# CSS-LM: A Contrastive Framework for Semi-supervised Fine-tuning of Pre-trained Language Models

Yusheng Su, Xu Han, Yankai Lin, Zhengyan Zhang, Zhiyuan Liu, Peng Li, Jie Zhou, Maosong Sun

Abstract—Fine-tuning pre-trained language models (PLMs) has demonstrated its effectiveness on various downstream NLP tasks recently. However, in many scenarios with limited supervised data, the conventional fine-tuning strategies cannot sufficiently capture the important semantic features for downstream tasks. To address this issue, we introduce a novel framework (named "CSS-LM") to improve the fine-tuning phase of PLMs via contrastive semi-supervised learning. Specifically, given a specific task, we retrieve positive and negative instances from large-scale unlabeled corpora according to their domain-level and class-level semantic relatedness to the task. We then perform contrastive semi-supervised learning on both the retrieved unlabeled instances and original labeled instances to help PLMs capture crucial task-related semantic features. The experimental results show that CSS-LM achieves better results than the conventional fine-tuning strategy on a series of downstream tasks with few-shot settings by up to 7.8%, and outperforms the latest supervised contrastive fine-tuning strategy by up to 7.1%. Our datasets and source code will be available to provide more details.

Index Terms—Pre-trained Language Model, Few-shot Learning, Contrastive Learning, Semi-supervised Learning, Fine-tuning

# **1** INTRODUCTION

**P**RE-TRAINED language models (PLMs) like BERT [1] and RoBERTa [2] can learn general language understanding abilities from large-scale unlabeled corpora, and provide informative contextual representations for downstream tasks. In recent years, instead of learning task-oriented models from scratch, it has gradually become a consensus to finetune PLMs for specific tasks, which has been demonstrated on various downstream NLP tasks, including dialogue [3], summarization [4], [5], question answering [6], [7], and relation extraction [8], [9].

Although fine-tuning PLMs has become a dominant paradigm in the NLP community, it still requires large amounts of supervised data to capture critical semantic features for downstream tasks [10], [11]. Without sufficient supervised data for downstream tasks, the conventional fine-tuning strategy might capture biased features for the downstream tasks that may cause errors or have decision boundary bias in Fig. 1. Therefore, fine-tuning PLMs is still challenging in those scenarios with limited data, and cannot be well generalized to many real-world applications whose labeled data is hard and expensive to obtain. Hence, a natural question to ask is: *How can we effectively capture crucial semantic features for downstream tasks with limited supervised data*?

To address this issue, some preliminary works have made some attempts to utilize semi-supervised methods

{liuzy, sms}@tsinghua.edu.cn

with unlabeled data for fine-tuning PLMs [12]. Nevertheless, these methods require extra efforts on high-quality labeling to start the semi-supervised learning process. Another way, which can capture crucial semantic features from the limited supervised data in the fine-tuning stage without labeling, is applying contrastive learning [13], by forcing the positive instances to be close to each other in the semantic space, and meanwhile forcing the negative instances to be far away from the positive ones. However, existing contrastive learning methods for enhancing PLMs still lack leveraging the rich large-scale open-domain corpora.

In order to make the use of large-scale open-domain corpora to capture crucial semantic features better, we introduce a novel Contrastive framework for Semi-Supervised Finetuning PLMs (named "CSS-LM"), which extends the conventional supervised fine-tuning strategy into a semi-supervised form enabling PLMs to leverage extra task-related data from large-scale open-domain corpora without annotating new labels. More specifically, CSS-LM separates the unlabeled data into the positive instances and the negative ones according to these instances' semantic relatedness to downstream tasks. Afterward, CSS-LM applies contrastive learning to let PLMs distinguish the nuances between these instances, so that PLMs can learn informative semantic features not expressed by the limited supervised data of the downstream task. We give an intuitive motivation description in Fig. 1 to show how our framework could better capture crucial semantics from unsupervised data to make final fine-tuned pre-trained language models more discriminative and have better decision boundaries.

As to use all unlabeled data is nearly impossible, we require to find the most informative positive and negative instances for our contrastive learning. Therefore, one key challenge of CSS-LM is how to measure the semantic relatedness of unlabeled instances to downstream tasks.

<sup>•</sup> Yusheng Su, Xu Han, Zhengyan Zhang, Zhiyuan Liu, Maosong Sun are with the Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China. E-mail: {sys19, hanxu17, zy-z19}@mails.tsinghua.edu.cn,

Yankai Lin, Peng Li and Jie Zhou are with the Pattern Recognition Center, WeChat AI Department, Tencent, Beijing 100080, China. Email: {yankailin, patrickpli, withtomzhou}@tencent.com

<sup>•</sup> *Zhiyuan Liu is the corresponding author.* 



Fig. 1. The overall framework of CSS-LM, which illustrates how we leverage specific task instances (red and blue dots) to retrieve task-related instances (light red and light blue dots in the green cluster) by measuring domain-related and class-related semantics. In the upper figure, by performing contrastive semi-supervised leaning fine-tuning, we can obtain a better decision boundary for the task.

In this paper, we consider both domain-level and classlevel semantic relatedness of instances to downstream tasks. More specifically, we encode unlabeled instances into the domain-level and class-level semantic spaces, respectively, and define the domain-level and class-level relatedness as the similarity of their representations in the corresponding semantic spaces. When performing contrastive semisupervised learning for the domain-level representations, all supervised instances of downstream tasks are used as the initialized positive instances. After that, we continue to retrieve unlabeled instances closest and farthest from the existing positive instances as new positive and negative instances. When performing contrastive semi-supervised learning for the class-level representations, we apply operations similar to learning domain-level representations; the only difference is that learning class-level representations requires considering fine-grained class-level semantics rather than coarse-grained domain-level semantics when selecting positive and negative instances.

In the experiments, the results on three typical classification tasks, including sentiment classification, intent classification, and relation classification, show that our proposed framework can outperform the conventional fine-tuning strategy, the conventional semi-supervised strategies, and the latest supervised contrastive fine-tuning strategies under the limited supervised data settings. These results explicitly indicate a promising direction of utilizing unlabeled data for fine-tuning PLMs based on contrastive semi-supervised learning.

To summarize, our major contributions are as follows:

- (1) We propose a contrastive semi-supervised framework CSS-LM, which can better leverage unlabeled instances from open-domain corpora to capture task-related features and enhance the model performance on the downstream tasks.
- (2) CSS-LM is free to the domain dependence, which can efficiently capture domain information from open-domain

corpora and retrieve domain-related instances.

(3) We conduct experiments on six classification datasets. The experimental results show that CSS-LM outperforms various typical fine-tuning models under the few-shot settings. Besides, sufficient empirical analyses of our retrieval mechanism and retrieval instances demonstrate that contrastive semi-supervised learning is more helpful than semi-supervised learning (pseudo labeling) to learn from unlabeled data.

# 2 RELATED WORK

### 2.1 Pre-trained Language Models

Various recent PLMs like BERT [1], RoBERTa [2] and XL-Net [14], provide a new perspective for NLP models to utilize a large amount of open-domain unlabeled data. Inspired by these works, a series of works have designed specific self-supervised learning objectives to help PLMs learn specific abilities in the pre-training phase, such as representing token spans [6] and entities [15], [16], [17], [17], [18], conferential reasoning [19], multi-lingualism [20], multi-modality [21], [22], [23], [24], [25], etc. Besides, some PLMs [26], [27], [28] are devoted to learning specific domain semantics by pre-training on the specific domain corpora, and these works demonstrate their effectiveness as well.

As PLMs are aimless concerning various downstream tasks in the pre-training stage. Hence, to adapt PLMs for a specific task requires a fine-tuning on extra supervised data of the tasks. Specifically, fine-tuning often replaces the top layers of PLMs with a specific task sub-network, and continues to update the parameters with the supervised data. Fine-tuning PLMs has also demonstrated its effectiveness on various downstream tasks, including dialogue [3], summarization [4], [5], question answering [7], and relation extraction [8], [9].

However, without sufficient supervised data, the conventional fine-tuning methods cannot effectively capture useful features for downstream tasks, which may lead to the side effect on performance. To address this issue, a series of works focus on exploring various heuristics during tuning the parameters of PLMs. Howard and Ruder et al. [29] gradually unfreeze the layers of PLMs with a heuristic learning rate schedule to enhance the fine-tuning performance of PLMs. Then, Peters et al. [30] study which layers of PLMs should be adapt or freeze during the fine-tuning stage. Houlsby et al. [31] and Stickland et al. [32] leverage some additional layers to PLMs and update parameters of specific additional layers during the fine-tuning phase.

Besides, some preliminary works have made some attempts to utilize unlabeled data for fine-tuning PLMs: Gu et al. [33] conduct selective language modeling with unlabeled data to focus on the semantic features related to the finetuning tasks. Gururangan et al. [34] propose a framework to retrieve task-related data from large-scale in-domain corpora to enhance fine-tuning PLMs. The in-domain instances often meet an individual margin probability distribution over instances, and they thus have correlated semantics in the feature space that is beneficial for specific tasks. Du et al. [12] further introduce a text retriever trained on a large amount of supervised data to retrieve task-specific in-domain data from large-scale open-domain corpora.

Existing works for enhancing fine-tuning rely on predefined in-domain corpora, massive supervised data, or extra efforts on labeling, limiting them to be applied to broad realworld applications. Unlike these works, our proposed contrastive semi-supervised framework can automatically and iteratively utilize large-scale open-domain data to improve fine-tuning under the few-shot settings without annotating extra data.

## 2.2 Contrastive Learning

Unlike conventional discriminative methods that learn a mapping to labels and generative methods that reconstruct input instances, contrastive learning is a learning paradigm based on comparing. Specifically, contrastive learning can be considered as learning by comparing among different instances instead of learning from individual instances one at a time. The comparison can be performed between positive pairs of "similar" instances and negative pairs of "dissimilar" instances. The early efforts for self-supervised contrastive learning have led to significant advances in NLP [9], [35], [36], [37], [38] and CV [39], [40], [41], [42], [43], [44] tasks. Nevertheless, self-supervised contrastive learning still has a limitation: it cannot utilize the supervised data of downstream tasks and sufficiently capture the fine-grained semantics to specific classes. Intuitively, self-supervised contrastive learning tends to distinguish instances rather than to classify instances.

Therefore, Khosla et al. [45] propose supervised contrastive learning to leverage the supervised instances of downstream tasks. More recently, Gunel et al. [13] verify the effectiveness of supervised contrastive learning in the finetuning stage of PLMs. However, it only considers the limited amounts of the supervised data, and ignores the potential information distributed in the unlabeled data.

In this paper, our approach is a general semi-supervised fine-tuning framework based on contrastive learning, which

could capture richer semantics from unlabeled data and align features with the supervised data of downstream tasks, simultaneously distinguishing and classifying instances.

# **3** METHODOLOGY

In this work, we mainly focus on fine-tuning PLMs for classification tasks. Unlike conventional fine-tuning methods aiming to learn features by classifying the limited labeled training data, CSS-LM leverages unlabeled data in opendomain corpora to capture better features. More specifically, as illustrated in Fig. 1, CSS-LM utilizes contrastive semisupervised learning in the fine-tuning phase of PLMs, aiming to identify all crucial semantic features to distinguish the instances of different domains and classes with both labeled and unlabeled data, and further obtain the better decision boundary.

CSS-LM consists of five important modules, including (1) contrastive semi-supervised learning, (2) informative instance retrieval, (3) semantic representation learning, (4) downstream task fine-tuning, and (5) efficient representation updating. In this section, we will first give some essential notations and then introduce these essential modules.

# 3.1 Notations

Given a specific downstream classification task, we denote its class set as  $\mathcal{Y}$ , and its training set as  $\mathcal{T} = \{(x_{\mathcal{T}}^1, y_{\mathcal{T}}^1), (x_{\mathcal{T}}^2, y_{\mathcal{T}}^2), \dots, (x_{\mathcal{T}}^N, y_{\mathcal{T}}^N)\}$ , where  $y_{\mathcal{T}}^i \in \mathcal{Y}$  is the supervised annotation of the instance  $x_{\mathcal{T}}^i$  and N is the instance number of the training set. Under the few-shot setting, for each class  $y \in \mathcal{Y}$ , we denote the number of instances whose class is y as  $K_y$ .

As our framework will utilize large-scale unsupervised open-domain corpora, we denote the open-domain unlabeled corpora as  $\mathcal{O} = \{x_{\mathcal{O}}^1, x_{\mathcal{O}}^2, \cdots, x_{\mathcal{O}}^M\}$ , where M is the instance number of unlabeled corpora. We use the bold face to indicate the representation of an instance computed by PLMs, e.g., the representation of  $x_{\mathcal{T}}^i$  is  $\mathbf{x}_{\mathcal{T}}^i$ . As  $\mathbf{x}_{\mathcal{T}}^i$  is computed by the PLM-based encoder, we denote the encoder as  $\text{Enc}(\cdot)$ , which will be introduced in 3.4.1.

Note that our framework is applicable to any PLMs. In experiments, we select BERT [1] and RoBERTa [2] as our encoders, considering they are the state-of-the-art and widely-used ones of existing PLMs.

#### 3.2 Contrastive Semi-supervised Learning

Given the training set of a specific downstream task T, we first introduce how to apply contrastive learning for supervised fine-tuning:

$$\mathcal{L}_{CS} = -\sum_{i=1}^{N} \frac{\sum_{j \in \mathcal{C}(i)} f(\mathbf{x}_{\mathcal{T}}^{i}, \mathbf{x}_{\mathcal{T}}^{j}) - \log Z^{i}}{|\mathcal{C}(i)|},$$

$$Z^{i} = \sum_{k=1}^{N} e^{f(\mathbf{x}_{\mathcal{T}}^{i}, \mathbf{x}_{\mathcal{T}}^{k})},$$
(1)

where  $C(i) = \{j | j \neq i, y_T^j = y_T^i\}$  is the index set of instances that share the same label with the instance  $x_T^i$ , and  $f(\cdot, \cdot)$  is the similarity function. Considering many downstream tasks only have limited supervised data, we apply contrastive semi-supervised learning in the fine-tuning phase of PLMs instead of Eq. (1). Specifically, we retrieve positive and negative instances from the open-domain corpora O for each supervised instance.

We denote the positive instance index set of  $x_{\mathcal{T}}^i$  as  $\mathcal{P}(i)$  and the negative one as  $\mathcal{N}(i)$ . Formally, the contrastive semisupervised objective for fine-tuning PLMs is:

$$\mathcal{L}_{CSS} = -\sum_{i=1}^{N} \left( \frac{\sum_{j \in \mathcal{C}(i)} f(\mathbf{x}_{\mathcal{T}}^{i}, \mathbf{x}_{\mathcal{T}}^{j}) - \log Z_{s}^{i}}{|\mathcal{C}(i)|} + \frac{\sum_{j \in \mathcal{P}(i)} f(\mathbf{x}_{\mathcal{T}}^{i}, \mathbf{x}_{\mathcal{O}}^{j}) - \log Z_{s}^{i}}{|\mathcal{P}(i)|} + \frac{\sum_{j \in \mathcal{P}(i)} \sum_{k \in \mathcal{P}(i), j \neq k} f(\mathbf{x}_{\mathcal{O}}^{j}, \mathbf{x}_{\mathcal{O}}^{k}) - \log Z_{u}^{j}}{(|\mathcal{C}(i)| + 1) \times |\mathcal{P}(i)| \times (|\mathcal{P}(i)| - 1)} \right),$$

$$(2)$$

where  $Z_s^i$  and  $Z_u^j$  are calculated as:

$$Z_{s}^{i} = \sum_{k=1}^{N} e^{f(\mathbf{x}_{\mathcal{T}}^{i}, \mathbf{x}_{\mathcal{T}}^{k})} + \sum_{k \in \mathcal{N}(i) \cup \mathcal{P}(i)} e^{f(\mathbf{x}_{\mathcal{T}}^{i}, \mathbf{x}_{\mathcal{O}}^{k})},$$

$$Z_{u}^{j} = \sum_{k \in \mathcal{P}(i), j \neq k} e^{f(\mathbf{x}_{\mathcal{O}}^{j}, \mathbf{x}_{\mathcal{O}}^{k})} + \sum_{k \in \mathcal{N}(i)} e^{f(\mathbf{x}_{\mathcal{O}}^{j}, \mathbf{x}_{\mathcal{O}}^{k})}.$$
(3)

Intuitively, by applying contrastive semi-supervised learning with Eq. (2), we can simultaneously consider the similarities among both supervised and retrieved instances, which can let PLMs well capture semantics from the unlabeled data of the task T.

## 3.3 Informative Instance Retrieval

Given a supervised instance  $x_{\mathcal{T}}^i$ , we retrieve the most informative instances from the open-domain corpora  $\mathcal{O}$  to build the positive instance index set  $\mathcal{P}(i)$  and the negative one  $\mathcal{N}(i)$ .

A straightforward retrieval solution is to regard the most similar instances as positive instances and the most dissimilar ones as negative instances, i.e., retrieve the instances  $x_{\mathcal{O}}^*$ according to  $f(\mathbf{x}_{\mathcal{T}}^i, \mathbf{x}_{\mathcal{O}}^*)$ . However, this straightforward solution is too coarse to select those most informative instances, since the instances of the open-domain corpora may belong to various quite different domains and classes. We cannot know which semantic levels would make the greater contribution with the coarse-grained function  $f(\mathbf{x}_{\mathcal{T}}^i, \mathbf{x}_{\mathcal{O}}^o)$  measuring.

Thus, instead of using  $f(\cdot)$  to retrieve instances, we introduce a similarity function  $f_T(\cdot)$  for retrieving the most informative instances by empirically considering the semantic relatedness between instances from two perspectives: domain-level and class-level similarities. More specifically, given instances  $x_T^i, x_O^*, f_T(\cdot)$  obtains the instance relatedness by calculating the summation of the domain relatedness and class relatedness as:

$$f_T(\mathbf{x}^i_{\mathcal{T}}, \mathbf{x}^*_{\mathcal{O}}) = f_D(\mathbf{d}^i_{\mathcal{T}}, \mathbf{d}^*_{\mathcal{O}}) + f_C(\mathbf{c}^i_{\mathcal{T}}, \mathbf{c}^*_{\mathcal{O}}), \quad (4)$$

where  $f_D(\cdot)$  is the domain similarity function,  $f_C(\cdot)$  is the class similarity function,  $\mathbf{d}_{\mathcal{T}}^i, \mathbf{d}_{\mathcal{T}}^*$  are the domain-level representations, and  $\mathbf{c}_{\mathcal{T}}^i, \mathbf{c}_{\mathcal{T}}^*$  are the class-level representations. In this way, we could build the positive instance index set  $\mathcal{P}(i)$  and the negative one  $\mathcal{N}(i)$ .

However, how to obtain domain-level and class-level representations is still a problem. Next, we will introduce



Fig. 2. Given an instance, we add two special tokens <code>[DOMAIN]</code> and <code>[CLASS]</code> into the sequence, and then input the sequence into the PLM-based encoder. d and c are the representations of <code>[DOMAIN]</code> and <code>[CLASS]</code> respectively. x is the concatenation of the two special token representations.

how to encode instances into domain-level and class-level spaces, and how to learn their representations respectively under our contrastive semi-supervised framework.

### 3.4 Semantic Representation Learning

This section gives the details of the text encoder and introduces how to learn domain-level and class-level instance representations.

### 3.4.1 Encoder Based on Pre-trained Language Models

Given an instance x (either a supervised instance or a unlabeled one), as show in Fig. 2, we first add two special tokens in front of the input sequence of x, i.e.,  $\bar{x} = [[DOMAIN], [CLASS], x]$ . Then, we input  $\bar{x}$  into the encoder, which is a multi-layer bidirectional Transformer encoder, such as BERT<sub>BASE</sub> [1] and RoBERTa<sub>BASE</sub> [2], and use the output representations of [DOMAIN] and [CLASS] as the domain-level representation **d** and the class-level representation **c** respectively. We then define the whole instance representation **x** as the concatenation of the domain-level and class-level representations as:

$$\mathbf{x} = [\mathbf{d}; \mathbf{c}] = \operatorname{Enc}(\bar{x}), \tag{5}$$

where d and c are the representations of [DOMAIN] and [CLASS] respectively,  $[\cdot; \cdot]$  is the concatenation of representations, and  $Enc(\cdot)$  indicates the PLM-based text encoder.

## 3.4.2 Domain-level Representations

We hope that domain-level representations can push instances with different domains far away, and meanwhile, cluster those instances with similar domains. Hence, the whole learning objective is:

$$\mathcal{L}_{D} = -\sum_{i=1}^{N} \Big( \frac{\sum_{j=1}^{N} f_{D}(\mathbf{d}_{\mathcal{T}}^{i}, \mathbf{d}_{\mathcal{T}}^{j}) - \log Z_{d}^{i}}{N} + \frac{\sum_{j \in \mathcal{P}_{D}(i)} f_{D}(\mathbf{d}_{\mathcal{T}}^{i}, \mathbf{d}_{\mathcal{O}}^{j}) - \log Z_{d}^{i}}{|\mathcal{P}_{D}(i)|} \Big),$$

$$Z_{d}^{i} = \sum_{k=1}^{N} e^{f_{D}(\mathbf{d}_{\mathcal{T}}^{i}, \mathbf{d}_{\mathcal{T}}^{k})} + \sum_{k=1}^{M} e^{f_{D}(\mathbf{d}_{\mathcal{T}}^{i}, \mathbf{d}_{\mathcal{O}}^{k})},$$
(6)

where  $\mathcal{P}_D(i)$  is the domain-related positive instance index set of  $x_{\mathcal{T}}^i$ .  $x_{\mathcal{T}}^i$  is retrieved from the open-domain corpora according to the similarity between domain-level representations computed by  $f_D(\cdot)$ . Besides the instances mentioned in  $x_{\mathcal{T}}^i$ , we regard all other instances of the open-domain corpora  $\mathcal{O}$  as domain-related negative instances.

## 3.4.3 Class-level Representations

Class-level representations aim to distinguish the semantic difference between the different classes of the task, and push away those instances not related to one specific class. Hence, the learning objective is formulated as:

$$\mathcal{L}_{C} = -\sum_{i=1}^{N} \left( \frac{\sum_{j=1}^{N} f_{C}(\mathbf{c}_{\mathcal{T}}^{i}, \mathbf{c}_{\mathcal{T}}^{j}) - \log Z_{c}^{i}}{N} + \frac{\sum_{j \in \mathcal{P}_{C}(i)} f_{C}(\mathbf{c}_{\mathcal{T}}^{i}, \mathbf{c}_{\mathcal{O}}^{j}) - \log Z_{c}^{i}}{|\mathcal{P}_{C}(i)|} \right),$$
(7)
$$Z_{c}^{i} = \sum_{k=1}^{N} e^{f_{C}(\mathbf{c}_{\mathcal{T}}^{i}, \mathbf{c}_{\mathcal{T}}^{k})} + \sum_{k \in \mathcal{P}_{C}(i) \cup \mathcal{N}_{C}(i)} e^{f_{C}(\mathbf{d}_{\mathcal{T}}^{i}, \mathbf{d}_{\mathcal{O}}^{k})},$$

where 
$$\mathcal{P}_C(i)$$
 is the class-related positive instance index set of  $x_{\mathcal{T}}^i$  retrieved from the open-domain corpora according to the similarity between the class-level representations computed by  $f_C(\cdot)$ . Note that, given the class of  $x_{\mathcal{T}}^i$ , we take the supervised instances belonging to all other classes and the retrieved positive instances of all other classes as the negative instances set of  $x_{\mathcal{T}}^i$ . The index set of these negative instances of  $x_{\mathcal{T}}^i$  is denoted as  $\mathcal{N}_C(i)$ .

#### 3.5 Efficient Representation Updating

Note that we set  $\mathcal{P}(i)$ ,  $\mathcal{P}_D(i)$  and  $\mathcal{P}_C(i)$  as empty sets at the beginning of the fine-tuning phase, since at that time, three kinds of representations are not well trained. Then, we will iteratively retrieve instances from the open-domain corpora to expand these sets. Since we found the PLM parameters will not change vastly during fine-tuning for the downstream task, for computational efficiency, we only update the parts of instance representations in the open-domain corpora  $\mathcal{O}$ , which are retrieved in every step, rather than updating the whole instance representations.

#### 3.6 Downstream Task Fine-tuning and Optimization

In classification tasks, we are devoted to finding effective decision boundaries with semantic representations. Given the training instances  $(x_{\mathcal{T}}^i, y_{\mathcal{T}}^i) \in \mathcal{T}$ , we encode  $x_{\mathcal{T}}^i$  to  $\mathbf{x}_{\mathcal{T}}^i$  with the text encoder  $\text{Enc}(\cdot)$ . Similar to the conventional

PLM fine-tuning, we apply the cross-entropy loss to learn task-specified classifier as:

$$\mathcal{L}_{CE} = -\sum_{i=1}^{N} \log f_{task}(\mathbf{x}_{\mathcal{T}}^{i}, y_{\mathcal{T}}^{i}), \qquad (8)$$

where  $f_{task}(\cdot)$  is the neural layers for specific tasks built on the PLMs to compute the probability  $p(y_{\mathcal{T}}^i|x_{\mathcal{T}}^i)$ .

## Algorithm 1: Contrastive Semi-supversied Learning

**Data:** Training set  $\mathcal{T}$ ; open-domain corpora  $\mathcal{O}$ ; development set  $\mathcal{E}$ ; **Result:** The optimal CSS-LM parameters  $\theta$ **Initialization:** Parameters of CSS-LM  $\theta$  [Refer to 4.2]; task-related sets  $\mathcal{P}(i)$ ,  $\mathcal{N}(i)$ ; domain-related set  $\mathcal{P}_D(i)$ ; class-related sets  $\mathcal{P}_C(i)$ ,  $\mathcal{N}_C(i)$ ; **for** *epoch*  $\leftarrow [0, ..., E]$  **do while**  $i \in N$  **do Encode** an instance of the training set:  $\mathbf{x}_{\mathcal{T}}^i = [\mathbf{d}_{\mathcal{T}}^i; \mathbf{c}_{\mathcal{T}}^i] = \operatorname{Enc}(\bar{x}_{\mathcal{T}}^i)$  [Refer to 3.4.1]; Retrieve the domain-related instances by  $f_D(\cdot)$ and update  $\mathcal{P}_D(i)$ ,  $\mathcal{N}_D(i)$ ,  $[\mathbf{d}_{\mathcal{O}}^j; \mathbf{c}_{\mathcal{O}}^j] = \operatorname{Enc}(\bar{x}_{\mathcal{O}}^j)$ , where  $\bar{x}_{\mathcal{T}}^i \in \{\mathcal{P}_D(i), \mathcal{N}_D(i)\}$ 

 $\bar{x}_{\mathcal{O}}^{j} \in \{\mathcal{P}_{D}(i), \mathcal{N}_{D}(i)\}, \\ \mathbf{d}_{\mathcal{O}}^{j} \leftarrow Select_{\mathbf{d}}([\mathbf{d}_{\mathcal{O}}^{j}; \mathbf{c}_{\mathcal{O}}^{j}]), \\ \mathbf{d}_{\mathcal{T}}^{i} \leftarrow Select_{\mathbf{d}}([\mathbf{d}_{\mathcal{T}}^{i}; \mathbf{c}_{\mathcal{T}}^{i}]), \\ \text{Perform } \mathcal{L}_{\mathcal{D}} [\text{Refer to } 3.4.2];$ 

Retrieve the class-related instances by  $f_C(\cdot)$ and update  $\mathcal{P}_C(i), \mathcal{N}_C(i),$  $[\mathbf{d}_{\mathcal{O}}^j; \mathbf{c}_{\mathcal{O}}^j] = \operatorname{Enc}(\bar{x}_{\mathcal{O}}^j),$  where  $\bar{x}_{\mathcal{O}}^j \in \{\mathcal{P}_C(i), \mathcal{N}_C(i)\},$  $\mathbf{c}_{\mathcal{O}}^j \leftarrow Select_{\mathbf{c}}([\mathbf{d}_{\mathcal{O}}^j; \mathbf{c}_{\mathcal{O}}^j]),$  $\mathbf{c}_{\mathcal{T}}^i \leftarrow Select_{\mathbf{c}}([\mathbf{d}_{\mathcal{T}}^i; \mathbf{c}_{\mathcal{T}}^i]),$ Perform  $\mathcal{L}_{\mathcal{C}}$  [Refer to 3.4.3];

Retrieve the task-related instances by  $f_T(\cdot)$ and update  $\mathcal{P}(i), \mathcal{N}(i)$  [Refer to 3.3],  $\mathbf{x}_{\mathcal{O}}^{j} = \text{Enc}(\bar{x}_{\mathcal{O}}^{j})$ , where  $\bar{x}_{\mathcal{O}}^{j} \in {\mathcal{P}(i), \mathcal{N}(i)}$ ,  $\mathbf{x}_{\mathcal{O}}^{j} \leftarrow \mathbf{x}_{\mathcal{O}}^{j}$ ,  $\mathbf{x}_{\mathcal{T}}^{i} \leftarrow \mathbf{x}_{\mathcal{T}}^{i}$ , Perform  $\mathcal{L}_{CSS}$  [Refer to 3.2];

Leverage  $Enc(\cdot)$  to update representations of  $\bar{x}_{\mathcal{O}}^{j} \in \{\mathcal{P}_{D}(i), \mathcal{N}_{D}(i), \mathcal{P}_{C}(i), \mathcal{N}_{C}(i), \mathcal{P}(i), \mathcal{N}(i)\}$  to the open-domain corpora representation set [Refer to 3.5];

Calculate the total loss:  

$$\mathcal{L} = \mathcal{L}_{CSS} + \mathcal{L}_D + \mathcal{L}_C + \mathcal{L}_{CE} \text{ [Refer to 3.6];}$$
Compute gradient  $\nabla_{\theta} \mathcal{L}(\theta; \mathcal{T})$ ;  
Update parameters  $\theta = \theta - \lambda \cdot \nabla_{\theta} \mathcal{L}(\theta; \mathcal{T})$ ;  
end

Save  $\theta$  as  $\theta_{epoch}$  in every epoch;

end

**Return**: the optimal  $\theta_{epoch}$  in the development set  $\mathcal{E}$ ;

We optimize the domain-level, class-level, and the whole instance representations jointly. Therefore, the overall learning objective is defined as the sum of four losses:

$$\mathcal{L} = \mathcal{L}_{CSS} + \mathcal{L}_D + \mathcal{L}_C + \mathcal{L}_{CE},\tag{9}$$

where  $\mathcal{L}_{CSS}$  is the contrastive semi-supervised loss, both  $\mathcal{L}_D$  and  $\mathcal{L}_C$  are the functions to let the encoder learn domainlevel and class-level representations respectively.  $\mathcal{L}_{CE}$  is the conventional fine-tuning cross-entropy loss.

The overall detail can refer to Algorithm 1. Given an instance in each step, first, CSS-LM retrieves domain-related, class-related, and task-related instances, then delivers to the corresponding positive and negative instance index set. Second, CSS-LM updates the representations of retrieved instances to the open-domain corpora representation set. Finally, CSS-LM leverage these instances to perform contrastive semi-supervised to obtain model parameters in every epoch. We will evaluate the CSS-LM on the development set every epoch and choose the optimal one as our model in downstream tasks.

## 4 EXPERIMENTS

In this section, we would first introduce the datasets, the experimental settings, and the details of the baseline models used in our experiments. After that, we give some empirical analyses to show the effectiveness of our contrastive semisupervised learning, indicating the promising results of leveraging unlabeled instances. Then, we perform some ablation studies to show which level of semantic relatedness mainly contributes to CSS-LM and the influence of the retrieved size. Finally, we perform visualization and case studies for a more intuitive observation.

## 4.1 Datasets and Tasks

We conduct our experiments on three typical text classification tasks including sentiment classification, intent classification, and relation extraction:

(1) **Sentiment classification** is the task of classifying the polarity of a given sentence. Sentiment classification is a core task of text classification. For sentiment classification, we select SemEval [46] and SST-5 [47] for our experiments.

(2) **Intent Classification** is the task of correctly labeling a natural language utterance from a predetermined set of intents. Similar to sentiment classification, intent classification is also a core task of text classification. For intent classification, we select SciCite [48] and ACL-ARC [48] for our experiments.

(3) **Relation Extraction** is the task of predicting attributes and relations for entities in a sentence. For example, given a sentence "Barack Obama was born in Honolulu, Hawaii.", a relation classifier aims at predicting the relation of "bornInCity". Relation extraction is the key component for building relational knowledge graphs, and it is of crucial significance to natural language understanding applications, such as structured search, question answering, and summarization. For relation extraction, we select SciERC [49] and ChemProt [50] for our experiments. More details of these datasets are shown in Table 1.

To build few-shot learning settings, we randomly sample a part of instances from the dataset as the training set 6

 $\mathcal{T}$ . Additionally, we prepare the open-domain corpora  $\mathcal{O}$  consisting of unused instances that share the same classes with  $\mathcal{T}$  from all downstream datasets. We also add English Wikipedia, which is used to train the original BERT<sub>BASE</sub> [1] and RoBERTa<sub>BASE</sub> [2], into  $\mathcal{O}$ . CSS-LM will leverage the open-domain corpora  $\mathcal{O}$  to perform contrastive semi-supervised learning; thus, to fairly compare baseline PLMs such as BERT<sub>BASE</sub> and RoBERTa<sub>BASE</sub>, the baseline PLMs are all previously pre-trained on the open-domain corpora  $\mathcal{O}$  in the experiments.

# 4.2 Experimental Settings

We choose BERT<sub>BASE</sub> [1] and RoBERTa<sub>BASE</sub> [2] as our encoders, using the official released parameters  $^{1,2}$ . The other parameters of CSS-LM are all initialized randomly.

For training, we set the learning rate as  $2 \times 10^{-5}$ , the batch size as 4. The remaining settings follow the original ones of BERT<sub>BASE</sub> and RoBERTa<sub>BASE</sub>. For the retrieved instance size, we perform a grid search over multiple hyper-parameters {16, 32, 48, 64}, and take the best one measured on the whole development set for our model.

The objective function of CSS-LM is  $\mathcal{L} = \mathcal{L}_{CE} + \mathcal{L}_{CSS} + \mathcal{L}_D + \mathcal{L}_C$ , which enables CSS-LM to retrieve task-related instances and enhance the performance of downstream tasks; however,  $\mathcal{L}_D$  and  $\mathcal{L}_C$  terms may make the performance drop on some downstream tasks during contrastive semi-supervised learning sometimes. Thus, we degrade the objective function  $\mathcal{L} = \mathcal{L}_{CE} + \mathcal{L}_{CSS} + \mathcal{L}_D + \mathcal{L}_C$  to  $\mathcal{L} = \mathcal{L}_{CE} + \mathcal{L}_{CSS}$  when the performance starts to drop in the development sets. Then, we continuously train CSS-LM with the new objective function.

Besides, to avoid the result instability [51], [52], we report the average performance across 5 different randomly sampled data splits.

# 4.3 Baselines

We compare our CSS-LM with the following effective finetuning strategies:

**Standard fine-tuning (Standard)**, which is the typical fine-tuning method in PLMs [1], [2] and its training objective is  $\mathcal{L}_{CE}$  mentioned in Eq. (8).

Supervised contrastive fine-tuning (SCF) [13], which only performs supervised contrastive learning with the supervised data of downsteam tasks, and its fine-tuning loss is  $\mathcal{L}_{CE} + \mathcal{L}_{CS}$ , referring to Eq. (1) and Eq. (8).

**CSS-LM-ST**, a variant of our CSS, which also retrieves task-related in-domain data with domain-related and classrelated relatedness (i.e., the positive instances) by contrastive semi-supervised learning as we mentioned in the section 3.2. The difference is CSS-LM-ST performing the standard finetuning method with pseudo labeling [53], which is the simple and efficient semi-supervised learning method for deep neural networks, instead of contrastive semi-supervised learning to capture critical features from the retrieved instances. We can denote the learning objective of CSS-LM-ST as  $\mathcal{L}_D + \mathcal{L}_C + \mathcal{L}'_{CE}$ , where  $\mathcal{L}'_{CE}$  is the same downstream task

<sup>1.</sup> https://storage.googleapis.com/bert\_models/2020\_02\_20/ uncased L-12 H-768\_A-12.zip

<sup>2.</sup> https://dl.fbaipublicfiles.com/fairseq/models/roberta.base.tar.gz

TABLE 1

The details of the datasets used in our experiments. To build few-shot settings, we sample  $N = K_y \times |\mathcal{Y}|$  instances from the original training set, where  $K_y$  is the sampled instance number for each class and  $|\mathcal{Y}|$  is the number of class types.

Task	Domain	Dataset	$ \mathcal{Y} $	#Train	#Test	#Dev	Class types
Sentiment	Review	SemEval	3	4,665	2,426	4,665	positive, neutral, negative
Classification	Review	SST-5	5	8,544	2,210	1,101	v. pos., positive, neutral, negative, v. neg.
Intent	Multi	Scicite	3	7,320	1,861	916	result, method, background
Classification	CS	ACL-ARC	6	1,688	139	114	background, uses, motivation, compareOrcontrast, extends, future
Relation	CS	SciERC	7	3,219	974	455	part-of, conjunction, hyponymy, used-for, feature-of, compare, evaluate-for
Classification	BIO	ChemProt	13	4,169	3,469	2,422	substrate, antagonist, indirect-upregulator, activator, indirect-downregulator, inhibitor, upregulator, downregulator, product-of, agonist, agonist-activator

#### TABLE 2

The results (%) of various fine-tuning methods on six classification datasets. All fine-tuning strategies are applied on RoBERTa<sub>BASE</sub> and BERT<sub>BASE</sub> models and set  $K_y = 16$  for the few-shot experimental settings.

Task	Sentiment Classification				Intent Classification				Relation Classification			
Dataset	SemEval SST		-5   SciC		ite	ACL-A	RC	SciEF	RC	ChemI	Prot	
Base Model	RoBERTa	BERT	RoBERTa	BERT	RoBERTa	BERT	RoBERTa	BERT	RoBERTa	BERT	RoBERTa	BERT
Fine-tune on the whole training set												
Standard [2] SCF [13]	89.1 88.9	87.4 86.9	56.8 57.4	54.1 53.6	86.0 <b>86.5</b>	84.8 <b>85.4</b>	81.3 <b>84.2</b>	77.5 77.8	88.7 87.5	88.6 87.7	82.3 81.8	80.6 81.0
CSS-LM-ST CSS-LM	89.2 <b>89.5</b>	87.0 <b>87.7</b>	57.5 57.5	53.5 <b>54.8</b>	86.0 86.0	85.0 <b>85.4</b>	82.0 <b>84.2</b>	78.2 <b>78.9</b>	88.7 <b>89.0</b>	88.4 <b>89.0</b>	<b>82.4</b> 82.3	81.0 <b>81.9</b>
Fine-tune on few-shot setting $(K_y = 16)$												
Standard [2] SCF [13]	68.9 69.1	63.9 65.3	39.4 <b>39.6</b>	<b>36.1</b> 35.5	74.3 75.8	75.0 75.0	45.7 <b>50.9</b>	43.7 47.1	46.9 52.0	40.5 40.3	47.6 46.9	44.1 45.4
CSS-LM-ST CSS-LM	71.1 <b>73.0</b>	70.0 70.0	39.5 39.5	36.0 <b>36.1</b>	75.9 <b>77.5</b>	76.9 <b>77.4</b>	50.4 48.8	46.6 <b>47.5</b>	53.0 <b>54.7</b>	46.1 <b>47.4</b>	48.2 49.0	47.6 <b>48.3</b>

fine-tuning loss in the section 3.6 but leverages extra pseudo labeled instances of open-domain corpora.

**CSS-LM** and **CSS-LM-ST** leverage the same method to retrieve task-related instances but apply different mechanisms to learn semantic features from the retrieved instances. CSS-LM learns semantic features by contrastive semi-supervised learning. CSS-LM-ST uses the conventional semi-supervised learning with pseudo labels to learn the features, which may be sensitive to the quality of labels. The key advantage of CSS-LM is free to pseudo labels. In the following parts, we will extensively study CSS-LM and CSS-LM-ST.

## 4.4 Overall Results

We conduct experiments on the three selected tasks under the standard and few-shot settings ( $K_y = 16$ ). The results are shown in Table 2. From the table, we can see that:

(1) Our CSS-LM framework achieves improvements on almost all six datasets compared to the baseline models (including the state-of-the-art SCF), especially under the few-shot settings. This demonstrates that our contrastive semi-supervised framework for fine-tuning could effectively capture the important semantic features for the task from the large-scale unlabeled data.

(2) Although utilizing the retrieved task-related indomain data can help fine-tuning, CSS-LM-ST outperforms standard fine-tuning yet does not obtain the same performance improvements as CSS-LM under the few-shot settings. It indicates that our framework can retrieve high-quality instances by contrastive semi-supervised learning; applying contrastive semi-supervised learning is better than assigning pseudo labels to the retrieved instances to train classifiers. In fact, under the few-shot settings, some retrieved instances will be linked to some classes close to their implicit golden labels rather than the golden labels. Therefore, directly annotating pseudo labels may lead to a biased model. In contrast, these instances corresponding to wrong classes can still provide correlation information for the contrastive learning of CSS-LM.

(3) We also compare CSS-LM under two different settings: fine-tuning with the few-shot instances and the whole training set instances. Although CSS-LM outperforms most baseline models under the few-shot setting, CSS-LM obtains slight improvements when fine-tuning with sufficient training data. This is an intuitive observation that directly fine-tuning PLMs will work well when the amount of data is sufficient. However, not all NLP tasks have enough data; low-resource tasks and long-tail classes are very common. Our framework can improve the model performance in few-shot learning scenarios without weakening the model performance.



Fig. 3. The results (%) of different fine-tuning strategies on the development sets of SemEval, Scicite and SciERC, with different numbers of instance per class.

#### 4.5 Effect of Supervised Data Size

In this part, we explore the effect of the supervised instance number  $K_y$  for each class. From Fig. 3, we have the following findings:

(1) Under the few-shot settings, both SCF and CSS-LM outperform the standard fine-tuning strategy on two datasets consistently. It indicates that performing contrastive learning between instances of different classes could help extract informative semantic features to distinguish them, and benefit downstream tasks in the fine-tuning stage.

(2) CSS-LM has better results than SCF on all the datasets of our experiments, especially when the supervised data size is small. It demonstrates that although SCF could discover the class relatedness from the supervised data to some extent, it may ignore semantic information that is trivial in the limited supervised data but crucial for the unseen task. In contrast, our CSS-LM could effectively take account of the unlabeled data to capture the informative semantic features not expressed by supervised data.

## 4.6 Effect of Retrieved Instance Size

In order to learn discriminative features, CSS-LM needs to leverage proper positive and negative pairs. As for CSS-LM-ST, pseudo labeling quality is essential for the performance. In this part, we explore the effect of different retrieved instance size on CSS-LM and CSS-LM-ST. As shown in Fig. 4, we have the following findings:

(1) When the size is small, CSS-LM retrieves too similar instances to make them close to each other, according to the contrastive learning paradigm, which cannot learn the essential information. As the size increases, the model

performance can gradually increase. Since more informative instances are retrieved from open-domain corpora, CSS-LM will take unrelated instances as similar instances to learn the false information when the size is too large. In future, it is meaningful to study how to denoise the retrieved instances for semi-supervised fine-tuning.

(2) When the size becomes large, the performance of CSS-LM-ST decays earlier than CSS-LM since CSS-LM-ST needs to learn features from high-quality labeled instances. Instead of learning the features from an individual instance, CSS-LM learns by comparing among different instances. Therefore, CSS-LM can better leverage sufficient unlabeled instances.

## 4.7 Difference between Retrieving Instances from Indomain Data and Open-domain Data

To show our framework can automatically extract in-domain instances from open-domain corpora, we perform CSS-LM<sub>[D+C]</sub> for the open-domain corpora, and perform CSS-LM<sub>[C]</sub> for the in-domain corpora. Specifically, the in-domain corpora contain some unused instances share same classes with  $\mathcal{T}$ . CSS-LM<sub>[D]</sub> is the domain-level term of our CSS-LM and CSS-LM<sub>[C]</sub> is the class-level term, CSS-LM<sub>[D+C]</sub> is the combination of <sub>[D]</sub> and <sub>[C]</sub>. The results are shown in Table 3.

From the table, we can see that: CSS-LM<sub>[D+C]</sub> can achieve comparable performance with CSS-LM<sub>[C]</sub> trained on the in-domain corpora. Although directly retrieving in-domain data is easier than retrieving open-domain data, we can see CSS-LM<sub>[C]</sub> is only slightly better than CSS-LM<sub>[D+C]</sub> in the well-defined domain such as SemEval (Restaurant review) and ChemProt (Biology); CSS-LM<sub>[D+C]</sub> even outperforms CSS-LM<sub>[C]</sub> in SciCite, which belongs to multiple domains



Fig. 4. The effect (%) of retrieved instance size on CSS-LM and CSS-LM-ST on the development sets of SemEval, Scicite and SciERC.

(not well-defined). It demonstrates that CSS-LM can efficiently learn to distinguish coarse-grained domains from a large amount of the open-domain corpora, free to domain dependence.

# 4.8 Contribution of Domain-Related and Class-Related Semantics to CSS-LM

As we empirically apply domain-related and class-related semantics for our framework, we wonder which semantic levels makes the greater contribution to our framework. In this part, we study the effect of the domain-level and classlevel terms of CSS-LM as shown in Table 4.

From the results, we can find that:  $CSS-LM_{[D+C]}$  is comparable to  $CSS-LM_{[C]}$  and outperforms  $CSS-LM_{[D]}$  in almost all datasets. Therefore, compared with the domain-level term, the class-level term is essential to retrieving high-quality instances to enhance CSS-LM performance.

Although this paper mainly focuses on applying the contrastive semi-supervised framework for fine-tuning rather than exploring retrieving informative instances, we think how to better consider domain-level semantic remains an interesting problem in the future.

TABLE 3

The results (%) of CSS-LM on the open-domain corpora and in-domain corpora. As in-domain corpora do not consider domain-level semantics, we perform CSS-LM<sub>[C]</sub> on in-domain corpora.

Fine-tune on few-shot $K_y = 16$						
Task	Sentiment Classification	Intent Classification	Relation Classification			
Dataset	SemEval	SciCite	SciERC			
RoBERTa <sub>BASE</sub>						
Open-domain corpora						
$CSS-LM_{[D+C]}$	73.0	77.5	54.7			
In-domain corpora						
$CSS-LM_{[C]}$	73.0	75.9	54.9			
BERT <sub>BASE</sub>						
Open-domain corpora						
$CSS-LM_{[D+C]}$	70.0	77.4	47.1			
In-domain corpora						
CSS-LM <sub>[C]</sub>	71.0	76.1	47.4			

TABLE 4 The results (%) of CSS-LM utilizing semantic relatedness at different levels.

Fine-tune on few-shot $K_y = 16$ in open-domain corpora							
Task	Sentiment Classification	Intent Classification	Relation Classification				
Dataset	SemEval	SciCite	SciERC				
RoBERTa <sub>Base</sub>							
CSS-LM <sub>[C]</sub>	73.6	75.1	53.6				
CSS-LM <sub>[D]</sub>	72.1	75.7	51.9				
$CSS-LM_{[D+C]}$	73.0	77.5	54.7				
BERT <sub>BASE</sub>							
CSS-LM <sub>[C]</sub>	69.8	75.8	46.0				
CSS-LM <sub>[D]</sub>	68.1	74.4	45.7				
$CSS-LM_{[D+C]}$	70.0	77.4	47.7				

#### 4.9 Visualization

We apply t-SNE [54] to visualize instance embeddings of the SemEval test set learned by standard, SCF, and CSS-LM fine-tuning strategies in Fig. 5. We give the two-dimensional points with different colors to represent its corresponding label in the downstream task.

From the figure, we can intuitively see that direct finetuning with sufficient data Fig. 5(a) has the best boundaries among all classes. Without sufficient data, all fine-tuning methods cannot detect neutral instances, but our CSS-LM can obtain the coarse-grained neutral cluster and achieve a better decision boundary between positive and negative instances.



Fig. 5. The tSNE plots of the embeddings learned by standard, SCF, and CSS-LM fine-tuning methods on the SemEval dataset. Green: Negative emotion; Blue: Positive emotion; Red: Neutral emotion. Except the embeddings in (a) are trained by the whole training set, the embeddings in (b) (c) (d) are only trained by  $K_y = 16$ .

TABLE 5 The analyses of some retrieved instances from the corpora by CSS-LM. Yellow: Class-related; Gray: Domain-related.

Class	Instances	Related Instances
Negative	Service was just ok, it is not what you'd expect for \$500.	I'm not even going to bother to describe it; speaks for itself. It's only \$1.95 for a regular slice and 4.00 for a slice with a mushroom , not mushrooms.
Neutral	A great way to make some money is to buy a case of snapple from costco and sell it right outside for only \$2.50.	Try the times square cocktail – ginger lemonade with vodka. I had the tuna tartare with sake, mushroom ravioli with pinot noir, and the chocolate sampler [].
Positive	The ambience was so fun, and the prices were great , on top of the fact that the food was really tasty .	Overall, I would highly recommend giving this one a try. [] Liverpool boss Klopp says win feels perfect. Oh yeah ever on the west side try there sister restaurant arties cafe.

# 4.10 Case Study

As shown in Table 5, we also give a case study to investigate whether CSS-LM can capture domain-level and class-level relatedness of instances. We color the sub-sequence containing human-annotated class-level information with yellow and domain-level information with gray.

From the table, we can see that CSS-LM can retrieve instances highly agreed with humans in most cases. Interestingly, the neutral class did not exist any apparent classlevel information recognized by humans, but CSS-LM can distinguish these instances in the sentence-level; however, CSS-LM takes the instance, ... *there sister restaurant arties cafe*, as the positive, which can be easily classified to the neutral class by humans with sentence-level meaning. To consider the sentence-level information during retrieving by contrastive learning is a future work we can study.

# 5 CONCLUSION AND FUTURE WORK

In this work, we introduce the CSS-LM framework to improve the fine-tuning phase of PLMs via contrastive semi-supervised learning. The experimental results on three typical text classification tasks show that CSS-LM could effectively capture crucial semantic features for downstream tasks with limited supervised data and achieve better performances than the conventional, supervised contrastive finetuning strategies. In future, we will explore the following promising directions:

(1) The CSS-LM framework makes an initial attempt to better fine-tune PLMs with limited supervised data of downstream tasks in text classification. Extending it to other NLP tasks, e.g., text generation, name entity recognition, and question answering, is a valuable direction.

(2) From our experimental results, we find that better methods to consider domain-level semantic remain a further exploration. Finding out a better strategy to retrieve and denoise instances is also a fascinating problem.

(3) CSS-LM is devoted to leveraging unannotated data from the open-domain corpora to capture crucial semantic features, which benefit the downstream tasks, instead of considering invariant features between different tasks or domains to perform transfer learning. How to better leverage these invariant features is a direction we can explore.

# ACKNOWLEDGMENTS

This work is supported by the National Key Research and Development Program of China (No. 2020AAA0106501) and the National Natural Science Foundation of China (NSFC No. 61772302).

#### REFERENCES

- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pretraining of deep bidirectional transformers for language understanding," in *Proceedings of NAACL*, 2019, pp. 4171–4186.
- [2] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," in *Proceedings of ICLR*, 2019.
- [3] Y. Zhang, S. Sun, M. Galley, Y.-C. Chen, C. Brockett, X. Gao, J. Gao, J. Liu, and B. Dolan, "Dialogpt: Large-scale generative pre-training for conversational response generation," in *Proceedings of ACL*, 2020, pp. 270–278.
- [4] J. Zhang, Y. Zhao, M. Saleh, and P. J. Liu, "Pegasus: Pre-training with extracted gap-sentences for abstractive summarization," in *Proceedings of ICML*, 2019, pp. 11 328–11 339.
- Proceedings of ICML, 2019, pp. 11 328–11 339.
  [5] Y. Liu and M. Lapata, "Text summarization with pretrained encoders," in *Proceedings of EMNLP*, 2019, pp. 3730–3740.
- [6] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, and O. Levy, "Spanbert: Improving pre-training by representing and predicting spans," in *arXiv*, 2019.
- [7] D. Adiwardana, M.-T. Luong, D. R. So, J. Hall, N. Fiedel, R. Thoppilan, Z. Yang, A. Kulshreshtha, G. Nemade, Y. Lu, and Q. V. Le, "Towards a human-like open-domain chatbot," in *arXiv*, 2020.
- [8] L. Baldini Soares, N. FitzGerald, J. Ling, and T. Kwiatkowski, "Matching the blanks: Distributional similarity for relation learning," in *Proceedings of ACL*, 2019, pp. 2895–2905.
- [9] H. Peng, T. Gao, X. Han, Y. Lin, P. Li, Z. Liu, M. Sun, and J. Zhou, "Learning from context or names? an empirical study on neural relation extraction," in *Proceedings of EMNLP*, 2020, pp. 3661–3672.
- [10] J.-C. Su, S. Maji, and B. Hariharan, "When does self-supervision improve few-shot learning?" in *Proceedings of ECCV*, 2020.
- [11] H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and T. Zhao, "Smart: Robust and efficient fine-tuning for pre-trained natural language models through principled regularized optimization," in *Proceedings of ACL*, 2020, pp. 2177–2190.
- [12] J. Du, E. Grave, B. Gunel, V. Chaudhary, O. Celebi, M. Auli, V. Stoyanov, and A. Conneau, "Self-training improves pre-training for natural language understanding," in *arXiv*, 2020.
- [13] G. Beliz, D. Jingfei, C. Alexis, and S. Veselin, "Supervised contrastive learning for pre-trained language model fine-tuning," in *Proceedings of ICLR*, 2021.
- [14] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, "Xlnet: Generalized autoregressive pretraining for language understanding," in *Proceedings of NeurIPS*, 2019.
- [15] Z. Zhang, X. Han, Z. Liu, X. Jiang, M. Sun, and Q. Liu, "ERNIE: Enhanced language representation with informative entities," in *Proceedings of ACL*, 2019, pp. 1441–1451.
- [16] Y. Su, X. Han, Z. Zhang, P. Li, Z. Liu, Y. Lin, J. Zhou, and M. Sun, "Cokebert: Contextual knowledge selection and embedding towards enhanced pre-trained language models," in *arXiv*, 2020.
- [17] M. E. Peters, M. Neumann, R. Logan, R. Schwartz, V. Joshi, S. Singh, and N. A. Smith, "Knowledge enhanced contextual word representations," in *Proceedings of EMNLP*, 2019, pp. 43–54.
- [18] Y. Sun, S. Wang, Y. Li, S. Feng, X. Chen, H. Zhang, X. Tian, D. Zhu, H. Tian, and H. Wu, "Ernie: Enhanced representation through knowledge integration," in *arXiv*, 2019.
- [19] D. Ye, Y. Lin, J. Du, Z. Liu, P. Li, M. Sun, and Z. Liu, "Coreferential reasoning learning for language representation," in *Proceedings of EMNLP*, 2020, pp. 7170–7186.
- [20] G. Lample and A. Conneau, "Cross-lingual language model pretraining," in *Proceedings of NeurIPS*, 2019.
- [21] L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang, "Visualbert: A simple and performant baseline for vision and language," in arXiv, 2019.
- [22] G. Li, N. Duan, Y. Fang, M. Gong, and D. Jiang, "Unicoder-vl: A universal encoder for vision and language by cross-modal pretraining," in *Proceedings of AAAI*, 2020.
- [23] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, "VI-bert: Pretraining of generic visual-linguistic representations," in *Proceedings* of ICLR, 2020.
- [24] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining taskagnostic visiolinguistic representations for vision-and-language tasks," in *Proceedings of NeurIPS*, 2019.
- [25] H. Tan and M. Bansal, "Lxmert: Learning cross-modality encoder representations from transformers," in *Proceedings of EMNLP*, 2019, pp. 5100–5111.

- [26] I. Beltagy, K. Lo, and A. Cohan, "Scibert: A pretrained language model for scientific text," in *Proceedings of EMNLP*, 2019, pp. 3615– 3620.
- [27] K. Huang, J. Altosaar, and R. Ranganath, "Clinicalbert: Modeling clinical notes and predicting hospital readmission," in arXiv, 2019.
- [28] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "Biobert: a pre-trained biomedical language representation model for biomedical text mining," in *Proceedings of Bioinformatics*, 2019, pp. 1234–1240.
- [29] J. Howard and S. Ruder, "Universal language model fine-tuning for text classification," in *Proceedings of ACL*, 2018, pp. 328–339.
- [30] M. E. Peters, S. Ruder, and N. A. Smith, "To tune or not to tune? adapting pretrained representations to diverse tasks," in *Proceedings* of ACL, 2019, pp. 7–14.
- [31] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, "Parameter-efficient transfer learning for NLP," in *Proceedings of ICML*, 2019, pp. 2790– 2799.
- [32] A. C. Stickland and I. Murray, "Bert and pals: Projected attention layers for efficient adaptation in multi-task learning," in *Proceedings* of ICML, 2019, pp. 5986–5995.
- [33] Y. Gu, Z. Zhang, X. Wang, Z. Liu, and M. Sun, "Train no evil: Selective masking for task-guided pre-training," in *Proceedings of EMNLP*, pp. 6966–6974.
- [34] S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith, "Don't stop pretraining: Adapt language models to domains and tasks," in *Proceedings of ACL*, 2020.
- [35] L. Logeswaran and H. Lee, "An efficient framework for learning sentence representations," in *Proceedings of ICLR*, 2018.
- [36] L. Kong, C. de Masson d'Autume, L. Yu, W. Ling, Z. Dai, and D. Yogatama, "A mutual information maximization perspective of language representation learning," in *Proceedings of ICLR*, 2019.
- [37] Y. Zhang, R. He, Z. Liu, K. H. Lim, and L. Bing, "An unsupervised sentence embedding method by mutual information maximization," in *Proceedings of EMNLP*, 2020, pp. 1601–1610.
- [38] Z. Wu, S. Wang, J. Gu, M. Khabsa, F. Sun, and H. Ma, "Clear: Contrastive learning for sentence representation," in *arXiv*, 2020.
- [39] T. Chen, K. Simon, N. Mohammad, and H. Geoffrey, "A simple framework for contrastive learning of visual representations," in *Proceedings of PMLR*, 2020, pp. 1597–1607.
- [40] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of CVPR*, 2020.
- [41] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," in *Proceedings of ICLR*, 2020.
- [42] A. Kolesnikov, X. Zhai, and L. Beyer, "Revisiting self-supervised visual representation learning," in *Proceedings of CVPR*, 2019, pp. 1920–1929.
- [43] A. Mnih and K. Kavukcuoglu, "Learning word embeddings efficiently with noise-contrastive estimation," in *Proceedings of NeurIPS*, 2013.
- [44] M. U. Gutmann and A. Hyvärinen, "Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics," in *Proceedings of JMLR*, 2012, pp. 307–361.
- [45] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," in *Proceedings of NeurIPS*, 2020, pp. 18661–18673.
- [46] W. Xue and T. Li, "Aspect based sentiment analysis with gated convolutional networks," in *Proceedings of ACL*, 2018, pp. 2514–2523.
- [47] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proceedings of EMNLP*, 2013, pp. 1631–1642.
- [48] A. Cohan, W. Ammar, V. Madeleine, Zuylen, and F. Cady, "Structural scaffolds for citation intent classification in scientific publications," in *Proceedings of NAACL*, 2019, pp. 3586–3596.
- [49] Y. Luan, L. He, M. Ostendorf, and H. Hajishirzi, "Multi-task identification of entities, relations, and coreferencefor scientific knowledge graph construction," in *Proceedings of EMNLP*, 2018, pp. 3219–3232.
- [50] I. Beltagy, K. Lo, and A. Cohan, "Scibert: Pretrained language model for scientific text," in *Proceedings of EMNLP*, 2019, pp. 3615–3620.
- [51] J. Dodge, G. Ilharco, R. Schwartz, A. Farhadi, H. Hajishirzi, and N. Smith, "Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping," in *arXiv*, 2020.

JOURNAL OF LATEX CLASS FILES, VOL. 14, NO. 8, AUGUST 2015

- [52] T. Zhang, F. Wu, A. Katiyar, K. Weinberger, and Y. Artzi, "Revisiting few-sample bert fine-tuning," in *Proceedings of ICLR*, 2021.
  [53] D.-H. Lee, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," in *Proceedings of ICML*, 2012. 2013.[54] D. Ulyanov, "Multicore-tsne," in *GitHub*, 2016.