# Nearest Neighbor Classifier with Margin Penalty for Active Learning

Yuan Cao[1,2,3], Zhiqiao Gao[2], Jie Hu[2], Mingchuan Yang[2], and Jinpeng Chen[1,3]✉

[1] School of Computer Science (National Pilot Software Engineering School), Beijing University of Posts and Telecommunications, Beijing, China
{caoyuanboy, jpchen}@bupt.edu.cn
[2] China Telecom Corporation Limited Research Institute, Beijing, China
{gaozhq6,hujie1,yangmch}@chinatelecom.cn
[3] Key Laboratory of Trustworthy Distributed Computing and Service (BUPT), Ministry of Education, Beijing, China

**Abstract.** As deep learning becomes the mainstream in the field of natural language processing, the need for suitable active learning method are becoming unprecedented urgent. Active Learning (AL) methods based on nearest neighbor classifier are proposed and demonstrated superior results. However, existing nearest neighbor classifiers are not suitable for classifying mutual exclusive classes because inter-class discrepancy cannot be assured. As a result, informative samples in the margin area can not be discovered and AL performance are damaged. To this end, we propose a novel **N**earest neighbor **C**lassifier with **M**argin penalty for **A**ctive **L**earning(NCMAL). Firstly, mandatory margin penalties are added between classes, therefore both inter-class discrepancy and intra-class compactness are both assured. Secondly, a novel sample selection strategy is proposed to discover informative samples within the margin area. To demonstrate the effectiveness of the methods, we conduct extensive experiments on three real-world datasets with other state-of-the-art methods. The experimental results demonstrate that our method achieves better results with fewer annotated samples than all baseline methods.

**Keywords:** Active Learning · Text Classification · Bert

## 1 Introduction

Recently, Deep Learning (DL) has shown unparalleled ability in many areas especially in the field of natural language processing (NLP). DL-based[4][11][12] text classification methods has changed the landscape of text classification and achieved state-of-the-art performance. However, DL's superb learn capabilities havily relies on large amount of labeled data. As a result, active learning (AL), which aims to maximize model performance while minimize labeling costs, is gradually receiving more attention[5][20][19][14][28], and may help ease the data shortage problems of DL.
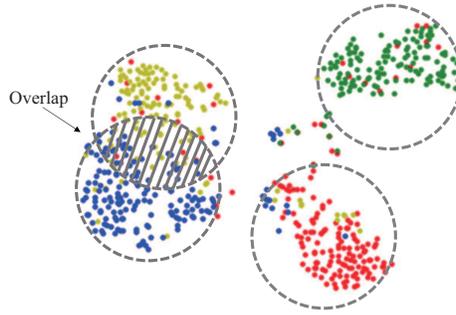
**Fig. 1.** Visualization of sample distribution results after 10 rounds of AL using NCENet on the **AGNEWS** dataset

So far, many AL methods analyse the output logits of the traditional softmax classifier for sample selection. The uncertainty-based method[21][6][5], a bunch of AL methods whose presence may date back to the era of machine learning, aims to calculate the uncertainty of the output logits to select the most uncertain samples for the model. Intuitively, those uncertainty-based methods are inherited in deep learning models. However, these ported methods didn't perform as well as they do on machine learning.

Fang et al.[26] pointed out that the problem encountered with deep models is actually a "false generalize" problem. DL models learn softmax classification boundaries form labeled samples, and incorrectly generalizes the classification boundaries to unlabeled samples. Fang then discards the traditional softmax classifier structure and utilizes a soft nearest neighbour classifier that classifies target samples by selecting prototype vectors on the feature space. NCENet[26] is proposed to avoid the "false generalize" problem by complete abandonment of the softmax classifier structure.

However, NCENet also has its shortcomings. The main structure of NCENet consists of $n$ sigmoid functions instead of one softmax function. Although the sigmoid structure can be used in multi-classification scenario, it may encounter difficulty with classes that are mutually exclusive. For example, Fig.1 shows a typical scene from a real AL training process on the **AGNEWS** dataset, visualised by t-SNE[16]. A clear overlap can be easily seen from the yellow class and blue class. The yellow class and the blue class are mixed together and no clear classification boundary can be established. In this way, two classes that were mutually exclusive become non-mutually exclusive. These are important indications that inter-class differences are not guaranteed by the sigmoid classifier. We define this phenomenon as the "non-exclusive problem". Solving the "non-exclusive problem" will enhance the performance of the model in classification and AL scenarios.

To this end, we propose **N**earest neighbor **C**lassifier with **M**argin penalty for **A**ctive **L**earning(NCMAL), which ensures class not overlapping by adding mandatory margins between classes so that the sigmoid classifier can be used in

classifying mutual exclusive classes. In other words, inter-class discrepancy can be assured with the mandatory inter-class margin added. And at the same time, as we project the whole feature space onto a n-dimensional hyper-sphere, higher inter-class discrepancy brings higher intra-class compactness. As a result, the classification accuracy can be improved. Meanwhile, with margin area added, we believe that unlabeled samples within the margin area has a relatively high uncertainty, and have a high priority when labeling. We proposed a sample selection strategy that focuses high priority samples in the margin area.

Our contributions are summarized as follows:
- The proposed NCMAL effectively increases the inter-class discrepancy and make sigmoid-based classifier suitable for classifying mutual exclusive classes.
- We prove samples in the margin area tend to be more informative and propose a special sample selection strategy, which gives high priority to samples in the margin area.
- NCMAL outperforms state-of-the-art AL methods for text-classification on three real-world datasets.

## 2   Related Work

We focus on AL in pool-based scenarios. Pool-based AL methods can be roughly divided into three categories: uncertainty-based, representation-based and fusion methods which combines uncertainty-based method and representation-based method.

**Uncertainty-based Method.** AL has been of interest to researchers since the days of machine learning, when one wanted to obtain better model performance with fewer labelled samples. In the most intuitive way of thinking, one determines whether a sample needs to be labelled by the uncertainty of the model on the sample. Different methods[18][17] have different measures of uncertainty, such as [2][25][24] based on least model confidence, [22] based on margin sampling, and [13][29] by measuring the entropy of the probability distribution to determine the uncertainty of this classification. In addition, [6] introduces Bayesian inference through the use of a Monte Carlo Dropout, which measures sample uncertainty more accurately by enabling the Dropout in the testing phase. However, the computational efficiency of this method is greatly limited by the need to perform multiple forward propagation.

**Representation-based Method.** Representation-based methods aims to select the most important samples for labelling by analyzing the distribution of the unlabelled samples. As in the DAL (Discriminative Active Learning)[7] method, a binary classifier is trained to discriminate whether a sample comes from the labelled set or the unlabelled set, so that samples that best represent the entire data set can be selected. The EGL(Expected Gradient Length)[9] method measures the impact of a sample on the model by calculating the EGL of the labeled sample, and selects the labeled sample based on this criterion. Coreset[23] is also an emerging and very effective method that models the entire AL process as a coreset problem. By solving the corresponding coreset problem

in the learned representation space, samples that can best represent the entire dataset are selected to be labeled.

**Fusion Method.** In addition, there are many methods that combine the uncertainty-based method with the representation-based method, e.g. the BADGE[1] method can be considered as a combination of the BALD[8] method and the Coreset[23] method, and has been experimented on several models. For example, LL4AL[27] uses an additional network structure to predict the "loss" of a sample, which gives a more accurate measure of the diversity and uncertainty of the sample, and the top $L$ labeled samples are obtained by sorting the loss values in descending order. The NCENet method uses a Nearest Neighbor Classifier to replace the traditional softmax classifier, thus solving the "false generalize" problem of the softmax classifier.

## 3    Methodology

In this section, we first introduce the NCMAL in detail. Then, the sample selection strategy, which aim to informative samples from the margin area, is described.
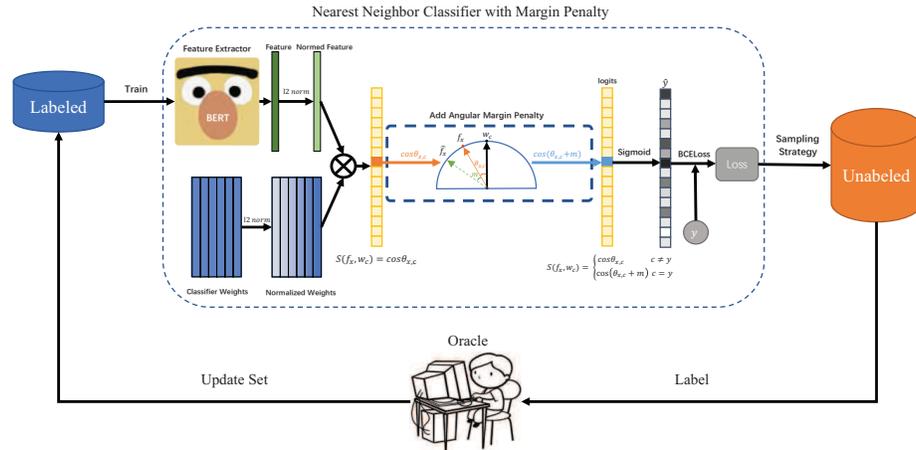


**Fig. 2.** One AL iteration of NCMAL

### 3.1    Overall Framework

Fig.2 shows the whole structure of NCMAL. This work is applied to pool-based active learning scenarios. Specifically, the algorithm is initialized with a small set of labeled samples $\mathcal{L}$ and a larger set of unlabeled samples $\mathcal{U}$. The samples $x_i \in \mathcal{L}$ all have corresponding labels $y_i$, while the unlabeled samples $x_i \in \mathcal{U}$ have

no labels. Using $\mathcal{L}$ as training data, a text classifier $g(x|\Theta) : \mathbf{X} \to \mathbf{Y}$ is trained. The goal of the sample selection strategy is to select $K$ samples from $\mathcal{U}$ by the classification result of the trained model $g(x|\Theta)$. The selected $K$ samples are then annotated and added to $\mathcal{L}$, and used as the training data of the next round of training. The whole algorithm can be summarized as Algorithm 1.

---

**Algorithm 1** *Nearest Neighbor Classifier with Margin Penalty for Active Learning(NCMAL)*

---

**Input:** Unlabeld set $\mathcal{U}$, initial budget $K_{init}$, budget $K$, margin factor $m$, deflation factor $s$, AL rounds $r$.
**Output:** Model parameters $\Theta$, labeled set $\mathcal{L}$
 1: Initialize $\Theta$ from Normal Distribution $\mathcal{N}(0, 0.01)$;
 2: $\mathcal{L} \longleftarrow Random\_Select\_K\_Sample\_From(\mathcal{U}, K_{init})$
 3: **for** $i = 1, 2, ..., r$ **do**
 4:     **for** $x_i, y_i \in \mathcal{L}$ **do**
 5:         Compute $o_{x_i,c}$ for every $c \in C$ according to Eq. 6;
 6:         Compute loss $L$ according to Eq. 7
 7:         Update parameter $\Theta$ by gradient decent optimization with loss $L$
 8:     **end for**
 9:     **for** $x \in \mathcal{U}$ **do**
10:         Compute Margin Confidence score $C_x^{Margin}$ according to Eq. 8
11:     **end for**
12:     $\mathcal{L}_i \longleftarrow Top\_K\_Sample\_Selection\_By\_Confidence\_Score(Conf(\mathcal{U}), K)$
13:     $\mathcal{L} \longleftarrow \mathcal{L} + \mathcal{L}_i$;
14:     $\mathcal{U} \longleftarrow \mathcal{U} - \mathcal{L}_i$
15: **end for**
16: **return** $\Theta, \mathcal{L}$;

---

### 3.2   Nearest Neighbor Classifier with Margin Penalty

In existing nearest neighbor classifier methods[26][10], take NCENet as an example, the classification result of an arbitrary sample mainly depends on the similarity between the feature vector $\boldsymbol{f}_x$ and the prototype vector $\boldsymbol{w}_c, c \in C$. The feature vector is extracted by a arbitrary feature extraction network (e.g. Bert, TextCNN). The output score function can be written as,

$$o_{x,c} = \sigma(\overline{S}(\boldsymbol{f}_x, \boldsymbol{w}_c)) \tag{1}$$

where $\overline{S}(\cdot, \cdot)$ is an arbitrary similarity function. As we can see from Eq.1, the main structure of the NCENet consists of $n$ sigmoid classifiers.

    As we mentioned before, the transformation from softmax to n-sigmoid brings the "non-exclusive" problem. In order to solve this problem, intuitively, we refer to the approach in [3] and add an angular margin penalty between classes during training, which can increase the inter-class discrepancy and the intra-class

compactness. In AL scenario, samples in the overlapping area contains much more information to better separate the overlapped classes. By adding mandatory margin, overlapping areas are now margin areas. Samples used to be in the overlapping area are now located in the created margin area, and can be better measured by the special sample selection strategy we describe later.

First, we utilize dot product to measure the similarity between vectors. We define the similarity between a feature vector $\boldsymbol{f}_x$ and prototype vector corresponding to class $c$ as

$$S(\boldsymbol{f}_x, \boldsymbol{w}_c) = \boldsymbol{f}_x^{\mathrm{T}} \boldsymbol{w}_c = ||\boldsymbol{f}_x||||\boldsymbol{w}_c|| cos\theta_{x,c} \tag{2}$$

where $\theta_{x,c}$ is the angle between $\boldsymbol{f}_x$ and $\boldsymbol{w}_c$. We apply $l_2$ regularization to $\boldsymbol{w}_c$ so that $||\boldsymbol{w}_c|| = 1$. We also regularise $\boldsymbol{f}_x$ and rescale to $s$. With $l_2$ regularisation, we project $\boldsymbol{f}_x$ and $\boldsymbol{w}_c$ onto a feature space shaped as a hypersphere with radius $s$, making the multi-classification prediction dependent only on the angle between the sample vector and the prototype vector. The similarity can be then described as

$$S(\boldsymbol{f}_x, \boldsymbol{w}_c) = s * cos\theta_{x,c}. \tag{3}$$

Since the sample features as well as the prototype vectors are projected onto the same hypersphere with radius $s$, adding a angular margin becomes possible. We add a angular margin penalty $m$ to $\theta_{x,y}$, where $y$ is the ground-truth label of sample $x$. The new similarity function of $\boldsymbol{f}_x$ and $\boldsymbol{w}_y$ can be written as

$$S(\boldsymbol{f}_x, \boldsymbol{w}_c) = s * cos(\theta_{x,y} + m) \tag{4}$$

The whole similarity function can be written as

$$S(\boldsymbol{f}_x, \boldsymbol{w}_c) = \begin{cases} s * cos\theta_{x,c} & c \neq y \\ s * cos(\theta_{x,c} + m) & c = y \end{cases} \tag{5}$$

We then apply sigmoid classifier to the similarity score in order to calculate the probability of sample $x$ belong to class $c$.

$$o_{x,c} = \sigma(S(\boldsymbol{f}_x, \boldsymbol{w}_c)) \tag{6}$$

A binary cross entropy loss function is applied. The loss function can be rewritten as

$$L = -\sum_c y \log o_{x,c} + (1 - y)\log(1 - o_{x,c})$$
$$= -\log \sigma(s * cos(\theta_{x,y} + m)) - \sum_{c \neq y} \log(1 - \sigma(s * cos\theta_{x,c})) \tag{7}$$

And in the testing phase, the sample will be predicted to be the class with maximum Eq. 3.

### 3.3   Sample Selection

In each round of active learning, we rely on the probability output $o_{x,c}$ to select the samples. As we mentioned in the previous section, the NCMAL creates a margin area between classes, and we believe that samples located in the margin area shall have higher priority when labeling. To best find the samples in the margin area, we proposed a confidence score function for NCMAL in order to give samples closer to the margin area a relatively high confidence score.

***Margin Confidence***

$$C_x^{Margin} = -|o_{x,c_0} - o_{x,c_1}|, \tag{8}$$

where $c_0, c_1$ are the classes with largest and second largest output probabilities respectively. It is worth noting that the Margin confidence score is closely related to the difference in hyperarc length from the sample point to the two nearest class centers after projected onto the feature hypersphere.

The higher the confidence score is, the higher priority the sample obtains when labeling. Samples with top-$k$ confidence score will be queried and manually labeled for the next AL iteration. The effects of different sampling strategies will be discussed in the experimental subsection.

## 4   Experiments

NCMAL is tested on several datasets on a text classification task. In this section, we describe the implementation and results of the experiments in detail.

### 4.1   Experimental Settings

***Datasets.*** Three different text classification datasets are used to prove the effectiveness of our method. The three datasets consists of two public datasets and one private dataset (**Telecom**). The **Telecom** dataset comprises of a total of 7,302 real-word messages from Chinese users. These messages were labeled and divided into 25 pre-defined mutually exclusive classes by a dedicated team. The dataset was divided into 5,841 training samples and 1,461 test samples. The Statistics for these three datasets are shown in Table 1.

**Table 1.** Statistics of datasets

| Dataset | AGNEWS | IMDb | Telecom |
|---|---|---|---|
| #class | 4 | 2 | 25 |
| #train | 120000 | 25000 | 5841 |
| #test | 7600 | 25000 | 1461 |
| #init budget | 50 | 100 | 500 |
| #budget | 10 | 20 | 20 |
| #round | 50 | 50 | 30 |

***Feature Extraction Network.*** The commonly used pre-trained Bert[4] were chosen as the Feature Extraction Network. For all AL methods, hyperparameters were chosen consistently for fairness considerations. Necessary changes are made on the original Bert structure due to the requirements of both NCENet and NCMAL, while the number of parameters remains the same for fairness considerations.

***Trainning Details.*** NCMAL is implemented on Pytorch and trained on 4* NVIDIA Tesla V100. Both init-budget and budget selection on different datasets is shown in Table 1. Batch size was set to 10 and the model trained on a learning rate of $2e^{-5}$ using the AdamW[15] optimizer. For each AL sampling method, 10 different random number seeds are used for testing and the final performance of the method was averaged over 10 experiments.

***Baselines.*** We compare our approach to the following baselines.

– **Random.** It is the most commonly used baseline in active learning. The samples added to the labeled set in each round are randomly selected from the unlabeled set.
– **DBAL(Deep Bayesian Active Learning).** Monte Carlo Dropout was used to provide a more accurate measure of the uncertainty of the classifier. Both Confidence and Entropy were used in the final sample selection stage and the best performing of the two methods was selected as the performance of this method in the end.
– **Coreset.** Samples that best cover the entire feature space are selected. We chose two implementations of Coreset as described in [23], the greedy version of Coreset are implemented.
– **BADGE[1].** It can be viewed as a combination of EGL and Coreset, and ensures diversity and uncertainty at the same time.
– **NCENet.** We implemented NCENet as described in [26].

The implementation[1] was based on the code[2] made available by [7].

### 4.2   Model Effect

**Table 2.** Accuracy performance on full training scenario. The best performing method in each row is boldfaced

| Classifier<br>Dataset | Softmax | NCENet | NCMAL |
|---|---|---|---|
| **AGNEWS** | 94.42% | 94.67% | **94.87**% |
| **IMDb** | 85.52% | 85.57% | **85.89**% |
| **Telecom** | 87.81% | 89.11% | **89.45**% |

---

[1] https://github.com/GhostAnderson/Nearest-Neighbor-Classifier-with-Margin-Penalty-for-Active-Learning

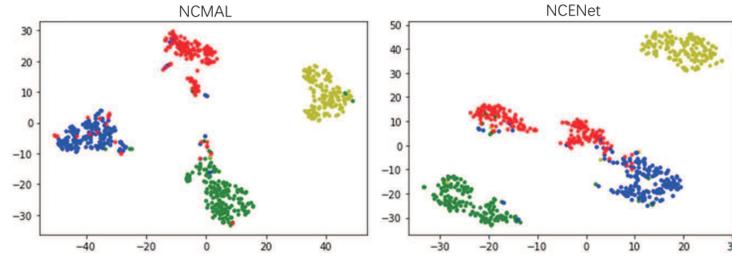[2] https://github.com/dsgissin/DiscriminativeActiveLearning

**Fig. 3.** Demonstration of sample distribution on **AGNEWS** dataset using NCMAL and NCENet

***Performance on full training.*** We first tested the performance of our model under full training with all samples labeled. This result is also equivalent to the test result under AL scenario at 100% sample labeled. Table 2 illustrates the accuracy performance of our NCMAL and baselines fully trained on three datasets. Our NCMAL outperforms all baselines on all three datasets. This demonstrates the structural advantage of our classifier over both traditional softmax classifier and NCENet. Fig.3 visualises the sample distribution of NCMAL and NCENet after full training on **AGNEWS** dataset. The classes formed by NCMAL are more compact compared to NCENet, which can be easily seen from the red class. At the same time, we can see from the distribution of blue class and red class that, the inter-class discrepancy are better assured by the NCMAL.



(a) AGNEWS              (b) IMDb              (c) Telecom

**Fig. 4.** Active learning performance and comparison with baseline methods

***Performance on AL.*** Fig.4 illustrates test accuracy curves of our NCMAL and baselines under AL scenario on three datasets. From Fig. 4, we can draw three following critic conclusions.

*1)* In most cases except for Coreset, AL methods outperforms random selection, which shows the importance of AL methods. Meanwhile, our NCMAL outperforms all other baseline methods. Significant performance gaps can be observed on **Telecom** and **AGNEWS** dataset, in which NCMAL has a large performance gap with other methods from beginning to end. Under **IMDb** dataset, though advantages of all active learning methods over random selection are not very clear, our NCMAL still shows comparable performance over other baseline methods.

*2)* From Fig.4 we can conclude that, the improvement are ascending as class number increases. Specifically, our method shows the greatest improvement over other methods on the **Telecom** dataset (25 classes) and marginal improvements on **IMDb** dataset (2 classes), which implies that, our methods is more suitable for classification with more classes. The reason may be that in classification with more classes, discrepancy between classes are even less guaranteed, and adding a margin term in such scenarios is more helpful in improving classification accuracy than in classification tasks where there are relatively few classes.

*3)* Compared with NCENet, our NCMAL shows constantly better performance especially in the front part of the learning curve. We believe that the addition of margin introduces a prior-knowledge to the model that classes are mutual exclusive, and thus the performance during early training period is improved.

### 4.3 Differenct Sample Selection Strategies



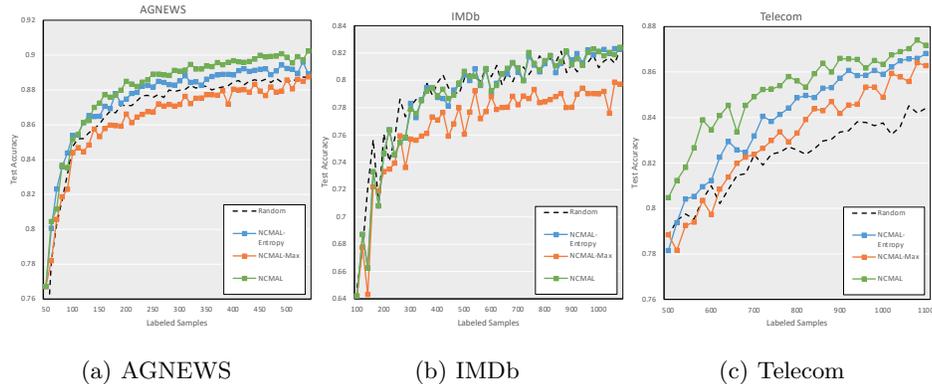| (a) AGNEWS | (b) IMDb | (c) Telecom |
|---|---|---|

**Fig. 5.** Active learning performance and comparison with different variants of NCMAL

With the same NCMAL network structure, the effect of sample selection strategies other than Margin Confidence is also studied. Except for Margin Confidence (NCMAL for simplicity consideration), we bring out two variants of NCMAL with different sample selection strategies as comparison.

highest confidence score to samples from the margin area, thus more samples from the margin area will be selected and labeled. And as a result, NCMAL obtains the best performance. The Entropy Confidence gives moderate priority to samples in the margin area while Max Confidence gives nearly no priority to samples in the margin area, and as a result, NCMAL-Entropy gains moderate performance and NCMAL-Max obtains the worst performance. It is easy to see that the more samples from the margin area are selected, the better the model performs. This finding supports the hypothesis we presented in the previous section that, samples from the margin area tend to be more informative than samples from other areas and should have high priority when labeling. The combination of our NCMAL and Margin Confidence can best discover informative samples from the margin area, thus gains significant performance improvement.

## 5    Conclusion

In this paper, we propose a novel nearest neighbour classifier with margin for active learning (NCMAL). We add angular margin penalties so that the inter-class discrepancy can be assured thus the sigmoid classifier structure can be applied to mutual exclusive classification scenarios. This solves the problem of overlapping class boundaries that can occur in [26] and achieves both better classification results and active learning results. We demonstrate the effectiveness of our method by comparing it with several baseline methods on different real datasets. We also proposed a special sample strategy in order to discover informative samples which lies in the margin area. The experimental results proves the correctness of our hypothesis and demonstrate the superiority of our method for text classification tasks.

## Acknowledgement

## References

1. Ash, J.T., Zhang, C., Krishnamurthy, A., Langford, J., Agarwal, A.: Deep batch active learning by diverse, uncertain gradient lower bounds. In: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net (2020), https://openreview.net/forum?id=ryghZJBKPS
2. Culotta, A., McCallum, A.: Reducing labeling effort for structured prediction tasks. In: AAAI. vol. 5, pp. 746–751 (2005)
3. Deng, J., Guo, J., Xue, N., Zafeiriou, S.: Arcface: Additive angular margin loss for deep face recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4690–4699 (2019)

4.  Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. In: Burstein, J., Doran, C., Solorio, T. (eds.) Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers). pp. 4171–4186. Association for Computational Linguistics (2019). https://doi.org/10.18653/v1/n19-1423, https://doi.org/10.18653/v1/n19-1423

5.  Dor, L.E., Halfon, A., Gera, A., Shnarch, E., Dankin, L., Choshen, L., Danilevsky, M., Aharonov, R., Katz, Y., Slonim, N.: Active learning for bert: An empirical study. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 7949–7962 (2020)

6.  Gal, Y., Islam, R., Ghahramani, Z.: Deep bayesian active learning with image data. In: International Conference on Machine Learning. pp. 1183–1192. PMLR (2017)

7.  Gissin, D., Shalev-Shwartz, S.: Discriminative active learning. arXiv preprint arXiv:1907.06347 (2019)

8.  Houlsby, N., Huszár, F., Ghahramani, Z., Lengyel, M.: Bayesian active learning for classification and preference learning. arXiv preprint arXiv:1112.5745 (2011)

9.  Huang, J., Child, R., Rao, V., Liu, H., Satheesh, S., Coates, A.: Active learning for speech recognition: the power of gradients. arXiv preprint arXiv:1612.03226 (2016)

10. Kontorovich, A., Sabato, S., Urner, R.: Active nearest-neighbor learning in metric spaces. J. Mach. Learn. Res. **18**, 195:1–195:38 (2017), http://jmlr.org/papers/v18/16-499.html

11. Lai, S., Xu, L., Liu, K., Zhao, J.: Recurrent convolutional neural networks for text classification. In: Twenty-ninth AAAI conference on artificial intelligence (2015)

12. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., Soricut, R.: ALBERT: A lite BERT for self-supervised learning of language representations. In: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net (2020), https://openreview.net/forum?id=H1eA7AEtvS

13. Lewis, D.D., Gale, W.A.: A sequential algorithm for training text classifiers. In: SIGIR'94. pp. 3–12. Springer (1994)

14. Li, C., Mao, K., Liang, L., Ren, D., Zhang, W., Yuan, Y., Wang, G.: Unsupervised active learning via subspace learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 8332–8339 (2021)

15. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net (2019), https://openreview.net/forum?id=Bkg6RiCqY7

16. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research **9**(11) (2008)

17. Nafa, Y., Chen, Q., Chen, Z., Lu, X., He, H., Duan, T., Li, Z.: Active deep learning on entity resolution by risk sampling. Knowledge-Based Systems **236**, 107729 (2022)

18. Nguyen, C.V., Ho, L.S.T., Xu, H., Dinh, V., Nguyen, B.T.: Bayesian active learning with abstention feedbacks. Neurocomputing **471**, 242–250 (2022)

19. Nguyen, Q.P., Low, B.K.H., Jaillet, P.: An information-theoretic framework for unifying active learning problems. In: Proc. AAAI. pp. 9126–9134 (2021)

20. Prabhu, S., Mohamed, M., Misra, H.: Multi-class text classification using bert-based active learning. In: Dragut, E.C., Li, Y., Popa, L., Vucetic, S. (eds.) 3rd Workshop on Data Science with Human in

the Loop, DaSH@KDD, Virtual Conference, August 15, 2021 (2021), `https://drive.google.com/file/d/1xVy4p29UPINmWl8Y7OospyQgHiYfH4wc/view`

21. Ren, P., Xiao, Y., Chang, X., Huang, P., Li, Z., Gupta, B.B., Chen, X., Wang, X.: A survey of deep active learning. ACM Comput. Surv. **54**(9), 180:1–180:40 (2022). https://doi.org/10.1145/3472291, `https://doi.org/10.1145/3472291`

22. Scheffer, T., Decomain, C., Wrobel, S.: Active hidden markov models for information extraction. In: International Symposium on Intelligent Data Analysis. pp. 309–318. Springer (2001)

23. Sener, O., Savarese, S.: Active learning for convolutional neural networks: A coreset approach. In: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net (2018), `https://openreview.net/forum?id=H1aIuk-RW`

24. Settles, B.: Active learning literature survey (2009)

25. Settles, B., Craven, M.: An analysis of active learning strategies for sequence labeling tasks. In: Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing. pp. 1070–1079 (2008)

26. Wana, F., Yuana, T., Fua, M., Jib, X., Yea, Q.H.Q.: Nearest neighbor classifier embedded network for active learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 10041–10048 (2021)

27. Yoo, D., Kweon, I.S.: Learning loss for active learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 93–102 (2019)

28. Zhou, B., Cai, X., Zhang, Y., Guo, W., Yuan, X.: Mtaal: Multi-task adversarial active learning for medical named entity recognition and normalization. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 14586–14593 (2021)

29. Zhu, J., Wang, H., Yao, T., Tsou, B.K.: Active learning with sampling by uncertainty and density for word sense disambiguation and text classification. In: Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008). pp. 1137–1144 (2008)