

Label Distribution Learning-Enhanced Dual-KNN for Text Classification

Bo Yuan^{*} Yulin Chen^{*} Zhen Tan[†] Wang Jinyan^{*} Huan Liu[†] Yin Zhang^{*‡}

Abstract

Many text classification methods usually introduce external information (e.g., label descriptions and knowledge bases) to improve the classification performance. Compared to external information, some internal information generated by the model itself during training, like text embeddings and predicted label probability distributions, are exploited poorly when predicting the outcomes of some texts. In this paper, we focus on leveraging this internal information, proposing a dual k nearest neighbor (DkNN) framework with two k NN modules, to retrieve several neighbors from the training set and augment the distribution of labels. For the k NN module, it is easily confused and may cause incorrect predictions when retrieving some nearest neighbors from noisy datasets (datasets with labeling errors) or similar datasets (datasets with similar labels). To address this issue, we also introduce a label distribution learning module that can learn label similarity, and generate a better label distribution to help models distinguish texts more effectively. This module eases model overfitting and improves final classification performance, hence enhancing the quality of the retrieved neighbors by k NN modules during inference. Extensive experiments on the benchmark datasets verify the effectiveness of our method.

Keywords- text classification, label distribution learning, k nearest neighbor, robust learning

1 Introduction.

Text classification is a fundamental task [1, 2] in NLP (natural language processing), in which a given text is classified into one or more of n categories. Text classification can be used for a variety of purposes, including relationship extraction [3], tag recommendation [4], and so on. Generally speaking, most text classification methods mainly consist of three steps: firstly feed the original text into a classification model to obtain text embedding; then feed the text embedding into a linear layer to obtain label prediction probability; finally, calculate loss with the distribution of predictions and real labels. However, such a standard process has its problems during inference: (1) The predicted results depend on the linear layer added directly on top of the text encoder, which may face the overfitting problem as pointed out by [5]. (2) Only the trained parameters are used, and other forms of information that can be directly obtained from historical training instances are

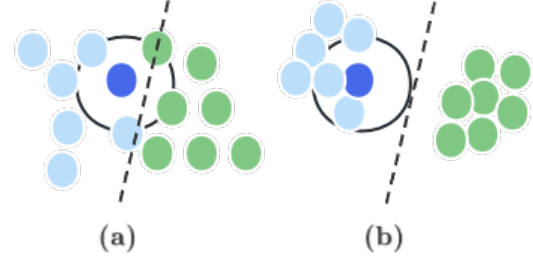


Figure 1: We use dots of different colors to denote different classes. (a) Previous work based on k NN may retrieve neighbors (the dots in the black circle) belonging to other classes (green dots) when the target text (deep blue dots) in similar datasets is relatively hard to distinguish (classification boundaries are close). (b) Our proposed label distribution learning can improve the performance of models (different class clusters get tighter and away from classification boundaries), thus the quality of retrieved neighbors in k NN is enhanced (the dots in the black circle are all blue).

ignored. To tackle these problems, it is intuitive to save more historical information during training, then make predictions based on them and not solely consider the ultimate output probability of the model.

Recently, retrieval-augmented methods have attracted lots of attention in the natural language processing field, such as language modeling [6], named entity recognition [7], multi-label text classification [8]. However, these methods retrieve neighboring instances based on text embedding or token embedding, and do not fully consider other information in the intermediate layers of the model. Motivated by these works, we add the k nearest neighbor mechanism based on the text embedding representations and predicted label probability representations to make predictions, instead of only using the ultimate output probability to make predictions.

Specifically, we propose a dual k NN (DkNN) framework for text classification tasks, which first constructs two representation stores to explicitly memorize the output representations of the model’s middle layers during training and then retrieves k nearest neighbors from

^{*}Zhejiang University. byuan, yulinchen, wangjy77, yinzh@zju.edu.cn

[†]Arizona State University. ztan36, huanliu@asu.edu

[‡]Corresponding Author

these cached representations during inference. When retrieving, we apply two k nearest neighbor modules, namely text- k NN and pro- k NN respectively. Intuitively, the retrieved neighbors may include noise (belongs to other classes) when the target text is relatively hard to distinguish (e.g., the texts with similar labels in the representation store have semantic similarity, or the relevant context is not enough in the representation store). Furthermore, we also find that the k NN mechanism is very sensitive to noise. Intuitively, as in the left part of Figure 1, the wrong retrieved results will make final prediction results poor robustness and generalization.

To address this problem, we propose to learn the label distribution based on label similarity and utilize contrastive learning to make the label distribution representation more distinctive. By doing this, both the quality of the retrieved neighbors and the classification performance of the original model are enhanced. We conducted extensive experiments to evaluate the effectiveness and robustness of our method. Experimental results show that our method can: (1) improve the classification performance of baseline models consistently on five benchmark datasets; (2) improve the defense ability against noise attacks significantly. Our contributions are summarized as follows:

- We propose a dual k NN (D k NN) framework for text classification tasks that utilize the representations (text embedding and predicted label probability) from training datasets in the inference stage.
- We introduce a label distribution learning module as an effective enhancement component for D k NN in the retrieval process. In addition, this module can also improve the performance of base classification models.
- Extensive experiments show that our proposed method produces significant improvements in baseline models. Furthermore, we conduct additional experiments to confirm the robustness of our method against noise attacks.

2 Related Work.

2.1 Text Classification. There have been many recent works that begin to exploit external information in text classification tasks, not only focusing on learning text representation. [9] treats conceptual knowledge as a type of knowledge and retrieves knowledge from an external knowledge source to enhance the semantic representation of short texts. [10] utilizes the statistical features of the corpus, such as word frequency and distribution over labels, to improve classification performance. [11] incorporates human rationales into attention-based

text classification models to improve the explainability of classification results. Despite all existing works, few studies pay attention to the internal information generated by the model during training. These data are based on referable historical instances, which can also aid the model in making predictions.

2.2 Nearest Neighbor Methods in NLP. The k NN (k nearest neighbors) classifier is a common basic machine learning method [12], and its basic idea is to determine the category of samples according to the surrounding limited adjacent samples (i.e., k nearest neighbors) [13]. Recent success on various NLP tasks has shown the effectiveness of k NN mechanism in improving the quality of NLP models. [14] extends a pre-trained neural language model by linearly interpolating with k NN mechanism. Based on k NN-MT proposed by [15], [16] introduces adaptive k NN-MT to determine the choice of k regarding each target token dynamically. However, we found that the k NN algorithm will retrieve some incorrect noise results on some datasets. [17] also realizes that the k NN method tends to introduce noise in simple applications, and they both introduce contrastive learning based on text representations to construct positive and negative samples. However, our method only focuses on label representations, which are different from them. These previous works utilize k NN algorithm based on the text embedding or the contextualized word embedding. In our paper, we will introduce another novelty k nearest neighbor (k NN) module based on the predicted label probability to consider more information generated by models during training.

2.3 Label Distribution Learning. Label Distribution Learning (LDL) [18] is a novel machine learning paradigm for applications where the overall distribution of labels matters. A label distribution covers a certain number of labels, representing the degree to which each label describes the instance. [18] also suggests LDL could be helpful for such scenarios: some labels are correlated with other labels. Based on it, [19] utilizes LDL to improve the performance of multi-label text classification tasks. However, for single-label text classification tasks where they only have a unique label for each sample, the true label distribution is hard to obtain. Thus, we explore utilizing label similarity to learn label distribution, which can make models more generalized.

3 Method.

The overview of our proposed method is shown in Figure 2, and this section describes our method in depth from its two primary parts. Firstly, we design a dual k nearest neighbor (D k NN) framework including two k NN

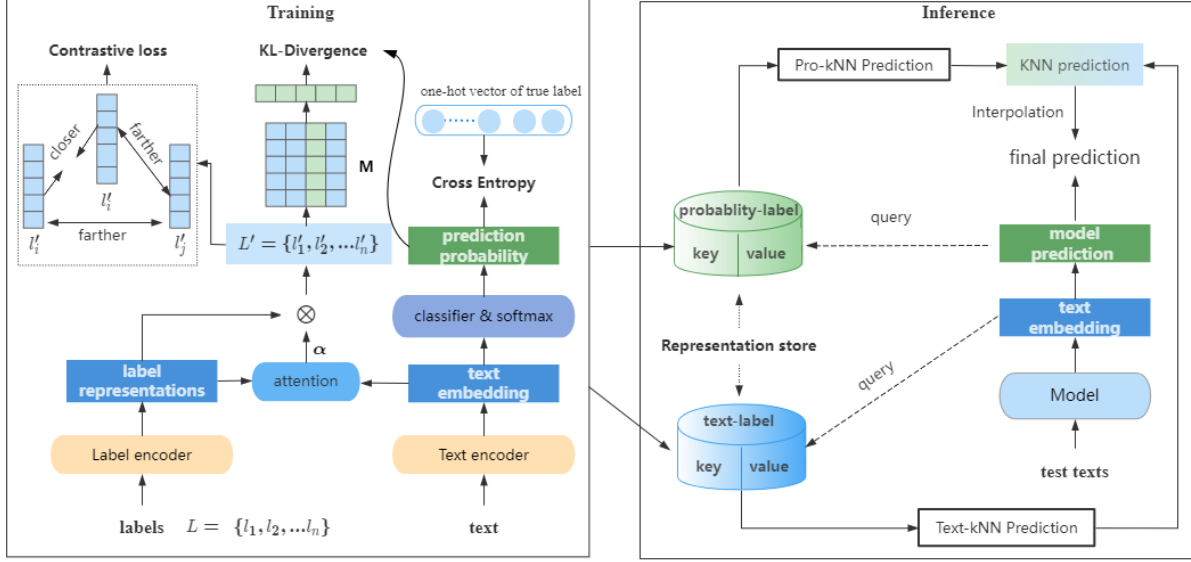


Figure 2: The overall framework of our proposed method. The representation store contains a set of representation-label pairs, which are extracted from the hidden states of the label distribution learning enhanced model. During inference, when querying k nearest neighbors from the representation store according to the similarity distance, the similarity distances are converted to k NN prediction distribution. Interpolating the k NN distribution with the vanilla model prediction distribution, we get the final distribution.

modules (text- k NN and pro- k NN) for the text classification model (single-label classification). Secondly, we propose a label distribution learning module based on label similarity to help the model distinguish texts more effectively. By training the model with this module, the Dk NN can also be enhanced.

3.1 Problem Formulation. Given an input text $\mathbf{x} = \{x_1, \dots, x_n\}$ with length n , $x_i (i = 1, \dots, n)$ denotes the i -th word token within this text. $D = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ is training datasets with N texts, where $(\mathbf{x}_i, \mathbf{y}_i)$ is a text sequence and corresponding one-hot label vector. For each text \mathbf{x} , most text classification models have a common paradigm to process it: use a deep neural network to obtain text embedding \mathbf{h} ; then feed it to a simple linear layer with softmax activation function to obtain the predicted label distribution \mathbf{p} ; finally calculate the cross entropy loss between \mathbf{p} and the true label vector \mathbf{y} .

3.2 k Nearest Neighbor in Text Classification. To comprehensively obtain the information generated by each part of the models during the inference stage, we propose the Dk NN framework, which includes the following steps: creating the key-value representation store, making two k NN module predictions based on it, and combining these two predicted distributions to get the final k NN prediction.

3.2.1 Representation store. The representation store consists of a set of key-value pairs, and we construct two representation stores to explicitly memorize texts. To be specific, for each text-label pair $(\mathbf{x}_i, \mathbf{y}_i)$ from training datasets D , we extract the text embedding representation \mathbf{h}_i and the predicted label probability representation \mathbf{p}_i by a single forward pass over \mathbf{x}_i . Then the text representation store can be constructed as follows: $S_{tr} = \{\mathcal{K}, \mathcal{V}\} = \{(\mathbf{h}_i, \mathbf{y}_i)\}_{i=1}^N$, and the predicted label probability representation store can be constructed as: $S_{pr} = \{\mathcal{K}, \mathcal{V}\} = \{(\mathbf{p}_i, \mathbf{y}_i)\}_{i=1}^N$, where \mathcal{K} is the key set and \mathcal{V} is the corresponding value set.

3.2.2 Inference. In the inference stage, given an input text \mathbf{x}_i , our model generates the predicted label probability distribution \mathbf{p}_i and the text embedding \mathbf{h}_i . For \mathbf{h}_i , the text- k NN queries the text representation store S_{tr} through the L^2 Euclidean distance $d_{L^2}(\cdot, \cdot)$ to obtain the k nearest neighbors $\mathcal{N}_{tr} = \{(k_i, v_i)\}_{i=1}^k \subseteq S_{tr}$. For \mathbf{p}_i , it is utilized by pro- k NN to query k nearest neighbors $\mathcal{N}_{pr} = \{(k_i, v_i)\}_{i=1}^k \subseteq S_{pr}$ from predicted label probability representation store S_{pr} with Kullback-Leibler divergence (KL-divergence) $d_{KL}(\cdot, \cdot)$ [20] as distance measure. Then, these two k NN predicted distributions over neighbors can be computed by applying a softmax to their negative distances and combining

multiple occurrences of the same label item.

$$(3.1) \quad \mathbf{p}_{text-knn}(\mathbf{y}_i|\mathbf{x}_i) \propto \sum_{(k_i, v_i) \in \mathcal{N}_{tr}} \mathbb{1}_{y=v_i} \exp(-d_{L^2}(\mathbf{k}_i, \mathbf{h}_i))$$

$$(3.2) \quad \mathbf{p}_{pro-knn}(\mathbf{y}_i|\mathbf{x}_i) \propto \sum_{(k_i, v_i) \in \mathcal{N}_{pr}} \mathbb{1}_{y=v_i} \exp(-d_{KL}(\mathbf{k}_i, \mathbf{p}_i))$$

Next, we improve high-confidence predictions while downgrading low-confidence ones by squaring and normalizing the two current predictions ($\mathbf{p}_{text-knn}$ and $\mathbf{p}_{pro-knn}$ can be viewed as \mathbf{p}_{ij}):

$$(3.3) \quad \mathbf{p}'_{ij} = \frac{\mathbf{f}_j}{\sum_{j'} \mathbf{f}_{j'}}, \mathbf{f}_j = \frac{\mathbf{p}_{ij}^2}{\sum_j \mathbf{p}_{ij}}$$

Finally, we interpolate the pure model prediction \mathbf{p}_i and final k NN prediction $\mathbf{p}_{knn}(\mathbf{y}_i|\mathbf{x}_i)$ with a tuned parameter λ to obtain final prediction results:

$$(3.4) \quad \mathbf{p}(\mathbf{y}_i|\mathbf{x}_i) = \lambda \mathbf{p}_{knn}(\mathbf{y}_i|\mathbf{x}_i) + (1 - \lambda) \mathbf{p}_i$$

$$(3.5) \quad \mathbf{p}_{knn}(\mathbf{y}_i|\mathbf{x}_i) = \frac{\mathbf{p}'_{pro-knn}(\mathbf{y}_i|\mathbf{x}_i) + \mathbf{p}'_{text-knn}(\mathbf{y}_i|\mathbf{x}_i)}{2}$$

where \mathbf{p}_{knn} is the average of the pro- k NN distribution $\mathbf{p}'_{pro-knn}$ and the text- k NN distribution $\mathbf{p}'_{text-knn}$.

3.3 Label Distribution Learning. In our experiments, we find that the traditional knn algorithm may cause prediction errors when the retrieved neighbors include noises, and cause over-fitting problems when the retrievals are most similar. The underlying reason is that original text classification models can not effectively distinguish semantic differences between texts, particularly for those different texts with semantically similar labels, thus their text embeddings are less discriminative and have a negative impact on the knn retrieval process. To address this problem, in the training stage, we propose a label distribution learning module based on label similarity to generate a distribution vector, which is a distinct representation with similarity simultaneously. In the following, we first define similarity measures among labels and then present how to equip the similarity measures to our training objective.

3.3.1 Label Similarity. Given texts and label names, we consider the similarities and differences semantically among various labels based on the texts. To utilize the semantic information of labels, we represent them by a label encoder, which is a deep neural network (DNN), to generate a trainable representation matrix $\mathbf{L} = [\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_c]$ with length c and dimension d , where

c is the number of label categories. For an input text \mathbf{x}_i , and its text embedding \mathbf{h}_i with dimension d , we measure the compatibility score for the entire label categories (i.e., the matching degree of texts and labels):

$$(3.6) \quad \boldsymbol{\alpha} = \text{softmax}(\mathbf{h}_i \mathbf{L})$$

where $\boldsymbol{\alpha} \in \mathbb{R}^c$. The label representation can be updated by weighting text-based attention score:

$$(3.7) \quad \mathbf{L}' = [\alpha_1 \mathbf{l}_1, \alpha_2 \mathbf{l}_2, \dots, \alpha_c \mathbf{l}_c]$$

However, focusing solely on the attention between texts and labels ignores the relationship between inter-classes.

3.3.2 Contrastive Learning with Label Similarity. Our next goal is to encourage the model to be aware of the relationship between inter-classes, which can help the model to better represent the text with different labels. Typically, contrastive learning methods aim to teach models to distinguish between target samples and other samples. They have been proven to be beneficial for various tasks and have caught our attention. So, we first propose a straightforward and efficient method for constructing a label similarity matrix to represent the relationship between different labels:

$$(3.8) \quad \mathbf{M} = \mathbf{L}' \mathbf{L}'^T$$

where $\mathbf{M} = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_c]$ is the matrix of size $c \times c$, each vector in \mathbf{M} represents the similarity value distribution between its corresponding label and all labels. Then, we introduce a contrastive objective \mathcal{L}_{CL} :

$$(3.9) \quad \mathcal{L}_{CL} = \frac{\sum_{i=1}^c \sum_{j=1, j \neq i}^c \max\{0, \rho - M_{i,i} + M_{i,j}\}}{c \times (c - 1)}$$

where $\rho \in [0, 1]$ is a pre-defined margin, $M_{i,j}$ is the value that denotes the similarity relationship between different labels and $M_{i,i}$ is the self-similarity value. Intuitively, the model learns to pull away the distances between representations of distinct labels and distinguish texts better by training with \mathcal{L}_{CL} .

3.3.3 Training Objective. For the ground-truth label, we first take out its corresponding column \mathbf{m} from \mathbf{M} and transfer the similarity value distribution to similarity probability distribution by $\mathbf{q} = \text{softmax}(\mathbf{m})$, then we utilize KL-divergence object \mathcal{L}_{KL} to measure the distribution difference of \mathbf{q} and \mathbf{p} . The \mathcal{L}_{KL} corresponds with Eq.3.2, which makes the model aware of the knn retrieval process. Denoting cross entropy loss as \mathcal{L}_{CE} , the overall training objective is defined as:

$$(3.10) \quad \mathcal{L} = \mathcal{L}_{CE} + \mathcal{L}_{KL} + \mathcal{L}_{CL}$$

Models	Accu. (%)				
	WOS-5736	20NG	AGNews	DBPedia	FDCNews
CNN [21]	65.15 \pm 0.0158	72.22 \pm 0.0046	89.73 \pm 0.0033	96.09 \pm 0.0011	80.78 \pm 0.0096
Ours	69.88 \pm 0.0164	77.12 \pm 0.0034	90.33 \pm 0.0013	96.92 \pm 0.0009	82.07 \pm 0.0056
LCM-CNN [22]	67.20 \pm 0.0037	76.44 \pm 0.0075	89.99 \pm 0.0023	96.96 \pm 0.0005	82.17 \pm 0.0039
Ours	69.20 \pm 0.0117	77.20 \pm 0.0046	90.33 \pm 0.0014	97.11 \pm 0.0011	82.39 \pm 0.0020
LSTM [23]	69.29 \pm 0.0160	69.84 \pm 0.0133	88.52 \pm 0.0040	93.44 \pm 0.0171	77.59 \pm 0.0073
Ours	71.71 \pm 0.0061	75.09 \pm 0.0102	89.27 \pm 0.0023	94.91 \pm 0.0018	80.44 \pm 0.0040
LCM-LSTM [22]	69.32 \pm 0.0157	74.00 \pm 0.0046	88.12 \pm 0.0049	94.09 \pm 0.0124	79.78 \pm 0.0096
Ours	71.25 \pm 0.0107	74.77 \pm 0.0102	89.10 \pm 0.0069	94.34 \pm 0.0098	80.64 \pm 0.0016
LCM-BERT [22]	78.33 \pm 0.0155	90.03 \pm 0.0031	91.30 \pm 0.0012	97.95 \pm 0.0009	96.15 \pm 0.0040
Ours	79.05 \pm 0.0142	90.38 \pm 0.0030	91.42 \pm 0.0004	98.32 \pm 0.0008	96.39 \pm 0.0064
ALBERT [24]	76.49 \pm 0.0223	88.27 \pm 0.0129	89.71 \pm 0.0014	96.74 \pm 0.0080	94.30 \pm 0.0029
Ours	77.27 \pm 0.0110	89.34 \pm 0.0053	92.04 \pm 0.0019	98.16 \pm 0.0026	96.15 \pm 0.0066
BERT [17]	77.79 \pm 0.0117	88.21 \pm 0.0005	90.24 \pm 0.0014	97.49 \pm 0.0011	96.63 \pm 0.0016
Ours	78.95 \pm 0.0113	90.53 \pm 0.0016	91.35 \pm 0.0014	98.31 \pm 0.0009	97.59 \pm 0.0028
Other baselines based on KNN or contrastive learning					
Base	77.79 \pm 0.0117	88.21 \pm 0.0005	90.24 \pm 0.0014	97.49 \pm 0.0011	96.63 \pm 0.0016
SCL [25]	77.90 \pm 0.6077	88.45 \pm 0.2592	90.28 \pm 0.0809	97.58 \pm 0.0722	96.65 \pm 0.1841
SCL-MoCo [26]	78.01 \pm 0.3389	88.46 \pm 0.1750	90.32 \pm 0.1153	97.64 \pm 0.0861	96.97 \pm 0.2922
KNN-BERT [17]	78.11 \pm 0.0093	88.97 \pm 0.0035	90.38 \pm 0.0013	97.71 \pm 0.0014	97.42 \pm 0.0022
Ours	78.95 \pm 0.0113	90.53 \pm 0.0016	91.35 \pm 0.0014	98.31 \pm 0.0009	97.59 \pm 0.0028

Table 1: Test accuracy on different datasets. We run all models and report their mean \pm standard deviation. The results with outstanding improvement over the base model are bolded, and these values indicate statistically significantly better (based on the student t -test, $p < 0.05$) performances across the board.

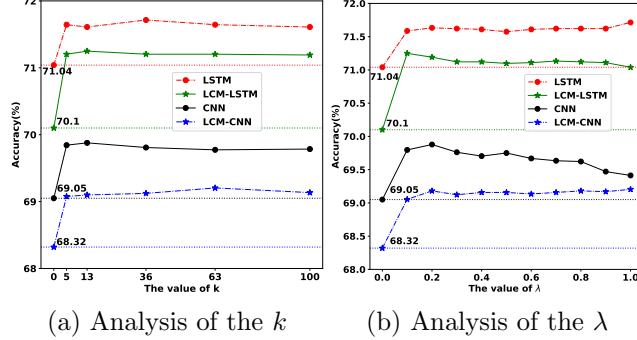


Figure 3: Hyperparameter settings of the DkNN.

4 Experiment.

4.1 Baselines. We first adopt the following models as our baselines and apply our method to all of them: CNN [21], LSTM [23], BERT [27], ALBERT [24], and LCM [22]. LCM is an enhancement component to current popular text classification models, which is similar to our proposed label distribution learning module. For LCM, we follow [22] to implement it based on LSTM, CNN, and BERT. Then, we compare our k NN-based approach with KNN-BERT [17]. Specifically, KNN-BERT utilizes the traditional k NN classifier in pretrained model fine-tuning and trains clustered representations based on a supervised contrastive learn-

ing framework to improve the performance of pretrained model fine-tuning. Further, we follow KNN-BERT to select a few related works SCL [25] and SCL-MoCo [26] as other baselines, and use the same backbone (BERT) for these baseline methods. For a fair comparison, we train them using the same parameters with ours.

4.2 Datasets and Evaluation Metric. We evaluate the effectiveness and robustness of our proposed model on the following several text classification datasets: 20NG dataset¹ [28] is an English news dataset categorized into 20 different categories, in total, it contains 18846 documents; AGNews dataset² [29] consists of news articles, which contains 127600 documents pertaining to the four largest classes; DBPedia³ [29] is an ontology classification dataset containing 630000 documents over 14 categories; FDCNews dataset⁴ is a Chinese dataset constructed by Fudan University, which has 9833 documents divided into 20 categories; WOS-5736⁵ [30] is an academic paper dataset, which includes 5736 documents classified into 3 coarse-grained cate-

¹<http://qwone.com/~jason/20Newsgroups/>

²<http://groups.di.unipi.it/gulli/AGcorpusofnewsarticles.html>

³<http://dbpedia.org>

⁴<http://www.nlpir.org>

⁵<http://archive.ics.uci.edu/index.php>

gories and 11 fine-grained categories. We conduct experiments using fine-grained categories. We follow the data processing of [22], then train the model with a train set and evaluate on development set after every epoch.

4.3 Parameter Settings. In our experiments, we set the batch size to 128. For non-pretrained models, we set the maximum sequence lengths as 100 tokens and the embedding size as 64. For CNN, we use two convolution blocks, the number of filters for each block is 64 and 32 respectively, and the filter sizes are both 3. For BERT and ALBERT, the maximum sequence length is 128 tokens and the hidden size is 768. For LCM, we use the default parameter settings in the original papers. All models are trained with the Adam Optimizer [31] on GTX 2080 Ti and RTX 3090.

4.4 Experiment Results. In our experiments, we split each dataset into different train and test sets 5 times with a 7:3 splitting ratio. Then we evaluate all models 5 times on these various train test splitting and report the average performance. We use accuracy as an evaluation metric following most previous works.

4.4.1 Overall Performance. The results of our method compared to other methods are listed in Table 1. The LCM-X means that the LCM model uses X as a basic predictor. As shown in Table 1, our method can improve the performance of baseline models on all datasets. For instance, our method improves the accuracy of CNN by 4.73% on WOS-5736 and 4.9% on 20NG. The performance improvements on LSTM are also quite obvious with 5.25% on 20NG and 2.85% on FDCNews. For pre-trained models, we also further compare our method with other baselines based on k NN and contrastive learning. As seen, our method performs better than them, which further verifies the effectiveness and applicability of our proposed method.

4.4.2 Ablation Test. As mentioned above, our method mainly includes the following parts: a label distribution learning module (denoted as LL) and a dual k nearest neighbor framework (denoted as Dk NN). We perform an ablation test to demonstrate the effectiveness of each part. As shown in Table 2, Dk NN and LL can both generally enhance the performance of the baseline models. Even though the retrieval effect of Dk NN remains consistent on some baseline models, when combined with our label distribution learning module, the improvements brought by the Dk NN have significantly improved. This demonstrates that the Dk NN is successfully enhanced by our LL.

5 Discussion.

5.1 Analysis for Dk NN. The k NN algorithm in the existing work makes predictions based on cached text embeddings and corresponding target tokens. We emphasize the importance of calculating and comparing the predicted label probability distribution distance between the samples to be classified and each historical sample, thus we introduce Dk NN consisting of text- k NN (denoted as t- k NN) and pro- k NN (denoted as p- k NN). Here, we conduct a parameter analysis of Dk NN on the WOS-5736 dataset, and an ablation study to explore the effectiveness of Dk NN.

5.1.1 Hyperparameters in Dk NN In our proposed Dk NN, k is a significant hyperparameter that returns top- k nearest neighbors for each query. As shown in Figure 3(a), the hyperperformance with Dk NN is always better than that when using only the model output ($k=0$), which further verifies the effectiveness of Dk NN to utilize the information from the most similar historical samples. Additionally, the performance initially increases quickly, but as the k increases, it tends to plateau or even decline. That can be summed up that a large k is not necessary and a small k can still have a considerable but not infallible result, which can save the computational cost of retrieval. Moreover, we use a hyperparameter λ to interpolate between the distribution from Dk NN search over the dataset and the pure model prediction distribution during inference. Figure 3(b) demonstrates the trend of model performance with λ . The general trend is similar to that of k , which further demonstrates the performance is poor when relying simply on the model prediction ($\lambda=0$). We also notice that LSTM and LCM-CNN can obtain the greatest results when only using the model predictions ($\lambda=1$).

5.1.2 Effectiveness of Dk NN. In our paper, we use text- k NN to measure text embedding similarities and pro- k NN to calculate the distance between probability distributions. To verify the effectiveness of these two modules, we further perform ablation experiments on the WOS-5736 dataset. From Table 3, we observe that pro- k NN (denoted as p- k NN) is as important as text- k NN (denoted as t- k NN) to obtain a better performance. After removing p- k NN or t- k NN, the performances will drop to a large degree. In some experiments, the improvements brought by p- k NN outperform t- k NN, which further proves that retrieving neighbors according to the prediction probability is also an appropriate way.

5.2 Analysis for LL. In our method, the LL part generates the learned label distribution, which explores

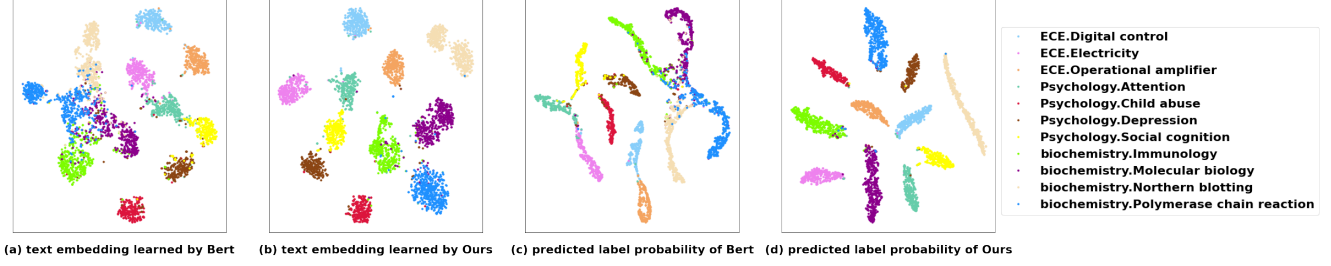


Figure 4: visualize the cached in the representation store by t-SNE tool.

Models	WOS-5736	20NG
CNN	65.15	72.22
CNN+DkNN	65.24	72.54
CNN+LL	68.54	76.59
CNN+LL+DkNN	69.88	77.12
LCM-CNN	67.20	76.44
LCM-CNN+DkNN	67.45	76.94
LCM-CNN+LL	68.32	76.57
LCM-CNN+LL+DkNN	69.20	77.20
LSTM	69.29	69.84
LSTM+DkNN	69.37	70.01
LSTM+LL	70.45	74.12
LSTM+LL+DkNN	71.71	75.09
LCM-LSTM	69.32	74.00
LCM-LSTM+DkNN	69.45	74.01
LCM-LSTM+LL	70.10	74.16
LCM-LSTM+LL+DkNN	71.25	74.77
BERT	77.79	88.21
BERT+DkNN	78.13	88.98
BERT+LL	78.45	90.13
BERT+LL+DkNN	78.95	90.53
LCM-BERT	78.33	90.03
LCM-BERT+DkNN	78.37	90.05
LCM-BERT+LL	78.93	90.23
LCM-BERT+LL+DkNN	79.05	90.38

Table 2: Accuracy (%) of the ablation tests for main components in our method. DkNN means the dual k nearest neighbor framework proposed in our method, and LL denotes the label distribution learning module.

the similarity relationship among labels by taking into account the complicated relationship between the input texts and labels. Then we introduce this label distribution representation vector into the calculation of our training objective. There are two reasons for this operation: (1) The k Nearest Neighbor mechanism relies heavily on the training sample set and has poor fault tolerance to the training data. If some samples in the training datasets are mislabeled and just right next to the text to be classified, they will directly lead to the inaccuracy of the predicted result. In this case, solely using cross entropy with one-hot label representation during training will make models more susceptible to

Models	DkNN	w/o p-kNN	w/o t-kNN
LSTM	69.37	69.27	69.04
LSTM+LL	71.71	71.66	71.28
LCM-LSTM	69.45	69.41	68.90
LCM-LSTM+LL	71.25	71.14	71.01
CNN	69.37	69.27	69.04
CNN+LL	69.88	69.03	69.06
LCM-CNN	67.45	67.37	67.40
LCM-CNN+LL	69.20	69.06	68.93
BERT	78.13	78.00	78.06
BERT+LL	78.95	78.58	78.66
LCM-BERT	78.37	78.16	78.08
LCM-BERT+LL	79.05	78.92	78.93

Table 3: Results of accuracy (%) on the WOS-5736 datasets. We experiment ablation tests for DkNN.

the effects of noisy data (mislabeled data). However, the distribution generated by LL is based on similarity among different labels, by measuring the distance between this distribution and the predicted label probability distribution with KL-divergence loss function (denoted as \mathcal{L}_{KL}), the value on the index of the wrong label will be alleviated. (2) During the inference stage, we also use KL-divergence as the distance measure to query k nearest neighbors from the predicted label probability representation store, which exactly corresponds to the \mathcal{L}_{KL} in the overall training objective. Intuitively, by using LL, the predicted label probability cached in the representation store will be better, and the quality of the retrieved neighbors will also be improved.

5.2.1 Impact of Training Objective. We further conduct experiments to analyze the impact of our training objective. As shown in Table 6, we equip \mathcal{L}_{KL} and contrastive loss objective \mathcal{L}_{CL} separately with the original cross entropy loss function (denoted as \mathcal{L}_{CE}). From the results, we can see both \mathcal{L}_{KL} objective and \mathcal{L}_{CL} objective can effectively improve the performance.

5.2.2 Representation Store Visualizations. To explore whether the text embeddings and the predicted

	biochemistry	ECE	Psychology
Labels	biochemistry.Polymerase chain reaction biochemistry.Molecular biology biochemistry.Northern blotting biochemistry.Immunology	ECE.Electricity ECE.Digital control ECE.Operational amplifier	Psychology.Attention Psychology.Child abuse Psychology.Social cognition Psychology.Depression

Table 4: Labels of the WOS-5736 dataset

Models	0%	3%	6%	30%	50%
Bert	77.79	77.11	75.78	69.74	57.86
+DkNN	78.13	77.26	75.99	70.35	59.00
Increase	0.34	0.15	0.21	0.61	1.14
Bert+LL	78.45	78.16	77.27	69.76	58.36
+DkNN	78.95	78.42	77.57	70.83	59.50
Increase	0.5	0.26	0.3	1.07	1.26
LCM-Bert	78.33	77.44	76.27	70.85	59.85
+DkNN	78.37	77.56	76.47	71.11	61.08
Increase	0.04	0.12	0.2	0.26	1.23
LCM-Bert+LL	78.93	78.37	77.26	69.94	58.06
+DkNN	79.05	78.55	77.56	70.75	59.31
Increase	0.12	0.18	0.3	0.81	1.25

Table 5: The experimental results (accuracy, %) on the noisy datasets, where the percentage (0%, 3%, 6%, 30%, 50%) represents the proportion (noisy ratio) of randomly mislabeled samples.

Bert	WOS-5736	20NG
\mathcal{L}_{CE}	77.79	89.51
$\mathcal{L}_{CE} + \mathcal{L}_{KL}$	78.23	90.10
$\mathcal{L}_{CE} + \mathcal{L}_{CL}$	78.13	89.86
$\mathcal{L}_{CE} + \mathcal{L}_{CL} + \mathcal{L}_{KL}$	78.45	90.13

Table 6: Ablation over the KL-divergence object and contrastive learning objective of the proposed LL module using the Bert model on WOS-5736 dataset and 20NG datasets.

label probability are actually better by our LL, we use the t-SNE tool [32] to visualize these representations cached in the representation store. Figure 4a and Figure 4b respectively show the visualization of text embeddings learned by the original Bert model and LL enhanced Bert model on the WOS-5735 datasets. We also plot the predicted label probability of the original Bert model and our model in Figure 4c and Figure 4d. As shown in Figure 4a and Figure 4c, we can find that the text embeddings and predicted label probability produced by the original model can be easily confused, especially those with similar labels. After using LL, the distance between different categories of texts becomes larger, which is shown in Figure 4b and Figure 4d. Even though some samples are still jumbled

in the incorrect cluster, the boundary is more distinct and separate from the other categories. In this case, the neighbors retrieved by DkNN during inference are more accurate. In Figure 4, we also notice that the representation of text embeddings is more horizontal and the predicted label probability representation is more vertical, which shows our DkNN can consider the distance among representations from more directions when retrieving in representation store.

5.3 Performance on Noisy Datasets. Current text classification models require a large number of human-labeled texts as training data. Label noise is unavoidable when collecting data. As we all know, noisy data will severely reduce classification performance. Besides, the kNN mechanism is also sensitive to noisy data that has labeling errors. In the WOS-5735 dataset, we found that there are 11 fine-grained labels but they belong to 3 coarse-grained categories. As shown in Table 4, the labels that belong to the same category are similar in some way and hard to distinguish, which makes labeling errors more likely to occur and makes models overfit. To verify the robustness of our method against noise attacks, we generated some noisy datasets from the WOS-5735 dataset with different percentages of noise, and the mislabeling samples are all drawn from the same category to make the noisy datasets closer to reality. Then we conduct experiments with Bert as the base model. The results are shown in Table 5, we can see that our LL module still enhances the kNN mechanism on the datasets with obvious label errors, and the increase of our method outperforms the improvement of the base model and LCM to a large degree. Even without DkNN, the LL module in our method still performs significantly better than LCM on noisy datasets. Meanwhile, this observation also proves that the proposed LL can make the model robust enough.

6 Conclusion.

In this paper, we propose a dual kNN framework (DkNN) with two kNN modules for current text classification models. Then we introduce a label distribution learning module based on label similarity with contrastive learning to improve the performance of original models directly, as well as indirectly improve DkNN re-

trieval quality. By conducting extensive experiments, we confirm the effectiveness of our method and analyze the performance improvement brought about by each module. Besides, our method has higher robustness in noisy environments. In the future, we would like to explore two avenues. Firstly, designing a better framework to improve retrieval efficiency and reduce the size of the representation store for $DkNN$. Secondly, we can generalize our method to other classification tasks.

Acknowledgments

This work is supported by the China Knowledge Centre for Engineering Sciences and Technology (CKCEST-2022-1-7).

References

- [1] Z. Tan, L. Cheng, S. Wang, Y. Bo, J. Li, H. Liu, Interpreting pretrained language models via concept bottlenecks, arXiv:2311.05014 (2023).
- [2] Y. Chen, B. Yuan, B. Liao, D. M. Gabbay, A self-explanatory contrastive logical knowledge learning method for sentiment analysis, *Knowl. Based Syst.* 278 (2023) 110863.
- [3] L. Lucas, D. Tomás, J. G. Rodríguez, Exploiting the relationship between visual and textual features in social networks for image classification with zero-shot deep learning, *CoRR abs/2107.03751* (2021).
- [4] K. Lei, Q. Fu, M. Yang, Y. Liang, Tag recommendation by text classification with attention-based capsule network, *Neurocomputing* 391 (2020) 65–73.
- [5] Y. Lin, Breaking the softmax bottleneck for sequential recommender systems with dropout and decoupling, *CoRR abs/2110.05409* (2021).
- [6] U. Khandelwal, O. Levy, D. Jurafsky, L. Zettlemoyer, M. Lewis, Generalization through memorization: Nearest neighbor language models, *CoRR abs/1911.00172* (2019).
- [7] S. Wang, X. Li, Y. Meng, T. Zhang, R. Ouyang, J. Li, G. Wang, knn-ner: Named entity recognition with nearest neighbor search, *CoRR abs/2203.17103* (2022).
- [8] R. Wang, X. Dai, et al., Contrastive learning-enhanced nearest neighbor mechanism for multi-label text classification, in: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2022, pp. 672–679.
- [9] J. Chen, Y. Hu, J. Liu, Y. Xiao, H. Jiang, Deep short text classification with knowledge powered attention, *AAAI Press*, 2019, pp. 6252–6259.
- [10] X. Li, Z. Li, H. Xie, Q. Li, Merging statistical feature via adaptive gate for improved text classification, 2021.
- [11] I. Arous, L. Dolamic, J. Yang, A. Bhardwaj, G. Cuccu, P. Cudré-Mauroux, MARTA: leveraging human rationales for explainable text classification, 2021.
- [12] P. Zhao, L. Lai, Efficient classification with adaptive KNN, in: *AAAI*, 2021, pp. 11007–11014.
- [13] Y. Wang, Z. Pan, J. Dong, A new two-layer nearest neighbor selection method for knn classifier, *Knowl. Based Syst.* 235 (2022) 107604.
- [14] U. Khandelwal, O. Levy, D. Jurafsky, L. Zettlemoyer, M. Lewis, Generalization through memorization: Nearest neighbor language models, in: *ICLR*, 2020.
- [15] U. Khandelwal, A. Fan, D. Jurafsky, L. Zettlemoyer, M. Lewis, Nearest neighbor machine translation, *CoRR abs/2010.00710* (2020).
- [16] X. Zheng, Z. Zhang, J. Guo, S. Huang, B. Chen, W. Luo, J. Chen, Adaptive nearest neighbor machine translation, arXiv preprint arXiv:2105.13022 (2021).
- [17] L. Li, D. Song, R. Ma, X. Qiu, X. Huang, KNN-BERT: fine-tuning pre-trained models with KNN classifier, *CoRR abs/2110.02523* (2021).
- [18] X. Geng, Label distribution learning, *IEEE Trans. Knowl. Data Eng.* 28 (7) (2016) 1734–1748.
- [19] F. Kang, R. Jin, R. Sukthankar, Correlated label propagation with application to multi-label learning, in: *CVPR*, 2006, pp. 1719–1726.
- [20] S. Kullback, R. A. Leibler, On information and sufficiency, *The annals of mathematical statistics* 22 (1) (1951) 79–86.
- [21] Y. Kim, Convolutional neural networks for sentence classification, in: *EMNLP, ACL*, 2014, pp. 1746–1751.
- [22] B. Guo, S. Han, X. Han, H. Huang, T. Lu, Label confusion learning to enhance text classification models, 2021, pp. 12929–12936.
- [23] P. Liu, X. Qiu, X. Huang, Recurrent neural network for text classification with multi-task learning, *IJCAI* (Jul 2016).
- [24] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, R. Soricut, ALBERT: A lite BERT for self-supervised learning of language representations, in: *ICLR*, 2020.
- [25] B. Gunel, J. Du, A. Conneau, V. Stoyanov, Supervised contrastive learning for pre-trained language model fine-tuning, *ICLR* (2021).
- [26] K. He, H. Fan, Y. Wu, S. Xie, R. Girshick, Momentum contrast for unsupervised visual representation learning, in: *CVPR*, 2020.
- [27] J. Devlin, M. Chang, K. Lee, K. Toutanova, BERT: pre-training of deep bidirectional transformers for language understanding, in: *Proceedings of the 2019 Conference of the North*, 2019, pp. 4171–4186.
- [28] K. Lang, Newsweeder: Learning to filter netnews, in: *Machine Learning Proceedings 1995*, Morgan Kaufmann, 1995, pp. 331–339.
- [29] X. Zhang, J. J. Zhao, Y. LeCun, Character-level convolutional networks for text classification, 2015, pp. 649–657.
- [30] K. Kowsari, D. E. Brown, M. Heidarysafa, K. J. Meimandi, M. S. Gerber, L. E. Barnes, Hdtex: Hierarchical deep learning for text classification, in: *ICMLA*, 2017, pp. 364–371.
- [31] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, 2015.
- [32] L. Van der Maaten, G. Hinton, Visualizing data using t-sne., *Journal of machine learning research* 9 (11)

(2008).