

Label-aware Document Representation via Hybrid Attention for Extreme Multi-Label Text Classification

Xin Huang, Boli Chen, Lin Xiao, and Liping Jing

Beijing Key Lab of Traffic Data Analysis and Mining, Beijing Jiaotong University,
Beijing 100044, China
{18120367, 18120345, 17112079, lpjing}@bjtu.edu.cn

Abstract. Extreme multi-label text classification (XMTC) aims at tagging a document with most relevant labels from an extremely large-scale label set. It is a challenging problem especially for the tail labels because there are only few training documents to build classifier. This paper is motivated to better explore the semantic relationship between each document and extreme labels by taking advantage of both document content and label correlation. Our objective is to establish an explicit label-aware representation for each document with a hybrid attention deep neural network model (*LAHA*). *LAHA* consists of three parts. The first part adopts a multi-label self-attention mechanism to detect the contribution of each word to labels. The second part exploits the label structure and document content to determine the semantic connection between words and labels in a same latent space. An adaptive fusion strategy is designed in the third part to obtain the final label-aware document representation so that the essence of previous two parts can be sufficiently integrated. Extensive experiments have been conducted on five benchmark datasets by comparing with the state-of-the-art methods. The results show the superiority of our proposed *LAHA* method, especially on the tail labels.

Keywords: Extreme Multi-label Text Classification · Deep Neural Network · Attention · Tail Label · Label-aware Document Representation.

1 Introduction

Extreme multi-label text classification (XMTC) aims at automatically tagging a document with most relevant labels from an extremely large label set. For instance, there are millions of categories on Wikipedia and one might wish to build a classifier that can annotate a given message with the subset of most relevant categories [8]. XMTC has become increasingly important due to the boom of big data, while it becomes significantly challenging because it has to simultaneously handle massive documents, features and labels. Thus it is emergency to develop effective extreme multi-label classifier for various real-applications such as product categorization in e-commerce, news annotation and etc.

Multi-label text classification, unlike the traditional multi-class classification, allows for the co-existence of more than one labels for a single document. Meanwhile, there may be a large number of 'tail labels' with very few positive documents in XMTC tasks. To tackle the aforementioned issues, researchers pay much attention on two facets: 1) how to represent label so that the correlation among labels can be accurately mined, and 2) how to represent document so that the dependency among text can be sufficiently captured. Recently, state-of-the-art extreme multi-label learning methods have been proposed in each facet. Among them, tree-based and embedding-based methods become popular to find the label correlation as they can obtain notable accuracy improvement by constructing a hierarchy structure [17] or learning a low-dimensional latent space [8]. Deep learning-based methods (e.g., convolutional neural network [5]) have achieved great success to represent text data. These methods usually characterize one document with the same representation on all labels. In this case, the probability of document belonging to a class is determined by their overall matching score regardless of the label-aware semantic information. Recent works, *AttentionXML* [6] and *EXAM* [7], turn attention to this issue with the aid of attentive neural network. However, they only focus on document or label content but ignoring the label structure among extreme labels which has been proved very important in extreme multi-label learning [8].

To solve the above-mentioned problems, we introduce a **Label-Aware** document representation model via a **Hybrid Attention** neural network (*LAHA*) by considering both document content and label structure. *LAHA* consists of three parts. The first part aims at detecting the importance of each word to all labels via a self-attention bidirectional LSTM neural network. The second part tries to explore the semantic connection between words and labels in a latent space. Here the word embedding is obtained by the bidirectional LSTM neural network. The label embedding is determined from the label co-exist graph so that the label structure can be sufficiently maintained in the same latent space with words'. Based on these two embeddings, we introduce an interaction-attention mechanism to explicitly compute the semantic relation between the words and labels. The last part is to represent each document along each label via an adaptive fusion strategy. The goal of fusion strategy is to adaptively extract proper information from the previous two parts so that the final document representation has discriminative ability to construct classifier.

The proposed XMTC model *LAHA* has been evaluated on five benchmark datasets and get competitive results, we summarize the major contributions.

- *LAHA* is the first work to construct label-aware document representation by simultaneously considering document content and label structure.
- The hybrid attention mechanism is firstly designed to adaptively extract the semantic relation between each document and all labels for XMTC.
- The performance of *LAHA* was thoroughly investigated on widely-used benchmark datasets, indicating the advantage over the baselines.

- The code and hyper-parameter settings are released¹ to facilitate other researchers.

2 Related Work

Significant progress has been made for XMTC. They can be roughly categorized into two categories: embedding-based and tree-based methods. Recently, due to the powerful ability of representation, deep learning technology has been introduced to effectively represent document for XMTC tasks. Next, we will briefly review them.

2.1 Embedding-based Methods

Embedding-based methods aim at reducing the huge label space to a low dimensional space while preserving the label correlation as much as possible, and then compressed label embedding are decompressed for prediction. Various approaches have been presented such as compressed sensing [18], output codes [19], Singular Value Decomposition [20], landmark labels [21], Bloom filters [22], etc. To efficiently handle large-scale label set, these embedding-based methods usually assume that the label matrix is low-rank. However, such methods have been proved unable to deliver high prediction accuracies as the low rank assumption is violated in most real world applications [8]. *SLEEC* [8] can be taken as the most representative embedding-based method due to its significant accuracy and computationally efficiency. Its main idea is to learn a small ensemble of local distance preserving embeddings. Specifically, *SLEEC* divides the training data set into several clusters, and in each cluster it detects embedding vectors by capturing non-linear label correlation and preserving the pairwise distance between labels. The k -nearest neighbors search is used to do prediction only in the cluster into which the test document is fallen. Later, Zhang et al. [14] adopted deep neural network for non-linear modeling the label embedding. Although these methods perform well, they play a heavy price in terms of prediction accuracy due to the loss of information during the compression and decompression phases.

2.2 Tree-based Methods

Tree-based methods introduce a tree structure to divide the documents recursively at each non-leaf node, so that documents in each leaf node share similar label distribution. The most representative method *FastXML* [17] implements this process by optimizing the normalized discounted cumulative gain (nDCG)-based ranking loss function. Then, a base classifier is trained at each leaf node which only focuses on a few active labels. To enhance the robustness of predictions, an ensemble of multiple induced trees are learned. The main advantage of tree-based methods is that the prediction time complexity is typically sublinear in the training-set size and would be logarithmic if the induced tree is balanced.

¹ https://github.com/HX-idiot/Hybrid_Attention_XML

A recent extension work of *FastXML* is *PfastreXML* [9], which adopted a propensity scored objective function instead of nDCG-based loss which is more friendly to tail labels. *Parabel* [10] is another tree-based method, which constructs balanced trees partitioning labels rather than instances. These tree-based methods represent document via bag-of-words, where the words are treated as independent features, which will ignore the semantic dependency among words.

2.3 Deep Learning-based Methods

To capture semantic dependency among words, researchers adopted deep learning models in text classification task due to its strong ability of representation. The popular deep models include CNN [26], GRU [25], RNN [12], LSTM [13], Bi-LSTM [27], BERT [29] and several combination networks [23,28]. Even though they have achieved great success in traditional NLP tasks, few work is designed for XMTC.

XML-CNN [5] can be taken as the first and most representative work using deep learning model in XMTC. It takes advantage of CNN, dynamic max-pooling and bottle-neck layer to build the deep model. Due to the limited window size, *XML-CNN* can not capture the long-distance dependency among text. Later, GRU and Bi-LSTM language models are adopted in *AttentionXML* [6] and *EXAM* [7] to effectively represent document for XMTC. Meanwhile, these two methods consider the difference of one document representation along different labels. *AttentionXML* [6] adopts self-attention mechanism [15], while *EXAM* [7] exploits the label content information to calculate the relations between words and classes. Although *AttentionXML* obtains promising performance, it ignores the label structure which has been proved very important in embedding-based and tree-based multi-label learning methods.

Therefore, in this paper, we propose a new XMTC deep model with hybrid attention to build label-aware document representation, which sufficiently exploits both document content and label structure.

3 LAHA model

In this section, we introduce the proposed deep model (*LAHA*) to handle XMTC tasks. The overall structure of *LAHA* is shown in Fig. 1). Our goal is to build a multi-label learning model from the training documents with a large-size label set. Let $D = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ be the given raw training document set containing total N documents and belonging to k labels. Each document has n tokens (or words) and each word is represented via a d -dimensional deep semantic dense vector acquired from word2vec technique, $\mathbf{e}_t \in \mathbb{R}^d$ ($t = 1, \dots, n$). $\mathbf{y}_i \subseteq \{0, 1\}^k$ is the corresponding label vector, and $y_{ij} = 1$ iff the j -th label is turned on for the i -th document $\mathbf{x}_i = (\mathbf{e}_1, \dots, \mathbf{e}_n)$.

3.1 Feature Embedding

To build the proposed *LAHA* multi-label text classifier, the raw text data is preprocessed via word embedding technique so that each word is represented

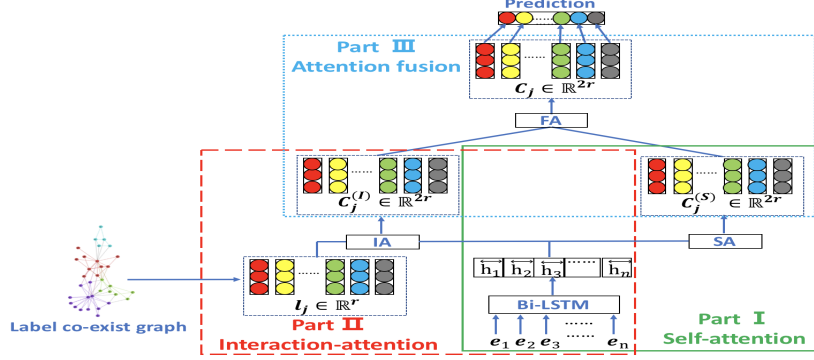


Fig. 1. The architecture of **LAHA**. The solid green box indicates the self-attention process, the dashed red box represents interaction-attention process, and the dotted blue box indicates attention fusion to integrate self-attention and interaction-attention.

as a low-dimensional dense vector. The extreme labels are embedded into dense vectors from the label co-exist graph so that the label correlation and local structure can be sufficiently captured.

Word Embedding Once having the d -dimensional word vector $\mathbf{e}_t \in \mathbb{R}^d$ for each word ($t = 1, \dots, n$), the whole document can be taken as a sequence of words $(\mathbf{e}_1, \dots, \mathbf{e}_n)$ as the input of *LAHA*. In order to capture the bi-directional contextual information, we adopt Bi-LSTM [27] to learn the word embedding for each input document. So the whole output of Bi-LSTM can be obtained by

$$H = (H^{(f)}; H^{(b)}) \text{ with } H^{(f)} = (\vec{h}_1, \dots, \vec{h}_n) \in \mathbb{R}^{r \times n}; \quad H^{(b)} = (\overleftarrow{h}_1, \dots, \overleftarrow{h}_n) \in \mathbb{R}^{r \times n} \quad (1)$$

where $\vec{h}_t \in \mathbb{R}^r$ and $\overleftarrow{h}_t \in \mathbb{R}^r$ are the forward and backward word context representations respectively. The whole document is taken as a matrix $H \in \mathbb{R}^{2r \times n}$.

Label Embedding To better extract label correlation information, we firstly build a label co-exist graph from the training data where each labels are represented by nodes. There will be an edge connecting the i -th label and the j -th label if they share at least one document [14]. Our goal is to represent the extreme labels in a low-dimensional latent space so that two nearby labels in the graph have similar representation, i.e., the local structure among labels are preserved as much as possible. Thus, the popular and powerful node2vec [16] is adopted here because it has ability to explore the labels' diverse neighborhoods by a flexible biased random walk procedure in a breadth-first sampling as well as depth-first sampling fashion. Finally, each label will be represented by a r -dimensional dense vector, i.e., $\mathbf{l}_i \in \mathbb{R}^r$ for the i -th label ($i = 1, \dots, k$) and the whole label set can be described by $L = (\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_k) \in \mathbb{R}^{r \times k}$.

3.2 Hybrid Attention Mechanism

Hybrid attention mechanism aims at better representing each document by taking advantage of both document content and label structure. It is composed of self-attention mechanism on document content and interaction-attention mechanism to exploit document content and label structure.

Self-attention (SA) has been successful used in text mining tasks such as relation extraction [30]. In multi-label data, since one document may be tagged by more than one labels, each document should have the most relative context to its corresponding labels. That is, the words in one document make different contributions to each label. To focus on different aspects of document, thus, we introduce self-attention mechanism (SA) [15] on the output of Bi-LSTM (H). The attention score $A^{(S)} \in \mathbb{R}^{n \times k}$ is calculated by

$$T = \tanh(W_{s1}H); \quad A^{(S)} = \text{softmax}(W_{s2}T) \quad (2)$$

where $W_{s1} \in \mathbb{R}^{d_a \times 2r}$ and $W_{s2} \in \mathbb{R}^{k \times d_a}$ are parameters to be trained. $A_j^{(S)} \in \mathbb{R}^n$ is the attention scores of words along the j -th label. To efficiently handle extreme multi-label data, we adopt negative sampling strategy [11] to update W_{s2} and computer $A_j^{(S)}$, so that all positive labels and a random small subset of negative labels are considered. Then, we can obtain the linear combination of context words for each label through self-attention mechanism as $C_j^{(S)} = HA_j^{(S)}$, which can be taken as the representation of the input document along the j -th lable. The whole matrix $C^{(S)} \in \mathbb{R}^{2r \times k}$ is the label-aware document representation under the self-attention mechanism.

Interaction-attention (IA) aims to determine the semantic connection between words and labels in a latent space. With the help of word embedding and label embedding technique, all words and labels are represented in the r -dimensional latent space as $H = (H^{(f)}; H^{(b)})$ and L respectively. To conveniently align the latent space of words and that of labels, a bridge mapping marix $W_q \in \mathbb{R}^{r \times r}$ is trained via $Q = W_q L$. Similar to SA, we can do negative sampling on L to produce $L^* \in \mathbb{R}^{r \times k^*}$ that is extracted from L according to sampled indices, and just use L^* for the following computation.

Inspired by the interaction mechanism [7], we take $Q \in \mathbb{R}^{r \times k}$ as the attention queries for each label, and use H to construct the key-value pairs in terms of forward and backward information for each word. Then, the interactive matching score $M^{(I)} \in \mathbb{R}^{n \times k}$

$$M^{(I)} = \begin{bmatrix} H^{(f)T} & H^{(b)T} \end{bmatrix} \begin{bmatrix} Q \\ Q \end{bmatrix} \quad (3)$$

To make sure the attention weight value fall into the range of $[0, 1]$, we normalize $M^{(I)}$ to obtain the interaction-attention weight $A^{(I)} = (A_{tj}^{(I)})_{t=\{1, \dots, n\}, j=\{1, \dots, k\}}$ as follows.

$$A_{tj}^{(I)} = e^{M_{tj}^{(I)}} / \sum_{i=1}^n e^{M_{ij}^{(I)}} \quad (4)$$

Similar to self-attention mechanism, the label-aware document representation can be calculated by linear combining the label’s context words as $C_j^{(I)} = HA_j^{(I)}$, which can be taken as the representation of the input document along the j -th label. The whole matrix $C^{(I)} \in \mathbb{R}^{2r \times k}$ is the label-aware document representation under the interaction-attention mechanism.

3.3 Attention Fusion (FA)

The above $C^{(S)}$ and $C^{(I)}$ are label-aware document representation. The former focuses on document content, while the latter prefers to the label structure. In order to take advantage of these two parts, an attention fusion strategy is designed here to adaptively extract proper information from these two components and build accurate label-aware document representation. More specifically, a fully connected layer is used to transform the input ($C^{(S)}$ and $C^{(I)}$) to weights $\alpha \in \mathbb{R}^{k \times 1}$ and $\beta \in \mathbb{R}^{k \times 1}$ via

$$\alpha = \sigma(F_1(C^{(S)})); \quad \beta = \sigma(F_2(C^{(I)})) \quad (5)$$

where σ is sigmoid function to ensure the weights falling into $(0, 1)$. Among them, α_j and β_j indicates the importances of self-attention and interaction-attention to final representation along the j -th label respectively. Therefore, we normalize them as $\alpha_j = \alpha_j / (\alpha_j + \beta_j)$ and $\beta_j = 1 - \alpha_j$. With the aid of fusion weights, we can get the final label-aware representation of input document along the j -th label

$$C_j = \alpha_j \times C_j^{(S)} + \beta_j \times C_j^{(I)}. \quad (6)$$

The whole matrix $C \in \mathbb{R}^{2r \times k}$ is the final label-aware document representation.

3.4 Prediction Layer

Once having $C \in \mathbb{R}^{2r \times k}$, we can build the classifier via a fully connected and output layer. The final predictions are obtained by $\hat{y} = \sigma(W_o(f(W_f C)))$ where $W_f \in \mathbb{R}^{r \times 2r}$, $W_o \in \mathbb{R}^{1 \times r}$, f is the activation function ReLU, and σ is adopted to ensure that the output value can be taken as a probability. In this case, the binary cross-entropy loss can be used as loss function which has been proved suitable for XMTC tasks [5].

$$L_{loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k [y_{ij} \log(\hat{y}_{ij}) + (1 - y_{ij}) \log(1 - \hat{y}_{ij})] \quad (7)$$

where N is the number of training documents. The ground truth $y_{ij} = 1$ if the i -th document belongs to the j -th class, otherwise $y_{ij} = 0$.

4 Experiments

In this section, we evaluate the proposed *LAHA* on five benchmark datasets by comparing with the state-of-the-art extreme multi-label learning methods in terms of widely used metrics.

Table 1. Summary of experimental datasets. N is the number of training documents, M is the number of testing documents, D is the number of features, L is the number of class labels, \hat{L} is the average number of labels per document, \hat{N} is the average number of documents per label.

datasets	N	M	D	L	\hat{L}	\hat{N}
<i>AAPD</i> [1]	54,840	1,000	69,399	54	2.41	2444.0
<i>Kan-Shan Cup</i> ²	2,799,967	200,000	411,721	1,999	2.3	3513.1
<i>EUR-Lex</i> [2]	11,585	3,865	171,120	3,956	5.3	15.6
<i>Amazon-12K</i> [4]	490,310	152,981	135,895	12,277	5.4	214.5
<i>Wiki-30K</i> [3]	12,959	5,992	100,819	29,947	18.7	8.1

4.1 Datasets

A series of experiments were carried out on five multi-label datasets with label sizes from 54 to 29,947. The dataset statistics are summarized in Table 1.

4.2 Methodology

Baseline Algorithms The proposed *LAHA* is a deep neural network model, thus the recent deep learning-based XMTC methods (*XML-CNN* [5] and *AttentionXML* [6]) are selected as baselines. Meanwhile, the existing powerful *SLEEC* [8] (an embedding-based method) and *PfastreXML* [9] (a tree-based method) are used as baselines because they obtained the best performance in each type as shown in the Extreme Classification Repository ².

Parameter Settings For all the five datasets, we adopt Glove(300-dimension) [11] as word embedding. The number of Bi-LSTM hidden units is set to $r = 256$. For the self-attention mechanism, $d_a = 256$. In the prediction layer, ReLU is adopted as non-linear activation function. The whole deep model is trained using Adam with the initial learning rate (0.001) and the batch size (64).

Evaluation Metrics In XMTC tasks, rank-based evaluation metrics are popular used to evaluate model performance, including Precision at τ ($P@_\tau$) and normalized Discounted Cumulative Gain at τ ($nDCG@_\tau$). Both of them have been widely used in XMTC tasks. They are defined as

$$P@_\tau = \frac{1}{\tau} \sum_{l \in r_\tau(\hat{\mathbf{y}})} \mathbf{y}_l; \quad nDCG@_\tau = \frac{\sum_{l \in r_\tau(\hat{\mathbf{y}})} \mathbf{y}_l / \log(l+1)}{\sum_{l=1}^{\min(\tau, \|\mathbf{y}\|_0)} 1 / \log(l+1)} \quad (8)$$

where $\mathbf{y} \in \{0, 1\}^k$ is the ground truth label vector of a document and $r_\tau(\hat{\mathbf{y}})$ is the label indexes of top τ highest scores of current prediction result. $\|\mathbf{y}\|_0$ counts the number of relevant labels in the ground truth label vector \mathbf{y} . Larger $P@_\tau$ and $nDCG@_\tau$ indicates better performance.

² <http://manikvarma.org/downloads/XC/XMLRepository.html>.

² <https://biendata.com/competition/zhihu/>

4.3 Ablation Test of *LAHA*

In this section, we firstly demonstrate the effect of each component on *LAHA*. To reach this goal, we do ablation test for self-attention mechanism (SA), interaction-attention mechanism (IA) and attention fusion mechanism (FA) respectively with two datasets: one sparse dataset *EUR-lex* and one dense dataset *AAPD*.

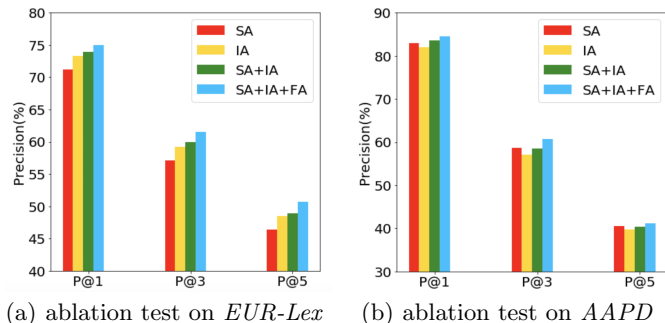


Fig. 2. Ablation test on *EUR-Lex* and *AAPD*. SA=*self-attention*, IA=*interaction-attention*, FA=*attention fusion*, *LAHA*=SA+IA+FA.

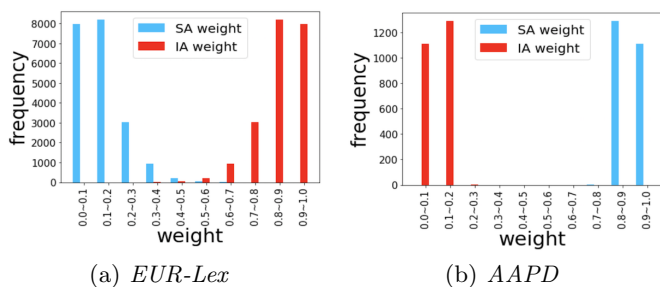


Fig. 3. Weight distributions for two components on *EUR-Lex* and *AAPD*. x -axis indicates the range of weight from 0 to 1 with 0.1 gap. y -axis indicates the frequency that the specific range occurs in current label group.

Fig.2 lists the results on these two datasets in terms of $P@_\tau$ ($\tau = \{1, 3, 5\}$). It can be seen that SA performs well on dense dataset (*AAPD*). However, neither SA nor IA can obtain good result on sparse dataset (*EUR-Lex*). Fortunately, combining SA and IA improves the prediction performance (SA+IA gets better performance than SA and IA). SA prefers to extract the useful content information when constructing the label-aware document representation, but SA ignores the label structure during the learning process. IA implements this by using the label embedding learnt from the label co-exist graph. However, in real application, such graph may contain noisy information (say in dense data). Therefore,

coupling with both attention components does really helpful for final performance because they can benefit each other on different datasets.

To adaptively extract proper information to learn the final label-aware document representation, the attention fusion mechanism is introduced in *LAHA*. Fig.3 lists the distribution of weights on SA and IA. It can be seen that for sparse data (*EUR-Lex*), the interaction-attention plays much more important role than self-attention on learning process, vice verse for dense dataset (*AAPD*). This result further clarifies that IA mechanism can leverage the label structure to improve the prediction performance for sparse data. On the other hand, in *AAPD*, each label has sufficient documents, i.e., SA mechanism can sufficiently capture the label-aware document information and perform well. That is why larger weights are assigned to SA on dense data. Similar trend can be found on other datasets, which are omitted due to the page limitation.

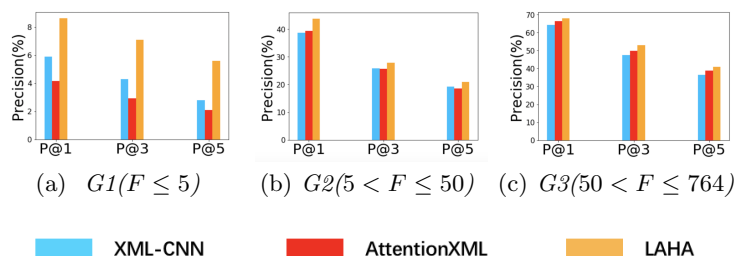


Fig. 4. Comparing *XML-CNN*, *AttentionXML* and *LAHA* on different label groups of sparse data(*EUR-Lex*) in terms of $P@_{\tau}$ ($\tau = \{1, 3, 5\}$). F is frequency of label occurring in training set.

4.4 Comparison with Deep Methods on Sparse Datasets

In order to explore the effect of *LAHA* on sparse datasets, we further divide labels into three groups according to their occurring frequencies. Fig.4 shows the prediction performance obtained by three deep methods. Obviously, label prediction in *G1* is much harder than in other two groups due to the lack of training documents. All methods become better from *G1* to *G3*, which is reasonable since *G3* contains more training documents than *G1*. *LAHA* has an overall improvement for all groups compared with two baselines. This result further demonstrates the superiority of the proposed hybrid attention mechanism on XMTC with large-scale tail labels. Similar phenomena can be found on other sparse datasets, which are omitted due to the page limitation.

To further investigate the attention-based methods, we visualize the attention weights on the original document using heat map, as shown in Fig.5. This example document belongs to 28 labels named as *autism*, *children*, *childhood*, *disease*, *asperger*, *social norm*, *health*, *neurology*, *abnormal* and etc. From the

autism explicitly cited references three reference anything but notable references subject tends attract lot controversy lod nine two nine nine zero zero omim two zero nineeight five zero zero zero one five two six med three two zero two emedicine mult zero zero one three two one one four four two autism overview autism disordercharacterized impaired social interaction verbal non verbal communication restricted repetitive behavior parents usually notice signs first two years child life signs oftendevelop gradually though children autism reach developmental milestones normal pace regress diagnostic criteria require sym ptoms become apparent early childhoodtypically age three autism highly heritable researchers suspect environmental genetic factors causes rare cases autism strongly associated agents cause birth defect sc ontroversies surround proposed environmental causes example vaccine hypotheses autism affects information processing brain altering nerve cells synapses connectorganize occurs well understood one thre e recognized disorders autism spectrum two asperger syndrome lacks delays cognitive development language pervasivedevelopmental disorder otherwise specified commonly abbreviated nos diagnosed full set criteria autism asperger syndrome met early speech behavioral interventions help children autism gain self care social communication skills although known cure reported cases children recovered many c hildren autism live independently reachingadulthood though become successful autistic culture developed individuals seeking cure others believing autism accepted difference treated disorder two zero one ze rounumber people affected estimated one two per one zero zero zero worldwide occurs four five times often boys girls one five children united states one six eight diagnosedthree zero increase one eight eight t wo zero one two rate autism among adults aged one eight years united kingdom one one number people diagnosed increasing dramatically since one nine eight zero partly due changes diagnostic practice go vurnment subsidized financial incentives named diagnoses question whether actual ratesincreased unresolved characteristics autism highly variable disorder first appears infancy childhood generally follows at eady course without remission overt symptomsgradually begin age six months become established age two three years tend continue adulthood although often muted form distinguished single symptom chara cteristictriad symptoms impairments social interaction communication restricted interests repetitive behavior aspects atypical eating also common essential diagnosisautism individual symptoms o our general population appear associate highly without sharp line separating severe common traits social development social deficitsdistinguish autism related autism spectrum disorders see classification de velopmental disorders people autism social impairments often lack intuition others many peopletake granted noted autistic temple described inability understand social communication people normal neural de velopment leaving feeling like anthropologist mars unusualsocial development becomes apparent early childhood autistic infants show less attention social stimuli smile look others less often respond less nam e autistic different strikingly social norms example less eye contact turn taking ability use simple movements express pointing things three five year old children autism less likely exhibit socialunderstanding appro ach others spontaneously imitate respond emotions communicate take turns others however form attachments primary caregivers autism display moderately less attachment security children although differen ce disappears children higher

(a) The words with largest label-aware attention weights output by *AttentionXML* (word→{labels}) are: *autism*→{autism}; *cure others*→{disease}; *disorder*→{abnormal}; *social communication, approach others*→{social norm}; *children*→{children, childhood}.

autism explicitly cited references three reference anything put notable references subject tends attract lot controversy lod nine two nine nine zero zero omim two zero nineeight five zero zero zero one five two six med three two zero two emedicine mult zero zero one three two one one four four two autism overview autism disordercharacterized impaired social interaction verbal non verbal communication restricted repetitive behavior parents usually notice signs first two years child life signs oftendevelop gradually though children autism reach developmental milestones normal pace regress diagnostic criteria require sym ptoms become apparent early childhoodtypically age three autism highly heritable researchers suspect environmental genetic factors causes rare cases autism strongly associated agents cause birth defect as ontroversies surround proposed environmental causes example vaccine hypotheses autism affects information processing brain altering nerve cells synapses connectorganize occurs well understood one thre e recognized disorders autism spectrum two asperger syndrome lacks delays cognitive development language pervasivedevelopmental disorder otherwise specified commonly abbreviated nos diagnosed full set criteria autism asperger syndrome met early speech behavioral interventions help children autism gain self care social communication skills although known cure reported cases children recovered many c hildren autism live independently reachingadulthood though become successful autistic culture developed individuals seeking cure others believing autism accepted difference treated disorder two zero one ze rounumber people affected estimated one two per one zero zero zero worldwide occurs four five times often boys girls one five children united states one six eight diagnosedthree zero increase one eight eight t wo zero one two rate autism among adults aged one eight years united kingdom one one number people diagnosed increasing dramatically since one nine eight zero partly due changes diagnostic practice go vurnment subsidized financial incentives named diagnoses question whether actual ratesincreased unresolved characteristics autism highly variable disorder first appears infancy childhood generally follows at eady course without remission overt symptomsgradually begin age six months become established age two three years tend continue adulthood although often muted form distinguished single symptom chara cteristictriad symptoms impairments social interaction communication restricted interests repetitive behavior aspects atypical eating also common essential diagnosisautism individual symptoms o our general population appear associate highly without sharp line separating severe common traits social development social deficitsdistinguish autism related autism spectrum disorders see classification de velopmental disorders people autism social impairments often lack intuition others many peopletake granted noted autistic temple described inability understand social communication people normal neural de velopment leaving feeling like anthropologist mars unusualsocial development becomes apparent early childhood autistic infants show less attention social stimuli smile look others less often respond less nam e autistic different strikingly social norms example less eye contact turn taking ability use simple movements express pointing things three five year old children autism less likely exhibit socialunderstanding appro ach others spontaneously imitate respond emotions communicate take turns others however form attachments primary caregivers autism display moderately less attachment security children although differen ce disappears children higher

(b) The words with largest label-aware attention weights output by *LAHA* (word→{labels}) are: *autism*→{autism}; *disorder*→{abnormal}; *child life, in-fancy childhood, children*→{children, childhood}; *diagnostic, genetic factor, lack intuition*→{disease}; *synapses connect organize*→{neurology}; *asperger syndrome*→{asperger}; *social communication*→{social norm}; *security*→{disease, health}.

Fig. 5. Heat map of label-aware attention weights obtained by (a) *AttentionXML* and (b) *LAHA* on an example document from *Wiki30K*.

attention weights, we can see that *AttentionXML* only captures few key words for few related labels. As expected, *LAHA* focuses on the related information as much as possible due to the capacity making full use of label structure and document content.

4.5 Comparison Results and Discussion

In this section, the proposed *LAHA* is evaluated on five benchmark datasets by comparing with four baselines in terms of $P@τ$ and $nDCG@τ$ ($τ = \{1, 3, 5\}$). Table 2 shows the averaged performance of all test documents. According to the formula (8), we know $P@1 = nDCG@1$, thus only $nDCG@3$ and $nDCG@5$ are listed. In each line, the best result is marked in bold, and the second best is underlined.

From Table 2, we can make a number of observations about these results. Firstly, *LAHA* outperforms the traditional powerful embedding-based and tree-based methods in most cases, while slightly underperforms the embedding-based method *SLEEC* on *EUR-Lex* and *Wiki-30K*. From Table 1, we can see there are only 11,585 and 12,959 training documents in these two datasets, in this case, the deep model may be not sufficiently trained. Second, *LAHA* is consistently superior to the state-of-the-art deep XMTC methods. The main reason is that

Table 2. Comparing *LAHA* with four baselines in terms of various metrics on five benchmark datasets.

Datasets	Metric	<i>SLEEC</i>	<i>PfastreXML</i>	<i>XML-CNN</i>	<i>AttentionXML</i>	<i>LAHA</i>
<i>AAPD</i>	<i>P@1</i>	81.96%	82.35%	76.25%	<u>83.02%</u>	84.48%
	<i>P@3</i>	57.48%	58.01%	54.34%	<u>58.72%</u>	60.72%
	<i>P@5</i>	38.99%	40.13%	37.84%	<u>40.56%</u>	41.19%
	<i>nDCG@3</i>	77.65%	<u>78.26%</u>	72.01%	78.01%	80.11%
	<i>nDCG@5</i>	81.59%	82.03%	76.40%	<u>82.31%</u>	83.70%
<i>Kan-Shan Cup</i>	<i>P@1</i>	51.41%	52.29%	49.68%	<u>53.69%</u>	54.38%
	<i>P@3</i>	32.81%	32.99%	32.27%	<u>34.10%</u>	34.60%
	<i>P@5</i>	24.29%	24.58%	24.17%	<u>25.16%</u>	25.88%
	<i>nDCG@3</i>	49.32%	49.96%	46.65%	<u>51.03%</u>	51.70%
	<i>nDCG@5</i>	49.74%	50.11%	49.60%	<u>53.96%</u>	54.65%
<i>EUR-Lex</i>	<i>P@1</i>	75.18%	73.03%	70.94%	71.89%	<u>74.95%</u>
	<i>P@3</i>	61.67%	60.39%	56.02%	57.74%	<u>61.48%</u>
	<i>P@5</i>	50.23%	49.69%	45.36%	47.35%	50.71%
	<i>nDCG@3</i>	<u>63.79%</u>	62.51%	59.68%	61.29%	64.89%
	<i>nDCG@5</i>	<u>58.03%</u>	57.72%	53.82%	56.71%	59.28%
<i>Amazon-12K</i>	<i>P@1</i>	93.49%	<u>93.95%</u>	93.15%	93.75%	94.87%
	<i>P@3</i>	78.01%	78.33%	76.11%	<u>78.36%</u>	79.16%
	<i>P@5</i>	62.09%	<u>62.77%</u>	60.51%	62.14%	63.16%
	<i>nDCG@3</i>	86.89%	<u>88.41%</u>	86.75%	87.62%	89.13%
	<i>nDCG@5</i>	84.53%	<u>86.23%</u>	84.01%	86.06%	87.57%
<i>Wiki-30K</i>	<i>P@1</i>	85.26%	82.81%	82.90%	81.98%	<u>84.18%</u>
	<i>P@3</i>	73.91%	68.48%	67.46%	67.27%	<u>73.14%</u>
	<i>P@5</i>	62.55%	59.93%	57.09%	56.43%	62.87%
	<i>nDCG@3</i>	76.01%	72.15%	71.04%	70.77%	<u>75.64%</u>
	<i>nDCG@5</i>	68.27%	63.83%	62.92%	62.35%	<u>67.82%</u>
<i>Win times</i>		6	0	0	0	19

LAHA has ability to sufficiently determine the label-aware document representation while *XML-CNN* does not. Even though *AttentionXML* tries to find the relation between each pair of document and label, it only focuses on document content, which will degrade its performance on tail labels due to lack of information. Fortunately, *LAHA* addresses this issue by simultaneously considering label structure via a hybrid attention mechanism.

5 Conclusions and Future Work

In this paper, a new XMTC method, *LAHA*, is proposed. *LAHA* utilizes self-attention and interaction-attention to extract the semantic relation between words and labels, and an attention fusion to construct the label-aware document representation. Extensive experiments on five benchmark datasets prove the superiority of *LAHA* by comparing with the state-of-the-art XMTC methods. In a nutshell, the novelty of *LAHA* lies in its providing a label-aware document representation that captures both document content and label structure, and has

better discriminative ability than baselines. In real applications, more contents can be collected such as label content, which is proved to be helpful in XMTC [7]. We therefore plan to extend the current model with such information.

References

1. Yang P, Sun X, Li W, Ma S, Wu W, Wang H. SGM: sequence generation model for multi-label classification. In: Proc. of COLING. 2018: 3915-3926.
2. Mencia E L, Frnkranz J. Efficient pairwise multilabel classification for large-scale problems in the legal domain. In: Proc. of ECML & PAKDD. Springer, 2008: 50-65.
3. Zubiaga A. Enhancing navigation on wikipedia with social tags. arXiv:1202.5469, 2012.
4. McAuley J, Leskovec J. Hidden factors and hidden topics: understanding rating dimensions with review text. In: Proc. of ACM RecSys. 2013: 165-172.
5. Liu J, Chang C, Wu Y, Yang Y. Deep learning for extreme multi-label text classification. In: Proc. of the 40th ACM SIGIR, 2017: 115-124.
6. You R, Dai S, Zhang Z, Mamitsuka H, Zhu S. AttentionXML: extreme multi-label text classification with multi-label attention based recurrent neural networks. arXiv:1811.01727, 2018.
7. Du C, Chin Z, Feng F, Zhu L, Gan T, Nie L. Explicit Interaction Model towards Text Classification. arXiv:1811.09386, 2018.
8. Bhatia K, Jain H, Kar P, Varma M, Jain P. Sparse local embeddings for extreme multi-label classification. In: Proc. of NIPS. 2015: 730-738.
9. Jain H, Prabhu Y, Varma M. Extreme multi-label loss functions for recommendation, tagging, ranking & other missing label applications. In: Proc. of ACM SIGKDD, 2016: 935-944.
10. Prabhu Y, Kag A, Harsola S, Agrawal R, Varma M. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In: Proc. of WWW, 2018: 993-1002.
11. Pennington J, Socher R, Manning C. Glove: Global vectors for word representation. In: Proc. of EMNLP, 2014: 1532-1543.
12. Sutskever I, Vinyals O, Le Q V. Sequence to sequence learning with neural networks. In: Proc. of NIPS. 2014: 3104-3112.
13. Wang S, Jiang J. Learning natural language inference with LSTM. arXiv:1512.08849, 2015.
14. Zhang W, Yan J, Wang X, Zha H. Deep extreme multi-label learning. In: Proc. of ACM ICMR, 2018: 100-107.
15. Lin Z, Feng M, Santos N, Yu M, Xiang B, Zhou B, Bengio Y. A structured self-attentive sentence embedding. arXiv:1703.03130, 2017.
16. Grover A, Leskovec J. node2vec: Scalable feature learning for networks. In: Proc. of ACM SIGKDD, 2016: 855-864.
17. Prabhu Y, Varma M. Fastxml: A fast, accurate and stable tree-classifier for extreme multi-label learning. In: Proc. of ACM SIGKDD, 2014: 263-272.
18. Hsu D J, Kakade S M, Langford J, Zhang T. Multi-label prediction via compressed sensing. In: Proc. of NIPS, 2009: 772-780.
19. Zhang Y, Schneider J. Multi-label output codes using canonical correlation analysis. In: Proc. of AISTATS. 2011: 873-882.
20. Tai F, Lin H T. Multilabel classification with principal label space transformation. Neural Computation, 2012, 24(9): 2508-2542.

21. Balasubramanian K, Lebanon G. The landmark selection method for multiple output prediction. arXiv:1206.6479, 2012.
22. Cisse M, Usunier N, Artieres T, Gallinari P. Robust bloom filters for large multilabel classification tasks. In: Proc. of NIPS. 2013: 1851-1859.
23. Zhou C, Sun C, Liu Z, Lau F. A C-LSTM neural network for text classification. arXiv:1511.08630, 2015.
24. Munkhdalai T, Yu H. Neural semantic encoders. In: Proc. of ACL, 2017, 1: 397.
25. Cho K, Van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, Bengio Y. Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv:1406.1078, 2014.
26. Zhang X, Zhao J, LeCun Y. Character-level convolutional networks for text classification. In: Proc. of NIPS. 2015: 649-657.
27. Zhou P, Qi Z, Zheng S, Xu J, Bao H, Xu B. Text classification improved by integrating bidirectional LSTM with two-dimensional max pooling. arXiv:1611.06639, 2016.
28. Wang X, Jiang W, Luo Z. Combination of convolutional and recurrent neural network for sentiment analysis of short texts. In: Proc. of COLING, 2016: 2428-2437.
29. Devlin J, Chang W, Lee K, Toutanova K. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
30. Wang L, Cao Z, De Melo G, Liu Z. Relation Classification via Multi-Level Attention CNNs. In: Proc. of ACL. 2016: 1298-1307.