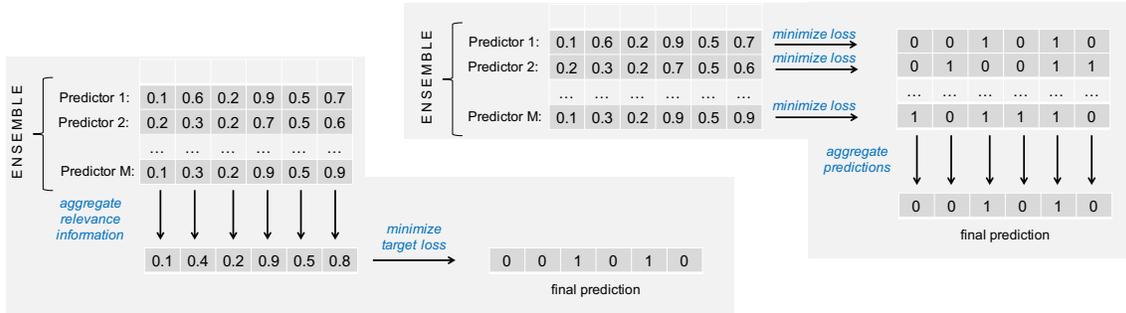


Fig. 1: Illustration of the “combine then predict” (red) and “predict then combine” (blue) approaches for the case where relevance information consists of marginal probabilities.

## 4 A Formal Framework

In the following, we define a formal framework for ensembles of multi-label classifiers.

### 4.1 Intermediate Relevance Information

Most MLC methods are two-step approaches in the sense that, prior to making a final prediction  $\hat{\mathbf{y}} \in \mathcal{Y}$ , intermediate results about the relevance of labels, their interdependencies, or similar information is compiled. We refer to such results as *relevance information*, which we distinguish from the final prediction. Important examples include the following:

- Estimates of *marginal probabilities* (1), which provide important information for the minimization of decomposable loss functions, or loss minimization in the case of label independence.
- The entire *joint distribution*  $\mathbf{p}(\cdot | \mathbf{x})$ , which might be needed for the minimization of non-decomposable losses in cases where the labels are not independent.
- *Probability estimates of a more general kind*. For example, [22] require the probabilities  $\mathbf{p}(y_k = 1, s_{\mathbf{y}} = s)$ ,  $k, s \in [K]$ , for loss minimization in the case of the F-measure.

In general, of course, the relevance information does not need to be probabilistic, but might be of a more general nature.

### 4.2 CTP versus PTC

In the context of ensemble learning, an important distinction between methods can be made depending on whether the relevance information provided by the different ensemble members is combined first, and a prediction is obtained afterwards, or whether individual predictions are produced first and then combined into an overall prediction (see Fig. 1 for an illustration). We refer to the former as “combine then predict” (CTP) and the latter as “predict then combine” (PTC).

In CTP, the relevance information  $\mathcal{R} = \{R_1, \dots, R_M\}$  provided by the individual ensemble members is first combined into a single condensed representation

$$R = \text{CTP-agg}(R_1, \dots, R_M). \quad (13)$$

Then, a final prediction  $\hat{\mathbf{y}}$  is produced on the basis of this representation, typically (though not necessarily) taking the underlying target loss  $\ell$  into account, i.e., minimizing expected loss with regard to  $\ell$  (cf. Section 2.2). Denoting the prediction step by  $\text{Pred}_\ell$ , this can be written compactly as follows:

$$\hat{\mathbf{y}} = \text{Pred}_\ell \left( \text{CTP-agg} (R_1, \dots, R_M) \right). \quad (14)$$

In PTC, each member of the ensemble first predicts a (loss minimizing) label combination  $\hat{\mathbf{y}}_j = \text{Pred}_\ell(R_j)$ . Then, in a second step, these predictions  $\hat{\mathbf{y}}_m, m \in [M]$ , are combined into an overall prediction  $\hat{\mathbf{y}}$  using a suitable aggregation function:

$$\hat{\mathbf{y}} = \text{PTC-agg} \left( \text{Pred}_\ell(R_1), \dots, \text{Pred}_\ell(R_M) \right).$$

Note that the commonly used techniques of weighted and binary voting as described in Section 3 can be seen as specific instantiations of CTP and PTC: Binary majority voting (BMV) first maps vectors of marginal label probabilities into label predictions, which are then combined via majority voting, and is thus an instance of PTC. Graded majority voting (GMV) first adds up the label probabilities into a single vector of marginal label probabilities, which are then thresholded for a final prediction, and is thus a special case of CTP. However, both voting methods are oblivious to specific loss functions.

### 4.3 Aggregation in CTP

The information that needs to be combined in both approaches, CTP and PTC, is of different nature. Thus, one may expect different types of aggregation functions to be suitable. In particular, relevance information to be combined in CTP is often *gradual* and represented in numerical form — probability estimates is again a typical example. Information of that kind is often reasonably combined through *averaging*. For instance, the arithmetic mean

$$\bar{p}_i = \frac{1}{M} \sum_{j=1}^M p_{i,j}, \quad (15)$$

produced by the ensemble members for the label  $\lambda_i$ , will be an improved estimate of the true marginal probability of that label. Of course, aggregation functions other than the arithmetic mean are also conceivable; for example, the median is known to be more robust toward outliers.

Moreover, aggregation does not necessarily need to be label-wise as in (15). Instead, it depends on what kind of relevance information is produced in the first place. Imagine, for example, that each ensemble member yields an estimate  $\hat{\mathbf{p}}_j(\cdot | \mathbf{x})$  of the joint label distribution on  $\mathcal{Y}$ . Aggregation should then be done at the same level, and averaging is again an obvious way for doing so:

$$\hat{\mathbf{p}}(\mathbf{y} | \mathbf{x}) = \frac{1}{M} \sum_{j=1}^M \hat{\mathbf{p}}_j(\mathbf{y} | \mathbf{x}), \forall \mathbf{y} \in \mathcal{Y}.$$

As already said, an approach of that kind might be advantageous in the case of non-decomposable losses, although it will not be tractable in general.

#### 4.4 Aggregation in PTC

In PTC, the problem is to combine (binary) predictions. More specifically, recalling the goal to minimize a given target loss  $\ell$ , the problem can be stated as follows: Given predictions  $\hat{y}_1, \dots, \hat{y}_M$ , which are all supposed to minimize  $\ell$  in expectation, what is a Bayes-optimal overall prediction  $\hat{y}$ ? The answer to this question is far from obvious and, to the best of our knowledge, has not been studied systematically in the literature so far. In fact, a formal analysis of this problem probably presupposes additional assumptions about how the predictions (15) may differ from the true Bayes-optimal prediction (obviously, they cannot all be Bayes-optimal at the same time, unless they all coincide).

In any case, it should be clear that averaging will be less suitable. First of all, binary predictions  $\hat{y}$  are *discrete* entities, and by averaging them one does not again end up with a discrete entity. This is to some extent comparable to the difference between ensemble *regression* (numerical case) and ensemble *classification* (categorical case): While arithmetic averaging is often used in the former, counting or “voting” techniques are more commonly applied in the latter. Second, even when solving this technical issue by turning an average into a discrete entity, for example by thresholding, undesirable effects might be produced, as shown by a simple example, in which the (conditional) ground-truth distribution  $p(\cdot | \mathbf{x})$  on the label space  $\mathcal{Y} = \{0, 1\}^3$  are given as follows:

$\mathbf{y}$	(0, 0, 0)	(1, 1, 1)	(0, 1, 1)	(1, 0, 1)	(1, 1, 0)
$p(\mathbf{y}   \mathbf{x})$	1/4	3/16	3/16	3/16	3/16

Obviously, the Bayes-optimal prediction for the subset 0/1 loss is (0, 0, 0), and ideally, this prediction is produced by each classifier in the ensemble. Now, since these classifiers are not perfect, suppose that the different label combinations are predicted in proportion to their conditional probabilities, i.e., (0, 0, 0) is predicted with probability 1/4, (1, 1, 1) with probability 3/16, etc. One easily verifies that, for each of the three labels, the probability of it being predicted as relevant (9/16) exceeds the probability for irrelevant (7/16). Therefore, by taking the arithmetic average over the ensemble members’ predictions, and then thresholding at 1/2, one will likely end up with the suboptimal prediction (1, 1, 1).

The reader may have noticed that this example is actually less problematic for the Hamming loss, for which the prediction (1, 1, 1) is indeed Bayes-optimal, and would be produced by the label-wise aggregation sketched above. More generally, it is plausible that a label-wise combination of predictions is indeed suitable for decomposable losses like Hamming, but suboptimal for non-decomposable losses.

Based on the discussion so far, we propose two aggregation functions for PTC, which can be seen as implementations of different types of voting, and will be used in our experimental study below:

- **Label-wise voting (PTC-lw):** For each individual label  $\lambda_i$ , the number of positive (relevant) and negative (irrelevant) votes in the predictions  $\hat{y}_1, \dots, \hat{y}_M$  is counted, and the majority is adopted.

Table 1: Datasets used in the experiments

#	Name	# Inst.	# Nom. Feat.	# Num. Feat.	# Lab.
1	Cal500	502	0	68	174
2	Emotions	593	0	72	6
3	Scene	2407	0	294	6
4	Yeast	2417	0	103	14
5	Mediamill	43907	0	120	101
6	Flags	194	9	10	7
7	Medical	978	1449	0	45
8	Bibtex	7395	1836	0	159

- **Mode (PTC-mode):** Counting is done at the level of the entire predictions, i.e., we predict the label combination  $\hat{\mathbf{y}}$  that occurs most frequently:

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \sum_{j=1}^M \llbracket \hat{\mathbf{y}}_i = \bar{\mathbf{y}}_j \rrbracket. \quad (16)$$

In case the maximum is not unique, ties are broken by choosing the maximal prediction with the highest score

$$s(\hat{\mathbf{y}}) = \sum_{k=1}^K \sum_{j=1}^M \llbracket \hat{\mathbf{y}}_k = \hat{\mathbf{y}}_{k,j} \rrbracket. \quad (17)$$

## 5 Experimental Evaluation

We perform experiments on eight standard benchmark datasets (cf. Table 1) from the MULAN repository<sup>5</sup>, following a 10-fold cross-validation procedure. Our primary goal is to confirm that the loss-based aggregation methods PTC and CTP outperform the commonly used voting techniques. Moreover, we conjecture that CTP performs better than PTC for decomposable losses, and PTC better than CTP for non-decomposable losses. This is because accurate marginal probabilities are of utmost importance for decomposable losses — which is exactly what CTP accomplishes through averaging label-wise predictions. Likewise, PTC is more apt at capturing label dependencies, which is important for non-decomposable losses, because it aggregates over several predictions tailored to the target loss (instead of producing only a single one, as CTP).

We conducted three series of experiments using the ensemble methods EMODT, EBR, and ECC with their cardinality set to 50. We employed logistic regression as the base classifiers for EBR and ECC and let them produce probabilistic predictions. Thus, each ensemble member provides a complete marginal probability vector.

The detailed results are shown in Table 2. The best way for getting an insight into the respective performances is to consider the averages of these ranks in the final column of the table. In particular, each column shows the results of one dataset, each line shows the results of a combination of ensemble technique, loss function, and aggregation technique. For each combination of ensemble, loss function, and dataset, we also report the respective ranks for the obtained losses over the aggregation approaches. The bold value indicates the best performance on each data set. According to

<sup>5</sup> <http://mulan.sourceforge.net/datasets.html>

Table 2: Predictive performance (in percent) and rank (small number) of aggregation methods with respect to the Hamming loss, the subset 0/1 loss and the F1-measure.

		Cal500	Emo- tions	Scene	Yeast	Flags	Medi- cal	Bibtex	Media- mill	Avg. ranks
<b>EMODT</b>										
Hamming loss ↓	GMV	<b>13.65</b> <sub>1</sub>	18.72 <sub>2</sub>	9.45 <sub>3</sub>	19.57 <sub>2</sub>	<b>23.48</b> <sub>1</sub>	1.55 <sub>2</sub>	1.27 <sub>3</sub>	2.67 <sub>3</sub>	2.13
	BMV	13.75 <sub>2</sub>	<b>18.52</b> <sub>1</sub>	9.11 <sub>2</sub>	<b>19.38</b> <sub>1</sub>	25.09 <sub>3</sub>	1.59 <sub>3</sub>	<b>1.25</b> <sub>1.5</sub>	<b>2.65</b> <sub>1.5</sub>	<b>1.88</b>
	CTP	equivalent to GMV								
	PTC-lw	equivalent to BMV								
	PTC-mode	14.41 <sub>3</sub>	19.66 <sub>3</sub>	<b>8.14</b> <sub>1</sub>	19.77 <sub>3</sub>	24.58 <sub>2</sub>	<b>1.52</b> <sub>1</sub>	<b>1.25</b> <sub>1.5</sub>	<b>2.65</b> <sub>1.5</sub>	2.00
Subset 0/1 loss ↓	GMV	<b>100</b> <sub>2</sub>	69.82 <sub>3</sub>	49.48 <sub>3</sub>	85.48 <sub>3</sub>	79.34 <sub>2</sub>	55.70 <sub>2</sub>	87.59 <sub>3</sub>	84.79 <sub>3</sub>	2.63
	BMV	<b>100</b> <sub>2</sub>	67.81 <sub>2</sub>	46.16 <sub>2</sub>	83.62 <sub>2</sub>	83.00 <sub>3</sub>	57.26 <sub>3</sub>	86.61 <sub>2</sub>	84.57 <sub>2</sub>	2.25
	CTP	equivalent to GMV								
	PTC-lw	equivalent to BMV								
	PTC-mode	<b>100</b> <sub>2</sub>	<b>64.45</b> <sub>1</sub>	<b>27.21</b> <sub>1</sub>	<b>74.43</b> <sub>1</sub>	<b>78.82</b> <sub>1</sub>	<b>50.96</b> <sub>1</sub>	<b>85.22</b> <sub>1</sub>	<b>79.21</b> <sub>1</sub>	<b>1.13</b>
F1-measure ↑	GMV	33.31 <sub>5</sub>	58.18 <sub>5</sub>	53.45 <sub>5</sub>	58.53 <sub>5</sub>	74.91 <sub>4</sub>	51.17 <sub>4</sub>	25.33 <sub>5</sub>	59.52 <sub>5</sub>	4.75
	BMV	37.33 <sub>4</sub>	61.79 <sub>4</sub>	57.02 <sub>4</sub>	60.99 <sub>4</sub>	74.33 <sub>5</sub>	50.72 <sub>5</sub>	27.69 <sub>4</sub>	60.74 <sub>4</sub>	4.25
	CTP	<b>48.31</b> <sub>1</sub>	<b>68.30</b> <sub>1</sub>	74.90 <sub>2</sub>	<b>66.61</b> <sub>1</sub>	75.89 <sub>3</sub>	<b>76.17</b> <sub>1</sub>	<b>48.77</b> <sub>1</sub>	<b>63.88</b> <sub>1</sub>	<b>1.38</b>
	PTC-lw	46.45 <sub>2</sub>	68.29 <sub>2</sub>	71.78 <sub>3</sub>	64.87 <sub>3</sub>	76.21 <sub>2</sub>	71.81 <sub>3</sub>	37.79 <sub>3</sub>	62.71 <sub>2</sub>	2.50
	PTC-mode	42.30 <sub>3</sub>	68.25 <sub>3</sub>	<b>77.55</b> <sub>1</sub>	65.38 <sub>2</sub>	<b>76.48</b> <sub>1</sub>	75.38 <sub>2</sub>	45.59 <sub>2</sub>	62.09 <sub>3</sub>	2.13
<b>ECC</b>										
Hamming loss ↓	GMV	14.08 <sub>2</sub>	<b>20.28</b> <sub>1</sub>	<b>8.63</b> <sub>1</sub>	20.24 <sub>2</sub>	23.26 <sub>2</sub>	0.89 <sub>3</sub>	<b>1.28</b> <sub>1</sub>	—	<b>1.71</b>
	BMV	<b>14.06</b> <sub>1</sub>	20.31 <sub>2</sub>	8.76 <sub>2</sub>	<b>20.11</b> <sub>1</sub>	23.48 <sub>3</sub>	<b>0.88</b> <sub>1.5</sub>	1.29 <sub>2.5</sub>	—	1.86
	CTP	equivalent to GMV								
	PTC-lw	equivalent to BMV								
	PTC-mode	14.24 <sub>3</sub>	20.48 <sub>3</sub>	8.77 <sub>3</sub>	20.39 <sub>3</sub>	<b>22.52</b> <sub>1</sub>	<b>0.88</b> <sub>1.5</sub>	1.29 <sub>2.5</sub>	—	2.43
Subset 0/1 loss ↓	GMV	<b>100</b> <sub>2</sub>	69.46 <sub>2</sub>	32.82 <sub>2</sub>	79.35 <sub>3</sub>	72.32 <sub>2</sub>	29.49 <sub>3</sub>	81.61 <sub>2</sub>	—	2.29
	BMV	<b>100</b> <sub>2</sub>	70.47 <sub>3</sub>	33.07 <sub>3</sub>	78.98 <sub>2</sub>	74.87 <sub>3</sub>	28.45 <sub>2</sub>	81.64 <sub>3</sub>	—	2.57
	CTP	equivalent to GMV								
	PTC-lw	equivalent to BMV								
	PTC-mode	<b>100</b> <sub>2</sub>	<b>68.45</b> <sub>1</sub>	<b>29.04</b> <sub>1</sub>	<b>77.28</b> <sub>1</sub>	<b>66.61</b> <sub>1</sub>	<b>28.35</b> <sub>1</sub>	<b>81.37</b> <sub>1</sub>	—	<b>1.14</b>
F1-measure ↑	GMV	32.13 <sub>5</sub>	63.20 <sub>5</sub>	72.84 <sub>4</sub>	62.63 <sub>5</sub>	72.80 <sub>4</sub>	81.74 <sub>4</sub>	40.08 <sub>5</sub>	—	4.57
	BMV	32.50 <sub>4</sub>	64.00 <sub>4</sub>	71.75 <sub>5</sub>	62.81 <sub>4</sub>	72.56 <sub>5</sub>	81.10 <sub>5</sub>	40.10 <sub>4</sub>	—	4.43
	CTP	<b>45.28</b> <sub>1</sub>	<b>67.65</b> <sub>1</sub>	<b>77.05</b> <sub>1</sub>	<b>64.85</b> <sub>1</sub>	74.21 <sub>2</sub>	<b>85.17</b> <sub>1</sub>	<b>49.30</b> <sub>1</sub>	—	<b>1.14</b>
	PTC-lw	42.05 <sub>2</sub>	67.58 <sub>2</sub>	75.52 <sub>3</sub>	64.69 <sub>2</sub>	<b>74.53</b> <sub>1</sub>	84.85 <sub>2</sub>	48.95 <sub>2</sub>	—	2.00
	PTC-mode	41.98 <sub>3</sub>	66.81 <sub>3</sub>	75.57 <sub>2</sub>	64.06 <sub>3</sub>	74.13 <sub>3</sub>	84.58 <sub>3</sub>	48.77 <sub>3</sub>	—	2.86
<b>EBR</b>										
Hamming loss ↓	GMV	<b>13.95</b> <sub>1</sub>	<b>20.15</b> <sub>1</sub>	9.82 <sub>2</sub>	19.89 <sub>3</sub>	24.05 <sub>2</sub>	<b>0.92</b> <sub>2</sub>	<b>1.22</b> <sub>1.5</sub>	—	<b>1.79</b>
	BMV	14.02 <sub>2</sub>	20.32 <sub>3</sub>	9.83 <sub>3</sub>	<b>19.87</b> <sub>1</sub>	<b>23.62</b> <sub>1</sub>	<b>0.92</b> <sub>2</sub>	<b>1.22</b> <sub>1.5</sub>	—	1.93
	CTP	equivalent to GMV								
	PTC-lw	equivalent to BMV								
	PTC-mode	14.10 <sub>3</sub>	20.21 <sub>2</sub>	<b>9.76</b> <sub>1</sub>	19.88 <sub>2</sub>	24.49 <sub>3</sub>	<b>0.92</b> <sub>2</sub>	1.23 <sub>3</sub>	—	2.29
Subset 0/1 loss ↓	GMV	<b>100</b> <sub>2</sub>	<b>72.86</b> <sub>1</sub>	46.11 <sub>3</sub>	84.78 <sub>3</sub>	82.45 <sub>3</sub>	30.71 <sub>3</sub>	82.12 <sub>3</sub>	—	2.57
	BMV	<b>100</b> <sub>2</sub>	73.53 <sub>3</sub>	46.03 <sub>2</sub>	84.49 <sub>2</sub>	81.92 <sub>2</sub>	30.39 <sub>2</sub>	81.88 <sub>2</sub>	—	2.14
	CTP	equivalent to GMV								
	PTC-lw	equivalent to BMV								
	PTC-mode	<b>100</b> <sub>2</sub>	73.37 <sub>2</sub>	<b>45.53</b> <sub>1</sub>	<b>84.32</b> <sub>1</sub>	<b>80.87</b> <sub>1</sub>	<b>30.09</b> <sub>1</sub>	<b>81.72</b> <sub>1</sub>	—	<b>1.29</b>
F1-measure ↑	GMV	34.17 <sub>5</sub>	58.38 <sub>4</sub>	61.99 <sub>5</sub>	61.37 <sub>5</sub>	72.60 <sub>4</sub>	79.18 <sub>5</sub>	39.75 <sub>5</sub>	—	4.71
	BMV	34.40 <sub>4</sub>	58.11 <sub>5</sub>	62.07 <sub>4</sub>	61.49 <sub>4</sub>	73.42 <sub>3</sub>	79.76 <sub>4</sub>	40.17 <sub>4</sub>	—	4.00
	CTP	<b>47.62</b> <sub>1</sub>	66.67 <sub>3</sub>	76.14 <sub>3</sub>	<b>65.07</b> <sub>1</sub>	<b>75.18</b> <sub>1</sub>	<b>85.27</b> <sub>1</sub>	<b>50.78</b> <sub>1</sub>	—	<b>1.57</b>
	PTC-lw	47.46 <sub>2</sub>	66.96 <sub>2</sub>	<b>76.29</b> <sub>1</sub>	65.01 <sub>2</sub>	74.80 <sub>2</sub>	84.63 <sub>2</sub>	49.02 <sub>3</sub>	—	2.00
	PTC-mode	47.33 <sub>3</sub>	<b>67.07</b> <sub>1</sub>	76.18 <sub>2</sub>	64.99 <sub>3</sub>	71.41 <sub>5</sub>	84.58 <sub>3</sub>	49.43 <sub>2</sub>	—	2.71

the Friedman/Nemenyi test, differences are statistically significant for a critical distance between the average ranks of 1.10/1.25 for  $\alpha=0.1/0.05$  for EMODT, and similarly 1.94/2.16 and 2.08/2.31 for ECC and EBR, respectively. The Friedman test fails for all Hamming loss comparisons.

- **Loss-based aggregation vs. voting.** Especially for F1 and subset 0/1 loss, there are (statistically significant) large differences between the voting-based decompositions on the one side, and PTC/CTP on the other side. This confirms our expectation that GMV and BMV are poorly suited for the case of non-decomposable performance measures. Only for Hamming loss, the voting-based techniques are in the same range, and, in fact, sometimes even better (yet, no significant difference).

This result is also expected because Hamming loss is decomposable, so that its performance primarily depends on accurate marginal probabilities. In fact, in this case, label-wise PTC and CTP are equivalent to binary and graded voting, respectively. For subset 0/1 loss, assuming label independence and marginal probability as the relevance information, our loss-based instantiations of CTP and PTC-lw are equivalent to GMV and BMV, respectively. As can be seen from the results, however, this assumption is most likely invalid for the investigated datasets, because PTC-mode, which addresses the problem of finding the mode of the joint label distribution, typically outperforms the alternatives.

- **PTC vs. CTP.** With respect to the two different approaches, the mode-based PTC decomposition performs significantly better for subset 0/1 loss, whereas CTP (or, in this case, equivalently GMV) seems to perform better for Hamming loss. These results provide clear evidence in favor of our conjecture. The results for F1 are a bit more difficult to interpret but also consistent. Given marginal probabilities, we derive the loss minimizer for F1 under the assumption of label independence, and in this case, accurate marginal probabilities are again crucial. This is probably the reason for why CTP has an advantage over PTC.

We also conducted a series of experiments using EMODT with the number of ensemble members varying from 1 to 100 ( $M \in \{1, 5, 10, 20, 30, \dots, 100\}$ ). Here, our interest was to study the influence of the ensemble size on the performance of the aggregation methods. For each value of the ensemble cardinality, we have run a 10 times 10-fold cross-validation, for which we report the average scores. As expected, the results shown in Figure 2 confirm that the MLC scores typically improve with an increasing size of the ensembles. This is in agreement with the observation on the performance of ECC reported in [9]. More importantly, we also see differences between the different aggregation methods, and that suitable instantiations of CTP and PTC can indeed reach better performance than standard voting techniques. In particular, the visible gaps for the subset 0/1 loss re-confirm the superiority of PTC-mode for non-decomposable losses. Finally, we note that the performances change rapidly in the beginning and tend to converge when the number of ensemble members reaches moderate values (i.e., 30 or 40), except for the subset 0/1 loss and PTC-mode. This is again in agreement with our expectations, because PTC-mode does voting at the level of the entire predictions, and the number of possible predictions increases exponentially with the number of labels, so that more iterations are necessary for convergence. A similar effect can be observed for PTC-lw/BMV, whose label-wise votings converge less rapidly to accurate marginal probability estimates than CTP/GMV, but are able to catch up with increasing number of votes. For EMODT, there seems to be even an advantage in the end for using the vote distributions,

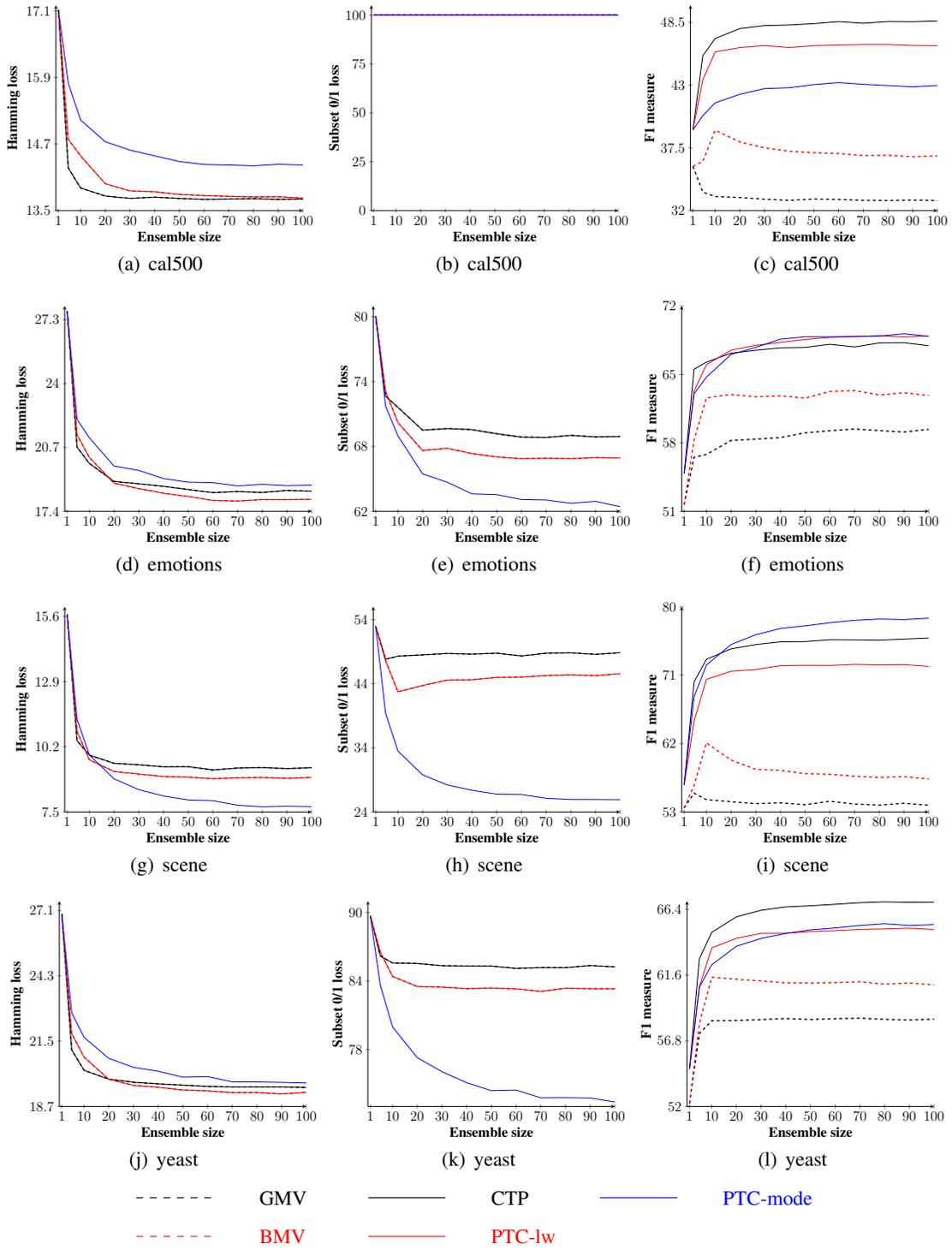


Fig. 2: Predictive performance ( $y$ -axis) of aggregation methods as a function of the cardinality of ensembles ( $x$ -axis) in terms of Hamming loss (left column), subset 0/1 loss (middle), and F1 (right column) for four datasets.

possibly due to less accurate probability estimates of the trees. Results similar to those shown in Figure 2 have been obtained for EBR and ECC.

## 6 Conclusion

This paper studied the question of how to aggregate the predictions of individual members of an ensemble of multilabel classifiers in a systematic way. We introduced a formal framework of ensemble multi-label classification, in which we distinguish two principal approaches, referred to as “predict then combine” (PTC) and “combine then predict” (CTP). Both approaches generalize voting techniques commonly used for EMLC, while allowing one to explicitly take the target performance measure into account. Our framework supports the analysis of existing EMLC methods as well as the systematic development of new ones. Besides, it suggests a number of interesting theoretical problems, like the question of how to combine predictions in PTC in a provably optimal way. Experimentally, we showed that standard voting techniques are indeed outperformed by suitable instantiations of CTP and PTC. Moreover, our results suggest that CTP performs well for decomposable loss functions, whereas PTC is the better choice for non-decomposable losses.

## References

1. Breiman, L.: Bagging predictors. *Machine Learning* **24**(2), 123–140 (1996)
2. Cheng, W., Hüllermeier, E., Dembczyński, K.J.: Bayes optimal multilabel classification via probabilistic classifier chains. In: *Proceedings of the 27th International Conference on Machine Learning (ICML)*. pp. 279–286 (2010)
3. Dembczyński, K., Waegeman, W., Cheng, W., Hüllermeier, E.: On label dependence and loss minimization in multi-label classification. *Machine Learning* **88**(1-2), 5–45 (2012)
4. Dembczyński, K., Waegeman, W., Hüllermeier, E.: An analysis of chaining in multi-label classification. In: *Proceedings of the 20th European Conference on Artificial Intelligence (ECAI)*. pp. 294–299. IOS Press (2012)
5. Dietterich, T.G.: Ensemble methods in machine learning. In: *Proceedings of the 1st International Workshop on Multiple Classifier Systems (MCS)*. pp. 1–15. Springer-Verlag (2000)
6. Gharroudi, O.: Ensemble Multi-label Learning in Supervised and Semi-supervised Settings. Ph.D. thesis, Université de Lyon (2017)
7. Gharroudi, O., Elghazel, H., Aussem, A.: Ensemble multi-label classification: A comparative study on threshold selection and voting methods. In: *Proceedings of the 27th IEEE International Conference on Tools with Artificial Intelligence (ICTAI)*. pp. 377–384. IEEE Computer Society (2015)
8. 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)*. pp. 624–631. Springer-Verlag (2007)
9. Li, N., Zhou, Z.H.: Selective ensemble of classifier chains. In: *Proceedings of the 11th International Workshop on Multiple Classifier Systems (MCS)*. vol. 7872, p. 146. Springer (2013)
10. 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)
11. Moyano, J.M., Gibaja, E.L., Cios, K.J., Ventura, S.: Review of ensembles of multi-label classifiers: models, experimental study and prospects. *Information Fusion* **44**, 33–45 (2018)
12. Murthy, S.K.: Automatic construction of decision trees from data: A multi-disciplinary survey. *Data Mining and Knowledge Discovery* **2**(4), 345–389 (1998)
13. Quinlan, J.R.: Induction of decision trees. *Machine Learning* **1**(1), 81–106 (1986)
14. Read, J., Pfahringer, B., Holmes, G., Frank, E.: Classifier chains for multi-label classification. In: *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases: Part II (ECML/PKDD)*. pp. 254–269. Springer-Verlag (2009)
15. Read, J., Pfahringer, B., Holmes, G., Frank, E.: Classifier chains for multi-label classification. *Machine Learning* **85**(3), 333 (2011)

16. Saha, S., Sarkar, D., Kramer, S.: Exploring multi-objective optimization for multi-label classifier ensembles. In: Proceedings of the IEEE Congress on Evolutionary Computation (CEC). pp. 2753–2760. IEEE, Wellington, New Zealand (2019)
17. Shi, C., Kong, X., Fu, D., Yu, P.S., Wu, B.: Multi-label classification based on multi-objective optimization. *ACM Transactions on Intelligent Systems and Technology* **5**(2), 35:1–35:22 (2014)
18. Shi, C., Kong, X., Yu, Philip, S., Wang, B.: Multi-label ensemble learning. In: Proceedings of the Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML-PKDD). pp. 223–239. Springer (2011)
19. Tsoumakas, G.: Random k-labelsets: An ensemble method for multilabel classification. In: Proceedings of the 18th European Conference on Machine Learning (ECML), pp. 406–417 (2007)
20. Tsoumakas, G., Katakis, I., Vlahavas, I.: Mining multi-label data. In: *Data Mining and Knowledge Discovery Handbook*, pp. 667–685. Springer (2009)
21. Tsoumakas, G., Katakis, I., Vlahavas, I.: Random k-labelsets for multilabel classification. *IEEE Transactions on Knowledge and Data Engineering* **23**(7), 1079–1089 (2010)
22. Waegeman, W., Dembczyński, K., Jachnik, A., Cheng, W., Hüllermeier, E.: On the Bayes-optimality of F-measure maximizers. *Journal of Machine Learning Research* **15**(1), 3333–3388 (2014)
23. 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)