

---

# A no-regret generalization of hierarchical softmax to extreme multi-label classification

---

**Marek Wydmuch**

Institute of Computing Science  
Poznan University of Technology, Poland  
mwydmuch@cs.put.poznan.pl

**Kalina Jasinska**

Institute of Computing Science  
Poznan University of Technology, Poland  
kjasinska@cs.put.poznan.pl

**Mikhail Kuznetsov**

Yahoo! Research  
New York, USA  
kuznetsov@oath.com

**Róbert Busa-Fekete**

Yahoo! Research  
New York, USA  
busafekete@oath.com

**Krzysztof Dembczyński**

Institute of Computing Science  
Poznan University of Technology, Poland  
kdembczynski@cs.put.poznan.pl

## Abstract

Extreme multi-label classification (XMLC) is a problem of tagging an instance with a small subset of relevant labels chosen from an extremely large pool of possible labels. Large label spaces can be efficiently handled by organizing labels as a tree, like in the hierarchical softmax (HSM) approach commonly used for multi-class problems. In this paper, we investigate probabilistic label trees (PLTs) that have been recently devised for tackling XMLC problems. We show that PLTs are a *no-regret* multi-label generalization of HSM when  $\text{precision}@k$  is used as a model evaluation metric. Critically, we prove that *pick-one-label* heuristic—a reduction technique from multi-label to multi-class that is routinely used along with HSM—is not consistent in general. We also show that our implementation of PLTs, referred to as EXTREME TEXT (XT), obtains significantly better results than HSM with the pick-one-label heuristic and XML-CNN, a deep network specifically designed for XMLC problems. Moreover, XT is competitive to many state-of-the-art approaches in terms of statistical performance, model size and prediction time which makes it amenable to deploy in an online system.

## 1 Introduction

In several machine learning applications, the label space can be enormous, containing even millions of different classes. Learning problems of this scale are often referred to as *extreme classification*. To name a few examples of such problems, consider image and video annotation for multimedia search (Deng et al., 2011), tagging of text documents for categorization of Wikipedia articles (Dekel & Shamir, 2010), recommendation of bid words for online ads (Prabhu & Varma, 2014), or prediction of the next word in a sentence (Mikolov et al., 2013).

To tackle extreme classification problems in an efficient way, one can organize the labels into a tree. A prominent example of such label tree model is *hierarchical softmax* (HSM) (Morin & Bengio, 2005), often used with neural networks to speed up computations in multi-class classification with large output spaces. For example, it is commonly applied in natural language processing problems such as language modeling (Mikolov et al., 2013). To adapt HSM to *extreme multi-label classification* (XMLC), several very popular tools, such as FASTTEXT (Joulin et al., 2016) and LEARNED TREE (Jernite et al., 2017), apply the *pick-one-label* heuristic. As the name suggests, this heuristic randomly picks one of the labels from a multi-label training example and treats the example as a multi-class one.

In this work, we exhaustively investigate the multi-label extensions of HSM. First, we show that the pick-one-label strategy does not lead to a proper generalization of HSM for multi-label setting. More precisely, we prove that using the pick-one-label reduction one cannot expect any multi-class learner to achieve zero regret in terms of marginal probability estimation and maximization of precision@ $k$ . As a remedy to this issue, we are going to revisit *probabilistic label trees* (PLTs) (Jasinska et al., 2016) that have been recently introduced for solving XMLC problems. We show that PLTs are a theoretically motivated generalization of HSM to multi-label classification, that is, 1) PLTs and HSM are identical in multi-class case, and 2) a PLT model can get *zero regret* (i.e., it is *consistent*) in terms of marginal probability estimation and precision@ $k$  in the multi-label setting.

Beside our theoretical findings, we provide an efficient implementation of PLTs, referred to as XT, that we build upon FASTTEXT. The comprehensive empirical evaluation shows that it gets significantly better results than the original FASTTEXT, LEARNED TREE, and XML-CNN, a specifically designed deep network for XMLC problems. XT also achieves competitive results to other state-of-the-art approaches, being very efficient in model size and prediction time, particularly in the online setting.

This paper is organized as follows. First we discuss the related work and situate our approach in the context. In Section 3 we formally state the XMLC problem and present some useful theoretical insights. Next, we briefly introduce the HSM approach, and in Section 5 we show theoretical results concerning the pick-one-label heuristic. Section 6 formally introduces PLTs and presents the main theoretical results concerning them and their relation to HSM. Section 7 provides implementation details of PLTs. The experimental results are presented in Section 8. Finally we make concluding remarks.

## 2 Related work

Historically, problems with a large number of labels were usually solved by nearest neighbor or decision tree methods. Some of today’s algorithms are still based on these classical approaches, significantly extending them by a number of new tricks. If the label space is of moderate size (like a few thousands of labels) then an independent model can be trained for each class. This is the so-called 1-VS-ALL approach. Unfortunately, it scales linearly with the number of labels, which is too costly for many applications. The extreme classification algorithms try to improve over this approach by following different paradigms such as sparsity of labels (Yen et al., 2017; Babbar & Schölkopf, 2017), low-rank approximation (Mineiro & Karampatziakis, 2015; Yu et al., 2014; Bhatia et al., 2015), tree-based search (Prabhu & Varma, 2014; Choromanska & Langford, 2015), or label filtering (Vijayanarasimhan et al., 2014; Shrivastava & Li, 2015; Niculescu-Mizil & Abbasnejad, 2017).

In this paper we focus on tree-based algorithms, therefore we discuss them here in more detail. There are two distinct types of these algorithms: decision trees and label trees. The former type follows the idea of classical decision trees. However, the direct use of the classic algorithms can be very costly (Agrawal et al., 2013). Therefore, the FASTXML algorithm (Prabhu & Varma, 2014) tackles the problem in a slightly different way. It uses sparse linear classifiers in internal tree nodes to split the feature space. Each linear classifier is trained on two classes that are formed in a random way first and then reshaped by optimizing the normalized discounted cumulative gain. To improve the overall accuracy FASTXML uses an ensemble of trees. This algorithm, like many other decision tree methods, works in a batch mode. Choromanska & Langford (2015) have succeeded to introduce a fully online decision tree algorithm that also uses linear classifiers in internal nodes of the tree.

In label trees each label corresponds to one and only one path from the root to a leaf. Besides PLTs and HSM, there exist several other instances of this approach, for example, filter trees (Beygelzimer et al., 2009b; Li & Lin, 2014) or label embedding trees (Bengio et al., 2010). It is also worth to underline that algorithms similar to HSM have been introduced independently in many different research fields, such as nested dichotomies (Fox, 1997) in statistics, conditional probability estimation trees (Beygelzimer et al., 2009a) in multi-class classification, multi-stage classifiers (Kurzynski, 1988) in pattern recognition, and probabilistic classifier chains (Dembczyński et al., 2010) in multi-label classification under the subset 0/1 loss. All these methods have been jointly analyzed in (Dembczyński et al., 2016).

A still open problem in label tree approaches is the tree structure learning. FASTTEXT (Joulin et al., 2016) uses HSM with a Huffman tree built on the label frequencies. Jernite et al. (2017)

have introduced a new algorithm, called `LEARNED TREE`, which combines HSM with a specific hierarchical clustering that reassigns labels to paths in the tree in a semi-online manner. Prabhu et al. (2018) follows another approach in which a model similar to PLTs is trained in a batch mode and a tree is built by using recursively balanced k-means over the label profiles. In Section 7 we discuss this approach in more detail.

The HSM model is often used as an output layer in neural networks. The `FASTTEXT` implementation can also be viewed as a shallow architecture with one hidden layer that represents instances as averaged feature (i.e., word) vectors. Another neural network-based model designed for XMLC has been introduced in (Liu et al., 2017). This model, referred to as XML-CNN, uses a complex convolutional deep network with a narrow last layer to make it work with large output spaces. As we show in the experimental part, this quite expensive architecture gets inferior results in comparison to our PLTs built upon `FASTTEXT`.

### 3 Problem statement

Let  $\mathcal{X}$  denote an instance space, and let  $\mathcal{L} = \{1, \dots, m\}$  be a finite set of  $m$  class labels. We assume that an instance  $\mathbf{x} \in \mathcal{X}$  is associated with a subset of labels  $\mathcal{L}_{\mathbf{x}} \in 2^{\mathcal{L}}$  (the subset can be empty); this subset is often called a set of relevant labels, while the complement  $\mathcal{L} \setminus \mathcal{L}_{\mathbf{x}}$  is considered as irrelevant for  $\mathbf{x}$ . We assume  $m$  to be a large number (e.g.,  $\geq 10^5$ ), but the size of the set of relevant labels  $\mathcal{L}_{\mathbf{x}}$  is much smaller than  $m$ , i.e.,  $|\mathcal{L}_{\mathbf{x}}| \ll m$ . We identify a set  $\mathcal{L}_{\mathbf{x}}$  of relevant labels with a binary (sparse) vector  $\mathbf{y} = (y_1, y_2, \dots, y_m)$ , in which  $y_j = 1 \Leftrightarrow j \in \mathcal{L}_{\mathbf{x}}$ . By  $\mathcal{Y} = \{0, 1\}^m$  we denote a set of all possible label vectors. We assume that observations  $(\mathbf{x}, \mathbf{y})$  are generated independently and identically according to the probability distribution  $\mathbf{P}(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y})$  (denoted later by  $\mathbf{P}(\mathbf{x}, \mathbf{y})$ ) defined on  $\mathcal{X} \times \mathcal{Y}$ .

The problem of XMLC can be defined as finding a *classifier*  $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_m(\mathbf{x}))$ , which in general can be defined as a mapping  $\mathcal{X} \rightarrow \mathcal{R}^m$ , that minimizes the *expected loss* (or *risk*):

$$L_{\ell}(\mathbf{h}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbf{P}(\mathbf{x}, \mathbf{y})}(\ell(\mathbf{y}, \mathbf{h}(\mathbf{x}))),$$

where  $\ell(\mathbf{y}, \hat{\mathbf{y}})$  is the (*task*) *loss*. The optimal classifier, the so-called *Bayes classifier*, for a given loss function  $\ell$  is:

$$\mathbf{h}_{\ell}^* = \arg \min_{\mathbf{h}} L_{\ell}(\mathbf{h}).$$

The *regret* of a classifier  $\mathbf{h}$  with respect to  $\ell$  is defined as:

$$\text{reg}_{\ell}(\mathbf{h}) = L_{\ell}(\mathbf{h}) - L_{\ell}(\mathbf{h}_{\ell}^*) = L_{\ell}(\mathbf{h}) - L_{\ell}^*.$$

The regret quantifies the suboptimality of  $\mathbf{h}$  compared to the optimal classifier  $\mathbf{h}^*$ . The goal could be then defined as finding  $\mathbf{h}$  with a small regret, ideally equal to zero.

In the following, we aim at estimating the marginal probabilities  $\eta_j(\mathbf{x}) = \mathbf{P}(y_j = 1 | \mathbf{x})$ . As we will show below, marginal probabilities are a key element to optimally solve extreme classification for many performance measures, like Hamming loss, macro-F measure, and precision@ $k$ . To obtain the marginal probability estimates one can use the label-wise log loss as a surrogate:

$$\ell_{\log}(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \sum_{j=1}^m \ell_{\log}(y_j, h_j(\mathbf{x})) = \sum_{j=1}^m (y_j \log(h_j(\mathbf{x})) + (1 - y_j) \log(1 - h_j(\mathbf{x}))).$$

Then the expected label-wise log loss for a single  $\mathbf{x}$  (i.e., the so-called *conditional risk*) is:

$$\mathbb{E}_{\mathbf{y}} \ell_{\log}(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \sum_{j=1}^m \mathbb{E}_{\mathbf{y}} \ell_{\log}(y_j, h_j(\mathbf{x})) = \sum_{j=1}^m L_{\log}(h_j(\mathbf{x}) | \mathbf{x}).$$

Therefore, it is easy to see that the pointwise optimal prediction for the  $j$ -th label is given by:

$$h_j^*(\mathbf{x}) = \arg \min_{h_j} L_{\log}(h_j(\mathbf{x}) | \mathbf{x}) = \eta_j(\mathbf{x}).$$

As shown in (Dembczyński et al., 2010), the Hamming loss is minimized by  $h_j^*(\mathbf{x}) = \mathbb{I}[\eta_j(\mathbf{x}) > 0.5]$ . For the macro F-measure it suffices in turn to find an optimal threshold on marginal probabilities for

each label separately as proven in (Ye et al., 2012; Narasimhan et al., 2014; Jasinska et al., 2016; Dembczyński et al., 2017). In the following, we will show a similar result for precision@ $k$  which has become a standard measure in extreme classification (although it is also often criticized, as it favors the most frequent labels).

Precision@ $k$  can be formally defined as:

$$\text{precision@}k(\mathbf{y}, \mathbf{x}, \mathbf{h}) = \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \mathbb{1}[y_j = 1], \quad (1)$$

where  $\hat{\mathcal{Y}}_k$  is a set of  $k$  labels predicted by  $\mathbf{h}$  for  $\mathbf{x}$ . To be consistent with the former discussion, let us define a loss function for precision@ $k$  as  $\ell_{p@k} = 1 - \text{precision@}k$ . The conditional risk is then:<sup>1</sup>

$$L_{p@k}(\mathbf{h} | \mathbf{x}) = \mathbb{E}_{\mathbf{y}} \ell_{p@k}(\mathbf{y}, \mathbf{x}, \mathbf{h}) = 1 - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \eta_j(\mathbf{x}).$$

The above result shows that the optimal strategy for precision@ $k$  is to predict  $k$  labels with the highest marginal probabilities  $\eta_j(\mathbf{x})$ . As the main theoretical result given in this paper is a regret bound for precision@ $k$ , let us define here the conditional regret for this metric:

$$\text{reg}_{p@k}(\mathbf{h} | \mathbf{x}) = \frac{1}{k} \sum_{i \in \mathcal{Y}_k} \eta_i(\mathbf{x}) - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \eta_j(\mathbf{x}),$$

where  $\mathcal{Y}_k$  is a set containing the top  $k$  labels with respect to the true marginal probabilities.

From the above results, we see that estimation of marginal probabilities is crucial for XMLC problems. To obtain these probabilities we can use the vanilla 1-VS-ALL approach trained with the label-wise log loss. Unfortunately, 1-VS-ALL is too costly in the extreme setting. In the following sections, we discuss an alternative approach based on the label trees that estimates the marginal probabilities with the competitive accuracy, but in a much more efficient way.

## 4 Hierarchical softmax approaches

Hierarchical softmax (HSM) is designed for multi-class classification. Using our notation, for multi-class problems we have  $\sum_{i=1}^m y_i = 1$ , i.e., there is one and only one label assigned to an instance  $(\mathbf{x}, \mathbf{y})$ . The marginal probabilities  $\eta_j(\mathbf{x})$  in this case sum up to 1.

The HSM classifier  $\mathbf{h}(\mathbf{x})$  takes a form of a label tree. We encode all labels from  $\mathcal{L}$  using a prefix code. Any such code can be given in a form of a tree in which a path from the root to a leaf node corresponds to a code word. Under the coding, each label  $y_j = 1$  is uniquely represented by a code word  $\mathbf{z} = (z_1, \dots, z_l) \in \mathcal{C}$ , where  $l$  is the length of the code word and  $\mathcal{C}$  is a set of all code words. For  $z_i \in \{0, 1\}$ , the code and the label tree are binary. In general, the code alphabet can contain more than two symbols. Furthermore,  $z_i$ s can take values from different sets of symbols depending on the previous values in the code word. In other words, the code can result with nodes of a different arity even in the same tree, like in (Grave et al., 2017) and (Prabhu et al., 2018). We will briefly discuss different tree structures in Section 7.

A tree node can be uniquely identified by the partial code word  $\mathbf{z}^i = (z_1, \dots, z_i)$ . We denote the root node by  $\mathbf{z}^0$ , which is an empty vector (without any elements). The probability of a given label is determined by a sequence of decisions made by node classifiers that predict subsequent values of the code word. By using the chain rule of probability, we obtain:

$$\eta_j(\mathbf{x}) = \mathbf{P}(y_j = 1 | \mathbf{x}) = \mathbf{P}(\mathbf{z} | \mathbf{x}) = \prod_{i=1}^l \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}).$$

By using logistic loss and a linear model  $f_{\mathbf{z}^i}(\mathbf{x})$  in each node  $\mathbf{z}^i$  for estimating  $\mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x})$ , we obtain the popular formulation of HSM. Let us notice that since we deal here with a multi-class distribution, we have that:

$$\sum_c \mathbf{P}(z_i = c | \mathbf{z}^{i-1}, \mathbf{x}) = 1. \quad (2)$$

<sup>1</sup>The derivation is given in Appendix A.

Because of this normalization, we can assume that a multi-class (or binary in the case of binary trees) classifier is situated in all internal nodes and there are no classifiers in the leaves of the tree. Alternatively, we can assume that each node, except the root, is associated with a binary classifier that estimates  $\mathbf{P}(z_i = c \mid z^{i-1}, \mathbf{x})$ , but then the additional normalization (2) has to be performed. This alternative formulation is important for the multi-label extension of HSM discussed in Section 6. In either way, learning of the node classifiers can be performed simultaneously as independent tasks.

Note that estimate  $\hat{\eta}_j(\mathbf{x})$  of the probability of label  $j$  can be easily obtained by traversing the tree along the path indicated by the code of the label. Unfortunately, the task of predicting top  $k$  labels is more involved as it requires searching over the tree. Popular solutions are beam search (Kumar et al., 2013; Prabhu et al., 2018), uniform-cost search (Joulin et al., 2016), and its approximate variant (Dembczyński et al., 2012; Dembczyński et al., 2016).

## 5 Suboptimality of HSM for multi-label classification

To deal with multi-label problems, some popular tools, such as FASTTEXT (Joulin et al., 2016) and its extension LEARNED TREE (Jernite et al., 2017), apply HSM with the pick-one-label heuristic which randomly picks one of the positive labels from a given training instance. The resulting instance is then treated as a multi-class instance. During prediction, the heuristic returns a multi-class distribution and the  $k$  most probable labels. We show below that this specific reduction of the multi-label problem to multi-class classification is not consistent in general.

Since the probability of picking a label  $j$  from  $\mathbf{y}$  is equal to  $y_j / \sum_{j'=1}^m y_{j'}$ , the pick-one-label heuristic maps the multi-label distribution to a multi-class distribution in the following way:

$$\eta'_j(\mathbf{x}) = \mathbf{P}'(y_j = 1 \mid \mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} \frac{y_j}{\sum_{j'=1}^m y_{j'}} \mathbf{P}(\mathbf{y} \mid \mathbf{x}) \quad (3)$$

It can be easily checked that the resulting  $\eta'_j(\mathbf{x})$  form a multi-class distribution as the probabilities sum up to 1. It is obvious that the heuristic changes the marginal probabilities of labels, unless the initial distribution is multi-class. Therefore this method cannot lead to consistent classifiers in terms of estimating  $\eta_j(\mathbf{x})$ . As we show below, it is also not consistent for precision@ $k$  in general.

**Proposition 1.** *A classifier  $\mathbf{h}$  such that  $h_j(\mathbf{x}) = \eta'_j(\mathbf{x})$  for all  $j \in \{1, \dots, m\}$  has in general a non-zero regret in terms of precision@ $k$ .*

*Proof.* We prove the proposition by giving a simple counterexample. Consider the following conditional distribution for some  $\mathbf{x}$ :

$$\mathbf{P}(\mathbf{y} = (1, 0, 0) \mid \mathbf{x}) = 0.1, \quad \mathbf{P}(\mathbf{y} = (1, 1, 0) \mid \mathbf{x}) = 0.5, \quad \mathbf{P}(\mathbf{y} = (0, 0, 1) \mid \mathbf{x}) = 0.4.$$

The optimal top 1 prediction for this example is obviously label 1, since the marginal probabilities are  $\eta_1(\mathbf{x}) = 0.6, \eta_2(\mathbf{x}) = 0.5, \eta_3(\mathbf{x}) = 0.4$ . However, the pick-one-label heuristic will transform the original distribution to the following one:  $\eta'_1(\mathbf{x}) = 0.35, \eta'_2(\mathbf{x}) = 0.25, \eta'_3(\mathbf{x}) = 0.4$ . The predicted top label will be then label 3, giving the regret of 0.2 for precision@1.  $\square$

The proposition shows that the heuristic is in general inconsistent for precision@ $k$ . Interestingly, the situation changes when the labels are conditionally independent, i.e.,  $\mathbf{P}(\mathbf{y} \mid \mathbf{x}) = \prod_{j=1}^m \mathbf{P}(y_j \mid \mathbf{x})$ .

**Proposition 2.** *Given conditionally independent labels, a classifier  $\mathbf{h}$  such that  $h_j(\mathbf{x}) = \eta'_j(\mathbf{x})$  for all  $j \in \{1, \dots, m\}$  has zero regret in terms of the precision@ $k$  loss.*

*Proof.* We show here only a sketch of the proof. The full proof is given in Appendix B. To prove the theorem, it is enough to show that in the case of conditionally independent labels the pick-one-label heuristic does not change the order of marginal probabilities. Let  $y_i$  and  $y_j$  be so that  $\mathbf{P}(y_i = 1 \mid \mathbf{x}) \geq \mathbf{P}(y_j = 1 \mid \mathbf{x})$ . Then in the summation over all  $\mathbf{y}$  in (3), we are interested in four different subsets of  $\mathcal{Y}$ ,  $S_{i,j}^{u,w} = \{\mathbf{y} \in \mathcal{Y} : y_i = u \wedge y_j = w\}$ , where  $u, w \in \{0, 1\}$ . Remark that during mapping none of  $\mathbf{y} \in S_{i,j}^{0,0}$  plays any role, and for each  $\mathbf{y} \in S_{i,j}^{1,1}$ , the value of  $y_t / (\sum_{t'=1}^m y_{t'}) \times \mathbf{P}(\mathbf{y} \mid \mathbf{x})$ , for  $t \in \{i, j\}$ , is the same for both  $y_i$  and  $y_j$ . Now, let  $\mathbf{y}' \in S_{i,j}^{1,0}$  and  $\mathbf{y}'' \in S_{i,j}^{0,1}$  be the same on all elements except the  $i$ -th and the  $j$ -th one. Then, because

of the label independence and the assumption that  $\mathbf{P}(y_i = 1 | \mathbf{x}) \geq \mathbf{P}(y_j = 1 | \mathbf{x})$ , we have  $\mathbf{P}(\mathbf{y}' | \mathbf{x}) \geq \mathbf{P}(\mathbf{y}'' | \mathbf{x})$ . Therefore, after mapping we obtain  $\eta'_i(\mathbf{x}) \geq \eta'_j(\mathbf{x})$ . Thus, for independent labels, the pick-one-label heuristic is consistent for precision@k.  $\square$

## 6 Probabilistic label trees

The section above has revealed that HSM cannot be properly adapted to multi-label problems by the pick-one-label heuristic. There is, however, a different way to generalize HSM to obtain no-regret estimates of marginal probabilities  $\eta_j(\mathbf{x})$ . The probabilistic label trees (PLTs) (Jasinska et al., 2016) can be derived in the following way. Let us encode  $y_j = 1$  by a slightly extended code  $\mathbf{z} = (1, z_1, \dots, z_l)$  in comparison to HSM. The new code gets 1 at the zero position what corresponds to a question whether there exists at least one label assigned to the example. As before, each node is uniquely identified by a partial code  $\mathbf{z}^i$  which says that there is at least one positive label in a subtree rooted in that node. It can be easily shown by the chain rule of probability that the marginal probabilities can be expressed in the following way:

$$\eta_j(\mathbf{x}) = \mathbf{P}(\mathbf{z} | \mathbf{x}) = \prod_{i=0}^l \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}). \quad (4)$$

The difference to HSM is the probability  $\mathbf{P}(z_0 = 1 | \mathbf{x})$  in the chain and a different normalization, i.e.:

$$\sum_c \mathbf{P}(z_i = c | \mathbf{z}^{i-1}, \mathbf{x}) \geq 1. \quad (5)$$

Only for  $z_0$  we have  $\mathbf{P}(z_0 = 1 | \mathbf{x}) + \mathbf{P}(z_0 = 0 | \mathbf{x}) = 1$ . Because of (5), the binary models that estimate  $\mathbf{P}(z_i = c | \mathbf{z}^{i-1}, \mathbf{x})$  (against  $\mathbf{P}(z_i \neq c | \mathbf{z}^{i-1}, \mathbf{x})$ ) are situated in all nodes of the tree (i.e., also in the leaves). The models can be trained independently as before for HSM. Only during prediction, one can re-calibrate the estimates when (5) is not satisfied, for example, by normalizing them to sum up to 1. It can be easily noticed that for a multi-class distribution, the resulting model of PLTs boils down to HSM, since  $\mathbf{P}(z_0 = 1 | \mathbf{x})$  is always equal 1, and in addition, normalization (5) will take the form of (2). In Appendix D we additionally present the pseudocode of training and predicting with PLTs.

Next, we show that the PLT model obeys strong theoretical guarantees. Let us first revise the result from (Jasinska et al., 2016) that relates the absolute difference between the true and the estimated marginal probability of label  $j$ ,  $|\eta_j(\mathbf{x}) - \hat{\eta}_j(\mathbf{x})|$ , to the surrogate loss  $\ell$  used to train node classifiers  $f_{\mathbf{z}^i}$ . It is assumed here that  $\ell$  is a strongly proper composite loss (e.g, logistic, exponential, or squared loss) characterized by a constant  $\lambda$  (e.g.  $\lambda = 4$  for logistic loss).<sup>2</sup>

**Theorem 1.** *For any distribution  $\mathbf{P}$  and internal node classifiers  $f_{\mathbf{z}^i}$ , the following holds:*

$$|\eta_j(\mathbf{x}) - \hat{\eta}_j(\mathbf{x})| \leq \sum_{i=0}^l \mathbf{P}(\mathbf{z}^{i-1} | \mathbf{x}) \sqrt{\frac{2}{\lambda}} \sqrt{\text{reg}_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x})},$$

where  $\text{reg}_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x})$  is a binary classification regret for a strongly proper composite loss  $\ell$  and  $\lambda$  is a constant specific for loss  $\ell$ .

Due to filtering of the distribution imposed by the PLT, the regret  $\text{reg}_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x})$  of a classifier  $f_{\mathbf{z}^i}$  exists only for  $\mathbf{x}$  such that  $\mathbf{P}(\mathbf{z}^{i-1} | \mathbf{x}) > 0$ , therefore we condition the regret not only on  $\mathbf{x}$ , but also on  $\mathbf{z}^{i-1}$ . The above result shows that the absolute error of estimating the marginal probability of label  $j$  can be upper bounded by the regret of the node classifiers on the corresponding path from the root to a leaf. The proof of Theorem 1 is given in Appendix A. Moreover, for zero-regret (i.e., optimal) node classifiers we obtain an optimal multi-label classifier in terms of estimation of marginal probabilities  $\eta_j(\mathbf{x})$ . This result can be further extended for precision@k.

**Theorem 2.** *For any distribution  $\mathbf{P}$  and classifier  $\mathbf{h}$  delivering estimates  $\hat{\eta}_j(\mathbf{x})$  of the marginal probabilities of labels, the following holds:*

$$\text{reg}_{p@k}(\mathbf{h} | \mathbf{x}) = \frac{1}{k} \sum_{i \in \mathcal{Y}_k} \eta_i(\mathbf{x}) - \frac{1}{k} \sum_{j \in \mathcal{Y}_k} \eta_j(\mathbf{x}) \leq 2 \max_i |\eta_i(\mathbf{x}) - \hat{\eta}_i(\mathbf{x})|$$

<sup>2</sup>For more detailed introduction to strongly proper composite losses, we refer the reader to (Agarwal, 2014).

The proof is based on adding and subtracting the following terms  $\frac{1}{k} \sum_{i \in \mathcal{Y}_k} \hat{\eta}_i(\mathbf{x})$  and  $\frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \hat{\eta}_j(\mathbf{x})$  to the regret (a detailed proof is given in Appendix A). By getting together both theorems we get an upper bound of the precision@ $k$  regret expressed in terms of the regret of the node classifiers. Again, for the zero-regret node classifiers, we get optimal solution in terms of precision@ $k$ .

## 7 Implementation details of PLTs

Given the tree structure, the node classifiers of PLTs can be trained as logistic regression either in online (Jasinska et al., 2016) or batch mode (Prabhu et al., 2018). Both training modes have their pros and cons, but the online implementation gives a possibility of learning more complex representation of input instances. The above cited implementations are both based on sparse representation, given either in a form of a bag-of-words or its TF-IDF variant. We opt here for training a PLT in the online mode along with the dense representation. We build our implementation upon FASTTEXT and refer to it as XT which stands for EXTREMETEXT.<sup>3</sup> In this way, we succeeded to obtain a very powerful and compressed model. The small dense models are important for fast online prediction as they do not need too much resources. The sparse models, in turn, can be slow and expensive in terms of memory usage as they need to decompress the node models to work fast. Remark also that, in general, PLTs can be used as an output layer of any neural network architecture (also that one used in XML-CNN (Yen et al., 2017)) to speed up training and prediction time.

In contrast to the original implementation of FASTTEXT, we use L2 regularization for all parameters of the model. To obtain representation of input instances we do not compute simple averages of the feature vectors, but use weights proportional to the TF-IDF scores of features. The competitive results can be obtained with feature and instance vectors of size 500. If a node classification task contains only positive instances, we use a constant classifier predicting 1 without any training. The training of PLT in either mode, online or batch, can be easily parallelized as each node classifier can be trained in isolation from the other classifiers. In our current implementation, however, we follow the parallelization on the level of training and test instances as in original FASTTEXT.

Our implementation, because of the additional use of the L2 regularization, has more parameters than original FASTTEXT. We have found, however, that our model is remarkably robust for the hyperparameter selection, since it achieves close to optimal performance for a large set of hyperparameters that is in the vicinity of the optimal one. Moreover, the optimal hyperparameters are close to each other across all datasets. We report more information about the hyperparameter selection in Appendix E.4.

The tree structure of a PLT is a crucial modeling decision. The vanishing regret for probability estimates and precision@ $k$  holds regardless of the tree structure (see Theorem 1 and 2), however, this theory requires the regret of the node classifiers also to vanish. In practice, we can only estimate the conditional probabilities in the nodes, therefore the tree structure does indeed matter as it affects the difficulty of the node learning problems. The original PLT paper (Jasinska et al., 2016) uses simple complete trees with labels assigned to leaves according to their frequencies. Another option, routinely used in HSM (Joulin et al., 2016), is the Huffman tree built over the label frequencies. Such tree takes into account the computational complexity by putting the most frequent labels close to the root. This approach has been further extended to optimize GPU operations in (Grave et al., 2017). Unfortunately, it ignores the statistical properties of the tree structure. Furthermore, for multi-label case the Huffman tree is no longer optimal even in terms of computational cost as we show it in Appendix C. There exist, however, other methods that focus on building a tree with high overall accuracy (Tagami, 2017; Prabhu et al., 2018). In our work, we follow the later approach, which performs a simple top-down hierarchical clustering. Each label in this approach is represented by a profile vector being an average of the training vectors tagged by this label. Then the profile vectors are clustered using balanced k-means which divides the labels into two or more clusters with approximately the same size. This procedure is then repeated recursively until the clusters are smaller than a given value (for example, 100). The nodes of the resulting tree are then of different arities. The internal nodes up to the leaves' parent nodes have  $k$  children, but the leaves' parent nodes are usually of higher arity. Thanks to this clustering, similar labels are close to each other in the tree. Moreover, the tree is balanced, so its depth is logarithmic in terms of the number of labels.

<sup>3</sup>Implementation of XT is available at <https://github.com/mwydmuch/extremeText>.

## 8 Empirical results

We carried out three sets of experiments. In the first, we compare exhaustively the performance of PLTs and HSM on synthetic and benchmark data. Due to lack of space, the results are deferred to Appendix E.1 and E.2. The results on synthetic data confirm our theoretical findings: the models are the same in the case of multi-class data, the performance of HSM and PLTs is on par using multi-label data with independent labels, and PLTs significantly outperform HSM on multi-label data with conditionally dependent labels. The results on the benchmark data clearly indicate the better performance of PLTs over HSM.

In the second experiment, we compare XT, the variant of PLTs discussed in the previous section, to the state-of-the-art algorithms on five benchmark datasets taken from XMLC repository,<sup>4</sup> and their text equivalents, by courtesy of Liu et al. (2017). We compare the models in terms of  $\text{precision}@\{1, 3, 5\}$ , model size, training and test time. The competitors for our XT are original FASTTEXT, its variant LEARNED TREE, a PLT-like batch learning algorithm PARABEL (we use the variant that uses a single tree instead of an ensemble), a XMLC-designed convolutional deep network XML-CNN, a decision tree ensemble FASTXML, and two 1-vs-All approaches tailored to XMLC problems, PPD-SPARSE and DISMEC. The hyperparameters of the models have been tuned using grid search. The range of the hyperparameters is reported in E.4.

The results presented in Table 1 demonstrate that XT outperforms the HSM approaches with the pick-one-label heuristic, namely FASTTEXT and LEARNED TREE, with a large margin. This proves the superiority of PLTs as the proper generalization of HSM to multi-label setting. In all the above methods we use vectors of length 500 and we tune the other hyperparameters appropriately for a fair comparison.

Moreover, XT scales well to extreme datasets achieving performance close to the state-of-the-art, being at the same time 1000x and 100x faster compared to DISMEC and PPDSPARSE during prediction. XT always responds below 2ms, what makes it a competitive alternative for an online setting. XT is also close to PARABEL in terms of performance. However, the reported times and model sizes of PARABEL are given for the batch prediction. The prediction times seem to be faster, but PARABEL needs to decompress the model during prediction, what makes it less suitable for online prediction. It is only efficient when the batches are sufficiently large. Finally, we would like to underline that XT outperforms XML-CNN, the more complex neural network, in terms of predictive performance with computational costs that are an order of magnitude smaller. Moreover, XML-CNN requires pretrained embedding vectors, whereas XT can be used with random initialization.

In the third experiment we perform an ablation analysis in which we compare different components of the XT algorithm. We analyze the influence of the Huffman tree vs. top-down clustering, the simple averaging of features vectors vs. the TF-IDF-based weighting, and no regularization vs. L2 regularization. Figure 1 clearly shows that the components need to be combined together to obtain the best results. The best combination uses top-down clustering, TF-IDF-based weighting, and L2 regularization, while top-down clustering alone gets worse results than Huffman trees with TF-IDF-based weighting and L2 regularization. In Appendix E.3 we give more detailed results of the ablation analysis performed on a larger spectrum of benchmark datasets.

## 9 Conclusions

In this paper we have proven that probabilistic label trees (PLTs) are no-regret generalization of HSM to the multi-label setting. Our main theoretical contribution is the  $\text{precision}@k$  regret bound for PLTs. Moreover, we have shown that the pick-one-label heuristic commonly-used with HSM in multi-label problems leads to inconsistent results in terms of marginal probability estimation and  $\text{precision}@k$ . Our implementation of PLTs referred to as XT, built upon FASTTEXT, gets state-of-the-art results, being significantly better than the original FASTTEXT, LEARNED TREE, and XML-CNN. The XT results are also close to the best known ones that are obtained by expensive 1-vs-All approaches, such as PPDSPARSE and DISMEC, and outperforms the other tree-based methods on many benchmarks. Our online variant has the advantage of producing very often much smaller models that can be efficiently used in fast online prediction.

---

<sup>4</sup>Additional statistics of these datasets are also included in Appendix F. Address of the XMLC repository: <http://manikvarma.org/downloads/XC/XMLRepository.html>



Table 1: Precision@ $k$  scores with  $k = \{1, 3, 5\}$  and statistics of FASTXML, PDPSPARSE, DISMEC, PARABEL (with 1 tree), FASTTEXT (FT), LEARNED TREE (LT), EXTREMETEXT (XT) and XML-CNN methods. Notation:  $N$  – number of samples,  $T$  – CPU time,  $m$  – number of labels,  $d$  – number of features, \* – result of offline prediction, \* – calculated on GPU, † – not reported by authors, ‡ – cannot be calculated due to lack of a text version of a dataset.

Dataset	Metrics	FASTXML	PPDSPARSE	DISMEC	FT	LT	XT	PARABEL	XML-CNN
<b>Wiki-30K</b> $N_{train} = 14146$ $N_{test} = 6616$ $d = 101938$ $m = 30938$	P@1	82.03	73.80	85.20	80.78	80.85	<b>85.23</b>	83.77	82.78
	P@3	67.47	60.90	<b>74.60</b>	50.46	50.59	73.18	71.96	66.34
	P@5	57.76	50.40	<b>65.90</b>	36.79	37.68	63.39	62.44	56.23
	$T_{train}$	16m	†	†	10m	12m	18m	<b>5m</b>	88m*
	$T_{test}/N_{test}$	3.00ms	†	†	1.88ms	10.09ms	<b>0.83ms</b>	1.63ms*	1.39ms*
	model size	354M	†	†	513M	513M	<b>259M</b>	<b>109M*</b>	*
<b>Delicious-200K</b> $N_{train} = 196606$ $N_{test} = 100095$ $d = 782585$ $m = 205443$	P@1	42.81	45.05	44.71	42.22	42.71	<b>47.85</b>	43.32	†
	P@3	38.76	38.34	38.08	37.90	36.27	<b>42.08</b>	38.49	†
	P@5	36.34	34.90	34.7	35.05	33.43	<b>39.13</b>	35.83	†
	$T_{train}$	458m	4781m	1080h	271m	563m	502m	105m	†
	$T_{test}/N_{test}$	4.86ms	275ms	5m	1.97ms	1.98ms	<b>1.41ms</b>	<b>1.31ms*</b>	†
	model size	15.4G	9.4G	18.0G	9.0G	9.0G	<b>1.9G</b>	<b>1.8G*</b>	†
<b>WikiLSHTC</b> $N_{train} = 1778351$ $N_{test} = 587084$ $d = 617899$ $m = 325056$	P@1	49.35	64.08	<b>64.94</b>	41.13	50.15	58.73	61.53	†
	P@3	32.69	41.26	<b>42.71</b>	24.09	31.95	39.24	40.07	†
	P@5	24.03	30.12	<b>31.5</b>	17.44	23.59	29.26	29.25	†
	$T_{train}$	724m	236m	750h	207	212m	550m	<b>34m</b>	†
	$T_{test}/N_{test}$	2.17ms	37.76ms	43m	1.25ms	4.76ms	<b>0.81ms</b>	0.92ms*	†
	model size	9.3G	5.2G	3.8G	6.5G	6.5G	<b>3.3G</b>	<b>1.1G*</b>	†
<b>Wiki-500K</b> $N_{train} = 1813391$ $N_{test} = 783743$ $d = 2381304$ $m = 501070$	P@1	54.10	70.16	<b>70.20</b>	32.73	37.18	64.48	66.12	59.85
	P@3	29.45	50.57	<b>50.60</b>	19.02	21.62	45.84	47.02	39.28
	P@5	21.21	39.66	<b>39.70</b>	14.46	16.01	35.46	36.45	29.81
	$T_{train}$	3214m	1771m	7495h	496m	531m	1253m	<b>168m</b>	7032m*
	$T_{test}/N_{test}$	8.03ms	113.70ms	155m	2.05ms	6.43ms	<b>1.07ms</b>	4.68ms*	21.06ms*
	model size	63G	3.4G	14.7G	11G	11G	<b>5.5G</b>	<b>2.0G*</b>	<b>3.7G*</b>
<b>Amazon-670K</b> $N_{train} = 490449$ $N_{test} = 153025$ $d = 135909$ $m = 670091$	P@1	34.24	45.32	<b>45.37</b>	25.47	27.67	39.90	41.59	35.39
	P@3	29.30	40.37	<b>40.40</b>	21.47	20.96	35.36	37.18	33.74
	P@5	26.12	36.92	<b>36.96</b>	18.61	17.72	32.04	33.85	32.64
	$T_{train}$	422m	102m	373h	162m	182m	241m	<b>8m</b>	3134m*
	$T_{test}/N_{test}$	3.39ms	66.09ms	23m	7.84ms	5.13ms	<b>1.72ms</b>	<b>0.68ms*</b>	16.18ms*
	model size	10G	6.0G	3.8G	3.2G	3.2G	<b>1.5G</b>	<b>0.7G*</b>	<b>1.5G*</b>

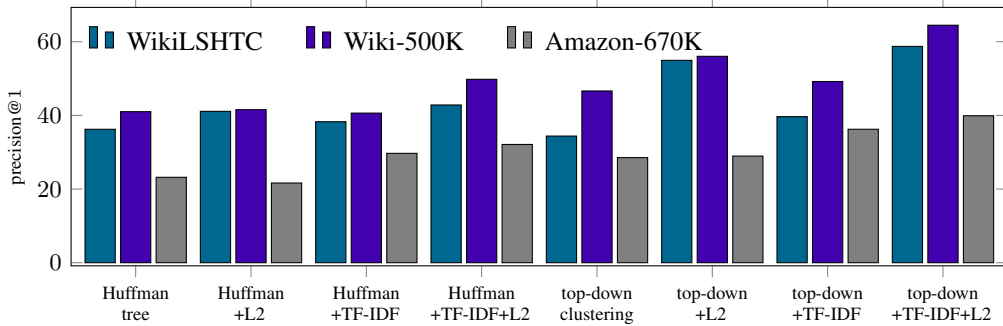


Figure 1: The ablation analysis of different variants of XT on WIKILSHTC, WIKI-500K, and AMAZON-670K.

## Acknowledgements

The work of Kalina Jasinska was supported by the Polish National Science Center under grant no. 2017/25/N/ST6/00747. The work of Krzysztof Dembczyński was supported by the Polish Ministry of Science and Higher Education under grant no. 09/91/DSPB/0651. Computational experiments have been performed in Poznan Supercomputing and Networking Center.

## References

- Agarwal, S. Surrogate regret bounds for bipartite ranking via strongly proper losses. *JMLR*, 15: 1653–1674, 2014.
- Agrawal, R., Gupta, A., Prabhu, Y., and Varma, M. Multi-label learning with millions of labels: Recommending advertiser bid phrases for web pages. In *WWW*, pp. 13–24. ACM, 2013.
- Babbar, Rohit and Schölkopf, Bernhard. Dismec: Distributed sparse machines for extreme multi-label classification. In *WSDM 2017*, pp. 721–729. ACM, 2017.
- Bengio, S., Weston, J., and Grangier, D. Label embedding trees for large multi-class tasks. In *NIPS 24*, pp. 163–171, 2010.
- Beygelzimer, A., Langford, J., Lifshits, Y., Sorkin, G., and Strehl, A. Conditional probability tree estimation analysis and algorithms. In *UAI*, pp. 51–58, 2009a.
- Beygelzimer, A., Langford, J., and Ravikumar, P. Error-correcting tournaments. In *ALT*, pp. 247–262. Springer, 2009b.
- Bhatia, K., Jain, H., Kar, P., Varma, M., and Jain, P. Sparse local embeddings for extreme multi-label classification. In *NIPS 29*, pp. 730–738, 2015.
- Choromanska, A. and Langford, J. Logarithmic time online multiclass prediction. In *NIPS 29*, pp. 55–63, 2015.
- Dekel, O. and Shamir, O. Multiclass-multilabel classification with more classes than examples. In *AISTATS*, pp. 137–144, 2010.
- Dembczyński, K., Cheng, W., and Hüllermeier, E. Bayes optimal multilabel classification via probabilistic classifier chains. In *ICML*, pp. 279–286, 2010.
- Dembczyński, K., Waegeman, W., Cheng, W., and Hüllermeier, E. An analysis of chaining in multi-label classification. In *ECAI*, 2012.
- Dembczyński, Krzysztof, Kotłowski, Wojciech, Waegeman, Willem, Busa-Fekete, Róbert, and Hüllermeier, Eyke. Consistency of probabilistic classifier trees. In *ECMLPKDD*. Springer, 2016.
- Dembczyński, Krzysztof, Kotłowski, Wojciech, Koyejo, Oluwasanmi, and Natarajan, Nagarajan. Consistency analysis for binary classification revisited. In Precup, Doina and Teh, Yee Whye (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 961–969. PMLR, 2017.
- Deng, J., Satheesh, S., Berg, A. C., and Fei-Fei, L. Fast and balanced: Efficient label tree learning for large scale object recognition. In *NIPS 24*, pp. 567–575, 2011.
- Fan, R.-E., Chang, K.-W., Hsieh, C.-J., Wang, X.-R., and Lin, C.-J. LIBLINEAR: A library for large linear classification. *JMLR*, 9:1871–1874, 2008.
- Fox, J. *Applied regression analysis, linear models, and related methods*. Sage, 1997.
- Grave, Edouard, Joulin, Armand, Cissé, Moustapha, Grangier, David, and Jégou, Hervé. Efficient softmax approximation for GPUs. In Precup, Doina and Teh, Yee Whye (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1302–1310, International Convention Centre, Sydney, Australia, 2017. PMLR.
- Jasinska, K., Dembczynski, K., Busa-Fekete, R., Pfannschmidt, K., Klerx, T., and Hüllermeier, E. Extreme F-measure maximization using sparse probability estimates. In *ICML*, 2016.
- Jernite, Yacine, Choromanska, Anna, and Sontag, David. Simultaneous learning of trees and representations for extreme classification and density estimation. In Precup, Doina and Teh, Yee Whye (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1665–1674, International Convention Centre, Sydney, Australia, 2017. PMLR.

- Joulin, Armand, Grave, Edouard, Bojanowski, Piotr, and Mikolov, Tomas. Bag of tricks for efficient text classification. *CoRR*, abs/1607.01759, 2016. URL <http://arxiv.org/abs/1607.01759>.
- Kumar, A., Vembu, S., Menon, A.K., and Elkan, C. Beam search algorithms for multilabel learning. In *Machine Learning*, 2013.
- Kurzynski, Marek. On the multistage bayes classifier. *Pattern Recognition*, 21(4):355–365, 1988.
- Langford, J., Strehl, A., and Li, L. Vowpal wabbit, 2007. <http://mloss.org/software/view/53/>.
- Li, Ch.-L. and Lin, H.-Ti. Condensed filter tree for cost-sensitive multi-label classification. In *ICML*, pp. 423–431, 2014.
- Liu, Jingzhou, Chang, Wei-Cheng, Wu, Yuexin, and Yang, Yiming. Deep learning for extreme multi-label text classification. In *SIGIR 2017*, pp. 115–124. ACM, 2017.
- Mikolov, Tomas, Sutskever, Ilya, Chen, Kai, Corrado, Greg S, and Dean, Jeff. Distributed representations of words and phrases and their compositionality. In *NIPS 26*, pp. 3111–3119. Curran Associates, Inc., 2013.
- Mineiro, P. and Karampatziakis, N. Fast label embeddings via randomized linear algebra. In *ECML/PKDD 2015*, pp. 37–51, 2015.
- Morin, F. and Bengio, Y. Hierarchical probabilistic neural network language model. In *AISTATS*, pp. 246–252, 2005.
- Narasimhan, H., Vaish, R., and S., Agarwal. On the statistical consistency of plug-in classifiers for non-decomposable performance measures. In *NIPS 27*, pp. 1493–1501, 2014.
- Niculescu-Mizil, Alexandru and Abbasnejad, Ehsan. Label Filters for Large Scale Multilabel Classification. In Singh, Aarti and Zhu, Jerry (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 1448–1457. PMLR, 2017.
- Prabhu, Y. and Varma, M. FastXML: A fast, accurate and stable tree-classifier for extreme multi-label learning. In *KDD*, pp. 263–272. ACM, 2014.
- Prabhu, Yashoteja, Kag, Anil, Harsola, Shrutendra, Agrawal, Rahul, and Varma, Manik. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In Champin, Pierre-Antoine, Gandon, Fabien L., Lalmas, Mounia, and Ipeirotis, Panagiotis G. (eds.), *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018*, pp. 993–1002. ACM, 2018.
- Shrivastava, A. and Li, P. Improved asymmetric locality sensitive hashing (ALSH) for maximum inner product search (mips). In *UAI*, 2015.
- Tagami, Yukihiro. Annexml: Approximate nearest neighbor search for extreme multi-label classification. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 455–464, New York, NY, USA, 2017. ACM.
- Vijayanarasimhan, Sudheendra, Shlens, Jonathon, Monga, Rajat, and Yagnik, Jay. Deep networks with large output spaces. *CoRR*, abs/1412.7479, 2014. URL <http://arxiv.org/abs/1412.7479>.
- Ye, N., Chai, A., Lee, W., and Chieu, H. Optimizing F-measures: A tale of two approaches. In *ICML*, 2012.
- Yen, Ian E.H., Huang, Xiangru, Dai, Wei, Ravikumar, Pradeep, Dhillon, Inderjit, and Xing, Eric. PPDsparse: A parallel primal-dual sparse method for extreme classification. In *KDD 2017*, pp. 545–553. ACM, 2017.
- Yu, Hsiang-Fu, Jain, Prateek, Kar, Purushottam, and Dhillon, Inderjit. Large-scale multi-label learning with missing labels. In *ICML 2014*, volume 32, pp. 593–601. PMLR, 2014.

## A Regret for precision@k

Precision@k is formally defined as:

$$\text{precision@}k(\mathbf{y}, \mathbf{x}, \mathbf{h}) = \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \mathbb{1}[y_j = 1],$$

where  $\hat{\mathcal{Y}}_k$  is a set of  $k$  labels predicted by  $\mathbf{h}$  for  $\mathbf{x}$ . The loss function for precision@k can be defined as  $\ell_{p@k} = 1 - \text{precision@}k$ . The conditional risk is then:

$$\begin{aligned} L_{p@k}(\mathbf{h} | \mathbf{x}) &= \mathbb{E}_{\mathbf{y}} \ell_{p@k}(\mathbf{y}, \mathbf{x}, \mathbf{h}) \\ &= 1 - \sum_{\mathbf{y} \in \mathcal{Y}} \mathbf{P}(\mathbf{y} | \mathbf{x}) \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \mathbb{1}[y_j = 1] \\ &= 1 - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \sum_{\mathbf{y} \in \mathcal{Y}} \mathbf{P}(\mathbf{y} | \mathbf{x}) \mathbb{1}[y_j = 1] \\ &= 1 - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \eta_j(\mathbf{x}). \end{aligned}$$

The above result shows that the optimal strategy for precision@k is to predict  $k$  labels with the highest marginal probabilities  $\eta_j(\mathbf{x})$ .

We show now that PLTs obey strong theoretical guarantees. We first recall the result from (Jasinska et al., 2016) that relates the absolute difference between the true and the estimated marginal probability of label  $j$ ,  $|\eta_j(\mathbf{x}) - \hat{\eta}_j(\mathbf{x})|$ , to the surrogate loss  $\ell$  used to train node classifiers  $f_{z^i}$ . It is assumed here that  $\ell$  is a strongly proper composite loss (e.g, logistic, exponential, or squared-error loss) characterized by a constant  $\lambda$  (e.g.  $\lambda = 4$  for logistic loss).<sup>5</sup>

**Theorem 1.** *For any distribution  $\mathbf{P}$  and internal node classifiers  $f_{z^i}$ , the following holds:*

$$|\eta_j(\mathbf{x}) - \hat{\eta}_j(\mathbf{x})| \leq \sum_{i=0}^l \mathbf{P}(z^{i-1} | \mathbf{x}) \sqrt{\frac{2}{\lambda}} \sqrt{\text{reg}_\ell(f_{z^i} | z^{i-1}, \mathbf{x})},$$

where  $\text{reg}_\ell(f_{z^i} | z^{i-1}, \mathbf{x})$  is a binary classification regret for a strongly proper composite loss  $\ell$  and  $\lambda$  is a constant specific for loss  $\ell$ .

*Proof.* From Equation (4) we have:

$$\eta_j(\mathbf{x}) = \mathbf{P}(z | \mathbf{x}) = \prod_{i=0}^l \mathbf{P}(z_i | z^{i-1}, \mathbf{x}) = \mathbf{P}(z^{n-1} | \mathbf{x}) \prod_{i=n}^l \mathbf{P}(z_i | z^{i-1}, \mathbf{x}),$$

for any  $1 \leq n \leq l$ . A similar equation holds for the estimates  $\hat{\eta}_j(\mathbf{x})$ ,  $\hat{\mathbf{P}}(z_i | z^{i-1}, \mathbf{x})$ , and  $\hat{\mathbf{P}}(z^{n-1} | \mathbf{x})$ .

By expressing  $\eta_j(\mathbf{x})$  and  $\hat{\eta}_j(\mathbf{x})$  in the aforementioned way we get:

$$\begin{aligned} |\eta_j(\mathbf{x}) - \hat{\eta}_j(\mathbf{x})| &= \left| \prod_{i=0}^l \mathbf{P}(z_i | z^{i-1}, \mathbf{x}) - \prod_{i=0}^l \hat{\mathbf{P}}(z_i | z^{i-1}, \mathbf{x}) \right| \\ &= \left| \mathbf{P}(z^{l-1} | \mathbf{x}) \mathbf{P}(z_l | z^{l-1}, \mathbf{x}) - \hat{\mathbf{P}}(z^{l-1} | \mathbf{x}) \hat{\mathbf{P}}(z_l | z^{l-1}, \mathbf{x}) \right| \end{aligned}$$

<sup>5</sup>For more detailed introduction to strongly proper composite losses, we refer the reader to (Agarwal, 2014)

By adding and subtracting  $\widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x})\mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x})$  and using the fact that  $|a - b| \leq |a - c| + |b - c|$ , and that probability values are in  $[0, 1]$ , we can write:

$$\begin{aligned}
|\eta_j(\mathbf{x}) - \hat{\eta}_j(\mathbf{x})| &= \left| \mathbf{P}(\mathbf{z}^l | \mathbf{x}) - \widehat{\mathbf{P}}(\mathbf{z}^l | \mathbf{x}) \right| \\
&= \left| \mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x})\mathbf{P}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) - \widehat{\mathbf{P}}(\mathbf{z}^{l-1} | \mathbf{x})\widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) \right| \\
&= \left| \mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x})\mathbf{P}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) - \widehat{\mathbf{P}}(\mathbf{z}^{l-1} | \mathbf{x})\widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) \right. \\
&\quad \left. + \widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x})\mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) - \widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x})\mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) \right| \\
&\leq \left| \mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x})\mathbf{P}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) - \widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x})\mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) \right| \\
&\quad + \left| \widehat{\mathbf{P}}(\mathbf{z}^{l-1} | \mathbf{x})\widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) - \widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x})\mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) \right| \\
&= \mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) \left| \mathbf{P}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) - \widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) \right| \\
&\quad + \widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) \left| \widehat{\mathbf{P}}(\mathbf{z}^{l-1} | \mathbf{x}) - \mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) \right| \\
&\leq \mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) \left| \mathbf{P}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) - \widehat{\mathbf{P}}(z_l | \mathbf{z}^{l-1}, \mathbf{x}) \right| \\
&\quad + \left| \mathbf{P}(\mathbf{z}^{l-1} | \mathbf{x}) - \widehat{\mathbf{P}}(\mathbf{z}^{l-1} | \mathbf{x}) \right|
\end{aligned}$$

We notice that rightmost term corresponds to the absolute value of the difference of probabilities corresponding to one-element shorter code  $\mathbf{z}^{l-1}$ . Therefore we can use recursion and write:

$$|\eta_j(\mathbf{x}) - \hat{\eta}_j(\mathbf{x})| \leq \sum_{i=0}^l \mathbf{P}(\mathbf{z}^{i-1} | \mathbf{x}) \left| \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}) - \widehat{\mathbf{P}}(z_i | \mathbf{z}^{i-1}, \mathbf{x}) \right|. \quad (6)$$

Next, we express the above bound in terms of the regret of the strongly proper composite losses. The  $(\mathbf{x}, z_i)$  pairs are generated i.i.d. according to  $\mathbf{P}(\mathbf{x}, z_i | \mathbf{z}^{i-1})$ . Assume that a node classifier has a form of a real-valued function  $f_{\mathbf{z}^i}$ . Moreover, there exists a strictly increasing (and therefore invertible) link function  $\psi : [0, 1] \rightarrow \mathbb{R}$  such that  $f_{\mathbf{z}^i}(\mathbf{x}) = \psi(\mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}))$ . Recall that the regret of  $f_{\mathbf{z}^i}$  in terms of a loss function  $\ell$  at point  $\mathbf{x}$  is defined as:

$$\text{reg}_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x}) = L_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x}) - L_\ell^*(\mathbf{z}^{i-1}, \mathbf{x}),$$

where  $L_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x})$  is the expected loss at point  $\mathbf{x}$ :

$$L_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x}) = \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x})\ell(1, f_{\mathbf{z}^i}(\mathbf{x})) + (1 - \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}))\ell(-1, f_{\mathbf{z}^i}(\mathbf{x})),$$

and  $L_\ell^*(\mathbf{x})$  is the minimum expected loss at point  $\mathbf{x}$ .

If a node classifier is trained by a learning algorithm that minimizes a strongly proper composite loss, then the bound (6) can be expressed in terms of the regret of this loss function (Agarwal, 2014):

$$\left| \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}) - \psi^{-1}(f_{\mathbf{z}^i}) \right| \leq \sqrt{\frac{2}{\lambda}} \sqrt{\text{reg}_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x})}.$$

By putting the above inequality into (6), we get

$$\begin{aligned}
|\eta(\mathbf{x}, j) - \hat{\eta}(\mathbf{x}, j)| &\leq \sum_{i=0}^l \mathbf{P}(\mathbf{z}^{i-1} | \mathbf{x}) \left| \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}) - \widehat{\mathbf{P}}(z_i | \mathbf{z}^{i-1}, \mathbf{x}) \right| \\
&= \sum_{i=0}^l \mathbf{P}(\mathbf{z}^{i-1} | \mathbf{x}) \left| \mathbf{P}(z_i | \mathbf{z}^{i-1}, \mathbf{x}) - \psi^{-1}(f_{\mathbf{z}^i}) \right| \\
&\leq \sum_{i=0}^l \mathbf{P}(\mathbf{z}^{i-1} | \mathbf{x}) \sqrt{\frac{2}{\lambda}} \sqrt{\text{reg}_\ell(f_{\mathbf{z}^i} | \mathbf{z}^{i-1}, \mathbf{x})}
\end{aligned}$$

□

The above result shows that the absolute error of estimating the marginal probability of label  $j$  can be upperbounded by the regret of the node classifiers on the corresponding path from the root to a leaf. Moreover, for zero-regret (i.e., optimal) node classifiers we obtain an optimal multi-label classifier in terms of estimation of marginal probabilities  $\eta_j(\mathbf{x})$ . This result can be further extended for precision@ $k$ .

Let us denote a set of the top  $k$  labels with respect to the true marginals by  $\mathcal{Y}_k$  and a set of the top  $k$  labels with respect to predicted marginals by  $\hat{\mathcal{Y}}_k$ . The conditional regret for precision@ $k$  is given then by:

$$\text{reg}_{p@k}(\mathbf{h} | \mathbf{x}) = \frac{1}{k} \sum_{i \in \mathcal{Y}_k} \eta_i(\mathbf{x}) - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \eta_j(\mathbf{x})$$

**Theorem 2.** *For any distribution  $\mathbf{P}$  and classifier  $\mathbf{h}$  delivering estimates  $\hat{\eta}_j(\mathbf{x})$  of the marginal probabilities of labels, the following holds:*

$$\text{reg}_{p@k}(\mathbf{h} | \mathbf{x}) = \frac{1}{k} \sum_{i \in \mathcal{Y}_k} \eta_i(\mathbf{x}) - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \eta_j(\mathbf{x}) \leq 2 \max_l |\eta_l(\mathbf{x}) - \hat{\eta}_l(\mathbf{x})|$$

*Proof.* Let us add and subtract the following two terms,  $\frac{1}{k} \sum_{i \in \mathcal{Y}_k} \hat{\eta}_i(\mathbf{x})$  and  $\frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \hat{\eta}_j(\mathbf{x})$ , to the regret and reorganize the expression in the following way:

$$\begin{aligned} \text{reg}_{p@k}(\mathbf{h} | \mathbf{x}) &= \frac{1}{k} \sum_{i \in \mathcal{Y}_k} \eta_i(\mathbf{x}) - \frac{1}{k} \sum_{i \in \mathcal{Y}_k} \hat{\eta}_i(\mathbf{x}) \\ &\quad \underbrace{\leq \frac{1}{k} \sum_{i \in \mathcal{Y}_k} |\eta_i(\mathbf{x}) - \hat{\eta}_i(\mathbf{x})|} \\ &+ \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \hat{\eta}_j(\mathbf{x}) - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \eta_j(\mathbf{x}) \\ &\quad \underbrace{\leq \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} |\hat{\eta}_j(\mathbf{x}) - \eta_j(\mathbf{x})|} \\ &+ \frac{1}{k} \sum_{i \in \mathcal{Y}_k} \hat{\eta}_i(\mathbf{x}) - \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \hat{\eta}_j(\mathbf{x}) \\ &\quad \underbrace{\leq 0} \end{aligned}$$

Because of the relations given under the braces, we finally get:

$$\text{reg}_{p@k}(\mathbf{h} | \mathbf{x}) \leq 2 \max_l |\eta_l(\mathbf{x}) - \hat{\eta}_l(\mathbf{x})| .$$

□

By getting together both theorems we get an upper bound of the precision@ $k$  regret expressed in terms of the regret of the node classifiers. Again, for the zero-regret node classifiers, we get optimal solution in terms of precision@ $k$ .

## B Hierarchical softmax with the pick-one-label heuristic

**Proposition 2.** *Given conditionally independent labels, a classifier  $\mathbf{h}$  such that  $h_j(\mathbf{x}) = \eta'_j(\mathbf{x})$  for all  $j \in \{1, \dots, m\}$  has zero regret in terms of the precision@ $k$  loss.*

*Proof.* To proof the proposition it suffices to show that for conditionally independent labels the order of labels induced by the marginal probabilities  $\eta_j(\mathbf{x})$  is the same as the order induced by the values of  $\eta'_j(\mathbf{x})$  obtained by the pick-one-label heuristic (3):

$$\eta'_j(\mathbf{x}) = \mathbf{P}'(y_j = 1 | \mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} \frac{y_j}{\sum_{j'=1}^m y_{j'}} \mathbf{P}(\mathbf{y} | \mathbf{x}).$$

In other words, for any two labels  $i, j \in \{1, \dots, m\}$ ,  $i \neq j$ ,  $\eta_i(\mathbf{x}) \geq \eta_j(\mathbf{x}) \Leftrightarrow \eta'_i(\mathbf{x}) \geq \eta'_j(\mathbf{x})$ .

Let  $\eta_i(\mathbf{x}) \geq \eta_j(\mathbf{x})$ . The summation over all  $\mathbf{y}$  in (3) can be written in the following way:

$$\eta'_j(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} y_j N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x}),$$

where  $N(\mathbf{y}) = (\sum_{i=1}^m y_i)^{-1}$  is a value that depends only on the number of positive labels in  $\mathbf{y}$ . In this summation we consider four subsets of  $\mathcal{Y}$ , creating a partition of this set:

$$\mathcal{S}_{i,j}^{u,w} = \{\mathbf{y} \in \mathcal{Y} : y_i = u \wedge y_j = w\}, \quad u, w \in \{0, 1\}.$$

The subset  $\mathcal{S}_{i,j}^{0,0}$  does not play any role because  $y_i = y_j = 0$  and therefore do not contribute to the final sum. Then (3) can be written in the following way for the  $i$ -th and  $j$ -th label:

$$\eta'_i(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{S}_{i,j}^{1,0}} N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x}) + \sum_{\mathbf{y} \in \mathcal{S}_{i,j}^{1,1}} N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x}) \quad (7)$$

$$\eta'_j(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{S}_{i,j}^{0,1}} N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x}) + \sum_{\mathbf{y} \in \mathcal{S}_{i,j}^{1,1}} N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x}) \quad (8)$$

The contribution of elements from  $\mathcal{S}_{i,j}^{1,1}$  is equal for both  $\eta'_i(\mathbf{x})$  and  $\eta'_j(\mathbf{x})$ . It is so because the value of  $N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x})$  is the same for all  $\mathbf{y} \in \mathcal{S}_{i,j}^{1,1}$ : the conditional joint probabilities  $\mathbf{P}(\mathbf{y}|\mathbf{x})$  are fixed and they are multiplied by the same factors  $N(\mathbf{y})$ .

Consider now the contributions of  $\mathcal{S}_{i,j}^{1,0}$  and  $\mathcal{S}_{i,j}^{0,1}$  to the relevant sums. By the definition of  $\mathcal{Y}$ ,  $\mathcal{S}_{i,j}^{1,0}$ , and  $\mathcal{S}_{i,j}^{0,1}$ , there exists bijection  $b_{i,j} : \mathcal{S}_{i,j}^{1,0} \rightarrow \mathcal{S}_{i,j}^{0,1}$ , such that for each  $\mathbf{y}' \in \mathcal{S}_{i,j}^{1,0}$  there exists  $\mathbf{y}'' \in \mathcal{S}_{i,j}^{0,1}$  equal to  $\mathbf{y}'$  except on the  $i$ -th and the  $j$ -th position.

Notice that because of the conditional independence assumption the joint probabilities of elements in  $\mathcal{S}_{i,j}^{1,0}$  and  $\mathcal{S}_{i,j}^{0,1}$  are related to each other. Let  $\mathbf{y}'' = b_{i,j}(\mathbf{y}')$ , where  $\mathbf{y}' \in \mathcal{S}_{i,j}^{1,0}$  and  $\mathbf{y}'' \in \mathcal{S}_{i,j}^{0,1}$ . The joint probabilities are:

$$\mathbf{P}(\mathbf{y}'|\mathbf{x}) = \eta_i(\mathbf{x})(1 - \eta_j(\mathbf{x})) \prod_{l \in \mathcal{L} \setminus \{i,j\}} \eta_l(\mathbf{x})^{y_l} (1 - \eta_l(\mathbf{x}))^{1-y_l}$$

and

$$\mathbf{P}(\mathbf{y}''|\mathbf{x}) = (1 - \eta_i(\mathbf{x}))\eta_j(\mathbf{x}) \prod_{l \in \mathcal{L} \setminus \{i,j\}} \eta_l(\mathbf{x})^{y_l} (1 - \eta_l(\mathbf{x}))^{1-y_l}.$$

One can easily notice the relation between these probabilities:

$$\mathbf{P}(\mathbf{y}'|\mathbf{x}) = \eta_i(\mathbf{x})(1 - \eta_j(\mathbf{x}))q_{i,j} \quad \text{and} \quad \mathbf{P}(\mathbf{y}''|\mathbf{x}) = (1 - \eta_i(\mathbf{x}))\eta_j(\mathbf{x})q_{i,j},$$

where  $q_{i,j} = \prod_{l \in \mathcal{L} \setminus \{i,j\}} \eta_l(\mathbf{x})^{y_l} (1 - \eta_l(\mathbf{x}))^{1-y_l} \geq 0$ . Consider now the difference of these two probabilities:

$$\begin{aligned} \mathbf{P}(\mathbf{y}'|\mathbf{x}) - \mathbf{P}(\mathbf{y}''|\mathbf{x}) &= \eta_i(\mathbf{x})(1 - \eta_j(\mathbf{x}))q_{i,j} - (1 - \eta_i(\mathbf{x}))\eta_j(\mathbf{x})q_{i,j} \\ &= q_{i,j}(\eta_i(\mathbf{x})(1 - \eta_j(\mathbf{x})) - (1 - \eta_i(\mathbf{x}))\eta_j(\mathbf{x})) \\ &= q_{i,j}(\eta_i(\mathbf{x}) - \eta_j(\mathbf{x})). \end{aligned}$$

From the above we see that  $\eta_i(\mathbf{x}) \geq \eta_j(\mathbf{x}) \Rightarrow \mathbf{P}(\mathbf{y}'|\mathbf{x}) \geq \mathbf{P}(\mathbf{y}''|\mathbf{x})$ . Due to the properties of the bijection  $b_{i,j}$ , the number of positive labels in  $\mathbf{y}'$  and  $\mathbf{y}''$  is the same and  $N(\mathbf{y}') = N(\mathbf{y}'')$ , therefore we also get  $\eta_i(\mathbf{x}) \geq \eta_j(\mathbf{x}) \Rightarrow \sum_{\mathbf{y} \in \mathcal{S}_{i,j}^{1,0}} N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x}) \geq \sum_{\mathbf{y} \in \mathcal{S}_{i,j}^{0,1}} N(\mathbf{y}) \mathbf{P}(\mathbf{y}|\mathbf{x})$ , which by (7) and (8) gives us finally  $\eta_i(\mathbf{x}) \geq \eta_j(\mathbf{x}) \Rightarrow \eta'_i(\mathbf{x}) \geq \eta'_j(\mathbf{x})$ .

The implication in the other side, i.e.,  $\eta_i(\mathbf{x}) \geq \eta_j(\mathbf{x}) \Leftarrow \mathbf{P}(\mathbf{y}'|\mathbf{x}) \geq \mathbf{P}(\mathbf{y}''|\mathbf{x})$  holds obviously for  $q_{i,j} > 0$ . For  $q_{i,j} = 0$ , we can notice, however, that  $\mathbf{P}(\mathbf{y}'|\mathbf{x})$  and  $\mathbf{P}(\mathbf{y}''|\mathbf{x})$  do not contribute to the appropriate sums as they are zero, and therefore we can follow a similar reasoning as above, concluding that  $\eta_i(\mathbf{x}) \geq \eta_j(\mathbf{x}) \Leftarrow \eta'_i(\mathbf{x}) \geq \eta'_j(\mathbf{x})$ .

Thus for conditionally independent labels, the order of labels induced by marginal probabilities  $\eta_j(\mathbf{x})$  is equal to the order induced by  $\eta'_j(\mathbf{x})$ . As the precision@ $k$  is optimized by  $k$  labels with the highest marginal probabilities, we have that prediction consisted of  $k$  labels with highest  $\eta'_j(\mathbf{x})$  has zero regret for precision@ $k$ .  $\square$

## C Huffman codes for PLTs

In this section we analyze computational cost of PLTs in binary case, i.e. every inner node has two children. We define the cost of a tree as the total expected fraction of instances which are used in the inner nodes and show that for the multi-class case minimization of this cost coincides with minimization of Huffman criteria. However, in the multi-label case, this does no longer hold.

We shall use prefix codes, as it is introduced in Section 4, to identify a path from the root to the leaf. Accordingly, a prefix code  $\mathbf{z} = (z_1, \dots, z_\ell) \in \mathcal{C}$  determines a path with length  $|\mathbf{z}| = \ell$ . The probability to observe a (possibly partial) prefix code is  $p_{\mathbf{z}^i} = \mathbf{P}(\mathbf{z}^i)$  with respect to the data distribution. If we are given a data that consists of  $n$  instances, then the expected number of positive instances that is used to train the node classifier in node  $\mathbf{z}^i$  is  $np_{\mathbf{z}^i}$ . Note that in the training process of PLT, only the positive instances shows up in the training data of a child node of  $\mathbf{z}^i$ , thus one can define the expected computational cost of a PLT with fixed structure, thus with fixed set of prefix code  $\mathcal{C}$ , as

$$\sum_{\mathbf{z} \in \mathcal{C}} \sum_{i=1}^{|\mathbf{z}|} p_{\mathbf{z}^{i-1}}. \quad (9)$$

The next proposition defines the relation of (9) and Huffman coding. Huffman coding can be naturally defined in the multi-class case using  $\mathbf{P}(y_i = 1)$  as weights.

**Proposition 3.** *Huffman code minimizes the expected computational cost that is given in (9) among binary codes if the data is multi-class.*

*Proof.* Recall that the Huffman code minimizes the following criterion:

$$\sum_{\mathbf{z} \in \mathcal{C}} p_{\mathbf{z}} |\mathbf{z}|.$$

In multi-class setting score  $p_{\mathbf{z}^i}$  of each inner node  $\mathbf{z}^i$  equals to the sum of scores of its children, and the objective function given in (9) can be rewritten as

$$\sum_{\mathbf{z} \in \mathcal{C}} \sum_{i=1}^{|\mathbf{z}|} p_{\mathbf{z}^{i-1}} = 1 + \sum_{\mathbf{z} \in \mathcal{C}} p_{\mathbf{z}} \sum_{i=1}^{|\mathbf{z}|} \deg(\mathbf{z}^i),$$

where  $\deg(\mathbf{z}^i) = \#\{c : \hat{z}_{i+1} = c, \mathbf{z}^i = \hat{\mathbf{z}}^i, \hat{\mathbf{z}} \in \mathcal{C}\}$  is the number of children of the node  $\mathbf{z}^i$ . In the case of binary codes  $\deg(\mathbf{z}^i) = 2$  and

$$1 + \sum_{\mathbf{z} \in \mathcal{C}} p_{\mathbf{z}} \sum_{i=1}^{|\mathbf{z}|} \deg(\mathbf{z}^i) = 1 + 2 \sum_{\mathbf{z} \in \mathcal{C}} p_{\mathbf{z}} |\mathbf{z}|$$

which completes the proof.  $\square$

Note that Grave et al. (2017) considered similar notion of computational cost and optimized in a restricted setup, assuming that the tree depth is at most two. Of course, in practice, the  $p_{\mathbf{z}^i}$  values are not known, but are estimated based on observations from the data distribution. Let us denote an estimate  $\hat{p}_{\mathbf{z}^i}$  of  $p_{\mathbf{z}^i}$ . An interesting question to address is that to what extent the empirical computational cost  $\sum_{\mathbf{z} \in \mathcal{C}} \sum_{i=1}^{|\mathbf{z}|} \hat{p}_{\mathbf{z}^{i-1}}$  concentrates around its expected values, and on what parameter of the code/tree structure does depend on.

For multi-label case score of an inner node cannot be represented by a sum of children scores, and this result does not hold.

## D Pseudocode of PLTs

The pseudocode below presents the training and prediction procedures of PLTs in detail. PLTs can be trained either in the online or batch mode. Algorithm 1 follows the former mode used, for example, in XT, the state-of-the-art variant of PLTs described in Section 7. Note that learning of feature embeddings used in XT is not included in the pseudocode.



In Algorithm 2 we present an inference algorithm, which uses the uniform-cost search for finding the top  $k$  labels. The algorithm searches the tree by starting from the root node  $z^0$ . It uses a priority queue  $Q$  to store pairs of the node  $z^i$  and its probability estimate  $\hat{\eta}_{z^i}$ . The algorithm pops the element with the highest probability estimate from the queue. If the element is a leaf we add the corresponding label to the final prediction. Otherwise, the algorithm estimates probabilities of the children nodes and adds them into the queue. Once we found the  $k$ -th label or we run out of the pairs in the queue, we stop the search procedure.

---

**Algorithm 1** Incremental learning of a PLT:

---

**Input:**  
 $T$ : a label tree with  $t$  nodes  
 $A_{\text{online}}$ : an incremental learning algorithm  
 $\mathcal{D}_N$ : a set of training examples  $(\mathbf{x}, \mathbf{y})$

**Output:**  
 $\mathcal{F}_t = \{f_{z^i}\}^t$ : a set of  $t$  node (binary) classifiers

- 1:  $\mathcal{F}_t = \emptyset$
- 2: **for** each node  $z^i \in T$  **do** ▷ Initialization of binary classifiers
- 3:      $\mathcal{F}_t \leftarrow \mathcal{F}_t \cup$  new classifier  $f_{z^i}$
- 4: **for** each training example  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_N$  **do**
- 5:     **if**  $\sum_{j=1}^m y_j = 0$  **then** ▷ Select nodes for the positive and negative update
- 6:          $\mathcal{Z}_{\text{positive}} \leftarrow \emptyset$
- 7:          $\mathcal{Z}_{\text{-positive}} \leftarrow z^0$
- 8:     **else**
- 9:         **for** each  $y_j = 1, j \in \{1, \dots, m\}$  **do**
- 10:              $z \leftarrow$  encode  $j$
- 11:              $\mathcal{Z}_{\text{positive}} \leftarrow \mathcal{Z}_{\text{positive}} \cup \left\{ \bigcup_{i=0}^l z^i \right\}$
- 12:             **for** each  $z^i \in \mathcal{Z}_{\text{positive}}$  **do**
- 13:                  $\mathcal{Z}_{\text{-positive}} \leftarrow \mathcal{Z}_{\text{-positive}} \cup \left\{ \bigcup_{z_i} (z^{i-1}, z_i) \right\} \setminus \mathcal{Z}_{\text{positive}}$
- 14:     **for** each node  $z^i \in \mathcal{Z}_{\text{positive}}$  **do** ▷ Update node classifiers
- 15:          $f_{z^i} \leftarrow A_{\text{online}}(f_{z^i}, \mathbf{x}, 1)$
- 16:     **for** each node  $z^i \in \mathcal{Z}_{\text{-positive}}$  **do**
- 17:          $f_{z^i} \leftarrow A_{\text{online}}(f_{z^i}, \mathbf{x}, 0)$
- 18: **return**  $\mathcal{F}_t$ .

---

## E Additional experimental results

### E.1 Comparison of PLTs and HSM on synthetic data

In this section we validate our theoretical results presented in Section 5, which show that HSM is not amenable to model the marginal probabilities in general for multi-label problems. In this case, HSM with pick-one-label heuristic should be outperformed in terms of precision@ $k$  by PLTs which are consistent for this performance measure. To validate this claim empirically, we compare the performance of PLTs and HSM on synthetic datasets of three types: multi-label data with independent labels, multi-label data with conditionally dependent labels and multi-class data.

All synthetic models are based on linear models parametrized by a weight vector  $w$  of size  $d$ . The values of the vector are sampled uniformly from a  $d$ -dimensional sphere of radius 1. Each instance  $\mathbf{x}$ , in turn, is represented as a vector sampled from a  $d$ -dimensional disc of the same radius.

---

**Algorithm 2** Top-k prediction with a PLT:

---

**Input:**

$T$ : a label tree with  $t$  nodes  
 $\mathcal{F}_t = \{f_{z^i}\}^t$ : a set of  $t$  node (binary) classifiers  
 $\mathbf{x}$ : a test example  
 $k$ : a size of prediction

**Output:**

$\hat{\mathbf{y}}$ : a vector with top  $k$  labels for the test example

```
1:  $\hat{\mathbf{y}} \leftarrow \{0\}^m$ 
2:  $Q \leftarrow \text{PRIORITYQUEUE}()$  ▷ Initialization of priority queue
3:  $\text{add}(Q, (1, z^0))$ 
4: while  $Q \neq \emptyset$  and  $\sum_{j=1}^m \hat{y}_j < k$  do ▷ Check if  $k$  labels not found
5:    $(\hat{\eta}_{z^i}, z^i) \leftarrow \text{pop}(Q)$ 
6:   if  $i \in \mathcal{C}$  then ▷ Check if leaf node reached
7:      $j \leftarrow \text{decode } z^i$ 
8:      $\hat{y}_j \leftarrow 1$ 
9:   else
10:    for each  $z_{i+1}$  do ▷ For each node's children
11:       $z^{i+1} \leftarrow (z^i, z_i)$ 
12:       $\text{add}(Q, (\hat{\eta}_{z^i} \cdot \hat{\mathbf{P}}_{f_{z^i}}(z_{i+1} | z^i, \mathbf{x}), z^{i+1}))$ 
13: return  $\hat{\mathbf{y}}$ .
```

---

**Multi-class distribution.** We associate a weigh vector  $\mathbf{w}_j$  with each label  $j \in \{1, \dots, m\}$ . The model assigns probabilities to labels at point  $\mathbf{x}$  based on the softmax schema

$$\eta_j(\mathbf{x}) = \frac{\exp(c\mathbf{w}_j^\top \mathbf{x})}{\sum_{j'=1}^m \exp(c\mathbf{w}_{j'}^\top \mathbf{x})} \quad (10)$$

and draws the positive label according to this probability distribution over labels. Scaling factor  $c$  is added to control noise in the model. Higher values of  $c$  give less noisy model.

**Multi-label distribution with conditionally independent labels.** The model is similar to the previous one used for the multi-class distribution. The difference lays is normalization as the marginal probabilities do not have to sum up to 1. To get a probability of the  $j$ -th label, we use the logistic transformation:

$$\eta_j(\mathbf{x}) = \frac{\exp(\mathbf{w}_j^\top \mathbf{x})}{1 + \exp(\mathbf{w}_j^\top \mathbf{x})}.$$

Then, we assign a label to an instance based on:

$$y_j = \llbracket r < \eta_j(\mathbf{x}) \rrbracket,$$

where the random value  $r$  is sampled uniformly and independently from range  $[0, 1]$  for each instance  $\mathbf{x}$  and label  $j \in \{1, \dots, m\}$ .

**Multi-label distribution with conditionally dependent labels.** To model conditionally dependent labels we use the mixing matrix model based on latent scoring functions generated by  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m)$ . The  $m \times m$  mixing matrix  $\mathbf{M}$  introduces dependencies between noise  $\epsilon$ , which stands for the source of randomness in the model. The models  $\mathbf{w}_j$  are sampled from a sphere of radius 1, as in previous cases. The values in the mixing matrix  $\mathbf{M}$  are sampled uniformly and independently from  $[-1, 1]$ . The random noise vector  $\epsilon$  is sampled from  $N(0, 0.25)$ . The label vector  $\mathbf{y}$  is then obtained by element-wise evaluation of the following expression:

$$\mathbf{y} = \llbracket \mathbf{M}(\mathbf{W}^\top \mathbf{x} + \epsilon) > 0 \rrbracket$$

Notice that if  $\mathbf{M}$  was an identity matrix the model would generate independent labels.

Table 2: Means and standard deviations of 50 runs of each experiment. The p-values on the right indicate the significance of the observed differences. The results of HSM and PLTs on multi-class data were always equal.

	HSM		PLT		t-test	p-values	
	mean	stdev	mean	stdev		sign	Wilcoxon
multi-class	21.90	2.74	21.90	2.74			
multi-label independent	32.57	0.34	32.58	0.33	0.4367	0.3222	0.5980
multi-label dependent	70.68	5.73	71.68	5.65	9.80E-14	3.71E-11	3.90E-09

**Experimental setting.** PLTs and HSM are usually implemented as online learning algorithms, i.e., the node classifiers are updated in an incremental way, example by example. To minimize the impact of the hyperparameter tuning of online algorithms, we have decided to implement batch versions of both algorithms using the LIBLINEAR-based (Fan et al., 2008) logistic regression. The sets of training instances for each node classifier are appropriately constructed by taking the corresponding conditioning of the probabilistic models into account. In the case of HSM with the pick-one-label heuristic, we first transform each multi-label example to  $s$  multi-class copies of it, one copy for each its label. Each such copy gets then a weight of  $1/s$ . Such transformation should be concordant with the theoretical model (3).

In the experiments we used the following parameters of the synthetic models:  $d = 3$  (i.e., the number of features),  $n = 100000$  instances (split 1 : 1 for training and test subsets), and  $m = 32$  labels or classes. In the case of the multi-class model we report results with the scaling factor  $c = 10$ . The choice of  $c$  does not change the interpretation of the results. To train logistic regression we use a fixed value of the regularization parameter, standing for a very weak regularization. For all experiments we report the results in terms of precision@1.

**Observations.** Table 2 presents the average results of all experiments besides with the standard deviation of obtained values. As expected the performance of PLTs and HSM on the multi-class data are exactly the same. For the other models, we additionally include the p-values of statistical tests run to verify, whether there is a significant difference in performance between PLTs and HSM with the pick-one-label heuristic. In the case of the multi-label data with conditionally independent labels the test shows that there is no evidence to reject the hypothesis that the performance of PLTs and HSM is the same. In the case of multi-label data with conditionally dependent labels, the statistical tests show that PLTs are significantly better than HSM with the pick-one-label heuristic.

## E.2 Comparison of PLTs and HSM on benchmark data

In Table 3 we present similar results, but obtained on the benchmark datasets. We use two models. The first one follows the sparse representation and uses an implementation of PLTs and HSM in VOWPAL WABBIT (Langford et al., 2007). The second one is based on FASTTEXT and produces the dense representation. In all models we use Huffman trees. The results clearly indicate the better performance of PLTs over HSM.

Table 3: Precision@ $k$  with  $k = \{1, 3, 5\}$  of a simple HSM and a simple PLT implementations.

Dataset	VOWPAL WABBIT		FASTTEXT		Dataset	VOWPAL WABBIT		FASTTEXT	
	HSM	PLT	HSM	PLT		HSM	PLT	HSM	PLT
EUR-Lex	56.98	<b>74.55</b>	66.39	<b>73.19</b>	AmazonCat-13K	86.69	<b>91.46</b>	90.18	<b>92.98</b>
	46.99	<b>60.60</b>	54.05	<b>57.79</b>		72.00	<b>76.00</b>	72.53	<b>75.75</b>
	39.09	<b>50.05</b>	44.73	<b>46.98</b>		57.97	<b>61.40</b>	56.20	<b>59.53</b>
Wiki-30K	70.20	<b>84.34</b>	83.02	<b>85.11</b>	Delicious-200K	41.58	<b>45.27</b>	42.17	<b>46.98</b>
	60.11	<b>72.34</b>	69.66	<b>73.12</b>		33.24	<b>38.95</b>	37.94	<b>40.99</b>
	53.17	<b>62.72</b>	59.50	<b>62.67</b>		28.04	<b>35.59</b>	35.77	<b>38.04</b>
WikiLSHTC-325K	36.90	<b>41.63</b>	41.28	<b>41.78</b>	Amazon-670K	33.64	<b>36.85</b>	25.04	<b>26.18</b>
	22.30	<b>26.78</b>	24.68	<b>24.96</b>		28.58	<b>32.48</b>	21.06	<b>22.76</b>
	16.60	<b>20.39</b>	18.08	<b>18.53</b>		25.01	<b>29.15</b>	18.28	<b>20.29</b>

### E.3 The ablation analysis on benchmark datasets

Table 4 contains results of the ablation analysis in which we compare different components of the XT algorithm. We analyze the influence of the Huffman tree vs. top-down clustering, the simple averaging of features vectors vs. the TF-IDF-based weighting, and no regularization vs. L2 regularization. For every configuration, we conducted a grid search of hyperparameters from ranges reported in Appendix E.4. The results clearly show that the components need to be combined together to obtain the best results. The best combination is usually the one that uses top-down clustering, TF-IDF-based weighting, and L2 regularization. It is worth to notice that top-down clustering alone gets worse results than Huffman trees with TF-IDF-based weighting and L2 regularization.

Table 4: Precision@ $k$  scores with  $k = \{1, 3, 5\}$  of different variants of PLT in FASTTEXT (EXTREMETEXT)

Dataset	Metrics	Huff.	Huff. + TF-IDF	Huff. + L2	Huff. + L2 + TF-IDF	Clus.	Clus. + TF-IDF	Clus. + L2	Clus. + L2 + TF-IDF
Eurlex	P@1	63.39	71.20	62.79	74.60	68.31	75.05	65.05	<b>77.68</b>
	P@3	50.48	57.26	48.99	60.36	54.62	61.65	50.87	<b>63.37</b>
	P@5	41.19	46.93	40.16	49.72	45.23	50.99	41.71	<b>52.85</b>
AmazonCat-13K	P@1	90.10	89.19	90.84	91.08	72.95	77.13	91.73	<b>92.43</b>
	P@3	72.67	72.81	73.10	75.27	61.66	65.76	75.07	<b>77.65</b>
	P@5	57.69	58.30	57.69	60.34	48.38	52.72	59.63	<b>62.74</b>
Wiki-30K	P@1	78.04	78.14	82.28	<b>85.23</b>	78.40	78.69	82.13	85.21
	P@3	63.31	67.47	68.75	<b>73.52</b>	65.62	66.40	69.22	73.18
	P@5	52.65	56.00	58.33	<b>63.71</b>	55.61	56.77	58.69	63.39
Delicious-200K	P@1	45.71	47.24	46.48	<b>47.85</b>	44.58	45.13	46.55	47.31
	P@3	40.69	40.88	41.38	<b>42.08</b>	40.42	40.63	41.08	41.63
	P@5	38.02	37.74	38.78	<b>39.13</b>	38.15	38.23	38.25	38.88
WikiLSHTC-325K	P@1	36.23	38.28	41.10	42.83	34.39	39.66	54.95	<b>58.73</b>
	P@3	20.60	22.43	25.62	26.33	21.17	24.79	36.42	<b>39.24</b>
	P@5	14.70	16.33	19.10	19.45	15.09	18.28	27.25	<b>29.26</b>
Wiki-500K	P@1	41.01	40.62	41.55	49.80	46.63	49.20	56.04	<b>64.48</b>
	P@3	24.79	25.68	25.94	32.39	31.59	34.52	38.87	<b>45.84</b>
	P@5	18.64	19.74	19.79	24.62	24.30	26.83	30.46	<b>35.46</b>
Amazon-670K	P@1	23.19	29.70	21.64	32.11	28.54	36.24	28.95	<b>39.90</b>
	P@3	19.80	25.84	18.11	27.77	25.52	32.15	25.74	<b>35.36</b>
	P@5	17.48	22.96	15.83	24.64	23.30	29.19	23.22	<b>32.04</b>

### E.4 Tuning of hyperparameters

The PLT has only one global hyperparameter which is the degree of the tree denoted by  $b$ . The other hyperparameters are associated with the node classifiers. The HSM and PLT in Vowpal Wabbit was tuned with the stochastic gradient descent with a step size  $\eta_t$  calculated separately for each node according to  $\eta_t = \eta \times (1/t)^p$  where  $t$  is number of node updates and  $\eta$  and  $p$  are hyperparameters. In FASTTEXT-based methods, HSM, LEARNED TREE and XT,  $\eta_t$  decreased linearly during training from  $\eta$  to 0.0. In the XT the optimization methods has been extended by L2 regularization, so it has one additional parameter. Balanced k-means clustering used to build a tree in XT has also a stopping parameter  $\epsilon$  set by default to 0.001.

Table 5: The hyperparameters of the HSM and PLT methods and their ranges used in hyperparameter optimization. Notation:  $b$  – tree arity,  $\eta$  – learning rate

Hyperparameter	Range
$b$	$\{2, \dots, 32\}$
$\eta$	$[0.0001 - 1.0]$
number of epochs	$\{20, 30, 40\}$

## F Information about the benchmark datasets

Table 8 contains the basic statistics of the benchmark datasets used in the experiments taken from the Extreme Classification Repository: <http://manikvarma.org/downloads/XC/XMLRepository.html>.

Table 6: The hyperparameters of the FASTTEXT and LEARNED TREE methods and their ranges used in hyperparameter optimization. Notation:  $b$  – tree arity,  $\eta$  – learning rate

Hyperparameter	Range
$b$	$\{2, \dots, 32\}$
$\eta$	$[0.0001 - 1.0]$
number of epochs	$\{20, 30, 40\}$
dim	$\{500\}$

Table 7: The hyperparameters of the XT method and their ranges used in hyperparameter optimization. Notation:  $b$  – tree arity (number of centroids used in k-means clustering),  $\epsilon$  – stopping condition of k-means clustering,  $\eta$  – learning rate

Hyperparameter	Range
$b$	$\{2\}$
$\epsilon$	$\{0.001\}$
$\eta$	$\{0.05, 0.1, 0.2, 0.3, 0.4, 0.5\}$
L2 regularization	$\{0.001, 0.002, 0.003\}$
number of epochs	$\{20, 30, 40\}$
dim	$\{500\}$
max leaves	$\{100\}$

Table 8: Statistics of the benchmark datasets.

Dataset	Number of features	Number of labels	Number of train points	Number of test points	Avg. points per label	Avg. labels per point
EURLex-4K	5000	3993	15539	3809	25.73	5.31
AmazonCat-13K	203882	13330	1186239	306782	448.57	5.04
Wiki10-31K	101938	30938	14146	6616	8.52	18.64
Delicious-200K	782585	205443	196606	100095	72.29	75.54
WikiLSHTC-325K	1617899	325056	1778351	587084	17.46	3.19
Wikipedia-500K	2381304	501070	1813391	783743	24.75	4.77
Amazon-670K	135909	670091	490449	153025	3.99	5.45