

Label Mask for Multi-Label Text Classification

Rui Song^a, Xingbing Chen^d, Zelong Liu^e, Haining An^e, Zhiqi Zhang^f,
Xiaoguang Wang^{b,*}, Hao Xu^d

^a*School of Artificial Intelligence, Jilin University*

^b*Public Computer Education and Research Center, Jilin University*

^c*College of Computer Science and Technology, Key Laboratory of Symbolic Computing and Knowledge Engineering of Ministry of Education, Jilin University*

^d*College of Electronic Science and Engineering, Jilin University*

^e*College of Construction Engineering, Jilin University*

^f*College of Software, Jilin University*

Abstract

One of the key problems in multi-label text classification is how to take advantage of the correlation among labels. However, it is very challenging to directly model the correlations among labels in a complex and unknown label space. In this paper, we propose a Label Mask multi-label text classification model (LM-MTC), which is inspired by the idea of cloze questions of language model. LM-MTC is able to capture implicit relationships among labels through the powerful ability of pre-train language models. On the basis, we assign a different token to each potential label, and randomly mask the token with a certain probability to build a label based Masked Language Model (MLM). We train the MTC and MLM together, further improving the generalization ability of the model. A large number of experiments on multiple datasets demonstrate the effectiveness of our method.

Keywords: Multi-label Text Classification, Bert, Cloze Questions, Masked Language Model

*Corresponding author

Email addresses: songrui20@mails.jlu.edu.cn (Rui Song), 1276402580@qq.com (Xingbing Chen), 18943698576@163.com (Zelong Liu), anhn2418@jlu.edu.cn (Haining An), 1005359144@qq.com (Zhiqi Zhang), wangxiaog@jlu.edu.cn (Xiaoguang Wang), xuhao@jlu.edu.cn (Hao Xu)

1. Introduction

Text categorization is a basic task in natural language processing, Multi-label text classification assigns multiple different labels to a document, rather than a document corresponding to a single category. In recent years, multi-label text classification has been widely used in sentiment analysis [1], topic classification [2], information retrieval [3], label recommendation [4]. How to fully capture semantic patterns from original documents, how to extract discriminant information related to corresponding labels from each document, and how to accurately mine the correlation between labels are the three aspects that researchers pay attention to [5]. Among them, the relevance modeling of complex and unknown label systems has always been the focus of scholars' efforts.

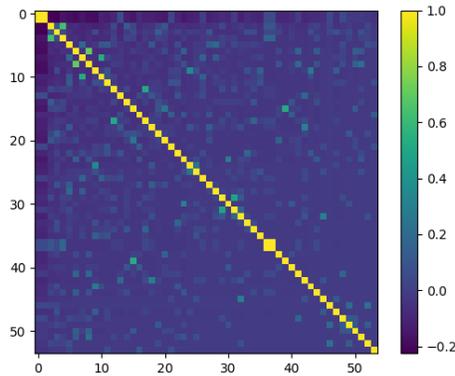


Figure 1: Pearson correlation coefficient between all label pairs in AAPD dataset. The lighter the color, the more relevant the label pairs are.

One of the most direct methods to solve multi-label text classification is to transform the multi-label text classification task into several binary classification tasks [6], however, this tends to ignore the relationship between multiple labels. Similarly, some deep learning approaches, such as CNN [7] and attention mechanism [2], can model the documents effectively, but they still ignore the relationship between labels. As shown in Figure 1, there are specific correla-

tions between different label pairs in the dataset AAPD calculated by Pearson correlation coefficient. For labels 0 and 1, the correlation is 1, which means that both appear together in all instances. Therefore, for some datasets with less label information or severe long-tail distribution, the association between labels can provide more important information [8].

The emergence of large-scale pre-trained language models, such as Bert (Bidirectional Encoder Representations from Transformers) [9] have made knowledge transfer in the field of natural language processing easier. Studies have confirmed that the intermediate layer of Bert encodes a wealth of linguistic information [10]. Inspired by Cloze Questions (CQ) methods based large-scale pre-trained language models, we propose a Label Mask multi-label text classification model (LM-MTC) to capture the potential semantic and association relations among labels [11, 12]. Specifically, we map different labels to different tokens and build a set of token prefix templates. During training, we spliced the token template with the sentences to be classified and input them to Bert. When predicting, we mask all the label tokens and predict them. The advantages of Bert can help the model adaptively capture the semantic relationship between labels and documents, as well as the association relationships. In addition, to make better use of the predictive ability of Bert, we constructed a multi-task framework, randomly masked label tokens, and used the Mask Language Model (MLM) to predict masked tokens to assist in optimizing the multi-label text classification learning task. Our contributions are as follows:

- We propose a Label Mask Multi-label Text Classification model (LM-MTC), which converts multi-label text classification into a Cloze Questions (CQ) task and captures the potential relationship between labels with the help of a pre-trained language model. We introduced MLM with LM-MTC for joint training which further improved the performance of the model.
- Through attention analysis, we confirm the LM-MTC’s ability to capture potential associations between labels.

- We carry out lots of experiments on different types of multi-label text classification tasks to prove the effectiveness of the proposed model.

2. Related Work

2.1. Multi-label Text Classification

Multi-label text classification is a basic task in NLP. There are some methods to solve it by transforming the multi-label text classification task into several binary classification tasks [6, 7, 2]. Some approaches take advantage of pair associations or mutexes between labels. Pairwise comparison (RPC) induces a binary preference relation using a natural extension of pairwise classification, which transforms the multi-label learning task into a label ranking task [13].

However, it is more efficient to assume that a label can be related to multiple labels and take advantage of the higher-order dependencies of the label. Classifier chains (CC) transforms the task of MTC into a chain of binary classification tasks [14]. K-labelsets (RAkEL) constructs small random subsets of labels and transforms MTC into a single-label classification task of random subsets [15]. With the development of deep learning in recent years, some studies have resorted to sequence learning models to solve MTC such as [16], and Sequence Generation Model (SGM) [17]. They generate a possible label sequence through the RNN decoder. However, sequence model requires searching for the optimal solution in the potential space, which is too time-consuming when there are too many labels.

Some approaches model the joint probability distribution of labels, rather than associations for specific labels, such as Bayesian networks [18, 19] and undirected graph models [20]. [21] strengthen the similarity between the joint distribution of multi-labels and the predicted multi-labels by means of adverse learning framework. In recent years, due to the effectiveness of Graph Neural Network (GNN) [22] in modeling non-Euclidean spatial data, some methods use GNN to capture the correlation of labels. Label-Specific Attention Network (LSAN) proposes a Label Attention Network model that considers both doc-

ument content and Label text, and uses self-attention mechanism to measure the contribution of each word to each Label [23]. MAGNET uses a feature matrix and a correlation matrix to capture and explore key dependencies between labels [24].

2.2. Cloze Questions

A natural way to gain knowledge from a pre-trained MLM is to set the task as fill-in-the-blanks. AUTOPROMPT creates prompts for a diverse set of tasks automatically and shows the inherent ability of MLM to perform emotion analysis and natural language reasoning [25]. [11] converts the input text into a cloze problem containing the task description combined with gradient-based optimization and proves that 'green' LM like Bert can still have competitive performance compared with GPT-3 [26]. Furthermore, PTE has successfully solved the problems of text classification and natural language reasoning in small samples by using cloze questions method [12]. In addition, [27] propose a similar P-tuning method to improve GPT's ability in natural language understanding tasks.

3. Preliminaries

First, we give some notations and describe the MTC task. For the given label space $\mathbf{L} = \{l_1, l_2, \dots, l_L\}$ and a text $x = \{w_1, w_2, \dots, w_m\}$, the MTC task aims to learn a mapping function $\chi : x \rightarrow \hat{L}$, where $\hat{L} \subseteq L$ and $|\hat{L}| \geq 1$. Unlike single-label tasks, a text can belong to several different categories. For ease of calculation, the output space of the labels is defined as a vector $L_O \in \mathbf{R}^{|\mathbf{L}|}$. For example, when $L = 3$ and $\hat{L} = \{l_1, l_2\}$, then L_O is $[1, 1, 0]$. Then we can rewrite the objective function as $\chi' : x \rightarrow L_O$. In addition, we define 1 as a positive label and 0 as a negative label.

Given a movie emotion dichotomy sentence, "The movie was so touching!", CQ usually generates a new sentence for input by a prefix/suffix template τ : "The movie was so touching! **I** $\#$ **it!**". $\#$ can be either "like" or "hate", indicating positive or negative emotions, respectively. The new input with a prefix

template can be expressed as:

$$x' = \tau || x \quad (1)$$

where $||$ denotes concatenation.

Consider an MLM M with a vocabulary V and a sentence x with mask m , the task of M is to predict the probability of the original word in the mask $p_M^m(w|x)$ where $w \in V$. The goal of CQ is to predict the true value of $\#$ by MLM, that is $p_M^m(\#|x')$. So the classification task can be transformed into the prediction task of MLM with the help of templates.

4. LM-MTC

In this section, we describe the proposed model in detail. The specific execution process is shown in Algorithm 1.

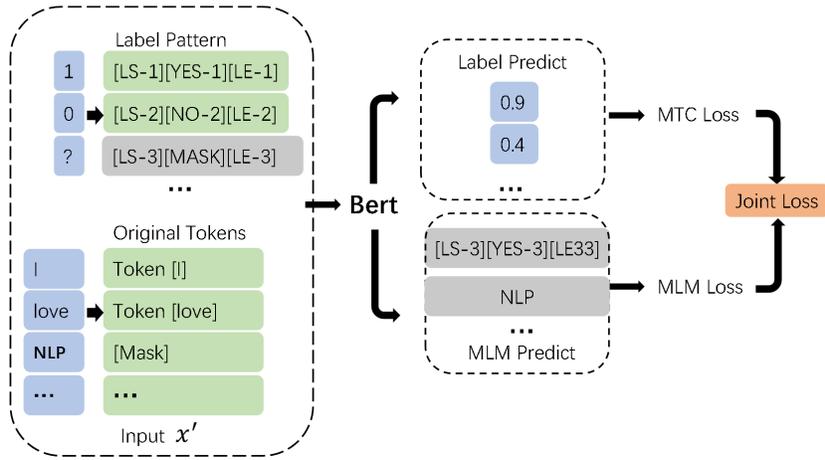


Figure 2: LM-MTC model structure.

4.1. Cloze Patterns

For cloze tasks, although some studies have demonstrated the advantages of a template approach, it is not clear whether the same template will work for every model, or what kind of template will fit the model better [25]. For

MTC, because different documents contain different numbers of true labels and the output is mapped to a single token, it is not possible to build a specific template for each label. To this end, we build a template system for the entire label space. First of all, considering the different states of labels, the labels in each position should have three different states: 0,1 and *mask*. We emphasize the order of the different labels, which is very important for label prediction [17]. In addition, we also introduce a location-based prompt to allow Bert to clearly recognize where the current Label is located. For label sequence with mask [1,0,*mask*], we generate the template as follows:

$$\begin{aligned}
& [LS - 1][YES - 1][LE - 1] \\
& [LS - 2][NO - 2][LE - 2] \\
& [LS - 3][MASK - 3][LE - 3]
\end{aligned} \tag{2}$$

where LS denotes *label start* and LE denotes *label end*.

4.2. Training and Inference

The forward propagation. After the template is generated, we treat it as a prefix to the original sentence and input x' to the pre-training model together. The training process has two main objectives: to predict the probability distribution of multiple labels in the label space, and to predict the *mask* by MLM. Assume that the output of Bert is $O \in \mathbf{R}^{|x'|_{max} * 768}$, then the prediction of the distribution of the labels and the prediction of the mask can be obtained by using one layer of full connection:

$$\begin{aligned}
O_l &= O * W_l + b_l \\
O_m &= O * W_m + b_m
\end{aligned} \tag{3}$$

where $|x'|_{max}$ denotes max token length, $W_l \in \mathbf{R}^{768 * |L|}$. In order for MLM to predict the mask, we need to extend V to V' . V' depends on the size of the label space, so that $W_m \in \mathbf{R}^{768 * |V'|}$.

Joint loss. We use the Binary Cross Entropy (BCE) as the loss function

for MTC and the Cross Entropy as the loss function for MLM. BCE loss can be written as follows:

$$\mathcal{L}_{mtc} = \frac{1}{|L|} \sum_{i=1}^{|L|} (y_{ti} \log(\sigma(y_{pi})) + (1 - y_{ti}) \log(1 - \sigma(y_{pi}))) \quad (4)$$

where σ denotes sigmoid activation function, y_t is the ground truth label and y_p denotes the predicted results. The final joint loss function is:

$$\mathcal{L} = \mathcal{L}_{mtc} + \lambda \mathcal{L}_{mlm} \quad (5)$$

Inference. When inferring, mask all the labels and calculate the probability of all the masked labels. As with training, we express the output of labels as O_l and use the logistic sigmoid function for probability normalization:

$$P_l = \sigma(O_l) \quad (6)$$

Then, all probability values greater than 0.5 are predicted to be positive labels, otherwise predicted to be negative labels.

Algorithm 1 Label Mask Multi-label Text Classification

Input: The original sentence X corresponds to the label Y_t .

Output: Predicted labels Y_p and a trained model M .

- 1: **if** train **then**
 - 2: generate templates T as Eq 2 and input tokens with *masks* X' as Eq 1.
 - 3: train M by optimizing the joint loss function as Eq 5.
 - 4: **else**
 - 5: generate templates T by masking all the labels and input tokens X' as Eq 1.
 - 6: calculate output O_l by feeding X' into M .
 - 7: predict Y_p by probability distribution as Eq 6.
 - 8: **end if**
-

Dataset	Label Sets	Train Size	Test Size
GAIC	17	24000	3000
AAPD	54	53840	1000
Reuters-21578	90	8630	2158
RMSC	22	5020	646
Emotion	28	43410	5427
Toxic	6	11357	4868

Table 1: Datasets statistics.

5. Experiments

5.1. Datasets

In view of the wide application of multi-label text classification, we applied our method on different types of data sets to verify the effectiveness of LM-MTC. The statistics of the dataset are shown in Table 1.

- GAIC¹. A competition dataset of desensitized medical texts which all words are replaced by numbers. We combined the dataset of the preliminary match and the semifinal match together, removed the 12 labels added in the semifinal match, and finally constituted a 17-label classified dataset of 30,000 samples. Then, we divide the training set, development set and test set according to the ratio of 8:1:1. We pre-trained a Bert model based on desensitization data and used the parameters of this Bert in the MTC task.
- AAPD [17]. A widely used large-scale classification dataset for multidisciplinary academic papers. The goal is to predict the subject by abstracts.
- Reuters-21578 [28]. Reuters news text dataset created in 1987 which has been a standard benchmark for MTC. We follow the classification criteria of [24] and use 90 categories.
- RMSC [29]. Collected from a Chinese popular music website². The goal is to distinguish musical styles based on different reviews. 22 styles are de-

¹<https://tianchi.aliyun.com/competition/entrance/531852>

²<https://music.douban.com>

fined. We divide the dataset into training/validation/test sets as described in the original paper. Note that since it is a Chinese dataset, we use the Chinese pre-trained Bert model³. In the process of data preprocessing, we remove non-Chinese and non-English special symbols.

- Emotion [30]. A largest manually annotated dataset of 58k English Reddit comments for fine-grained sentiment classification labeled for 27 emotion categories and neutral, 28 categories totally.
- Toxic Comments⁴. A dataset from Toxic Comment Classification Challenge Competition contains text that may be considered profane, vulgar, or offensive. We remove comments that don't carry any negative sentiment and keep only 16,225 tagged records as our dataset. We split the train/test set randomly in a 7:3 ratio.

5.2. Evaluation Metrics

The same as [15, 17], Hamming Loss and micro-F1 Score are used for the main evaluation metrics. Besides, we also use Accuracy and Micro-Jaccard for further evaluation.

- Accuracy. The strict accuracy of multi-label classification. For the given prediction result $Y_p \in \mathbf{R}^{|\Gamma|*|L|}$ and ground truth $Y_t \in \mathbf{R}^{|\Gamma|*|L|}$, Accuracy is calculated as follows:

$$Accuracy = \sum_i \frac{\Xi(Y_{ti}, Y_{pi})}{|\Gamma|} \quad (7)$$

where $|\Gamma|$ is the test set size, $\Xi(\cdot)$ is an indicator function. If the corresponding elements at all positions are equal in Y_{ti} and Y_{pi} , $\Xi(Y_{ti}, Y_{pi}) = 1$, else $\Xi(Y_{ti}, Y_{pi}) = 0$.

³<https://huggingface.co/bert-base-chinese>

⁴<https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge>

- Micro-F1 [17]. It can be interpreted as a weighted average of accuracy and recall. It calculates indicators globally by calculating total true positives, false negatives and false positives.
- Micro-Jaccard [31]. The Jaccard similarity coefficient is defined as the size of the intersection divided by the size of the union of two label sets.
- Hamming loss (HL) [32]. It directly calculates the proportion of misclassified labels. A value of 0 means that all labels for each sample have been assigned the correct label.

5.3. Baselines

We compare LM-MTC with the widely available baselines:

- Binary Relevance (BR) [6]. Categorize each label separately, regardless of the correlation between the labels.
- Classifier Chains (CC) [14]. Consider the high order correlation between labels and convert it into a binary classification chain.
- CNN [33]. Convolutional neural network is used to extract text features and outputs the distribution of labels in the label space.
- CNN-RNN [16]. CNN and RNN are combined for local and global text modeling.
- Hierarchical Attention Network (HAN) [2]. Model text with hierarchical attention to words and sentences.
- HAN+LG [29]. Introduce Label Graph (LG) on the basis of HAN, and the soft association relationship between labels is trained by LG.
- SGM [17]. View the MCT as a sequence generation problem, and apply a sequence generation model with a novel decoder structure to solve it.
- Bert [9]. Self-attention based pretrained language model. Make different fine-tuning for different downstream tasks.

- Bert+MLM. On the basic Bert classification, additional MLM tasks are added.
- MEGNET [24]. A graph attention network-based model to capture the attentive dependency structure among the labels.
- Label-Wise (LW) LSTM with PT and FT [34]. A document representation with label-aware information is obtained through a pre-training model and fine-tuned for different downstream tasks. PT denotes the pre-training method. FT denotes the fine-tuned method on downstream task.

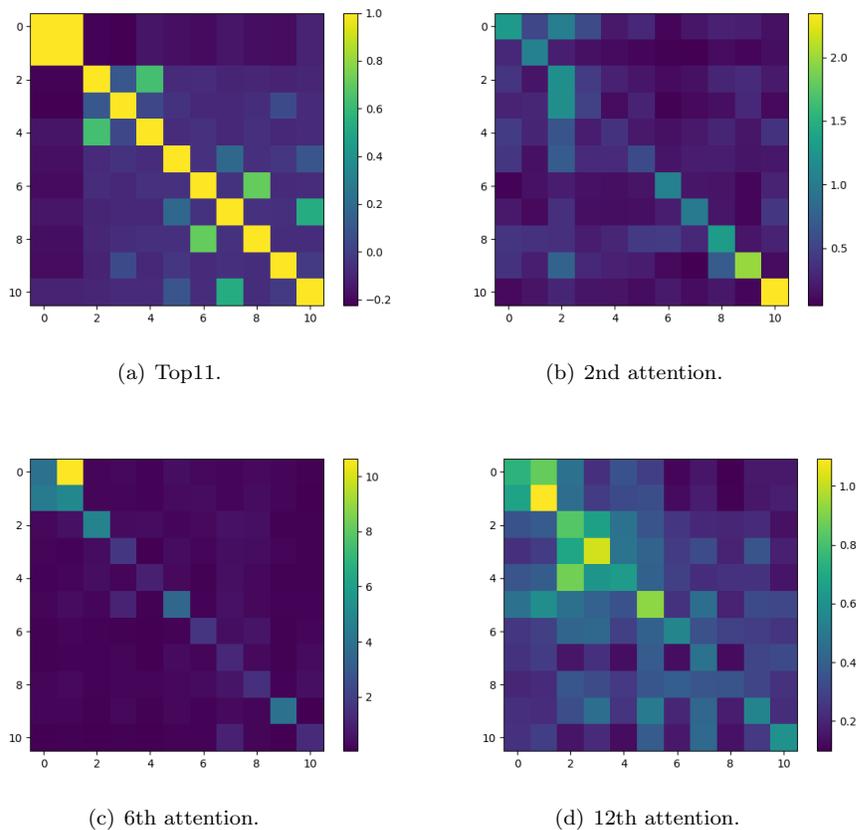


Figure 3: Spearman correlation coefficient of top11 labels and attention visualization of different Bert layers of AAPD test set.

Dataset	Model	Accuracy(+,%)	Micro-F1(+,%)	Micro-Jaccard(+,%)	HL(-)
# GAIC	Bert	89	92	86	0.02
	Bert+MLM	87	91	84	0.02
	LM	9	93	86	0.02
	LM+MLM	92	95	9	0.01
# AAPD	Bert	42.2	72.57	56.95	0.0226
	Bert+MLM	40.5	72.8	57.23	0.0227
	LM	41.3	73.73	58.39	0.0228
	LM+MLM	42.3	73.61	58.23	0.0224
# Reuters-21578	Bert	85.25	89.09	80.33	0.0029
	Bert+MLM	85.03	89.4	80.83	0.00274
	LM	84.15	89.6	81.15	0.0028
	LM+MLM	85.59	89.97	81.77	0.00264
# RMSC	Bert	37.46	72.21	56.51	0.0507
	Bert+MLM	39.63	73.18	57.71	0.0499
	LM	36.53	72.81	57.24	0.0499
	LM+MLM	38.39	73.8	58.47	0.0517
# Emotion	Bert	47.08	57.55	40.4	0.031
	Bert+MLM	43.97	54.7	37.7	0.033
	LM	45.35	58.02	40.86	0.0306
	LM+MLM	47.63	59.07	41.92	0.0304
# Toxic Comments	Bert	52.67	85.76	75.07	0.1043
	Bert+MLM	51.97	84.96	73.86	0.1087
	LM	52.71	85.86	75.23	0.1048
	LM+MLM	53.53	85.77	75.08	0.1041

Table 2: Performance over different datasets. The bold represents the optimal results.

5.4. Details

We set the learning rate as 5e-5, the batch size as 16, and the running epoch as 40. We set the warm up epoches ratio to 0.1, set the mask probability of MLM to 0.15. We use AdamW as the optimizer [35]. All of the code is written using PyTorch and runs on NVIDIA RTX 3090.

5.5. Overall Results

We report the experimental results of our method on all datasets in Table 2 and compare them with Bert and Bert+MLM. We calculate the Accuracy, Micro-F1, Micro-Jaccard and Hamming-Loss. Compared with BERT, in most cases, LM has significant performance improvement in all evaluation metrics, which indicates that the model performance can be effectively improved after

Model	AAPD		Reuters-21578		RMSC	
	Micro-F1(+,%)	HL(-)	Micro-F1(+,%)	HL(-)	Micro-F1(+)	HL(-)
BR	64.6	0.0316	87.8	0.0316	41.8	0.083
CC	65.4	0.0306	87.9	0.0306	44.3	0.107
CNN	65	0.0264	86.3	0.0287	59.1	0.0702
CNN-RNN	66.4	0.0278	85.5	0.0282	-	-
HAN	70.81	0.0236	-	-	66.75	0.059
HAN+LG	71.19	0.0235	-	-	68.21	0.058
SGM	71	0.0245	-	-	-	-
Bert	72.57	0.0226	89.09	0.0029	72.21	0.0507
Bert+MLM	72.23	0.0231	89.4	0.00274	73.18	0.0499
MEGNET	69.6	0.0252	89.9	0.0252	-	-
LW-LSTM+PT	71.21	0.0239	-	-	67.04	0.0583
LW-LSTM+FT	71.31	0.0241	-	-	72.18	0.0537
LM	73.73	0.0228	89.6	0.0028	72.81	0.0499
LM+MLM	73.61	0.0224	89.97	0.00264	73.8	0.0517

Table 3: Comparison with other baseline methods. The bold represents the optimal results.

transforming MTC into template populated tasks. We also notice that adding MLM can further improve the performance of ML and Bert, which illustrates the effectiveness of joint training.

We explain this phenomenon from the essence of Bert. When we input the label and the original sentence together into Bert, it is equivalent to constructing the context for the label, and self-attention can sensitively capture such context relations that do not exist in the original sentence. In this way, we introduce the association among labels which can increase the model’s ability to understand the label context. In addition, since Bert is essentially an MLM, allowing Bert to continue learning the mask for different downstream tasks can improve its performance.

We also compared LM-MTC with some other commonly used baseline models. The results are shown in Table 3. We notice that approaches with labels relationships is often superior to methods that do not consider the relationship between labels. For example, CC is better than BR, and HAN+LG is better than HAN. This demonstrates the necessary of label relevance. We also notice that our method outperforms all baselines, because our approach takes advantage of Bert’s powerful context-capture capabilities to capture semantic

associations between labels. This approach is more efficient than some methods of explicit label association modeling, such as graph-based or chain-based methods.

5.6. Analysis

5.6.1. Attention Visualization

The middle layer of Bert has been shown to adequately capture semantic relationships between words [10]. In LM-MTC, each potential label can be treated as a word, so we verify how the LM-MTC captures the correlation of labels by visualizing the attention of each layer.

Figure 3(a) shows the Spearman correlation between different labels of AAPD testset. To facilitate the observation, we chose the top11 labels with high relevance. Similar to Figure 1, test labels have a similar correlation distribution with train dataset.

After that, we take the attention output parameters of different Bert layers. We average all the attention heads and select the attention scores between all the label pairs. We add up all the batch data to get a global score matrix of attention on testset. We select the attention matrix of the 2nd layer (Figure 3(b)), the 6th layer (Figure 3(c)) and the last layer (Figure 3(d)) for visualization. Bert in the shallow layer learns some rough information, the 6th layer pays more attention to the local correlation, and the attention of the last layer is closer to the original label correlation distribution. This shows that deep Bert can capture the correlation between labels, which also provides a valid explanation for the advantages of LM-MTC.

5.6.2. MLM Loss Ratio

For multi-tasking learning, different task Loss should be of similar magnitude [36]. According to the changes in MTC and MLM loss as shown in Figure 4, we select different λ to investigate the effect of MLM task weight on model performance.

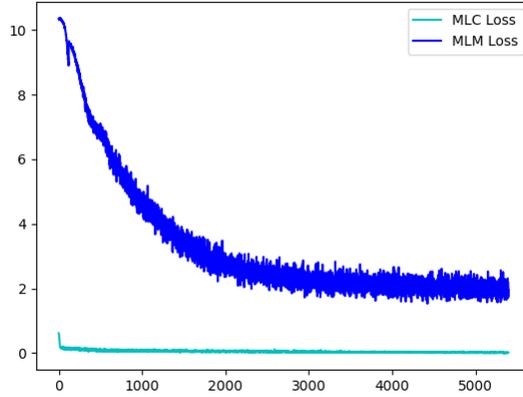


Figure 4: Changes of every 10 batch loss function during training.

As is shown in Figure 5(b), for two different datasets, RMSC and Emotion, the best performance occurs at $\lambda = 0.05$, λ is too large or too small will have a negative impact on the model. Therefore, So to keep the MTC task dominant, and to make the MLM task have enough impact, we set $\lambda = 0.05$ for all datasets.

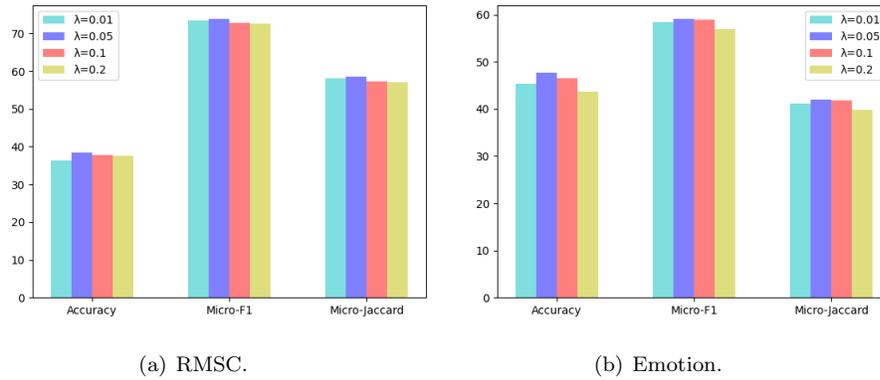


Figure 5: Model performance under different λ .

	RMSC		Emotion	
	Diff	Same	Diff	Same
Accuracy	38.39	21.67	47.63	37.31
Micro-F1	73.8	60.63	59.07	51.38
Micro-Jaccard	58.47	43.5	41.92	34.57
HL	0.0517	0.065	0.0304	0.0323

Table 4: The influence of different mask strategies on the results.

5.6.3. Different Mask Strategies

In fact, different templates may have different effects on the same task [25]. As a comparison, we assign the same token to each label in different positions, that is, we modify the template in Eq 2 as follows:

$$\begin{aligned}
& [LS][YES][LE] \\
& [LS][NO][LE] \\
& [LS][MASK][LE]
\end{aligned} \tag{8}$$

We repeat the experiment on RMSC and Emotion, and the results are shown in Table 4. When all labels use the same templates, there is a significant performance drop. This means that different label templates can better provide identification information for Bert, but same template does not.

6. Conclusion

In this paper, we propose LM-MTC model for multi-label text classification. We build prefix templates for multiple labels, transform MTC into cloze task, and combine training with MLM to improve the performance of the model under a variety of evaluation metrics. Further, we explain the ability of LM-MTC to capture the potential association between labels and LM-MTC can perform well in tests against multiple types of datasets. In future work, we will explore the ability of different templates to combine with different pre-trained language models.

7. Acknowledgements

This work was supported by National Natural Science Foundation of China (NSFC), "From Learning Outcome to Proactive Learning: Towards a Human-centered AI Based Approach to Intervention on Learning Motivation" (No.62077027).

References

- [1] E. Cambria, D. Olsher, D. Rajagopal, Senticnet 3: A common and common-sense knowledge base for cognition-driven sentiment analysis, AAAI (2014) 1515–1521.
- [2] Z. Yang, D. Yang, C. Dyer, X. He, J. A. Smola, H. E. Hovy, Hierarchical attention networks for document classification, HLT-NAACL (2016) 1480–1489.
- [3] S. Gopal, Y. Yang, Multilabel classification with meta-level features, SIGIR (2010) 315–322.
- [4] I. Katakis, I. Vlahavas, G. Tsoumakas, Multilabel text classification for automated tag suggestion, european conference on principles of data mining and knowledge discovery.
- [5] L. Xiao, X. Huang, B. Chen, L. Jing, Label-specific document representation for multi-label text classification, EMNLP/IJCNLP (1) (2019) 466–475.
- [6] R. M. Boutell, J. Luo, X. Shen, M. C. Brown, Learning multi-label scene classification, Pattern Recognition (2004) 1757–1771.
- [7] J. Liu, W.-C. Chang, Y. Wu, Y. Yang, Deep learning for extreme multi-label text classification, SIGIR (2017) 115–124.
- [8] L. Xiao, X. Zhang, C. Huang, M. Song, L. Jing, Does head label help for long-tailed multi-label text classification, AAAI 2021.

- [9] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, north american chapter of the association for computational linguistics.
- [10] G. Jawahar, B. Sagot, D. Seddah, What does bert learn about the structure of language, *ACL* (1) (2019) 3651–3657.
- [11] T. Schick, H. Schütze, It’s not just size that matters: Small language models are also few-shot learners, *arxiv*.
- [12] T. Schick, H. Schütze, Exploiting cloze-questions for few-shot text classification and natural language inference, *EACL* (2021) 255–269.
- [13] E. Hüllermeier, J. Fürnkranz, W. Cheng, K. Brinker, Label ranking by learning pairwise preferences, *Artif. Intell.* (2008) 1897–1916.
- [14] J. Read, B. Pfahringer, G. Holmes, E. Frank, Classifier chains for multi-label classification, *Machine Learning* (2011) 333–359.
- [15] G. Tsoumakas, I. Vlahavas, Random k-labelsets: An ensemble method for multilabel classification, *ECML* (2007) 406–417.
- [16] G. Chen, D. Ye, Z. Xing, J. Chen, E. Cambria, Ensemble application of convolutional and recurrent neural networks for multi-label text categorization, *IJCNN* (2017) 2377–2383.
- [17] P. Yang, X. Sun, W. Li, S. Ma, W. Wu, H. Wang, Sgm: Sequence generation model for multi-label classification, *COLING* (2018) 3915–3926.
- [18] Z. Barutcuoglu, E. R. Schapire, G. O. Troyanskaya, Hierarchical multi-label prediction of gene function, *Bioinformatics* (2006) 830–836.
- [19] M.-L. Zhang, K. Zhang, Multi-label learning by exploiting label dependency, *KDD* (2010) 999–1008.
- [20] S. Wang, J. Wang, Z. Wang, Q. Ji, Multiple emotion tagging for multimedia data by exploiting high-order dependencies among emotions, *IEEE Trans. Multimedia* (2015) 2185–2197.

- [21] S. Wang, G. Peng, Z. Zheng, Capturing joint label distribution for multi-label classification through adversarial learning, *IEEE Trans. Knowl. Data Eng.* (2020) 2310–2321.
- [22] F. Scarselli, M. Gori, C. A. Tsoi, M. Hagenbuchner, G. Monfardini, The graph neural network model, *IEEE Transactions on Neural Networks* (2009) 61–80.
- [23] L. Xiao, X. Huang, B. Chen, L. Jing, Label-specific document representation for multi-label text classification, *EMNLP/IJCNLP* (1) (2019) 466–475.
- [24] P. Ankit, S. Muru, S. Malaikannan, Multi-label text classification using attention-based graph neural network, *ICAART: PROCEEDINGS OF THE 12TH INTERNATIONAL CONFERENCE ON AGENTS AND ARTIFICIAL INTELLIGENCE, VOL 2* (2020) 494–505.
- [25] T. Shin, Y. Razeghi, L. L. R. IV, E. Wallace, S. Singh, Autoprompt: Eliciting knowledge from language models with automatically generated prompts, *empirical methods in natural language processing* (2020) 4222–4235.
- [26] T. B. B., M. Benjamin, R. Nick, S. Melanie, K. Jared, D. Prafulla, N. Arvind, S. Pranav, S. Girish, A. Amanda, A. Sandhini, H.-V. Ariel, K. Gretchen, H. Tom, C. Rewon, R. Aditya, D. Z. M., W. Jeffrey, W. Clemens, H. Christopher, C. Mark, S. Eric, L. Mateusz, G. Scott, C. Benjamin, C. Jack, B. Christopher, M. Sam, R. Alec, S. Ilya, A. Dario, Language models are few-shot learners, *NIPS 2020*.
- [27] X. Liu, Y. Zheng, Z. Du, M. Ding, Y. Qian, Z. Yang, J. Tang, Gpt understands, too.
- [28] F. Debole, F. Sebastiani, An analysis of the relative hardness of reuters-21578 subsets: Research articles, *Journal of the American Society for Information Science and Technology* (2005) 584–596.

- [29] G. Zhao, J. Xu, Q. Zeng, X. Ren, Review-driven multi-label music style classification by exploiting style correlations, arXiv: Computation and Language.
- [30] D. Dorottya, M.-A. Dana, K. Jeongwoo, C. Alan, N. Gaurav, R. Sujith, Goemotions: A dataset of fine-grained emotions, ACL (2020) 4040–4054.
- [31] P. Baldi, S. Brunak, Y. Chauvin, A. C. Andersen, H. Nielsen, Assessing the accuracy of prediction algorithms for classification: an overview, BIOINFORMATICS (2000) 412–424.
- [32] E. R. Schapire, Y. Singer, Improved boosting algorithms using confidence-rated predictions, Machine Learning (1998) 80–91.
- [33] Y. Kim, Convolutional neural networks for sentence classification, EMNLP (2014) 1746–1751.
- [34] H. Liu, C. Yuan, X. Wang, Label-wise document pre-training for multi-label text classification, international conference natural language processing (2020) 641–653.
- [35] I. Loshchilov, F. Hutter, Fixing weight decay regularization in adam, arXiv: Learning.
- [36] Z. Chen, V. Badrinarayanan, C.-Y. Lee, A. Rabinovich, Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks, international conference on machine learning.