

# Consistency of Probabilistic Classifier Trees

Krzysztof Dembczyński<sup>1</sup>✉, Wojciech Kotłowski<sup>1</sup>, Willem Waegeman<sup>2</sup>,  
Róbert Busa-Fekete<sup>3</sup>, and Eyke Hüllermeier<sup>3</sup>

<sup>1</sup> Poznan University of Technology, Poland  
[kdembczynski,wkotlowski]@cs.put.poznan.pl

<sup>2</sup> Ghent University, Belgium  
willem.waegeman@ugent.be

<sup>3</sup> Paderborn University, Germany  
[busarobi,eyke]@upb.de

**Abstract.** Label tree classifiers are commonly used for efficient multi-class and multi-label classification. They represent a predictive model in the form of a tree-like hierarchy of (internal) classifiers, each of which is trained on a simpler (often binary) subproblem, and predictions are made by (greedily) following these classifiers’ decisions from the root to a leaf of the tree. Unfortunately, this approach does normally not assure consistency for different losses on the original prediction task, even if the internal classifiers are consistent for their subtask. In this paper, we thoroughly analyze a class of methods referred to as probabilistic classifier trees (PCTs). Thanks to training probabilistic classifiers at internal nodes of the hierarchy, these methods allow for searching the tree-structure in a more sophisticated manner, thereby producing predictions of a less greedy nature. Our main result is a regret bound for 0/1 loss, which can easily be extended to ranking-based losses. In this regard, PCTs nicely complement a related approach called filter trees (FTs), and can indeed be seen as a natural alternative thereof. We compare the two approaches both theoretically and empirically.

## 1 Introduction

Multi-class and multi-label classification problems are nowadays characterized not only by large sample sizes and feature spaces, but also by a large number of labels. In application fields like image classification [12], text classification [8], online advertising [3], and video recommendation [23], it is not uncommon to deal with tens or hundreds of thousands [11], or even millions of labels [20].

*Label tree classifiers* belong to the most efficient approaches for problems at this scale [2]. In this approach, a solution to the original problem is represented in the form of a hierarchy of classifiers, each of which is trained on a simpler subproblem. A prediction for a new example is then derived from the predictions of these (internal) classifiers, each of which corresponds to a node in the tree-like hierarchical structure; typically, each label in the original classification problem is uniquely represented by a path from the root to a leaf of that tree.

However, combining conventional training of the internal classifiers with greedy inference, namely, following a single root-to-leaf path in the tree, does not guarantee consistency of this approach [4,10]. Thus, even perfect (zero regret) classifiers in each node of the tree do not imply a perfect (global) classification of new examples. There are two ways to remedy this problem: adjusting training and adjusting inference. The first idea is to modify the training of the internal classifiers so as to assure the consistency of greedy inference later on. The second approach, while training more conventionally, guarantees consistency by searching the tree-structure for an optimal prediction in a less greedy way.

The first idea is realized by the *filter tree* (FT) approach [4]. By constructing label trees in a bottom-up manner, an internal classifier can anticipate the decisions of its successor classifiers, and exploit this information to properly condition its own behavior to these classifiers. In the case of 0/1 loss, this is accomplished thanks to a specific filter technique, which removes examples from the training data on which successor classifiers made incorrect predictions. For this training procedure, a regret bound connecting the global performance with the average performance of node classifiers can be proved [4]. This bound can be generalized from 0/1 loss to any cost-based loss function, albeit at the price of a more expensive training procedure; ranking-based losses, which require the ordering of labels, cannot be tackled by FTs. Since inference can be done in a greedy way, the complexity of prediction is only logarithmic in the number of labels. More recently, the training of FTs has been further improved in the context of multi-label classification [17].

The second approach ensures consistency thanks to more sophisticated search of label trees in the inference phase [10,16,18]. To this end, probabilistic classifiers in each node of the tree are required, which allow for assessing the usefulness of different search directions. Label trees with probabilistic classifiers have already been considered in multi-class classification under the name of conditional probability trees [3] and nested dichotomies [14]. In multi-label classification, a similar approach has been referred to as probabilistic classifier chains [9]. The same concept also appears in neural networks and natural language processing under the name of hierarchical softmax [19]. In the following we unify all these approaches and jointly refer to them as *probabilistic classifier trees* (PCTs).

We restrict to binary label trees, which are especially natural for multi-label classification; here, each level of the binary tree directly corresponds to one label. Higher order trees (including nodes with more than two children) are often used in multi-class classification. This usually improves the predictive performance at the cost of an increase in prediction time. We also assume the tree structure to be given beforehand, or to have been induced using any of the methods developed for this purpose [3,2,23], and focus on the (orthogonal) problem of how training and prediction should be performed to ensure consistency (given the tree structure).

The main contribution of the paper is a regret bound for PCT in the case of 0/1 loss, which is expressed in terms of the search error and the Kullback-Leibler (KL) divergence (i.e., log-loss regret) of the internal classifiers. The regret bound

implies the consistency of the method, a good "sanity check" for any learning algorithm. Its form quantifies a trade-off between the computational complexity and the statistical accuracy. Moreover, we show that under log-loss we do not theoretically pay any price in terms of performance for representing the joint distribution over classes by a tree structure. Our regret analysis significantly extends and improves the results of [3] for the estimation error of conditional probability trees expressed in terms of squared error loss. We also point out that the bound can be further generalized to ranking-based losses, e.g., recall at  $k$ . We also generalize the tree search algorithms of [10] and [18] to get an anytime  $A^*$ -like algorithm and study its theoretical guarantees, extending the previous results given in [10]. Our theoretical contributions are complemented by a comparison of PCTs with filter trees, both conceptually and experimentally.

The paper is organized as follows. We formally state the problem in Section 2. Section 3 describes PCTs and gives a theoretical analysis of the generalized tree search algorithm. In Section 4, we prove the regret bound for 0/1 loss. Section 5 compares PCTs with other label tree approaches, particularly with conditional probability and filter trees. Section 6 discusses the use of PCTs for predicting top- $k$  labels and its extension to multi-label classification. Section 7 presents experimental results, prior to concluding the paper in Section 8.

## 2 Problem statement

We formalize our problem in the setting of multi-class classification. Let  $(\mathbf{x}, y)$  be an example coming from a probability distribution  $P(\mathbf{X} = \mathbf{x}, Y = y)$  (later denoted  $P(\mathbf{x}, y)$ ) on  $\mathcal{X} \times \mathcal{Y}$ , where  $\mathbf{x} \in \mathcal{X} = \mathbb{R}^d$  and  $y \in \mathcal{Y} = \{1, \dots, m\}$ . A classifier  $h$  predicts a label  $\hat{y} = h(\mathbf{x}) \in \mathcal{Y}$  for each  $\mathbf{x} \in \mathcal{X}$ . The prediction accuracy of  $h$  can be measured in terms of 0/1 loss:<sup>4</sup>

$$\ell_{0/1}(y, h(\mathbf{x})) = \llbracket y \neq h(\mathbf{x}) \rrbracket$$

We are interested in minimizing the expected loss, also referred to as the *risk*:

$$L_{0/1}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell_{0/1}(y, h(\mathbf{x}))] = \int_{\mathcal{X} \times \mathcal{Y}} \llbracket y \neq h(\mathbf{x}) \rrbracket dP(\mathbf{x}, y)$$

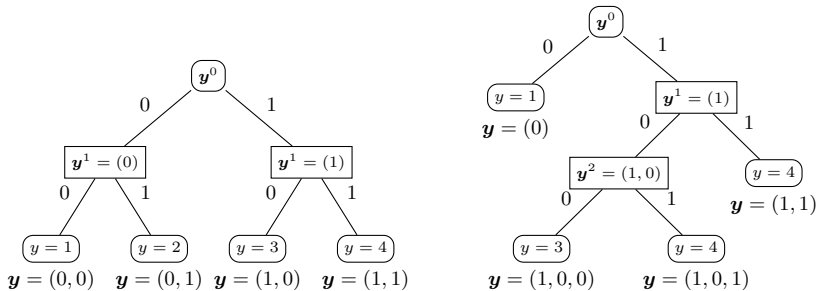
The *Bayes classifier*

$$h^* = \arg \min_h L_{0/1}(h)$$

minimizes the risk among all possible classifiers. While  $h^*$  may not be unique in general, the risk of  $h^*$ , denoted  $L_{0/1}^*$ , is unique, and is called the *Bayes risk*. Decomposing the risk over classes, i.e., writing  $L_{0/1}(h)$  in the form

$$L_{0/1}(h) = \int_{\mathcal{X}} \left( \underbrace{\sum_{y \in \mathcal{Y}} \llbracket y \neq h(\mathbf{x}) \rrbracket P(y|\mathbf{x})}_{=1 - P(h(\mathbf{x})|\mathbf{x})} \right) dP(\mathbf{x}),$$

<sup>4</sup> We use  $\llbracket P \rrbracket$  to denote a number that is 1 if condition  $P$  is satisfied, and 0 otherwise.



**Fig. 1.** Different binary codes in multi-class classification.

reveals that  $h^*$  minimizes risk in a pointwise manner, i.e., for every  $\mathbf{x}$ ,

$$h^*(\mathbf{x}) = \arg \min_{y \in \mathcal{Y}} \{1 - P(y|\mathbf{x})\} = \arg \max_{y \in \mathcal{Y}} P(y|\mathbf{x}).$$

Given a classifier  $h$ , the *regret* of  $h$  is defined as

$$\text{reg}_{0/1}(h) = L_{0/1}(h) - L_{0/1}^* = \int_{\mathcal{X}} \left( P(h^*(\mathbf{x})|\mathbf{x}) - P(h(\mathbf{x})|\mathbf{x}) \right) dP(\mathbf{x}). \quad (1)$$

The regret quantifies the suboptimality of  $h$  compared to the optimal classifier  $h^*$ . The goal is to train a classifier  $h$  with a small regret, ideally equal to zero.

In the following, we assume  $h$  to be represented as a label tree classifier. To this end, we encode the labels  $\{1, \dots, m\}$  using a prefix code. Any such code can be represented by a tree with 0/1 splits. Each path from the root to a leaf node then corresponds to a code word. Recall that codes of fixed length are also prefix codes. Figure 1 shows two examples of coding trees for multi-class classification with 4 classes. Under the coding, we represent each label  $y$  by a binary vector  $\mathbf{y} = (y_1, \dots, y_l)$ , where  $l$  is the maximum length of the code. The set of all code words we denote by  $\mathcal{C}$ . As another special case, consider the problem of multi-label (instead of multi-class) classification, where the goal is to predict the set of labels assigned to a given instance  $\mathbf{x}$ . Such a set can be represented by a binary vector  $\mathbf{y} = (y_1, \dots, y_m)$ , which in turn can be used as a prefix code.

In the label tree approach, we put a binary classifier in each non-leaf node of the tree. An internal node can be uniquely identified by the partial code word  $\mathbf{y}^i = (y_1, \dots, y_i)$ . We denote the root node by  $\mathbf{y}^0$ , which is an empty vector (without any elements). The final prediction is determined by a sequence of decisions of internal classifiers. In the next section, we present a specific instance of the label tree approach that uses probabilistic classifiers in internal nodes of the tree.

### 3 Probabilistic classifier trees

Probabilistic classifier trees (PCTs) are designed to estimate probabilities  $P(y|\mathbf{x})$  by following a path from the root to a leaf node, which corresponds to a code

word  $\mathbf{y} = (y_1, \dots, y_l)$  assigned to label  $y \in \mathcal{Y}$ . Recalling the chain rule of probability, the process corresponds to computing

$$P(y|\mathbf{x}) = P(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^l P(y_i|\mathbf{y}^{i-1}, \mathbf{x}), \quad (2)$$

where  $P(y_i|\mathbf{y}^{i-1}, \mathbf{x})$  are probabilities of  $y_i \in \{0, 1\}$ , estimated in non-leaf nodes  $\mathbf{y}^{i-1}$ . In the next two subsections, training and inference (classification of new examples) for PCT will be discussed in more detail.

### 3.1 Training

Training of PCT naturally decomposes into learning problems over non-leaf nodes of the tree. In each node  $\mathbf{y}^{i-1}$ , the task is to train a probabilistic classifier (e.g., logistic regression) to estimate  $P(y_i|\mathbf{y}^{i-1}, \mathbf{x})$ .

Looking at PCTs as a reduction technique, it is worth mentioning that its training complexity could be much lower than that of the 1-vs-all approach, since each example  $(\mathbf{x}, y)$  is used in only  $l$  instead of  $m$  binary problems, where  $l$  is the height of the tree (i.e.,  $l = \lceil \log_2 m \rceil$  if the tree is balanced). To further improve the training time complexity, one can use online learning methods, such as stochastic gradient descent [5]. Moreover, internal classifiers in PCT can be trained independently of each other, thereby allowing for a massive parallelization of the training procedure. Let us also remark that the learning process can be defined as a single task; this is the so-called one-classifier trick [4], in which a node indicator is used as an additional feature. Alternatively, one can use a separate task for each level of the tree. This approach is used in multi-label classification, as will be discussed in Section 6.

### 3.2 Inference

The classification procedure in PCTs is more involved. To begin with, note that a probability estimate  $Q(y|\mathbf{x})$  for any label  $y$  (given instance  $\mathbf{x}$ ) is obtained quite easily, simply by following the corresponding path in the tree and applying the chain rule:

$$Q(y|\mathbf{x}) = Q(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^l Q(y_i|\mathbf{y}^{i-1}, \mathbf{x})$$

However, being interested in minimization of 0/1 loss, we actually seek to find

$$\hat{\mathbf{y}}^* = \arg \max_{\mathbf{y} \in \mathcal{C}} Q(\mathbf{y}|\mathbf{x}), \quad (3)$$

preferably without computing the probability of each label first. A simple idea is to follow a single path in the tree, starting in the root and always choosing the branch  $y_i \in \{0, 1\}$  for which  $Q(y_i|\mathbf{y}^{i-1}, \mathbf{x}) > 0.5$ . However, while being efficient, this approach is not guaranteed to find the optimal solution [4,10].

---

**Algorithm 1** Inference with  $\epsilon$ -approximate  $A^*$ 


---

```

1: input:  $\mathbf{x}$  (test example)
2: priority list  $\mathcal{Q} \leftarrow \{\mathbf{y}_0\}$  (contains root node initially)
3: priority list  $\mathcal{K} \leftarrow \{\}$  (contains nodes whose both children were not inserted to  $\mathcal{Q}$ )
4:  $\epsilon \leftarrow 2^{-c}$  with  $1 \leq c \leq m$ 
5: while  $\mathcal{Q} \neq \emptyset$  do
6:    $\mathbf{v} \leftarrow$  pop first element in  $\mathcal{Q}$ 
7:   if  $\mathbf{v}$  is a leaf then delete all elements in  $\mathcal{K}$  and break the while loop
8:    $\mathbf{v}_1 \leftarrow (\mathbf{v}, 1)$  (left child of  $\mathbf{v}$ ) and  $\mathbf{v}_0 \leftarrow (\mathbf{v}, 0)$  (right child of  $\mathbf{v}$ )
9:   compute  $E(\mathbf{v}_1 | \mathbf{x})$  and  $E(\mathbf{v}_0 | \mathbf{x})$  recursively from  $E(\mathbf{v} | \mathbf{x})$  using Eqn. (4)
10:  if  $E(\mathbf{v}_1 | \mathbf{x}) \geq \epsilon$  then add  $(\mathbf{v}_1, E(\mathbf{v}_1 | \mathbf{x}))$  to  $\mathcal{Q}$  sorted in descending order of  $E$ 
11:  if  $E(\mathbf{v}_0 | \mathbf{x}) \geq \epsilon$  then add  $(\mathbf{v}_0, E(\mathbf{v}_0 | \mathbf{x}))$  to  $\mathcal{Q}$  sorted in descending order of  $E$ 
12:  if  $\mathbf{v}_1$  and  $\mathbf{v}_0$  are not inserted to  $\mathcal{Q}$  then add  $\mathbf{v}$  to  $\mathcal{K}$  in descending order of  $E$ 
13:  $\theta \leftarrow 0$ 
14: while  $\mathcal{K} \neq \emptyset$  do
15:    $\mathbf{v}' \leftarrow$  pop first element in  $\mathcal{K}$ 
16:    $\mathbf{v}' \leftarrow$  apply greedy search downward on  $\mathbf{v}'$ 
17:   if  $Q(\mathbf{v}' | \mathbf{x}) \geq \theta$  then  $\mathbf{v} \leftarrow \mathbf{v}'$  and  $\theta \leftarrow Q(\mathbf{v}' | \mathbf{x})$ 
18: return  $h_\epsilon(\mathbf{x}) = \hat{\mathbf{y}}_\epsilon = \mathbf{v}$ 

```

---

Better inference methods have been presented in recent years, based on search algorithms such as uniform-cost search [10], beam search [16], and  $A^*$  [18].

All three approaches allow for trading complexity against optimality, and hence for using PCTs in an anytime fashion, thanks to a hyper-parameter  $\epsilon$ . This parameter controls the degree of optimality, i.e., of finding the true loss minimizer (3), as a function of the runtime (it finds a solution  $\hat{\mathbf{y}}_\epsilon$  the conditional probability  $Q(\hat{\mathbf{y}}_\epsilon | \mathbf{x})$  of which is not much worse than the probability of the optimal solution  $\hat{\mathbf{y}}^*$  defined in Eqn. 3). In the analysis that follows, we will use this property to give a formal bound on the error made by such inference algorithms, with a particular focus on uniform-cost and  $A^*$  search. An extension of the analysis to beam search is straightforward and omitted due to lack of space. The pseudo code in Algorithm 1 unifies the approaches of [18] and [10]. This general algorithm, which we denote  $h_\epsilon(\mathbf{x})$ , is a variant of  $A^*$ . It fulfills the anytime property, i.e., the search can be stopped at any time and the algorithm will deliver a valid though possibly suboptimal solution.

Recall that each node in the tree is uniquely defined by a path from the root to this node, i.e., by the partial code word  $\mathbf{y}^i$ . We use  $\mathbf{v}$  to denote the node currently visited by the algorithm, and associate with this node the following value:

$$E(\mathbf{v} | \mathbf{x}) = E(\mathbf{y}^i | \mathbf{x}) = Q(\mathbf{y}^i | \mathbf{x}) \times H(\mathbf{y}^i | \mathbf{x})$$

This value can be interpreted as an approximation of the maximal value of  $Q(\mathbf{y} | \mathbf{x})$ , in which  $Q(\mathbf{y}^i | \mathbf{x})$  is the part of the path that can be computed when moving from the root to node  $\mathbf{v}$ , and  $H(\mathbf{y}^i | \mathbf{x})$  is a heuristic that optimistically guesses the part of the path that has not yet been computed (in the considered case,  $E(\mathbf{y}^i | \mathbf{x})$  has to overestimate or to be the same as the maximal value of

$Q(\mathbf{y} | \mathbf{x})$ ).  $Q(\mathbf{y}^i | \mathbf{x})$  can be computed recursively as follows:  $Q(\mathbf{y}^0 | \mathbf{x}) = 1$  and

$$Q(\mathbf{y}^i | \mathbf{x}) = Q(y_i = 1 | \mathbf{y}^{i-1}, \mathbf{x}) \times Q(\mathbf{y}^{i-1} | \mathbf{x}). \quad (4)$$

In [18], a procedure for computing  $H(\mathbf{y}^i | \mathbf{x})$  is proposed for the specific case of logistic regression as a base learner, whereas the heuristic is simply  $H(\mathbf{y}^i | \mathbf{x}) = 1$  in uniform-cost search used in [10]. The former approach has the advantage of providing a more accurate estimation of maximal  $Q(\mathbf{y} | \mathbf{x})$ , albeit with an additional computing cost, while the latter approach makes a more rough estimation without any additional cost. Interestingly, as shown in experiments in [18], the former approach is still more expensive in terms of the total search cost than the latter.

In a nutshell, Algorithm 1 starts from the root of the label tree, which is the single element of priority list  $\mathcal{Q}$ , sorted in descending order of  $E$ . In every iteration, the top element of the list is popped and the children  $\mathbf{v}_0$  and  $\mathbf{v}_1$  of the corresponding node  $\mathbf{v}$  are visited.  $E(\mathbf{y}^i | \mathbf{x})$  is then recursively computed for the children of node  $\mathbf{v}$ , which are added to the list if this quantity exceeds the threshold  $\epsilon = 2^{-c}$  with  $1 \leq c \leq l$ , where  $l$  is the maximal length of the path in the tree. Basically, they are inserted into the list at the appropriate position, so that the order imposed by  $E(\mathbf{y}^i | \mathbf{x})$  is respected. The first while-loop of the algorithm stops in two situations: (i) when the element popped from the list  $\mathcal{Q}$  corresponds to a leaf of the tree, or (ii) when the list  $\mathcal{Q}$  is empty. The label corresponding to the leaf is then returned in the former case, while in the latter case, inference by greedy search is applied to define a path from all nodes from the list  $\mathcal{K}$ . This list, also sorted in descending order of  $E$ , contains nodes for which none of their children has been added to  $\mathcal{Q}$ . The use of list  $\mathcal{K}$  ensures that by decreasing the value of  $\epsilon$ , the algorithm will always find a solution that is not worse than a solution that would be found with greater  $\epsilon$ .

Algorithm 1 enjoys strong theoretical guarantees. Assuming the cost for computing  $H(\mathbf{y}^i | \mathbf{x})$  to be constant, the following result immediately follows from a theorem proved in [10].

**Theorem 1.** *Let  $1 \leq c \leq l$ . Algorithm 1 with  $\epsilon = 2^{-c}$  needs at most  $\mathcal{O}(l\epsilon^{-1})$  iterations to find a prediction  $h_\epsilon(\mathbf{x}) = \hat{\mathbf{y}}_\epsilon$  such that*

$$Q(\hat{\mathbf{y}}^* | \mathbf{x}) - Q(\hat{\mathbf{y}}_\epsilon | \mathbf{x}) \leq \epsilon - 2^{-l}.$$

From the theorem, we see that the quality of the solution found by the algorithm improves with the length of the running time. Consequently, the algorithm will always find the optimal solution, provided its probability mass is greater than  $\epsilon$ . Reformulating the above, we can say that the algorithm finds the solution in time linear in  $1/q_{\max}$ , where  $q_{\max}$  is the probability mass of the best solution in the estimated distribution  $Q$ . For problems with low noise (high values of  $q_{\max}$ ), this method should work very fast.

The theorem also implies that the greedy search, which corresponds to the algorithm with  $\epsilon = 0.5$ , has very poor guarantees that approach the bound of 0.5 with  $m \rightarrow \infty$ .

## 4 Regret bounds for PCT

In this section, we are concerned with the generalization ability of the PCT classifier, measured by means of the regret (1). Assume for a moment that  $Q(\cdot|\mathbf{x})$ , the label distribution produced by PCT, coincides with the true conditional distribution  $P(\cdot|\mathbf{x})$  for every  $\mathbf{x}$ . Then, if the  $\epsilon$ -approximate inference algorithm is used for classification, Theorem 1 implies the regret of the PCT classifier is at most  $\epsilon$ , i.e., the expected classification error of PCT is at most  $\epsilon$  larger than the expected classification error of the Bayes classifier.

It is, however, unrealistic to assume that PCT is able to perfectly match the true data distribution, hence  $Q(\cdot|\mathbf{x})$  and  $P(\cdot|\mathbf{x})$  will differ in general. Thus, the question arises whether the expected classification error of PCT is still not much worse than the expected classification error of the Bayes classifier if  $Q(\cdot|\mathbf{x})$  and  $P(\cdot|\mathbf{x})$  do not coincide, but are *close* to each other in some sense. This section presents an affirmative answer to this question, delivering a regret bound on the classification error that takes into account the predictive performance of the internal classifiers. More precisely, we bound the PCT regret for 0/1 loss in terms of the difference between  $Q$  and  $P$ , quantified in terms of *log-loss regret*.

We start with a general definition of the log-loss. Consider a problem of estimating a probability distribution on some outcome space  $\mathcal{S}$ . The log-loss of probability estimate  $Q(\cdot)$  on  $\mathcal{S}$  when the observed outcome is  $y \in \mathcal{S}$  is given by

$$\ell_{\log}(y, Q) = -\log Q(y).$$

The log-loss is by far the most popular measure for quantifying the accuracy of probabilistic predictions, and plays an important role in information theory, data compression, and statistics [7] (we briefly analyze the other loss function, squared loss, in Section 5). The *log-loss risk* is the expected log-loss of  $Q(\cdot)$ :

$$L_{\log}(Q) = \mathbb{E}_{y \sim P}[\ell_{\log}(y, Q)],$$

where  $P(\cdot)$  is the true distribution of  $y$ . The log-loss is a *strictly proper loss*, which means that the unique minimizer of the risk is achieved at  $Q(\cdot) \equiv P(\cdot)$  (see, e.g., [21]). We thus define the *log-loss regret* as:

$$\text{reg}_{\log}(Q) = L_{\log}(Q) - L_{\log}(P) = \mathbb{E}_{y \sim P} \left[ \log \frac{P(y)}{Q(y)} \right] = D(P\|Q),$$

where  $D(\cdot\|\cdot)$  is the Kullback-Leibler (KL) divergence.

We now turn back to PCTs. Let us first fix an instance  $\mathbf{x} \in \mathcal{X}$  and consider the distribution over code words  $\mathbf{y} \in \mathcal{C}$ . There are two ways in which log-loss can be used in this setting:

- To measure the quality of the estimate of the joint distribution of labels given  $\mathbf{x}$ ,  $Q(\mathbf{y}|\mathbf{x})$ , i.e., the outcome space is  $\mathcal{S} = \mathcal{C}$ , and the log-loss is  $\ell_{\log}(\mathbf{y}, Q(\cdot|\mathbf{x})) = -\log Q(\mathbf{y}|\mathbf{x})$ . The log-loss regret is then the KL divergence between true joint conditional distribution  $P(\mathbf{y}|\mathbf{x})$  and its estimate  $Q(\mathbf{y}|\mathbf{x})$ ,  $\text{reg}_{\log}(Q(\cdot|\mathbf{x})) = D(P(\cdot|\mathbf{x})\|Q(\cdot|\mathbf{x}))$ .



- To measure the quality of individual classifiers in each node of the tree. Given a node  $\mathbf{y}^{i-1} = (y_1, \dots, y_{i-1})$ , the probability estimate for label  $y_i \in \{0, 1\}$  at this node is  $Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})$ . Thus, the outcome space is  $\mathcal{S} = \{0, 1\}$ , and  $\ell_{\log}(y_i, Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})) = -\log Q(y_i|\mathbf{y}^{i-1}, \mathbf{x})$ . The log-loss regret is then  $\text{reg}_{\log}(Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})) = D(P(\cdot|\mathbf{y}^{i-1}, \mathbf{x})\|Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x}))$ .

Both ways described above turn out to be equivalent. Indeed, we have

$$\begin{aligned} \ell_{\log}(\mathbf{y}, Q(\cdot|\mathbf{x})) &= -\log Q(\mathbf{y}|\mathbf{x}) = \sum_{i=1}^l -\log Q(y_i|\mathbf{y}^{i-1}, \mathbf{x}) \\ &= \sum_{i=1}^l \ell_{\log}(y_i, Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})), \end{aligned}$$

so that the log-loss of the joint distribution is equal to the sum of log-losses of individual node classifiers along the path from the root to leaf  $\mathbf{y}$ . Similarly,

$$\begin{aligned} \text{reg}_{\log}(Q(\cdot|\mathbf{x})) &= \mathbb{E}_{\mathbf{y} \sim P(\cdot|\mathbf{x})} \left[ \log \frac{P(\mathbf{y}|\mathbf{x})}{Q(\mathbf{y}|\mathbf{x})} \right] = \mathbb{E}_{\mathbf{y} \sim P(\cdot|\mathbf{x})} \left[ \sum_{i=1}^l \log \frac{P(y_i|\mathbf{y}^{i-1}, \mathbf{x})}{Q(y_i|\mathbf{y}^{i-1}, \mathbf{x})} \right] \\ &= \mathbb{E}_{\mathbf{y} \sim P(\cdot|\mathbf{x})} \left[ \sum_{i=1}^l \text{reg}_{\log}(Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})) \right], \end{aligned} \quad (5)$$

i.e., the log-loss regret of the joint distribution is equal to the sum of the regrets of node classifiers along the random path from the root to leaf  $\mathbf{y}$ , where  $\mathbf{y}$  is drawn from  $P(\cdot|\mathbf{x})$ . This basically expresses the chain rule for KL divergence [7]. The consequence of the above is that under log-loss we theoretically do not pay any price in terms of performance for representing the joint distribution by a tree structure.

We are now ready to present the main result of this section, which states that the 0/1-regret of the PCT classifier is bounded by means of the sum of log-loss regrets along a random path from the root to the leaf (or, equivalently, by the log-loss regret of the joint distribution) and the search error  $\epsilon$  of the inference procedure.

**Theorem 2.** *Consider PCT, which estimates the probability  $Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})$  in each non-leaf node  $\mathbf{y}^{i-1}$ , and let  $h_\epsilon$  be the classifier which for any  $\mathbf{x}$ , outputs  $\hat{\mathbf{y}}_\epsilon$  found by the  $\epsilon$ -approximate inference procedure (Algorithm 1). Then, for any distribution  $P$ ,*

$$\text{reg}_{0/1}(h_\epsilon) \leq \sqrt{2\text{reg}_{\log}(Q)} + \epsilon - 2^{-l},$$

where  $\text{reg}_{\log}(Q) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} \left[ \sum_{i=1}^l \text{reg}_{\log}(Q(\cdot|\mathbf{y}^i, \mathbf{x})) \right]$  is the expected sum of regrets at internal classifiers along a path from the root to the leaf.

*Proof.* We first condition everything on a fixed  $\mathbf{x}$ . Let  $\mathbf{y}^* = \arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x})$  be the mode of  $P(\cdot|\mathbf{x})$ , and let  $\hat{\mathbf{y}}_\epsilon = h_\epsilon(\mathbf{x})$  be the output of Algorithm 1 for input

$\mathbf{x}$ . Moreover, we let  $\hat{\mathbf{y}}^* = \arg \max_{\mathbf{y}} Q(\mathbf{y}|\mathbf{x})$  denote the mode of  $Q(\cdot|\mathbf{x})$ , and note that from Theorem 1,

$$Q(\hat{\mathbf{y}}^*|\mathbf{x}) - Q(\hat{\mathbf{y}}_\epsilon|\mathbf{x}) \leq \epsilon - 2^{-l}. \quad (6)$$

According to (1), the 0/1-regret of  $\hat{\mathbf{y}}_\epsilon$  conditioned at  $\mathbf{x}$  is given by

$$\text{reg}_{0/1}(\hat{\mathbf{y}}_\epsilon) = P(\mathbf{y}^*|\mathbf{x}) - P(\hat{\mathbf{y}}_\epsilon|\mathbf{x}).$$

Note that the regret is 0 if  $\mathbf{y}^* = \hat{\mathbf{y}}_\epsilon$ , hence we assume  $\mathbf{y}^* \neq \hat{\mathbf{y}}_\epsilon$  in what follows. From the definition of  $\hat{\mathbf{y}}^*$ ,  $Q(\hat{\mathbf{y}}^*|\mathbf{x}) - Q(\mathbf{y}^*|\mathbf{x}) \geq 0$ , which together with (6) gives  $Q(\hat{\mathbf{y}}_\epsilon|\mathbf{x}) - Q(\mathbf{y}^*|\mathbf{x}) + \epsilon - 2^{-l} \geq 0$ . Hence, we obtain the upper bound

$$\begin{aligned} \text{reg}_{0/1}(\hat{\mathbf{y}}_\epsilon) &\leq \left( P(\mathbf{y}^*|\mathbf{x}) - Q(\mathbf{y}^*|\mathbf{x}) \right) + \left( Q(\hat{\mathbf{y}}_\epsilon|\mathbf{x}) - P(\hat{\mathbf{y}}_\epsilon|\mathbf{x}) \right) + \epsilon - 2^{-l} \\ &\leq |P(\mathbf{y}^*|\mathbf{x}) - Q(\mathbf{y}^*|\mathbf{x})| + |Q(\hat{\mathbf{y}}_\epsilon|\mathbf{x}) - P(\hat{\mathbf{y}}_\epsilon|\mathbf{x})| + \epsilon - 2^{-l} \\ &\leq \sum_{\mathbf{y} \in \mathcal{C}} |P(\mathbf{y}|\mathbf{x}) - Q(\mathbf{y}|\mathbf{x})| + \epsilon - 2^{-l}, \end{aligned}$$

where the last inequality is from  $\mathbf{y}^* \neq \hat{\mathbf{y}}_\epsilon$ . We now make use of Pinsker's inequality

$$\frac{1}{2} \sum_{\mathbf{y} \in \mathcal{C}} |P(\mathbf{y}|\mathbf{x}) - Q(\mathbf{y}|\mathbf{x})| \leq \sqrt{\frac{1}{2} D(P(\cdot|\mathbf{x}) \| Q(\cdot|\mathbf{x}))},$$

which together with (5) implies

$$\text{reg}_{0/1}(\hat{\mathbf{y}}_\epsilon) \leq \sqrt{2\mathbb{E}_{\mathbf{y} \sim P(\cdot|\mathbf{x})} \left[ \sum_{i=1}^l \text{reg}_{\log}(Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})) \right]} + \epsilon - 2^{-l}. \quad (7)$$

Note that the 0/1-regret of  $h_\epsilon$ ,  $\text{reg}_{0/1}(h_\epsilon)$ , is just the expectation of the left-hand side of (7) with respect to  $\mathbf{x}$ . Thus, taking expectation on both sides of (7), and using  $\mathbb{E}[\sqrt{\cdot}] \leq \sqrt{\mathbb{E}[\cdot]}$  on the right-hand side (which is Jensen's inequality applied to the concave function  $x \mapsto \sqrt{x}$ ) gives

$$\begin{aligned} \text{reg}_{0/1}(h_\epsilon) &\leq \sqrt{2\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} \left[ \sum_{i=1}^l \text{reg}_{\log}(Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})) \right]} + \epsilon - 2^{-l} \\ &= \sqrt{2\text{reg}_{\log}(Q)} + \epsilon - 2^{-l}. \end{aligned}$$

□

Theorem 2 states that if the log-loss regret of node classifiers is small, the resulting  $\epsilon$ -approximate classifier will be close to the Bayes classifier in terms of 0/1 loss. This suggests to use node classifiers which minimize log-loss on the training sample, examples of which include logistic regression, Gradient Boosting Machines, deep neural networks,<sup>5</sup> and many others. One can show that the square-root dependence in the bound of Theorem 2 cannot be improved in general, since when the tree consists only of the root node, our bound essentially specializes to the bound in [1], which also exhibits square-root dependence.

<sup>5</sup> In this case, the log-loss is often referred to as “soft-max” function.

## 5 Relation to other label tree approaches

### 5.1 Conditional probability trees

Conditional probability trees (CPTs) [4] estimate a conditional probability distribution  $P(y|\mathbf{x})$  in the multiclass setting and have the same structure as PCTs. What distinguishes this approach from ours is that CPTs are used for probability estimation, with squared loss  $\ell_{\text{sq}}(y_i, Q(\cdot|\mathbf{y}^{i-1}, \mathbf{x})) = (y_i - Q(y_i|\mathbf{y}^{i-1}, \mathbf{x}))^2$  as a performance measure, whence there is no inference phase to determine the mode of the conditional distribution. The main result in [4] relates the squared loss regret on the joint distribution to the expected squared loss over the nodes of the tree. This result is analogous to the identity (5), except that an additional  $O(\sqrt{l})$  factor appears in the squared loss bound. Moreover, no result analogous to Theorem 2 is given, which would relate expected squared loss regret to the 0/1 classification regret.

In fact, we can show a *lower bound* on the 0/1 regret in terms of expected squared loss, which is at least a factor of  $\Omega(\sqrt{l})$  worse than our bound. To be more precise, one can show that for any  $l > 2$ , there exists a true distribution  $P$  and an estimate  $Q$  with the following property: even when assuming that the inference algorithm can identify the mode of the distribution exactly, it holds that  $\text{reg}_{0/1}(h_\epsilon) > \sqrt{l \text{reg}_{\text{sq}}(Q)}$ , where  $\text{reg}_{\text{sq}}(Q)$  is the corresponding regret with log-loss replaced by squared loss.<sup>6</sup> In other words, using squared loss yields a bound for classification error that is at least a factor  $\Omega(\sqrt{l})$  worse than the bound we obtained for log-loss.

### 5.2 Filter trees

The filter tree (FT) approach [3] is the first label tree algorithm for which a regret bound for the classification error has been proved. Interestingly, the specific training procedure used in FTs ensures that the greedy classification procedure is sufficient for obtaining consistent predictions.

FT uses the same tree structure as PCT, but with binary classifiers instead of class probability estimators in the non-leaf nodes of the tree. The method follows a bottom-up strategy, which can be interpreted as a single elimination tournament on the set of labels. A classifier in node  $\mathbf{y}^{i-1}$  is trained to predict  $y_i$ , but FT implicitly transforms the underlying distribution of examples in the node. The transformation for 0/1 loss relies on filtering out all training examples that have been misclassified by successor classifiers on a path to a leaf. The learning algorithm starts with classifiers on the lowest non-leaf level of the tree. The correctly classified examples are then moved upward to nodes one level above. This process is repeated until the root node is reached.

In [3], a regret bound for 0/1 loss has been proved that is conceptually similar to the one given in Theorem 2. The difference is that the right side of the

<sup>6</sup> We skip the details of the construction of  $P$  and  $Q$  due to the space limit.

bound is expressed in terms of 0/1 loss of the binary classifiers in non-leaf nodes. Therefore, these two bounds are not directly comparable.

Another advantage of FTs is that they can be used with any cost-based loss function. An appropriate bound has also been proved in [3]. The classification procedure still follows a greedy search, but training is more demanding. It requires weighting of examples, the use of cost-sensitive learners, and each training example generally occurs in each internal classifier.

## 6 Extensions of PCTs

Since PCT estimates the entire conditional distribution over labels, it can be used with any loss function. This comes with no additional cost during training, but may lead to very costly inference. Actually, inference can be performed efficiently only for certain losses, such as 0/1 loss as discussed in Section 3.2, but also some ranking-based loss functions. As an example, consider recall at  $k$ th position defined as

$$R_{@k}(\mathbf{y}, \mathbf{x}, \mathcal{Y}_k) = \mathbb{I}[\mathbf{y} \in \mathcal{Y}_k],$$

where  $\mathcal{Y}_k$  is a set of  $k$  labels predicted for  $\mathbf{x}$ . One can easily verify that an optimal  $\mathcal{Y}_k$  should contain  $k$  top-labels with largest  $P(y|\mathbf{x})$ . This can be approximated by  $k$  top-labels with largest  $Q(y|\mathbf{x})$ , which are easily obtained by PCT and a small extension of the  $\epsilon$ -approximate algorithm: it is enough to continue the search procedure until  $k$  leaves are visited. Moreover, the bound in Theorem 2 can be easily extended to this case.

As already mentioned, PCTs can also be used in multi-label classification. In this case, the tree is of height  $m$  and is fully balanced. Each path from the root to a leaf corresponds to one of possible label combinations. In principle, PCT contains a single classifier in each non-leaf node. In multi-label case, storing  $2^m - 1$  classifiers for large  $m$  is not feasible. One can, however, follow a trick used in probabilistic classifier chains [9] and condensed filter trees [17], which relies on using one binary classifier per tree level. In other words, prediction of the  $i$ th label corresponds to the prediction made by the classifier on level  $i$  with additional features that indicate a given node of the tree.

## 7 Experimental Results

We empirically evaluate PCTs and FTs in two scenarios: multi-label classification (MLC) and multi-class classification (MCC). We test the algorithms in terms of 0/1 loss and the computational costs of their training and testing procedures. For PCTs, we additionally report  $R_{@k}$ .

We conduct experiments on 3 multi-class and 3 multi-label datasets.<sup>7</sup> Table 1 provides a summary of basic statistics of the datasets. Notice that the number of

<sup>7</sup> Taken from the libsvm repository <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets> and the image net competition webpage <http://www.image-net.org/challenges/LSVRC/2010>.

**Table 1.** Multi-class (MCC) and multi-label (MLC) datasets and their properties: the number of training (#train) and test (#test) examples, the number of labels ( $m$ ) and features ( $d$ ).

Dataset	MCC					Dataset	MLC			
	#train	#test	$m$	$d$			#train	#test	$m$	$d$
Sector	6412	3207	105	55197		Yeast	1500	917	14	103
Aloi	97200	10800	1000	128		TMC	21519	7077	22	30438
ILSVR2010	1261406	150000	1000	1000		Mediamill	30993	12914	101	120

leaf nodes is equal to  $m$  (the number of labels) in case of multi-class problems, and  $2^m$  (the number of all possible label combinations) in case of multi-label problems. We therefore use multi-label datasets up to around 100 labels. For datasets with a greater number of labels, the 0/1 loss is usually very close to 1. We use the original split into a training and test set if available; otherwise, we use 90/10 train/test splits. For the ILSVR2010 dataset, we use the visual code words (sbow) vectors provided by the organizers of the challenge. Features were generated on the basis of the guidance contained in the ILSVR development kit.

## 7.1 Implementation

We carefully implemented PCTs and FTs in Java. As internal classifiers, we use  $L_2$  linear logistic regression trained by a variant of stochastic gradient descent (SGD) introduced in [13]. To deal with a large number of weights, we use feature hashing [22] shared over all tree nodes using hashes up to size of  $2^{24}$ . We use a random complete binary tree to code class labels in the MCC scenarios and train a classifier in each node of the tree. For MLC problems we take the original order of the labels to obtain the code words. We use one classifier per tree level. We tune the hyper-parameters of SGD in a 80/20 simple validation on the training set. We applied an off-the-shelf hyper-parameter optimizer [15] with a wide range of parameters. We tune PCTs to optimize the log-loss as suggested by our theoretical analysis. FTs are tuned to perform well on 0/1 loss.

We use PCTs with the  $\epsilon$ -approximate inference algorithm with different values of  $\epsilon \in \{0, 0.25, 0.5\}$ . The variant with  $\epsilon = 0.5$  corresponds to greedy search, while the algorithm with  $\epsilon = 0$  will always find the optimal solution, but may visit all nodes of the tree in the worst case (in fact,  $\epsilon$  should be set to  $2^{-l}$  instead to 0 to be concordant with the description of the algorithm; to keep the notation simple, we use 0 to indicate the smallest possible value of  $\epsilon$  for a given dataset).

## 7.2 Results

The results are given in Table 2. We can observe that PCTs improve with decreasing value of  $\epsilon$ . PCT with  $\epsilon = 0.5$  gets worse results than FT, which confirms

**Table 2.** Experimental results for 0/1 loss and  $1-R_{@5}$  (both in %), train ( $t_{trn}$ ) and test ( $t_{test}$ ) running times (in seconds), and the average (A) number of inner products per a test example. The *Top 1* column indicates the results for top-1 prediction, while column *Top 5* the results for top-5 prediction (only for PCT with  $\epsilon < 0.5$ ). The best results are indicated in bold (except for wall-clock times which can be affected by many factors). The value in subscript of PCT corresponds to the value of  $\epsilon$ .

	MCC						MLC							
	$t_{trn}$	0/1	Top 1 $t_{test}$	A	1- $R_{@5}$	Top 5 $t_{test}$	A	$t_{trn}$	0/1	Top 1 $t_{test}$	A	1- $R_{@5}$	Top 5 $t_{test}$	A
	Sector, $m = 105$						Yeast, $m = 14$							
FT	11.75	13.43	0.144	6.81	–	–	2.49	78.73	0.07	<b>14</b>	–	–	–	–
PCT <sub>.5</sub>	11.56	17.18	0.154	6.81	–	–	3.12	80.15	0.04	<b>14</b>	–	–	–	–
PCT <sub>.25</sub>	11.56	13.68	0.16	7.04	12.61	<b>0.24</b>	<b>7.5</b>	3.12	79.28	0.05	17.15	76.22	<b>0.12</b>	<b>21.3</b>
PCT <sub>0</sub>	11.56	<b>13.28</b>	0.198	7.13	<b>7.23</b>	0.48	18.2	3.12	<b>78.62</b>	0.09	23.82	<b>58.77</b>	0.17	64.6
	Aloi, $m = 105$						TMC, $m = 22$							
FT	15.11	88.98	0.14	<b>9.97</b>	–	–	30.7	77.06	0.47	<b>22</b>	–	–	–	–
PCT <sub>.5</sub>	13.43	88.99	0.14	<b>9.97</b>	–	–	34.3	75.06	0.39	<b>22</b>	–	–	–	–
PCT <sub>.25</sub>	13.43	<b>88.95</b>	0.15	9.98	88.64	<b>0.21</b>	<b>10.2</b>	34.3	73.74	0.45	27.97	68.50	<b>0.57</b>	<b>34.0</b>
PCT <sub>0</sub>	13.43	<b>88.95</b>	0.21	9.98	<b>76.19</b>	0.55	26.1	34.3	<b>73.18</b>	0.73	33.50	<b>41.18</b>	1.29	87.9
	ILSVR2010, $m = 1000$						Mediamill, $m = 101$							
FT	1710	95.10	10.12	<b>8.39</b>	–	–	220	90.79	2.24	<b>101</b>	–	–	–	–
PCT <sub>.5</sub>	1825	99.96	10.13	<b>8.39</b>	–	–	274	90.78	2.22	<b>101</b>	–	–	–	–
PCT <sub>.25</sub>	1825	95.30	13.23	10.03	95.30	<b>20.10</b>	<b>14.4</b>	274	90.06	2.79	107	89.14	<b>3.02</b>	<b>129</b>
PCT <sub>0</sub>	1825	<b>94.76</b>	15.20	10.57	<b>92.33</b>	44.31	34.3	274	<b>89.65</b>	5.23	274	<b>74.22</b>	9.50	529

the theoretical results, i.e., filtering of misclassified examples during training in FT improves for the greedy inference. For  $\epsilon = 0.25$ , the results are already very competitive to FT. For  $\epsilon = 0$ , PCT consistently outperforms FT, but the difference is not always large.

From a computational perspective, FTs achieve better performance. The training time of both approaches is very similar, but the testing time is in favor of FTs (and PCTs with  $\epsilon = 0.5$ ). To give a deeper insight into the time costs we also report the average number of inner products computed by internal classifiers per test example. Interestingly, PCT with  $\epsilon = 0$  always finds the solution in a reasonable time. Its testing time is never longer than three times that of FT. Similarly, the number of inner products is only up to three times greater than that of FT or PCT with  $\epsilon = 0.5$ .

Recall at  $k$ th position ( $R_{@k}$ ) can be measured only for PCTs. There is no way to deliver top- $k$  predictions in FTs, since this algorithm uses binary decisions in non-leaf nodes, so the search process results only in a single path from the root to a leaf node. From the results we observe that PCT efficiently finds topmost results. The positive label appears more often in the top-5 predictions than in the top-1. Similarly as for 0/1 loss,  $R_{@5}$  improves with decreasing value of  $\epsilon$ . Unfortunately, predicting top- $k$  labels increases test time. Therefore, the label tree search for  $\epsilon = 0$  requires about 2-3 times more steps to find top-5 labels.

## 8 Conclusions

In this paper, we analyzed probabilistic classifier trees for efficient multi-class and multi-label classification. In particular, we proved a regret bound for 0/1 loss, which provides a strong theoretical foundation of PCTs, and which can also be extended to ranking-based losses. Moreover, we compared PCTs with the closely related filter tree method. We conclude the paper by summarizing the main theoretical and empirical results of FTs and PCTs, pointing out advantages and disadvantages of both approaches.

An unquestionable advantage of FTs is their prediction time, which is logarithmic in the number of classes or possible label combinations. FT can be used with any type of binary classifier as base learner and relies on simple 0/1 predictions. However, to guarantee the consistency of greedy inference, it requires more demanding training. In the naïve implementation, classifiers are trained sequentially in a bottom-up manner. The most important disadvantage is a significant reduction of the number of training examples in the top levels of the tree, which is caused by filtering examples in each level from bottom to top. This sparsity of training data may deteriorate predictive performance. However, thanks to filtering, an internal classifier is aware of errors of the successor classifiers. FT can be used with any cost-based loss function, but it is not able to predict top- $k$  labels.

Prediction with PCTs requires search techniques, whence it is usually more demanding than FTs (yet significantly faster than 1-vs-all). Moreover, anytime algorithms can be used for searching the tree. The time complexity of PCT strongly depends on the noise contained in the data. If the signal-to-noise ratio is high, we can expect prediction time to be small. However, learning is much simpler for PCT than for FT, and can be easily parallelized. There is no filtering of training examples, so all examples are used for training on each level of the tree. The probabilistic nature of PCTs allows for delivering a list of top-labels and to work efficiently for  $R_{@k}$ .

The results we obtained for FTs are comparable with those reported in [6]. We stress that better results can be obtained by other algorithms, for example LomTrees introduced in the same paper. This is mainly because LomTrees train the tree structure online, along with the internal classifiers, whereas PCTs and FTs use random trees/coding. Interestingly, LomTrees are not consistent. Thus, an important challenge for future research is to find an algorithm that is able to train the tree structure online while ensuring consistency.

*Acknowledgments.* The work of Krzysztof Dembczyński and Wojciech Kotłowski has been supported by the Polish National Science Centre under grant no. 2013/09/D/ST6/03917 and 2013/11/D/ST6/03050, respectively.

## References

1. Bartlett, P.L., Jordan, M.I., McAuliffe, J.D.: Convexity, classification, and risk bounds. *Journal of the American Statistical Association* 101(473), 138–156 (2006)

2. Bengio, S., Weston, J., Grangier, D.: Label embedding trees for large multi-class tasks. In: NIPS 23. pp. 163–171. Curran Associates, Inc. (2010)
3. Beygelzimer, A., Langford, J., Lifshits, Y., Sorkin, G.B., Strehl, A.L.: Conditional probability tree estimation analysis and algorithms. In: UAI. pp. 51–58 (2009)
4. Beygelzimer, A., Langford, J., Ravikumar, P.D.: Error-correcting tournaments. In: ALT. pp. 247–262 (2009)
5. Bottou, L.: Large-scale machine learning with stochastic gradient descent. In: COMPSTAT. pp. 177–187. Springer (2010)
6. Choromanska, A., Langford, J.: Logarithmic time online multiclass prediction. In: NIPS 29 (2015)
7. Cover, T., Thomas, J.: Elements of Information Theory. Wiley (1991)
8. Dekel, O., Shamir, O.: Multiclass-multilabel learning when the label set grows with the number of examples. In: AISTATS (2010)
9. Dembczyński, K., Cheng, W., Hüllermeier, E.: Bayes optimal multilabel classification via probabilistic classifier chains. In: ICML. pp. 279–286. Omnipress (2010)
10. Dembczyński, K., Waegeman, W., Cheng, W., Hüllermeier, E.: An analysis of chaining in multi-label classification. In: ECAI (2012)
11. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Li, F.F.: ImageNet: A large-scale hierarchical image database. In: CVPR. pp. 248–255 (2009)
12. Deng, J., Satheesh, S., Berg, A.C., Li, F.F.: Fast and balanced: Efficient label tree learning for large scale object recognition. In: NIPS 24. pp. 567–575 (2011)
13. Duchi, J., Singer, Y.: Efficient online and batch learning using forward backward splitting. JMLR 10, 2899–2934 (2009)
14. Fox, J.: Applied regression analysis, linear models, and related methods. Sage (1997)
15. Hutter, F., Hoos, H.H., Leyton-Brown, K.: Sequential model-based optimization for general algorithm configuration. In: Learning and Intelligent Optimization. Springer (2011)
16. Kumar, A., Vembu, S., Menon, A.K., Elkan, C.: Beam search algorithms for multilabel learning. Machine Learning 92(1), 65–89 (2013)
17. Li, C.L., Lin, H.T.: Condensed filter tree for cost-sensitive multi-label classification. In: ICML. pp. 423–431 (2014)
18. Mena, D., Montañés, E., Quevedo, J.R., del Coz, J.J.: Using A\* for inference in probabilistic classifier chains. In: IJCAI. pp. 3707–3713 (2015)
19. Morin, F., Bengio, Y.: Hierarchical probabilistic neural network language model. In: AISTATS. pp. 246–252 (2005)
20. Prabhu, Y., Varma, M.: Fastxml: A fast, accurate and stable tree-classifier for extreme multi-label learning. In: KDD. pp. 263–272. ACM (2014)
21. Reid, M.D., Williamson, R.C.: Composite binary losses. JMLR 11, 2387–2422 (2010)
22. Weinberger, K., Dasgupta, A., Langford, J., Smola, A., Attenberg, J.: Feature hashing for large scale multitask learning. In: ICML. pp. 1113–1120. ACM (2009)
23. Weston, J., Makadia, A., Yee, H.: Label partitioning for sublinear ranking. In: ICML. pp. 181–189 (2013)