# Random Label Forests: An Ensemble Method with Label Subsampling For Extreme Multi-Label Problems

Sheng-Wei Chen National Taiwan University d09944003@ntu.edu.tw

### Abstract

Text classification is one of the essential topics in natural language processing, and each text is often associated with multiple labels. Recently, the number of labels has become larger and larger, especially in the applications of ecommerce, so the large space required by existing multi-label learning methods may no longer be affordable. To address the space concern, utilizing a distributed system to share that large memory requirement is a possible solution. We propose "random label forests," a distributed ensemble method with label subsampling, for handling extremely large-scale labels. Random label forests can reduce the memory usage per computer while keeping competitive performances over real-world data sets.

### 1 Introduction

Text classification is one of the essential topics in natural language processing fields. There are many valuable applications, such as product categorization for e-commerce (Shen et al., 2011; Agrawal et al., 2013; McAuley and Leskovec, 2013), coding diagnosis and procedures in medical records (Nuthakki et al., 2019), and document tagging (Zubiaga, 2009). Usually, the prediction in the text classification can be multi-labeled. Hence, a text classification problem falls into the category of multi-label classification, which is used to find the relevant labels of a data instance. For example, we can set the document contents and tags in the document tagging problem as the features and labels in a multi-label problem.

Recently, the number of labels has become larger and larger, especially in the applications of e-commerce. Thus, an emerging topic is extreme multi-label learning (XML), which focuses on large-scale candidate labels, input instances, and input features. Because of these three large-scale components, an XML method should further consider the model training time and memory usage in Chih-Jen Lin National Taiwan University / MBZUAI cjlin@csie.ntu.edu.tw

addition to the performance.

A simple and classic way to solve a multi-label classification problem is by the one-versus-rest (OVR) method with linear models. However, the time complexity and model size directly depend on the number of labels and features. To handle problems with many labels, two existing extensions of OVR were developed:

- DiSMEC (Babbar and Schölkopf, 2017) splits the label set into several subsets and lets each machine in a distributed system handle one subset. This way, the training time and model size per machine are reduced.
- Tree-based linear methods (Tsoumakas et al.; Prabhu et al., 2018; Khandagale et al., 2020; Yu et al., 2022) utilize the divide-and-conquer paradigm on labels via clustering methods such as *K*-means and then apply the OVR method to train the smaller problems on the clusters. We discuss the details in Section 2.2.

Besides the linear methods, we can also use lowrank embedding on features and labels to reduce the training time and memory usage, e.g., (Bhatia et al., 2015; Yu et al., 2014). However, some works (Khandagale et al., 2020; Babbar and Schölkopf, 2017) report that the distribution of positive instances over labels is long-tail in most XML data sets, so the label space cannot be wellapproximated to a low-rank embedding space. For the deep-learning methods, some earlier works (Kim, 2014; Liu et al., 2017) only show competitive performance on short-text XML problems (Yu et al., 2022, Section 5.2). However, several treebased deep-learning methods (You et al., 2019; Jiang et al., 2021; Zhang et al., 2021) can rank better on the public XML benchmark (Bhatia et al., 2016). Nevertheless, although neural network models perform better in many fields, Lin et al. (2023) point out that the linear model is still a strong baseline for certain multi-label classification data. Furthermore, the linear models are easy to understand

and more explainable, so we focus on linear models in this work.

Let us go back to discuss the linear methods. Although DiSMEC and tree-based linear methods are reasonable solutions for XML problems, each method still has disadvantages. DiSMEC can only handle a small number of labels, but not large-scale labels, on each machine. Tree-based linear methods require large memory space to handle all the labels during the training. Hence, we hope a method can take the advantages of DiSMEC and tree-based methods without having their disadvantages.

Label subsampling is a possible way to divide an XML problem into smaller subproblems. RAKEL (Tsoumakas and Vlahavas, 2007) is a pioneer in using label subsampling with the ensemble method. After the label subsampling, RAKEL converts each small-scale multi-label subproblem to a multi-class one by considering every label combination as a new class label. This setting, referred to as "label powerset" in multi-label learning, is not scalable to XML because covering the predictions of highly large-scale labels by the powerset method is almost impossible.

In this work, we utilize label subsampling and the distributed system to reduce the impact of the large-scale labels. Specifically, each computer can solve a smaller XML subproblem via some existing XML methods, such as tree-based linear methods. Hence, handling XML problems with a large number of labels becomes practical. Since we use the label subsampling technique with the tree-based linear method, we call this method "random label forests." Let us list our contributions as follows.

- We propose a natively parallel framework, random label forests, which is an ensemble method with label subsampling for the XML problem.
- Our experiments show that random label forests are competitive with the standard tree-based methods applied to all labels.
- We analyze the model size of tree-based methods and explain why random label forests can reduce memory usage in each computer of the distributed system.
- We also analyze the time complexity of treebased methods. The training time of random label forests in a distributed environment is shorter than the tree-based methods with all labels in a single computer.

Section 2 discusses OVR and tree-based methods for XML problems. Section 3 focuses on distributing a tree-based model and then presents random label forests, including discussions on the data processing, time complexity, and model size. Section 4 shows comparison results on performance, training time, and model size. Section 5 concludes this work.

# 2 Multi-Label Problems

A multi-label classification problem aims to find a function f with the parameter  $\theta$  that can predict whether a given instance x is associated with the label-j for j = 1, ..., m. We assume that x is a feature vector in  $\mathbb{R}^n$ , and n and m are the feature dimension and the number of labels, respectively. We use 0/1 to indicate if an instance is associated without/with a label. Hence, we can denote  $y \in$  $\{0, 1\}^m$  as the label vector of the instance  $x \in \mathbb{R}^n$ so that the prediction  $f(x; \theta)$  is close to y. We have mentioned in Section 1 that in this work we focus on using linear models in constructing  $f(x; \theta)$ .

### 2.1 One-Versus-Rest Method

The OVR method involves training a single model per label, with the instances of that label as positives and all other instances as negatives. Thus, when training an OVR linear model on a multilabel classification problem with the training set

$$D = \{ (\boldsymbol{y}_i, \boldsymbol{x}_i) \in (\{0, 1\}^m, \mathbb{R}^n) \mid i = 1, \dots, l \},\$$

where l is the number of the training instances, we solve m subproblems

$$\min_{\boldsymbol{w}_j \in \mathbb{R}^n} \frac{\lambda}{2} \boldsymbol{w}_j^T \boldsymbol{w}_j + \sum_{i=1}^l \xi(\boldsymbol{w}_j^T \boldsymbol{x}_i, [\boldsymbol{y}_i]_j), \quad (1)$$

for j = 1, ..., m. The model parameter  $\theta$  includes  $w_1, ..., w_m$ . In each subproblem,  $\lambda$  is the hyperparameter,  $[a]_j$  denotes the *j*th component of the vector a, and  $\xi$  is the loss function for binary classification. Subproblems in (1) can be easily handled by some mature binary classification libraries such as LIBLINEAR (Fan et al., 2008). After the training procedure, we can use the OVR model to get the scores  $[w_1^T x \cdots w_m^T x]$  for any given instance x. Moreover, a 0-1 function  $\delta$  can map the scores to a label vector  $[\delta(w_1^T x) \cdots \delta(w_m^T x)] \in \{0, 1\}^m$ as the prediction.

Many works (Babbar and Schölkopf, 2017; Lin et al., 2023) show that OVR linear models are use-ful, but

(i) the space requirement for the model parameter and

#### (ii) the training time of (1)

increase as m becomes larger. For large XML problems, the issues above become essential.

### 2.2 Tree-Based Methods

To overcome the training time issue (ii), past works (e.g., Prabhu et al., 2018; Khandagale et al., 2020; Yu et al., 2022) utilize the tree structure to reduce the training time. The structure is constructed based on recursively clustering labels by methods such as K-means. For clustering, we need label information. If such information is not directly available, some works (Prabhu et al., 2018; Khandagale et al., 2020; Yu et al., 2022) construct the label representations by averaging all instances of label-j:

$$\frac{\sum_{i} [\boldsymbol{y}_{i}]_{j} \boldsymbol{x}_{i}}{\|\sum_{i} [\boldsymbol{y}_{i}]_{j} \boldsymbol{x}_{i}\|_{2}}.$$
(2)

Next, we discuss a two-level tree as an example. The K-means procedure divides the set of labels  $\{1, \ldots, m\}$  into K partitions  $I_1, \ldots, I_K$ . Then, we can train a smaller OVR model of K weight vectors by solving

$$\min_{\tilde{\boldsymbol{w}}_{\tilde{j}} \in \mathbb{R}^n} \frac{\lambda}{2} \tilde{\boldsymbol{w}}_{\tilde{j}}^T \tilde{\boldsymbol{w}}_{\tilde{j}} + \sum_{i=1}^l \xi(\tilde{\boldsymbol{w}}_{\tilde{j}}^T \boldsymbol{x}_i, [\boldsymbol{z}_i]_{\tilde{j}}), \quad (3)$$

for  $\tilde{j} = 1, ..., K$ . For any  $\boldsymbol{x}$ , the model can determine if  $\boldsymbol{x}$  has pseudo-label- $\tilde{j}$ , where the pseudo-label vector  $\boldsymbol{z}$  is defined by

$$[\boldsymbol{z}]_{\tilde{j}} = \begin{cases} 1 & \boldsymbol{x} \text{ has any label in the partition } I_{\tilde{j}}, \\ 0 & \text{otherwise,} \end{cases}$$

for all  $\tilde{j} = 1, ..., K$ . Prabhu et al. (2018); Yu et al. (2022) point out that the model trained by (3) can estimate the probabilities

$$P(\boldsymbol{x} \text{ has pseudo-label-} j \mid \boldsymbol{x}; \tilde{\boldsymbol{w}}_{\tilde{j}}),$$
 (4)

for  $\tilde{j} = 1, \ldots, K$ , via the transform function

$$\sigma(\tilde{\boldsymbol{w}}_{\tilde{j}}^T \boldsymbol{x}) = \exp\left(-\max(1-\tilde{\boldsymbol{w}}_{\tilde{j}}^T \boldsymbol{x}, 0)^2\right) \quad (5)$$

if  $\xi$  is the square-hinge loss. However, we are interested in

$$P(\boldsymbol{x} \text{ has label-} j \mid \boldsymbol{x}; \boldsymbol{\theta}), \forall j = 1, \dots, m,$$

where  $\theta$  includes all the parameters (i.e.,  $\tilde{w}_{\tilde{j}}, w_j$ ,  $\forall \tilde{j}, j$ ) in the model. Therefore, Prabhu et al. (2018)



Figure 1: A two-level tree-based model.

utilize the property that

$$P(\boldsymbol{x} \text{ has label-} j \mid \boldsymbol{x}; \boldsymbol{\theta})$$
  
=  $P(\boldsymbol{x} \text{ has label-} j, j \in I_{\tilde{j}}$   
 $\mid \boldsymbol{x}, \boldsymbol{x} \text{ has pseudo-label-} \tilde{j}; \boldsymbol{w}_{j}) \cdot (4), \quad (6)$ 

so we can rely on a data subset

$$D_{\tilde{j}} = \{ (\boldsymbol{y}_i, \boldsymbol{x}_i) \mid \boldsymbol{x}_i \text{ has any label in } I_{\tilde{j}} \}$$
(7)

to train another OVR model

$$\min_{\boldsymbol{w}_j \in \mathbb{R}^n} \frac{\lambda}{2} \boldsymbol{w}_j^T \boldsymbol{w}_j + \sum_{(\boldsymbol{y}_i, \boldsymbol{x}_i) \in D_{\tilde{j}}} \xi(\boldsymbol{w}_j^T \boldsymbol{x}_i, [\boldsymbol{y}_i]_j), \quad (8)$$

for estimating

$$P(\boldsymbol{x} \text{ has label-} j, j \in I_{\tilde{j}}$$

$$\mid \boldsymbol{x}, \boldsymbol{x} \text{ has pseudo-label-} \tilde{j}; \boldsymbol{w}_{j}),$$
(9)

for all  $j \in I_{\tilde{j}}$ . The estimation of (9) corresponds to the training process at node  $I_{\tilde{j}}$  in Figure 1. Specifically, at node  $I_{\tilde{j}}$ , we have a smaller multi-label problem with labels in  $I_{\tilde{j}}$  and data points in  $D_{\tilde{j}}$ . We still use the OVR setting for training. Thus, we have K linear models trained by the full data set D in level-1 of the tree and m linear models trained by K smaller data subsets  $D_1, \ldots, D_K$  in level-2 of the tree. Furthermore, the transform function (5) and the property (6) estimate

$$P(\boldsymbol{x} \text{ has label-} j \mid \boldsymbol{x}; \boldsymbol{\theta}) \\\approx \sigma(\tilde{\boldsymbol{w}}_{\tilde{j}}^T \boldsymbol{x}) \cdot \sigma(\boldsymbol{w}_{j}^T \boldsymbol{x}), \forall j \in I_{\tilde{j}},$$
(10)

for all  $\tilde{j} = 1, ..., K$ , for predicting all labels with a given x. Note that for simplicity, we only show a two-level example here for describing the tree-based model.

For the training time of linear models, a popular method (Hsieh et al., 2008) costs  $O(\dot{n}l)$  per iteration, where  $\dot{n}$  is the average number of non-zero feature values per training instance. Because the number of iterations is usually not large, we can treat it as a constant in the time complexity analysis.

Therefore, training the standard OVR that involves m binary problems costs  $O(m\dot{n}l)$ . For a two-level tree-based model, in Appendix C, we derive the training time as

$$O\left(Km\tilde{c}\dot{n}R + (K + \frac{cm}{K})\dot{n}l\right),\qquad(11)$$

in which the first term corresponds to the cost of K-means. In (11),  $\tilde{c}\dot{n}$  with  $\tilde{c} > 1$  is the average number of non-zeros in each label representation,  $c \ge 1$  is a constant upper-bounded by the maximal label number of an instance, and R is the average number of K-means iterations. For data with many instances, we generally have

$$l \gg KR\tilde{c},$$

so the second part in (11) is the dominant term. If we compare it with the  $O(m\dot{n}l)$  cost of OVR, when m is enormous, a tree-based model costs much less.

We have discussed a tree-based model to solve the training time issue (ii). Let us check the model size. To begin, we assume that zero features have been removed before training any binary problem. Due to the use of  $l_2$  regularization (i.e.,  $\boldsymbol{w}_j^T \boldsymbol{w}_j$ ) in (1)), the resulting  $w_i$  is generally a dense vector, including many non-zero components. For the OVR method, we need O(mn) space to store m linear models. On the other hand, the two-level tree-based model has (K + m) linear models, so the model size is O(Kn + mn), which is larger than the OVR model. Fortunately, while we use a training data subset (7) to train the model for estimating (8), many instances may be removed. Thus, some features may not be used by any instance of the training data subset so we can reduce the dimensionality of the feature space. Assume the reduced dimensions are  $n_1, \ldots, n_K$  that correspond to the training subsets  $D_1, \ldots, D_K$ , and  $\bar{n}$ is the average. The two-level tree-based model size becomes  $O(Kn + m\bar{n})$ , less than OVR's model size O(mn) if  $\bar{n}$  is small enough. A detailed study on the size of tree-based models is in (Lin et al., 2024).

# **3** Distributed Settings to Address the Memory Issue

Using multiple computers to store a model can reduce the memory usage in each computer. In the past, DiSMEC distributed the training and the storage of an OVR model to multiple computers. However, the setting is still not scalable to data with

Data set	Max. labels	Ratio of total labels
eur-lex-4k	422	0.11
wiki10-31k	3289	0.11
amazoncat-13k	2854	0.21
amazon-670k	106963	0.16
amazon-3m	352094	0.13
wiki-500k	88769	0.18

Table 1: The maximal label number of K-means partitions over six data sets as K = 100.

extremely many labels. Hence, separating a treebased model into several computers is a possible solution because of the faster training and smaller model size than the OVR setting. Therefore, we discuss the distributed training of a tree-based model in the following subsections.

# 3.1 Distributing a Tree-Based Model

We discussed in Section 2.2 that a tree-based model separates the labels to the partitions  $I_1, \ldots, I_K$ , so distributing the sub-tree construction of these partitions is possible. Assume we have K machines. For the machine- $\tilde{j}$ , to simulate the task at node  $I_{\tilde{j}}$  in Figure 1, we must train a binary classifier on

data with labels in  $I_{\tilde{i}}$  versus without. (12)

To this end, the machine- $\tilde{j}$  must have the whole training set. After that, the machine- $\tilde{j}$  continues to construct the sub-tree. Therefore, the performance will be the same as the serial setting. However, there are two issues:

- (i) Unless we implement distributed *K*-means, the *K*-means procedure must be done on one computer to get the partitions.
- (ii) The partitions are imbalanced. Table 1 shows that the largest partition among the 100 clusters  $I_1, \ldots, I_{100}$  contains 10% or more of the whole data set. Thus, the model sizes of the partitions are hard to estimate, so the specifications of each computer are challenging to decide.

# 3.2 Random Label Forests

One possible solution to overcome these two issues is to omit K-means, uniformly split labels to the partitions  $I_1, \ldots, I_K$  in level-1, and everything else is the same. However, our experiment in Section 4.1 shows that the performance of a tree-based model using K-means partitions is better than using random partitions. The inferior performance seems to be from the poor estimations of level-1 probabilities in (4). We know that each label subset



(a) Without negative sampling.





(c) Negative sampling by labels.

Figure 2: An example of different negative samplings.

(b) Negative sampling by instances.

should contain instances with similar or even identical feature values. For a random split of labels, these similar instances may end up being on both positive and negative sides of the problem (12), an ambiguous situation that may result in a poor model. In contrast, K-means helps to put these labels into the same partition, so the issue may not occur. From the discussion so far, the question becomes how to alleviate the performance issue while controlling the model size per computer. We propose the following settings.

- We let each machine handle a random label subset. This way allows us to control the model size.
- We propose bypassing the level-1 probability estimate. Instead, we run the standard tree-based method on the label subset, applying *K*-means on every level.

For the random label subset in each computer, instead of one partition from splitting the whole label space, we can be general so that label subsets overlap. The remaining task is to let each machine produce a suitable probability estimate and ensemble results from different machines. Because each label subset corresponds to an independent tree, our method is an ensemble method with label subsampling. We call our method "random label forests" due to the similar idea from random forests (Breiman, 2001).

Let us formally discuss random label forests in detail. Suppose we have N computers and use the sample rate r on subsampling the labels  $\{1, \ldots, m\}$  to the subsets  $\hat{I}_t$ , for all  $t = 1, \ldots, N$ . Without loss of generality, let us focus on t = 1 as an example. Since we consider a label subset, the label space is modified from  $\{0, 1\}^m$  to  $\{0, 1\}^{|\hat{I}_1|}$ . By a function  $\phi_1$  to describe this label mapping, our training data subset becomes

$$\hat{D}_1 = \{(\phi_1(\boldsymbol{y}_i), \boldsymbol{x}_i) \mid i = 1, \dots, l\}$$

in the computer-1. We note that the smaller the sample rate r we set, the more instances are empty-labeled in  $\hat{D}_1$ , i.e.,

$$\phi_1(\mathbf{y}) = \mathbf{0} \in \{0, 1\}^{|\hat{I}_1|}$$

Thus, it seems that we have a choice of removing empty-labeled instances.

The training time of the models can be reduced if we remove the empty-labeled instances. However, the probability estimates from (4) may become inaccurate. Let us explain this issue by an example.

- Consider a multi-class problem, a particular case of multi-label problems, with three labels {red circle, green triangle, blue cross} in Figure 2a. If the red circle is the positive label and the others are negative labels, a linear model can be trained as the dark blue line in Figure 2a.
- If we uniformly sample the negative instances in Figure 2b, the linear model may not be affected. However, suppose we sample the negative instances by the labels in Figure 2c. In that case, the linear model can be impacted incorrectly. The reason is that we completely remove the data from some labels.
- Thus, a non-uniform negative sampling can affect (4) because the model does not estimate the probability well anymore.

The example above shows a critical point in a treebased model training on a subset of labels. If we train a two-level tree-based model on the data subset  $\hat{D}_1$ , the training problem in level-1 will be changed from (3) to

$$\min_{\tilde{\boldsymbol{w}}_{\tilde{j}} \in \mathbb{R}^n} \frac{\lambda}{2} \tilde{\boldsymbol{w}}_{\tilde{j}}^T \tilde{\boldsymbol{w}}_{\tilde{j}} + \sum_{i=1}^l \hat{\xi}(\tilde{\boldsymbol{w}}_{\tilde{j}}^T \boldsymbol{x}_i, [\hat{\boldsymbol{z}}_i]_{\tilde{j}}), \quad (13)$$

for  $\tilde{j} = 1, ..., K$ . In contrast to z in (3), a pseudolabel vector  $\hat{z}$  of the data subset  $\hat{D}_1$  is defined as

$$[\hat{z}]_{\tilde{j}} = \begin{cases} 1 & x \text{ has any label in } I_{\tilde{j}} \text{ of } \hat{D}_1, \\ 0 & \text{otherwise.} \end{cases}$$

#### Algorithm 1 Training random label forests.

Require: '	Training set $D$ , # of submodels $N$ , sample rate $r$
distribu	ted for $t = 1, \ldots, N$ do
$\hat{I}_t \leftarrow$	- subsample the label indices with the rate $r$ on the
full inde	$x \text{ set } \{1,\ldots,m\}.$
$\hat{D}_t \leftrightarrow$	– update the label part of D with $\hat{I}_t$ .
$\boldsymbol{\theta}_t \leftarrow$	- train a tree-based submodel with the subset $\hat{D}_t$ .
end dist	ributed for

Many  $\hat{z}_i = 0$  if  $x_i$ 's labels do not appear in the set  $\hat{I}_1$ . If these  $x_i$  are removed, it is similar to doing a non-uniform negative sampling, as in Figure 2c. The model trained by (13) may not estimate the probability (4) reasonably. Therefore, we should not remove those empty-labeled instances in  $\hat{D}_1$ . Besides this crucial point, all other details in the tree construction are the same as those shown in Section 2.2 for the tree using all labels.

With the label subsets  $\hat{I}_1, \ldots, \hat{I}_N$  and the training subsets  $\hat{D}_1, \ldots, \hat{D}_N$ , we can parallelly train the submodels  $\theta_1, \ldots, \theta_N$  in N computers. Algorithm 1 shows the whole training procedure.

Next, let us discuss the prediction procedure. We still assume that two-level trees are used. Since the prediction probabilities can be estimated by (10), we can estimate  $P(x \text{ has label-}j \mid x; \theta_t)$  for all label-j in the subsampled subset  $\hat{I}_t$  by the tth tree-based submodel  $\theta_t$ . Because label-j may appear in several label subsets, a natural setting is to average the several probability estimations.

$$= \frac{P(\boldsymbol{x} \text{ has label-} j \mid \boldsymbol{x}; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N)}{\left| \{t \mid j \in \hat{I}_t, \forall t = 1, \dots, N\} \right|}.$$
(14)

Nevertheless, an issue exists because a tree-based submodel  $\theta_t$  can only predict labels in the subsampled subset  $\hat{I}_t$ . Hence, if some labels are never subsampled, our model can never predict those labels. This situation occurs for some rare labels, so the performance may not be affected much. Our competitive performance, as shown in Section 4.2, seems to support this point. However, other ways of label subsampling can be a future study to mitigate the issue.

### 3.3 The Benefits of Random Label Forests

Section 2.2 discusses the time complexity and model size of a two-level tree-based model. Now, let us check the complexities in a computer when applying random label forests with two-level treebased submodels in the distributed system.

- Space complexity. Since we set the sample rate as r, the number of a subsampled label set becomes rm. Therefore, the model size changes from  $O(Kn + m\bar{n})$  to  $O(Kn + rm\bar{n})$ .
- Time complexity. Similarly, the time complexity changes from (11) to

$$O\left(Krm\tilde{c}\dot{n}R + (K + \frac{crm}{K})\dot{n}l\right)$$

Therefore, we can roughly control the model size and training time by the rate r in random label forests if each computer handles a tree-based submodel. The experiments for comparing a treebased model with all labels and a tree-based submodel of random label forests are discussed in Section 4.3 and Section 4.4.

# 4 Experiments

Section 3.2 discusses different distributed settings on tree-based models, so we first analyze which setting is the better choice over four smaller XML data sets in Section 4.1. After deciding on the distributed setting, we compare three linear models: OVR, a tree-based model with all labels, and random label forests on six XML data sets in Section 4.2.

Throughout this section, the label number m is shown after the name in each data set. We leave the details of data sets in Appendix A and discuss the hyper-parameters in Appendix B. To measure the performance of an XML model, we follow the works (Prabhu et al., 2018; Khandagale et al., 2020; Yu et al., 2022; You et al., 2019; Jiang et al., 2021) to use precision at 1, 3, and 5 as the metrics of the predictions. To mitigate randomization issues in tree-based methods due to K-means or the random label selection, we execute training and prediction procedures ten times on the data sets "eur-lex-4k," "wiki10-31k," "amazoncat-13k," and "amazon-670k." We only execute the procedures once for the data sets "amazon-3m" and "wiki-500k" because the training time is too long.

We discuss the model size in Section 4.3, and the comparison of the training time is in Section 4.4.

#### 4.1 Analysis of Different Distributed Settings

Section 3.2 discusses that using a uniform split of labels is a possible solution for avoiding the imbalanced K-means clusters under the distributed setting. Now, we check the performance between these two partition methods. They follow the standard setting of tree-based models but intend to con-

	Random partitions	Random label forests
level-1	data with $I_{\tilde{i}}$ versus without	nothing
level-2	data with $I_{\hat{j}}$ versus those with $I_{\tilde{j}}$ but not $I_{\hat{j}}$	data with $I_{\hat{j}}$ versus without

Table 2: Comparison between random partitions and random label forests.

tinue the tree construction after level-1 in parallel. On the other hand, we have the proposed random label forests, which independently generate tree-based submodels on label subsets. Thus, we compare these three methods with the following settings.

- Tree with all labels. The standard tree-based method. Note that we set K = 100.
- **Random 100 partitions**. The tree-based method that replaces *K*-means in level-1 with a uniform split of labels.
- Random 10 partitions and Random 10 partitions×10. In the previous setting, we consider 100 random partitions because of following the K = 100 in K-means. While K-means may require a careful selection of K, when using random splits, we can instead control each cluster's size according to the capacity of a machine. Thus, we try 10 partitions of the data set. Another reason for doing so is that later, for the proposed random label forests, we let each machine handle 10% of labels. Thus, we need the setting of 10 random partitions as a comparison.

In Section 3.2, we discussed how to ensemble the prediction results of various trees. Therefore, an extension is to run the label partition several times to generate more trees. Here, we run the "Random 10 partitions" 10 times to generate 100 label subsets. We call this setting "Random 10 partitions  $\times 10$ ."

• Random label forests-10P and Random label forests-10P×10. We consider label subsets generated in the previous setting to run the proposed label forests. Note that the current setting has some subtle differences from the previous one. To illustrate the differences, we extend Figure 1 to a tree with more than two levels.



For both the current and the previous settings,  $I_1, \ldots, I_K$  with K = 10 correspond to the 10 subsets from a random partition. However, they differ in binary problems solved at each level, as illustrated in Table 2. Note that "Random 10 partitions" is the standard setting for training a tree model, though we run the procedure related to  $I_{\tilde{i}}$  in one particular computer. In the prediction of "Random 10 partitions," we use the model obtained in level-1 to estimate the probability in (4). In contrast, for "Random label forests-10P," we can say that there is no level-1. Using  $I_{\tilde{i}}$  and all data, we construct an independent tree. It is important to note that, as explained in Section 3.2, for the binary problem involving those with  $I_{\hat{i}}$  as positive, we need all other data as negative. We cannot just consider those in  $I_{\tilde{j}}$  but not  $I_{\hat{j}}$ ; see the row "level-2" in Table 2. In prediction, each tree obtains its own probability estimation, and we calculate the average in (14).

Random label forests-100U. In the previous settings, I<sub>j</sub> is one partition of a random split because our generation of label sets extends from the K-means procedure of splitting all labels. Instead, we check the possibility of uniform sampling on {1,...,m} to have label subsets. Here, we randomly draw labels with replacements to obtain 100 label subsets, each of which has 10% of the labels. Thus, these subsets overlap with each other.

Ransom label forests are applicable to use these overlapped label subsets because, first, the  $\tilde{j}$ th tree gives

 $P(\boldsymbol{x} \text{ has label-} j \mid \boldsymbol{x}, \boldsymbol{x} \text{ has pseudo-label-} \tilde{j}; \boldsymbol{\theta}_{\tilde{j}}),$ 

for all  $j \in I_{\tilde{j}}$ , and second, we can apply (14) to get the ensembled prediction. We call this setting "Random label forests-100U." It differs from "Random label forests  $10P \times 10$ " only in how label subsets are generated.

We present comparisons in Table 3 and have the following observations:

• A comparison between "tree with all labels" and "Random 100 partitions" shows that using *K*means for the label partitions is better than using

Method	P@1	P@3	P@5	P@3	P@5	P@1	P@3	P@5	P@1	P@3	P@5
	e	ur-lex-4	k	wiki1	0-31k	ama	azoncat-	13k	am	azon-67	0k
Tree with all labels	82.29	69.35	57.91	74.72	65.86	93.19	79.55	64.61	44.58	39.44	35.64
Random 100 partitions	80.96	67.70	56.40	73.17	63.87	92.53	77.84	62.71	38.39	32.47	27.86
Random 10 partitions	80.23	66.88	55.94	73.41	64.52	92.47	78.01	63.14	41.74	36.46	32.34
Random label forests-10P	82.21	69.08	57.68	73.56	65.01	93.31	79.22	64.18	42.79	37.69	33.97
Random 10 partitions×10	82.23	69.06	58.05	74.17	65.31	94.05	79.95	65.09	44.84	39.89	36.25
Random label forests-10P×10	83.17	69.96	58.73	74.36	65.75	94.21	80.18	65.19	45.23	40.27	36.69
Random label forests-100U	83.08	69.92	58.74	74.35	65.68	94.19	80.17	65.19	45.19	40.23	36.64

Table 3: Comparison of different distributed settings in precisions. For wiki10-31k, since half of the instances are associated with a unique label, precision at 1 is not a discriminable metric on this data set. Therefore, we do not show the precision at 1 results in wiki10-31k for the space limitation.

random partitions. We explained this observation at the beginning of Section 3.2.

- The ensemble method enhances the performance. This result can be seen by comparing
  - "Random 10 partitions" and "Random 10 partitions × 10," and
  - "Random label forests-10P" and "Random label forests-10P×10."
- In Table 2, we illustrated an important difference between random partitions and the proposed random label forests. We organize Table 3 in a way so that, except for the first two rows, the following parts serve as a comparison between the two settings. We can see that random label forests are better than random partitions, no matter whether the ensemble method is considered. Results confirm that random label forests mitigate the issue mentioned in Section 3.2 about the poor estimations of level-1 probabilities in (4).
- Regarding the generation of label subsets, for the proposed random label forests, we should compare "Random label forests-10P×10" and "Random label forests-100U," where the former considers partitions of the label space and the latter applies uniform sampling. Results seem to indicate that the impact of the generation of label subsets is minor.

Experiments in this section fully show that distributing a tree-based method is not trivial because of different considerations and options. Although both the settings "Random label forests- $10P \times 10$ " and "Random label forests-100U" perform well in Table 2, generating uniform sampled label subsets is easier in pratice. Therefore, we use "Random label forests-100U" as our distributed solution in subsequent experiments.

# 4.2 Performance Comparison with Existing Multi-Label Methods

We compare the following methods in the rest of this section.

- OVR: The standard one-vs-rest method. This method, though not scalable even in a distributed setting, gives the performance on small data sets to serve as the target to be achieved by tree-based methods.
- Tree with all labels: This standard tree-based method was evaluated in Section 4.1. Though we discussed its possible parallelization in Section 4.1, here we assume that it is run on a single computer for checking the effectiveness of our proposed distributed setting.
- Random label forests: This is the "Random label forests-100U" setting in Section 4.1.

Table 4 shows the comparison results, and we have the following observations.

- OVR is slightly better than tree-based methods. This result is reasonable because tree-based methods are a kind of "hierarchical approximation" of OVR to address the scalability issue.
- The proposed random label forests are consistently better than the tree-based method with all labels except for "wiki-500k." The performance is close to that of OVR.

The worse results on "wiki-500k" are an example of our method's limitations, and we discuss this issue in Section 6.1.

### 4.3 Comparison on the Model Size

In Section 3.3, we have analyzed the model size of random label forests that use two-level tree-based submodels. However, because

- the sparsity of data points can affect  $\bar{n}$ , and
- a tree-based model may have more than two layers,

we conduct experiments to check the model size in practice. Table 5 compares

Method	P@1	P@3	P@5	P@1	P@3	P@5
		eur-lex-4k		amazon-670k		
One-versus-rest	83.47	70.62	59.05	45.41	40.41	36.97
Tree with all labels	$82.29\pm0.30$	$69.35\pm0.09$	$57.91 \pm 0.12$	$44.58\pm0.07$	$39.44\pm0.04$	$35.64\pm0.03$
Random label forests	$83.08 \pm 0.19$	$69.92\pm0.06$	$58.74 \pm 0.07$	$45.19\pm0.04$	$40.23\pm0.01$	$36.64\pm0.02$
		wiki10-31k			amazon-3m	
One-versus-rest	85.23	75.80	67.11	-	-	-
Tree with all labels	$84.72\pm0.08$	$74.72\pm0.17$	$65.86 \pm 0.09$	47.48	44.74	42.63
Random label forests	$84.80\pm0.11$	$74.35\pm0.12$	$65.68 \pm 0.06$	48.69	45.67	43.49
		amazoncat-13k			wiki-500k	
One-versus-rest	94.14	79.71	64.69	-	-	-
Tree with all labels	$93.19\pm0.03$	$79.55\pm0.04$	$64.61\pm0.03$	68.39	48.90	38.00
Random label forests	$94.19\pm0.02$	$80.17\pm0.02$	$65.19\pm0.02$	64.37	45.83	36.09

Table 4: Comparison of random label forests and other linear methods in precision at 1, 3, and 5 over six data sets. OVR results on "amazon-3m" and "wiki-500k" are unavailable due to lengthy running time.

Data set	Tree with all labels	Random label forests
eur-lex-4k	561.02 MB	49.25 MB
amazoncat-13k	1.72 GB	193.13 MB
wiki10-31k	6.38 GB	749.07 MB
amazon-670k	20.56 GB	2.92 GB
amazon-3m	135.89 GB	12.69 GB
wiki-500k	161.37 GB	13.65 GB

Table 5: The model size comparison between a tree with all labels and a tree of random label forests.

Data set	Tree with all labels	Random label forests
eur-lex-4k	224.42 s	32.18 s
wiki10-31k	4220.77 s	543.33 s
amazoncat-13k	5311.23 s	910.06 s
amazon-670k	32068.89 s	2120.25 s
amazon-3m	503575.79 s	23023.76 s
wiki-500k	202261.66 s	22987.20 s

Table 6: Training time comparison between a tree with all labels and a tree of random label forests.

- the model size of a tree-based method with all labels in a single computer, and
- the model size of a tree in random label forests, which corresponds to the needed space in each computer of the distributed environment.

We can see that the ratio of space reduction is close to the sample rate r = 0.1. Hence, random label forests can reduce the model size in each computer of the distributed system, even though the whole model may be larger than that of a tree-based model with all labels. Therefore, random label forests are effective in addressing the memory difficulty of extreme multi-label classification.

### 4.4 Comparison on the Training Time

Table 6 shows the training time comparison between a tree-based model with all labels and a tree of random label forests. If, for the random label forests,

• each machine handles one tree-based submodel,

and

• all machines used have similar configurations,

then Table 6 gives the comparison of total running time between the standard tree-based method run in one computer and the proposed random label forests run in a distributed environment. We observe significant time reduction by random label forests, especially for problems with a large number of labels.

# 5 Conclusion

This work proposes random label forests, a distributed ensemble method with label subsampling, and tree-based linear models as the backbone. Random label forests give competitive performances using much less training time and memory usage in a machine. Hence, handling a problem with extremely many labels becomes practical. For all methods considered and evaluated in this work, the backbone is a linear classifier. In the future, we plan to consider more sophisticated techniques, such as neural networks.

# **6** Limitations

### 6.1 Breaking Label Relationships

In Section 4.2, the performance of random label forests is worse than that of a tree-based method with all labels over the "wiki-500k" data set. After the investigation, we think the correlation between the labels in "wiki-500k" may be higher than other data sets because the labels are the tags of the documents in Wikipedia. For example, the labels of the

	Tree with all labels		Random labe	forests
Labels	Prediction	rank	Prediction	rank
Andorra	0.0000	-	0.0001	-
Constitutional monarchies	0.0011	13	0.0001	-
Countries in Europe	0.0021	4	0.0003	-
Diarchies	0.0000	-	0.0000	-
Landlocked countries	0.0000	-	0.0001	-
Liberal democracies	0.0018	6	0.0005	-
Member states of La Francophonie	0.0011	15	0.0001	-
Member states of the United Nations	0.0040	1	0.0022	32
Catalans Countries	0.0000	-	0.0001	-
Principalities	0.0000	-	0.0000	-
Pyrenees	0.0000	-	0.0010	45
Romance countries and territories	0.0000	-	0.0001	-
States and territories established in 1278	0.0000	-	0.0000	-

Table 7: The comparison of the 16th instance's prediction between tree with all labels and random label forests. We only give the rank of labels in the top-50 predictions, so the notation '-' indicates that labels are out of the top-50 predictions.

16th instance in the test data of "wiki-500k" are

'Andorra,' 'Constitutional monarchies,' 'Countries in Europe,' 'Diarchies,' 'Landlocked countries,' 'Liberal democracies,' 'Member states of La Francophonie,' 'Member states of the United Nations,' 'Catalans Countries,' 'Principalities,' 'Pyrenees,' 'Romance countries and territories,' 'States and territories established in 1278.'

We can see labels are strongly related. For example, "Andorra" is one of the "Countries in Europe," one of the "Landlocked countries," and one of the "Catalans Countries."

Table 7 shows the ranking of these labels in the top-50 predictions. The ranking result of random label forests is less accurate than that of tree with all labels. Based on the above result, we think that uniform sampling in labels breaks the relations so that the performance may be hurt. If we increase the sample rate from 0.1 to 0.15, the performance of "wiki-500k" will become better to

P@1	P@3	P@5
65.45	46.84	36.92

However, that performance is still much lower than the tree-based model with all labels. Therefore, the label subsampling technique seems unsuitable in the data sets that contain high correlations between labels.

# References

Rahul Agrawal, Archit Gupta, Yashoteja Prabhu, and Manik Varma. 2013. Multi-label learning with millions of labels: Recommending advertiser bid phrases for web pages. In *Proceedings of the 22nd International Conference on World Wide Web (WWW)*, pages 13–24.

- Rohit Babbar and Bernhard Schölkopf. 2017. DiSMEC: Distributed sparse machines for extreme multi-label classification. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining* (WSDM), pages 721–729.
- Kush Bhatia, Kunal Dahiya, Himanshu Jain, Purushottam Kar, Anshul Mittal, Yashoteja Prabhu, and Manik Varma. 2016. The extreme classification repository: Multi-label datasets and code.
- Kush Bhatia, Himanshu Jain, Purushottam Kar, Manik Varma, and Prateek Jain. 2015. Sparse local embeddings for extreme multi-label classification. In *Advances in Neural Information Processing Systems* 28, pages 730–738.
- Leo Breiman. 2001. Random forests. *Machine Learn-ing*, 45(1):5–32.
- Philip A. Etter, Kai Zhong, Hsiang-Fu Yu, Lexing Ying, and Inderjit Dhillon. 2022. Enterprise-scale search: Accelerating inference for sparse extreme multi-label ranking trees. In *Proceedings of the ACM Web Conference*, pages 452–461.
- Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. 2008. LIBLINEAR: a library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874.
- Cho-Jui Hsieh, Kai-Wei Chang, Chih-Jen Lin, S. Sathiya Keerthi, and Sellamanickam Sundararajan. 2008. A dual coordinate descent method for largescale linear SVM. In *Proceedings of the Twenty Fifth International Conference on Machine Learning* (*ICML*).
- Ting Jiang, Deqing Wang, Leilei Sun, Huayi Yang, Zhengyang Zhao, and Fuzhen Zhuang. 2021. LightXML: Transformer with dynamic negative sampling for high-performance extreme multi-label text classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 7987– 7994.

- Sujay Khandagale, Han Xiao, and Rohit Babbar. 2020. Bonsai: diverse and shallow trees for extreme multilabel classification. *Machine Learning*, 109:2099– 2119.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751.
- He-Zhe Lin, Cheng-Hung Liu, and Chih-Jen Lin. 2024. Exploring space efficiency in a tree-based linear model for extreme multi-label classification. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP).*
- Yu-Chen Lin, Si-An Chen, Jie-Jyun Liu, and Chih-Jen Lin. 2023. Linear classifier: An often-forgotten baseline for text classification. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL), pages 1876–1888. Short paper.
- Jingzhou Liu, Wei-Cheng Chang, Yuexin Wu, and Yiming Yang. 2017. Deep learning for extreme multilabel text classification. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 115–124.
- Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of the 7th* ACM Conference on Recommender Systems, pages 165–172.
- Siddhartha Nuthakki, Sunil Neela, Judy W. Gichoya, and Saptarshi Purkayastha. 2019. Natural language processing of MIMIC-III clinical notes for identifying diagnosis and procedures with neural networks. *Preprint*, arXiv:1912.12397.
- Yashoteja Prabhu, Anil Kag, Shrutendra Harsola, Rahul Agrawal, and Manik Varma. 2018. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In *Proceedings of the 2018 World Wide Web Conference* (WWW), pages 993–1002.
- Dan Shen, Jean David Ruvini, Manas Somaiya, and Neel Sundaresan. 2011. Item categorization in the e-commerce domain. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management (CIKM)*, pages 1921–1924.
- Grigorios Tsoumakas, Ioannis Katakis, and Ioannis Vlahavas. Effective and efficient multilabel classification in domains with large number of labels. In *Proceedings of ECML/PKDD 2008 Workshop on Mining Multidimensional Data*.
- Grigorios Tsoumakas and Ioannis Vlahavas. 2007. Random k-labelsets: An ensemble method for multilabel classification. In *European Conference on Machine Learning*, pages 406–417.

- Ronghui You, Zihan Zhang, Ziye Wang, Suyang Dai, Hiroshi Mamitsuka, and Shanfeng Zhu. 2019. AttentionXML: Label tree-based attention-aware deep model for high-performance extreme multi-label text classification. In *Advances in Neural Information Processing Systems*, volume 32.
- Hsiang-Fu Yu, Prateek Jain, Purushottam Kar, and Inderjit S. Dhillon. 2014. Large-scale multi-label learning with missing labels. In *Proceedings of the Thirty First International Conference on Machine Learning* (*ICML*), pages 593–601.
- Hsiang-Fu Yu, Kai Zhong, Jiong Zhang, Wei-Cheng Chang, and Inderjit S. Dhillon. 2022. PECOS: Prediction for enormous and correlated output spaces. *Journal of Machine Learning Research*, 23(98):1–32.
- Jiong Zhang, Wei-Cheng Chang, Hsiang-Fu Yu, and Inderjit Dhillon. 2021. Fast multi-resolution transformer fine-tuning for extreme multi-label text classification. In *Advances in Neural Information Processing Systems*, volume 34, pages 7267–7280.
- Arkaitz Zubiaga. 2009. Enhancing navigation on Wikipedia with social tags. In Proceedings of Wikimania.

	l	n	m
Data Set	instances	features	labels
eur-lex-4k	15,449	186,104	3,956
wiki10-31k	14,146	104,374	30,938
amazoncat-13k	1,186,239	203,882	13,330
amazon-670k	490,449	135,909	670,091
amazon-3m	1,717,899	337,067	2,812,281
wiki-500k	1,779,881	2,381,304	501,070

Table 8: Statistics of data sets

# A Data Sets

We show the statistics of data sets in Table 8. We list the specific link of each set as follows.

- eur-lex-4k:
  - Training: https://www.csie.ntu.edu. tw/~cjlin/libsvmtools/datasets/ multilabel/eurlex\_tfidf\_train.svm. bz2
  - Test: https://www.csie.ntu.edu. tw/~cjlin/libsvmtools/datasets/ multilabel/eurlex\_tfidf\_test.svm.bz2
- wiki10-31k:
  - Training: Training:https://www.csie. ntu.edu.tw/~cjlin/libsvmtools/ datasets/multilabel/wiki10\_31k\_ tfidf\_train.svm.bz2
  - Test: Testing:https://www.csie.ntu. edu.tw/~cjlin/libsvmtools/datasets/ multilabel/wiki10\_31k\_tfidf\_test.svm. bz2
- amazoncat-13k:
  - Training: https://www.csie.ntu.edu. tw/~cjlin/libsvmtools/datasets/ multilabel/AmazonCat-13K\_tfidf\_ train\_ver1.svm.bz2
  - Test: https://www.csie.ntu.edu. tw/~cjlin/libsvmtools/datasets/ multilabel/AmazonCat-13K\_tfidf\_test\_ ver1.svm.bz2
- amazon-670k, wiki-500k and amazon-3m are downloaded from the link
  - https://archive.org/download/ pecos-dataset/inference-models/

which is supported by the following GitHub repository provided in Etter et al. (2022).

- https://github.com/ UniqueUpToPermutation/pecos/tree/ benchmark

Note that those data sets have already been preprocessed from documents to a popular sparse feature representation, "TF-IDF." Moreover, every data set has further been split into training and test subsets.

# **B** Details of Experimental Settings

In our experiments, we utilize the library LibMulti-Label<sup>1</sup> to handle all of the model training and evaluation. Moreover, the linear classifiers are trained by LIBLINEAR (Fan et al., 2008) in LibMultiLabel. In using LIBLINEAR, we consider the settings from the works (Khandagale et al., 2020; Yu et al., 2022):

- using squared hinge loss (L2-SVM) and
- taking  $\lambda = 1$  in the training problems.

In the implementations of (Khandagale et al., 2020; Yu et al., 2022), they decided to stop the training process early in each linear model training. However, we chose to spend more time in the model training to get a tight solution that is closer to the optimal solution of the training problem. The treebased methods in LibMultiLabel follow the implementation<sup>2</sup> in Khandagale et al. (2020), and we follow past works to consider the following commonly used hyper-parameters.

• K = 100 for K-means and the max depth of the tree is 10.

# C Time Complexity of Two-Level Tree-Based Models

Because the level-1 of a two-level tree-based model is a smaller OVR model, the level-1 only costs  $O(K\dot{n}l)$ . For the level-2, since the instances can belong to several label partitions, we assume that the average number of instances in each of the subsets  $D_1, \ldots, D_K$  is

$$\frac{cl}{K},$$

where  $c \ge 1$  is a small positive number. Moreover, c is bounded by the maximal label number of an instance<sup>3</sup>. Hence, the training cost of a subset in level-2 is about

$$O\left(\frac{\dot{n}cl}{K}\right)$$

and the two-level tree-based model then costs

$$O\left((K+\frac{cm}{K})\dot{n}l\right),$$

<sup>1</sup>https://www.csie.ntu.edu.tw/~cjlin/ libmultilabel

<sup>2</sup>The work (Khandagale et al., 2020) ensembles the predictions of three tree-based models in their experiments, but LibMultiLabel only considers the single tree-based model.

<sup>3</sup>Prabhu et al. (2018) assume  $O(\log(m))$  is the maximal label number of an instance, so *c* is bounded by  $O(\log(m))$ .

which is different from the complexity of OVR  $O(m\dot{n}l)$ . Besides the training time of linear models, we must check the cost of running K-means. The process involves several iterations, in each of which we calculate the distance between each label representation (2) and K centers of the current clusters. If the label representations are still sparse and the average number of non-zeros is  $\ddot{c}\dot{n}$ , where  $\ddot{c} > 1$  is a positive constant, checking the distance requires

$$O(\tilde{c}\dot{n}).$$

Thereby, one iteration of K-means requires

$$O(Km\tilde{c}\dot{n}).$$

If we set a constant R as the maximum iteration, the time complexity of K-means is

$$O(Km\tilde{c}\dot{n}R).$$

Hence, the total training time of a two-level treebased model is

$$O\left(Km\tilde{c}\dot{n}R + (K + \frac{cm}{K})\dot{n}l\right).$$