

InterCSE: Sentences Representation with Interactive Framework of Supervised Sentence Representation

Anonymous EMNLP submission

Abstract

Contrastive learning has achieved remarkable results in sentence representation, but its semantic representation remains independent in the process of training and inference, and could not pay attention to the interactive information of sentence pairs. Therefore, this paper proposes to introduce a multi-task contrastive learning method, which not only focuses on the sorting of sentence pairs embedding similarity, but also increases the interaction information as a supplement to the sentence embedding representation. Meanwhile, we propose a loss decaying strategy to balance multiple loss function. We evaluate the performance of InterCSE on standard semantic textual similarity (STS) tasks, and experiments show that our model using $BERT_{base}$ and $BERT_{large}$ achieve 82.11% and 82.88% spearman’s correlation, 0.54% and 0.43% improvement compared to SimCSE respectively. We also conduct experiments compared to Sentence-Transformers adding interactive network, which get 85.18%(+0.88% $BERT_{base}$) spearman’s correlation. Hence, adding interactive features to the traditional siamese network performed very well, and achieved the effect of state-of-the-art on sentence representation tasks. Our code is available at <https://github.com/2hip3ng/InterCSE>.

1 Introduction

Sentence representation learning is a vital component of natural language processing tasks (Cer et al., 2017). The rapid development of sentence representation technology has made a wide range of downstream tasks more intelligent, especially information retrieval and text clustering.

Recently, the pre-trained language model has become the cornerstone of natural language processing technology, such as BERT (Devlin et al., 2019; Liu et al., 2019), GPT (Radford et al., 2018, 2019; Brown et al., 2020), ERNIE (Sun et al., 2019a,b, 2021), which greatly affects the development of

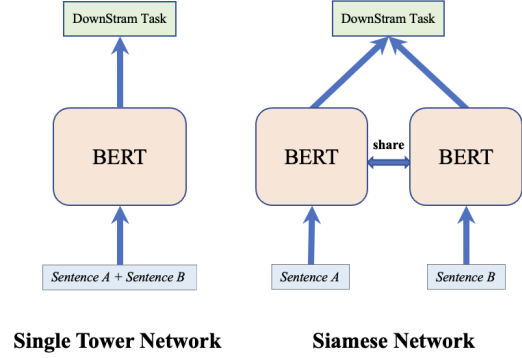


Figure 1: The architecture of single tower network(left) and siamese network(right).

various downstream tasks. Sentence representation could not get high performance directly with pre-trained language model because of anisotropic phenomenon (Gao et al., 2019), but contrastive learning play the role of a bridge. Many sentence representation approaches are transformed into point-wise classification tasks (Reimers and Gurevych, 2019), and there is a gap between the training optimization objectives and good sentence representation. Comparative learning applies the sorting method with InfoNCE (van den Oord et al., 2019) contrastive loss function, which is more suitable for the optimization goal of sentence representation.

The sentence representation model (Gao et al., 2022; Yan et al., 2021) of contrastive learning combined with pre-trained model is a siamese network, which encodes all sentences independently. The structure of the siamese network makes it possible to quickly produce the vector representation of the sentences and calculate the similarity during training and prediction, so as to complete the retrieval or sentence clustering task in massive data. However, the structure of independent encoding makes the sentences pair lose the interactive information, which reduces the accuracy rate.

There are obvious differences between single

| Network Type | Model | Spearman |
|----------------------|------------------------------|----------|
| Siamese Network | SBERT-base ♡ | 84.7 |
| | SBERT-large ♡ | 84.5 |
| | SimCSE-BERT-base ♠ | 84.3 |
| | SimCSE-BERT-large(reproduce) | 85.4 |
| | SROBERTa-base ♡ | 84.9 |
| | SROBERTa-large ♡ | 85.0 |
| | SimCSE-RoBERTa-base ♠ | 85.8 |
| | SimCSE-RoBERTa-large ♠ | 86.7 |
| Single Tower Network | BERT-base ♦ | 85.8 |
| | BERT-large ♦ | 86.5 |
| | RoBERTa-base ◇ | 87.2 |
| | RoBERTa-large ◇ | 88.1 |

Table 1: Comparison of single tower network and siamese network on STS-Benchmark(Cer et al., 2017) test set with spearman’s correlation. ♡: results from Sentence-Transformers(Reimers and Gurevych, 2019), ♠: results from SimCSE(Gao et al., 2022), ♦ : results from BERT (Devlin et al., 2019), ◇: the performance of RoBERTa-base and RoBERTa-large in single tower network are reproduce through fairseq(Ott et al., 2019). SimCSE models are trained on NLI datasets(Bowman et al., 2015), and the other are trained on STS-Benchmark train set. All results are rounded to one decimal place for comparison.

tower network and siamese network. The network structure is depicted in Figure 1. The single tower network is like the original BERT(Devlin et al., 2019) model structure. After concating the sentence pairs, the semantic features of the sentence pairs are extracted and input to the downstream classification or regression network. The siamese network(Koch et al., 2015) has two encoders that encode each sentence individually. After encoding, it can perform similarity calculation or put into the downstream network using sentence embedding, such as Sentence-Transformers(Reimers and Gurevych, 2019). During training, the two encoders share model parameters.

The performance of the single-tower network and the siamese network on the STS Benchmark test dataset is shown in Table 1. On the whole, the single tower network has a good improvement in spearman’s correlation index compared with the siamese network. We consider that when the single tower network encodes a sentence pair, the sentence is not an independent individual. It will refer to its counterparts for encoding and use the attention mechanism to extract features, which makes the single tower network in the sentence pair similarity task has a higher correlation coefficient. Recently, most sentence representation tasks use siamese networks. Although it could bring computational advantages on massive data, it inevitably reduces the accuracy of correlation.

In order to take full advantage of the accuracy

advantages of single tower networks and the inference speed advantages of siamese networks, this paper proposes a multi-task contrastive learning method. During the training process, on the basis of SimCSE, we added a single tower network as a supplement to form a framework for multi-task learning. This method can fully obtain the interaction information between sentence pairs during the process of training sentence representation.

Our contributions can be summarized as follows:

1. We propose an effective framework of multi-task contrastive learning for sentence representation tasks which could introduce sentence interactive semantic information.
2. We design a loss function for the proposed framework. When the interactive network introduce information increment, try to minimize the hurt to the original sentence representation.
3. Experiments show that our approach achieves new state-of-the-art performance on STS tasks.

2 Related Work

2.1 Large Language Model

Recently, transformer(Vaswani et al., 2017) structure shines in the field of deep learning and lays the foundation for large models. The attention mechanism is good at capturing the semantic relationship

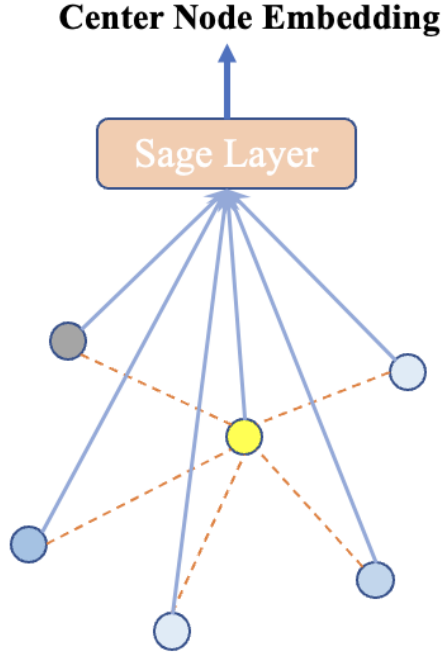


Figure 2: The architecture of GraphSage. Sage layer aggregates all neighbor nodes’ embedding for center node.

between sequences. Using transformer’s encoder, many language models have been born, such as GPT(Radford et al., 2018), BERT(Devlin et al., 2019), RoBERTa(Liu et al., 2019), XLNet(Yang et al., 2019), Ernie(Sun et al., 2019a), etc. These models mainly have some differences in masking mechanism, pre-training data, and pre-training methods. Among them, the BERT model mainly masks 15% words in the sentence and predicts the origin words as a unsupervised pre-training task. After obtaining the pre-trained model, NLP downstream tasks only need to complete fine-tuning on the model to achieve high performances and set new state-of-the-art results, including sentence classification, sequence tagging, question answering, machine reading and comprehension and sentence-pair classification or regression.

2.2 Contrastive Learning and Sentence Representation

Contrastive learning was first proposed in the field of computer vision. It mainly wants to optimize the similarity of sample pairs in a representation space, and make similar sample pairs gather and dissimilar sample pairs stay away. The core idea is that in a batch of sample pairs, there is only one pair

| Type of NLI datasets | Numbers |
|--------------------------|---------|
| Entailment | 314k |
| Neutral | 314k |
| Contradiction | 314k |
| Entailment+Contradiction | 270k |

Table 2: Statistics of different type of SNLI+MNLI datasets.

of samples that is positive samples, and the others are negative samples. During the training process, after the sample representation is obtained using siamese network embedding, the sample similarity is calculated to maximize the similarity of the positive samples. This optimization scenario is very suitable for unsupervised scenarios. Usually, positive samples can be obtained through simple data enhancement, and a large number of negative samples can be obtained through negative sampling.

In terms of sentence representation, the idea of contrastive learning continues to be used, leading to many research results, such as ConSERT(Yan et al., 2021), SimCSE(Gao et al., 2022) and ES-imCSE(Wu et al., 2022). ConSERT propose four different data augmentation strategies to generate views for contrastive learning, including adversarial attack, token shuffling, cutoff and dropout. SimCSE first describe an unsupervised approach, which takes an input sentence and predicts itself in a contrastive objective, with only standard dropout used as noise. Only using dropout as data augmentation become a popular method for contrastive learning. ESImCSE introduce two modifications based on SimCSE. It apply a simple repetition operation to modify the input sentence, and then pass the input sentence and its modified counterpart to the pre-trained Transformer encoder, respectively, to get the positive pair. And it introduce a momentum contrast, enlarging the number of negative pairs.

2.3 Multi-Task Learning

Multi-Task Learning (MTL) is an important research topic in machine learning, which aims to learn multiple tasks simultaneously. In MTL, multiple tasks are divided into multiple learning units, and a learner can learn multiple tasks by making multiple small models.

One of the important research directions in MTL is to develop efficient algorithms and theoretical models. In recent years, many works have been

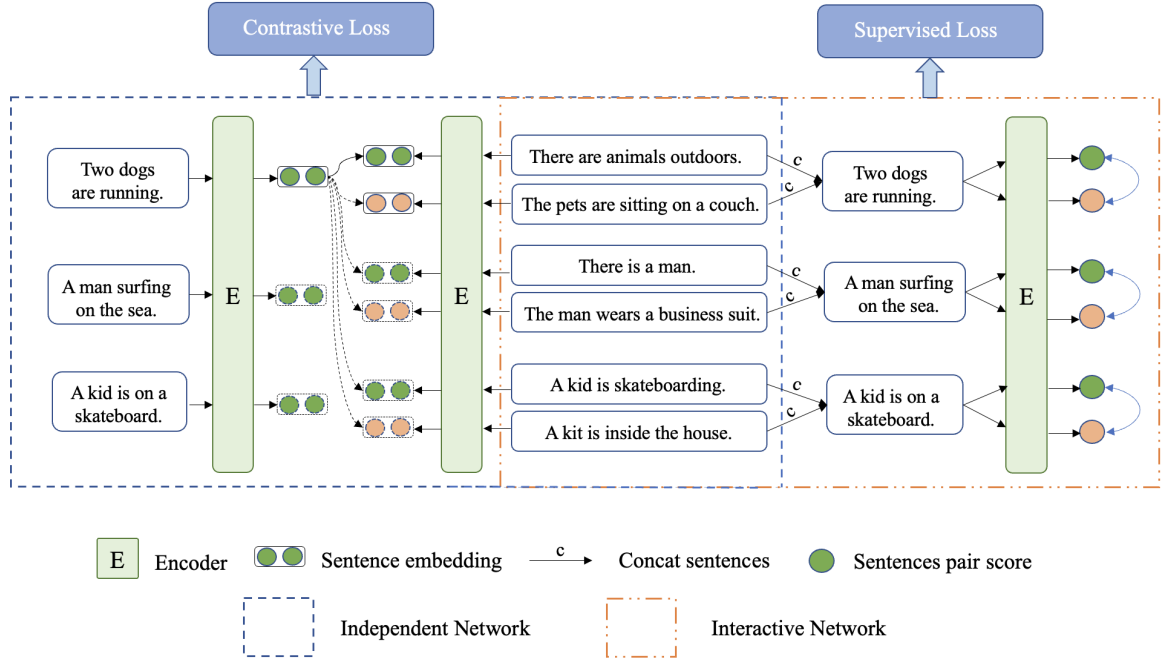


Figure 3: The architecture of InterCSE. Independent network learn the representation of sentence, and interactive network learn sentences interactive semantic.

carried out in this direction, including deep neural networks, attention mechanism, joint attention, graph attention, and so on. MTL has been successfully used in natural language processing applications, including text classification(Liu et al., 2017), machine translation(Luong et al., 2016), sequence labeling(Rei, 2017) and sentence representations(Ahmad et al., 2018). These algorithms and theoretical models can help MTL learn multiple tasks with better performance.

3 Approach

In this section, we present InterCSE (introducing **Interactive** semantic information in **Contrastive Sentences Embedding**) for sentence representation task. Firstly, we present the overall framework of our approach. Then, we introduce the training dataset for our model. Finally, we talk about the combination strategies of loss function for multiple objective.

3.1 Framework

Our approach is mainly inspired by GraphSage(Hamilton et al., 2018). Transformer-Sage layer aggregates the embedding of neighbor nodes, and realizes an information interaction between neighbors, which enhances the representation ability of center nodes, as shown in Figure 2.

The main idea of our approach is to introduce interactive information while encoding a sentence. Hence, a single tower network which we call interactive network in InterCSE is added on the basis of SimCSE(Gao et al., 2022) as shown in Figure 3. There are two major components in our framework: independent network and interactive network. During the training process, the interactive network completes the interactive feature extraction of sentence pairs. While inference, only independent network works. All the transformer encoders share same parameters.

Independent Network The goal of the independent network is to quickly produce semantic representation of single sentence without relying on other sentences. The independent network is composed of a transformer encoder and cosine network. Transformer encoder is a siamese network based on BERT-like model, and cosine network is generally cosine similarity calculation. During the training process, the input is exactly two sentence, such as sentence A and sentence A^+ (or sentence A^-)¹, and independent network encodes the sentence pairs independently by $[CLS]$ representation, then calculates the cosine similarity. There

¹For convenience of description, sentence A refers *Two dogs are running.*, sentence A^+ refers *There are animals outdoors.*, and sentence A^- refers *The pets sitting on a couch.* in Figure 3.

are a high semantic score between sentence A and sentence A^+ , and a low semantic score between sentence A and sentence A^- .

Interactive Network The goal of the interactive network is to obtain non-independent sentence semantic information through the attention interaction of words between sentence pairs. Non-independent sentences semantic information is an important correlation signal, which can keenly perceive the semantics of sentence to details. The interactive network consists of an transformer encoder and a feed forward network. Transformer encoder is a BERT-like model, and the input is the concatenation of two sentence and $[SEP]$ token. The $[CLS]$ embedding of transformer encoder is input into the feed forward network to evaluate the similarity of sentence pairs or classification. For example, concatenation of sentence A and sentence A^+ is positive, concatenation of sentence A and sentence A^- is negative.

3.2 Training Dataset

The model we proposed is suitable for supervised tasks, and the in-batch-negative sampling method of the independent network requires positive label data, so we introduce natural language inference (NLI) datasets to train our model, including the SNLI(Bowman et al., 2015) and MNLI(Williams et al., 2018) datasets. NLI datasets consist of high-quality pairs, and given a premise, human annotators generate three types of sentences: entailment(is absolutely true), neutral(might be true), and contradiction(is definitely false).

Following the work of SimCSE(Gao et al., 2022), we adopt the hard negative strategy, using entailment and contradiction to represent positive and negative samples, respectively. Therefore, a triplet is generated, (x_i, x_i^+, x_i^-) , where x_i is the premise, x_i^+ and x_i^- are entailment and contradiction hypotheses. Statistics of training datasets are shown in Table 2.

3.3 Loss Expression

The model we proposed is a multi-task model, including two modules of interactive network and independent network. The modeling objectives of the two modules are different, and the corresponding loss functions are also different. The following describes the loss functions corresponding to the two modules in detail.

For the independent encoding module, the input is a batch of data, one pair of sentences is a positive

sample, and the rest of the sentences are negative samples, which is in-batch negative sampling. We use hard negative in our model, and the contrastive loss $loss_{cl}$ can be expressed as:

$$-\log \frac{e^{\text{sim}(h_i, h_i^+)/\tau}}{\sum_{j=1}^N (e^{\text{sim}(h_i, h_i^+)/\tau} + e^{\text{sim}(h_i, h_i^-)/\tau})} \quad (1)$$

where $\text{sim}(\cdot)$ indicates cosine similarity function, τ controls the temperature, h_i , h_i^+ and h_i^- is the embedding of x_i , x_i^+ and x_i^- respectively.

For the interaction module, the input is a batch data where there are N pairs (x, x^+) and (x, x^-) , and the label of (x, x^+) and (x, x^-) are 1 and 0 respectively. We consider two kind of loss functions to perceive the interactive semantic information.

The first is classification loss for all sentences in a batch, and we use cross entropy loss $loss_{ce}$ to express as:

$$-\sum_{i=1}^{2N} (y_i \log(p_i) + (1 - y_i) \log(1 - p_i)) \quad (2)$$

where y_i is the label, p_i is the predicted score of feed forward network.

The second is margin ranking loss for the ranking relation between positive and negative samples in a batch. Margin ranking loss $loss_{mr}$ could enhance the pairwise distinction in semantic space, and it can be express as:

$$\sum_{i=1}^N \max(0, \text{margin} - (p_i^+ - p_i^-)) \quad (3)$$

where p_i^+ and p_i^- is the predicted score of feed forward network of (x_i, x_i^+) and (x_i, x_i^-) , margin controls the boundary.

With classification loss and margin ranking loss, we get supervised loss for interactive network as:

$$loss_{sl} = loss_{ce} + loss_{mr} \quad (4)$$

Multi-task learning needs to integrate different loss functions. We propose the weighted decaying methods to regularize three different loss functions:

While capturing the interactive semantic information, we need to minimize the hurt to siamese encoder. We propose the loss decaying methods to regularize three different loss functions as:

$$loss = loss_{cl} + \alpha * loss_{sl} \quad (5)$$

where α is piecewise constant decaying when training. For example, we evenly divide the

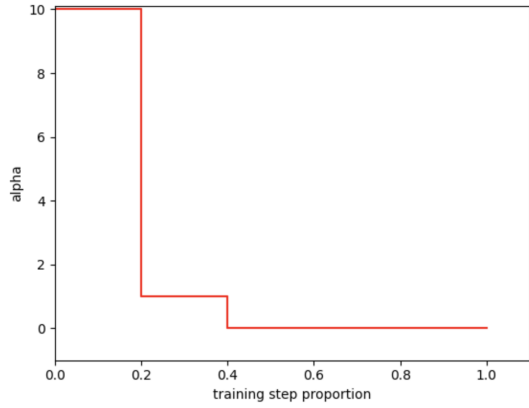


Figure 4: α piecewise constant decaying when training.

training process into 5 stages, each stage corresponds to a different α . When the value of α is $[10, 1, 0.1, 0.01, 0.001]$, the decaying is shown on Figure 4.

4 Experiments

Our approach is mainly proposed for supervised tasks, and we conducted multiple experiments on Semantic Textual Similarity (STS) task to verify the effectiveness of this approach.

4.1 Setups

Datasets Following previous works(Reimers and Gurevych, 2019; Gao et al., 2022; Yan et al., 2021), we evaluate our approach on 7 STS tasks: STS 2012-2016(Agirre et al., 2012, 2013, 2014, 2015, 2016), STS Benchmark(Cer et al., 2017) and SICK-Relatedness(Marelli et al., 2014). A pair of sentences in those datasets has a gold score between 0 and 5 to indicate their semantic similarity. The higher the score, the higher the similarity of sentences pair.

Since the gold score label of the STS datasets is not suitable for the training process of the contrastive learning model, we introduce the SNLI(Bowman et al., 2015) and MNLI(Williams et al., 2018) (NLI) datasets as supervised data to train our model. The details of NLI datasets have described in Section 3.2.

Baselines To show the effectiveness of our approach on supervised sentence representation, we select many state-of-the-art methods as comparison recently, including InferSent(Conneau et al., 2017), Universal Sentence Encoder(Cer et al., 2018), Sentence-BERT(Reimers and Gurevych, 2019), ConSERT(Yan et al., 2021), SimCSE(Gao et al., 2022).

Evaluation When evaluating the trained model, we first obtain the representation of sentences by $[CLS]$ token embeddings, then we report the spearman correlation between the cosine similarity scores of sentence representations and the human-annotated gold scores. When calculating spearman correlation, we merge all sentences together in a STS task, and calculate the mean of spearman correlation.

4.2 Training Details

Our implementation is based on the SimCSE. We start from pre-trained checkpoints of BERT (uncased) and take the $[CLS]$ representation as the sentence embedding, and train model on the combination of MNLI and SNLI datasets on 4 Tesla V100 gpus for 4 epochs. Hyper-parameters τ and $margin$ is set to 0.04 and 0.5. Since α is used to learn interactive information and not hurt performance of sentence embedding, the value is chose from $[10, 1, 0.1, 0.01, 0.001]$ following with the training process.

4.3 Main Results

We compare our approach to previous state-of-the-art sentence embedding methods on STS tasks, including InferSent(Conneau et al., 2017), Universal Sentence Encoder(Cer et al., 2018) Sentence-Transformer(Reimers and Gurevych, 2019), ConSERT(Yan et al., 2021) and SimCSE(Gao et al., 2022). Table 5 shows the evaluation results on 7 STS tasks. InterCSE using $BERT_{base}$ and $BERT_{large}$ achieve 82.11% and 82.88% spearman’s correlation, 0.54% and 0.43% improvement compared to SimCSE respectively, and achieves new state-of-the-art performance on STS tasks.

5 Compared with SBERT

We propose to add interactive network to the Sentence-Transformers(Reimers and Gurevych, 2019) model, and conduct some experiments to verify the performance on the STS-Benchmark dataset. Our proposed model is called Inter-Sentence-Transformers². The label of the STS-Benchmark dataset is a number between 0 and 5, so we divide it by 5.0 as the two sentences semantic similarity.

The independent network is following Sentence-Transformers with siamese transformers, and use

²<https://github.com/2hip3ng/InterCSE/tree/main/Inter-Sentence-Transformers>

| Model | STS12 | STS13 | STS14 | STS15 | STS16 | STS-B | SICK-R | Avg |
|------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| InferSent-Glove ♣ | 52.86 | 66.75 | 62.15 | 72.77 | 66.87 | 68.03 | 65.65 | 65.01 |
| Universal Sentence Encoder ♣ | 64.49 | 67.80 | 64.61 | 76.83 | 73.18 | 74.92 | 76.69 | 71.22 |
| SBERT(base) ♣ | 70.97 | 76.53 | 73.19 | 79.09 | 74.30 | 77.03 | 72.91 | 74.89 |
| SBERT-flow(base) ♠ | 69.78 | 77.27 | 74.35 | 82.01 | 77.46 | 79.12 | 76.21 | 76.60 |
| SBERT-whitening(base) ♠ | 69.65 | 77.57 | 74.66 | 82.27 | 78.39 | 79.52 | 76.91 | 77.00 |
| ConSERT-joint(base) ◇ | 70.53 | 79.96 | 74.85 | 81.45 | 76.72 | 78.82 | 77.53 | 77.12 |
| SimCSE-BERT(base) ♠ | 75.30 | 84.67 | 80.19 | 85.40 | 80.82 | 84.25 | 80.39 | 81.57 |
| PromCSE-BERT(base) △ | 75.58 | 84.33 | 79.67 | 85.79 | 81.24 | 84.25 | 80.79 | 81.81 |
| InterCSE-BERT(base) ♥ | 75.83 | 85.05 | 80.82 | 86.00 | 81.07 | 84.86 | 81.16 | 82.11 |
| SBERT(large) ♣ | 72.27 | 78.46 | 74.90 | 80.99 | 76.25 | 79.23 | 73.75 | 76.55 |
| SimCSE-BERT(large) ◆ | 76.15 | 86.29 | 80.80 | 86.27 | 81.29 | 85.39 | 80.98 | 82.45 |
| InterCSE-BERT(large) ♥ | 76.59 | 86.69 | 81.83 | 86.32 | 81.85 | 85.78 | 81.12 | 82.88 |

Table 3: Sentence embedding performance on STS tasks (Spearman’s correlation, "all" setting). ♥: results of our approach, ♣: results from sentence-transformer(Reimers and Gurevych, 2019), ◇: results from ConSERT(Yan et al., 2021), ♠: results from SimCSE(Gao et al., 2022), △: results from PromCSE(Jiang et al., 2022), ◆: results from our reproduced.

the mean of all tokens embedding as sentences representation. The predicted score produced by cosine-similarity on sentences representation is called *indep_score*.

The interactive network extracts the $[CLS]$ token representation as input of a layer of regression network which is made up of fully-connected layer and sigmoid function. The predicted score produced by sigmoid function is called *inter_score*.

Interactive and independent network both use mean-square error as loss function, and we combine these to generate a final loss with piecewise-constant decay as shown in equation (6) to (8).

$$L_{inter} = \frac{1}{N} \sum_{i=1}^N (y_i - inter_score_i)^2 \quad (6)$$

$$L_{indep} = \frac{1}{N} \sum_{i=1}^N (y_i - indep_score_i)^2 \quad (7)$$

$$Loss = L_{indep} + \alpha * L_{inter} \quad (8)$$

where y_i is the true score of sentences pair, α is decaying as training.

The piecewise constant value α is chose from $[10, 1, 0.1, 0.01, 0.001]$ in BERT-like base models and $[1, 0.1, 0.01, 0.001, 0.0001]$ in BERT-like large models. The training epochs is 32, and the others parameters is the same as Sentence-Transformer. The final performance is shown in the table 4. Inter-Sentence-Transformers get a state-of-the-art performances. It improves 0.88% and 1.07% based on BERT-base and RoBERTa-base respectively compared with Sentence-Transformers. And it slight increases based on BERT-large and RoBERTa-large.

| Model | Spearman |
|-----------------------------|--------------|
| SBERT-STSB-base ♣ | 84.30 |
| SBERT-STSB-large ♣ | 84.28 |
| SRoBERTa-STSB-base ♣ | 84.62 |
| SRoBERTa-STSB-large ♣ | 84.41 |
| Inter-SBERT-STSB-base ◇ | 85.18 |
| Inter-SBERT-STSB-large ◇ | 84.91 |
| Inter-SRoBERTa-STSB-base ◇ | 85.69 |
| Inter-SRoBERTa-STSB-large ◇ | 84.82 |

Table 4: Evaluation on the STS Benchmark test set with spearman’s correlation. ♣: reproduced on Sentence-Transformers(Reimers and Gurevych, 2019); ◇: adding interactive network based on Sentence-Transformers. All results are trained on STS Benchmark train set.

6 Ablation Experiment

In this section, we conduct further analyses to understand the role of interactive network in our approach.

Training Strategy We believe that there are three training strategies that can help the model to consider the interaction information of sentence pairs when capturing sentence independent semantics:(1) Firstly we train single tower network, and then train siamese network; (2) Firstly we train siamese network, and then train single tower network; (3) We use joint training strategy as InterCSE. The performance of three training strategy are shown as Table 5.

Loss Fuction From the experimental results of InterCSE and Inter-Sentence-Transformers, we can see the effectiveness of the interaction network in improving the overall performance. At the same

| Model | STS-B Spearman |
|---------------------|----------------|
| Single-Siamese-BERT | 37.74 |
| Siamese-Single-BERT | 84.76 |
| InterCSE-BERT | 86.25 |

Table 5: Evaluation on the STS Benchmark dev set with spearman’s correlation of training strategy.

| Model | STS-B | STS Avg. |
|-------------------|--------------|--------------|
| <i>Loss</i> | | |
| w/o $loss_{ce}$ | 86.19 | 82.02 |
| w/o $loss_{mr}$ | 86.16 | 81.97 |
| constant α | 86.08 | 81.75 |
| square penalty | 86.09 | 81.76 |
| decay α | 86.25 | 82.11 |

Table 6: Evaluation on the STS Benchmark dev set with spearman’s correlation of loss function.

time, in the section *Training Strategy*, the benefits of multi-task methods for training two sub-networks at the same time are also confirmed. However, there are still many methods how to combine the loss of multi-tasking. Therefore, we have carried out a variety of combination strategies:

1. remove $loss_{ce}$ from equation (4),
2. remove $loss_{mr}$ from equation (4),
3. α is a constant, such as 1.0,
4. add square penalty and show in equation (9),

$$loss = loss_{cl}^2 + loss_{sl}^2 \quad (9)$$

5. use loss decay.

The performance of different loss is shown in Table 6.

Loss Decay In our approach, we use a dynamically changing loss to express the multi-task learning goal, which does not affect the sentence semantic representation while taking more semantic information. Hence, we propose an attenuation strategy for the loss of the interactive network. As shown in the equation (5), α gradually decreases with the training process.

The different descent methods we compared, the results are shown in Table 7. In the early stage of training, the interactive network has a large loss, which makes the model perceive the interactive information brought by the interactive network and update the model parameters. In the later stage of training, the interaction network is weakened, and

| α decay list | STS-B Spearman |
|-------------------------------------|----------------|
| [100, 10, 1, 0.1, 0.01] | 86.19 |
| [10, 1, 0.1, 0.01, 0.001] | 86.25 |
| [1, 0.1, 0.01, 0.001, 0.0001] | 86.09 |
| [0.1, 0.01, 0.001, 0.0001, 0.00001] | 86.02 |

Table 7: Evaluation on the STS Benchmark dev set with spearman’s correlation of α decay list.

the model fully learns the semantic expression of the sentence.

7 Conclusion

In this paper, we propose multiple task method joint contrastive learning for sentence representation which is termed InterCSE. Experiments shows that InterCSE achieves considerable performance on 7 semantic text similarity tasks. Through InterCSE uses a simple framework, it makes a perfect combination of the advantages of single tower network and siamese network by decaying loss function. When perceiving fine-grained word information, try to minimize damage to sentence semantics. Therefore, in the future, we will focus on designing effective network which improve interactive information gain.

Limitations

One limitation of our work is that we can not experiment on unsupervised setting, though the interactive network need labeled sentence pairs.

Acknowledgements

TODO

References

- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, German Rigau, Larraitz Uria, and Janyce Wiebe. 2015. [SemEval-2015 task 2: Semantic textual similarity, English, Spanish and pilot on interpretability](#). In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 252–263, Denver, Colorado. Association for Computational Linguistics.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2014. [SemEval-2014 task 10: Multilingual semantic textual similarity](#). In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 81–91, Dublin, Ireland. Association for Computational Linguistics.

| | | | |
|-----|---|--|-----|
| 534 | Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, | Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, | 592 |
| 535 | Aitor Gonzalez-Agirre, Rada Mihalcea, German | Brian Strope, and Ray Kurzweil. 2018. Universal | 593 |
| 536 | Rigau, and Janyce Wiebe. 2016. SemEval-2016 | sentence encoder for English . In <i>Proceedings of</i> | 594 |
| 537 | task 1: Semantic textual similarity, monolingual | <i>the 2018 Conference on Empirical Methods in Nat-</i> | 595 |
| 538 | and cross-lingual evaluation . In <i>Proceedings of the</i> | <i>tural Language Processing: System Demonstrations</i> , | 596 |
| 539 | <i>10th International Workshop on Semantic Evaluation</i> | pages 169–174, Brussels, Belgium. Association for | 597 |
| 540 | <i>(SemEval-2016)</i> , pages 497–511, San Diego, Califor- | Computational Linguistics. | 598 |
| 541 | nia. Association for Computational Linguistics. | | |
| 542 | Eneko Agirre, Daniel Cer, Mona Diab, and Aitor | Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc | 599 |
| 543 | Gonzalez-Agirre. 2012. SemEval-2012 task 6: A | Barraut, and Antoine Bordes. 2017. Supervised | 600 |
| 544 | pilot on semantic textual similarity . In <i>*SEM 2012:</i> | learning of universal sentence representations from | 601 |
| 545 | <i>The First Joint Conference on Lexical and Compu-</i> | natural language inference data . In <i>Proceedings of</i> | 602 |
| 546 | <i>tational Semantics – Volume 1: Proceedings of the</i> | <i>the 2017 Conference on Empirical Methods in Nat-</i> | 603 |
| 547 | <i>main conference and the shared task, and Volume</i> | <i>tural Language Processing</i> , pages 670–680, Copen- | 604 |
| 548 | <i>2: Proceedings of the Sixth International Workshop</i> | hagen, Denmark. Association for Computational Lin- | 605 |
| 549 | <i>on Semantic Evaluation (SemEval 2012)</i> , pages 385– | guistics. | 606 |
| 550 | 393, Montréal, Canada. Association for Computa- | Jacob Devlin, Ming-Wei Chang, Kenton Lee, and | 607 |
| 551 | tional Linguistics. | Kristina Toutanova. 2019. Bert: Pre-training of deep | 608 |
| 552 | Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez- | bidirectional transformers for language understand- | 609 |
| 553 | Agirre, and Weiwei Guo. 2013. *SEM 2013 shared | ing . | 610 |
| 554 | task: Semantic textual similarity . In <i>Second Joint</i> | Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie- | 611 |
| 555 | <i>Conference on Lexical and Computational Semantics</i> | Yan Liu. 2019. Representation degeneration problem | 612 |
| 556 | <i>(*SEM), Volume 1: Proceedings of the Main Confer-</i> | in training natural language generation models . | 613 |
| 557 | <i>ence and the Shared Task: Semantic Textual Similar-</i> | Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2022. | 614 |
| 558 | <i>ity</i> , pages 32–43, Atlanta, Georgia, USA. Association | Simcse: Simple contrastive learning of sentence em- | 615 |
| 559 | for Computational Linguistics. | beddings . | 616 |
| 560 | Wasi Uddin Ahmad, Xueying Bai, Zhechao Huang, | William L. Hamilton, Rex Ying, and Jure Leskovec. | 617 |
| 561 | Chao Jiang, Nanyun Peng, and Kai-Wei Chang. 2018. | 2018. Inductive representation learning on large | 618 |
| 562 | Multi-task learning for universal sentence embed- | graphs . | 619 |
| 563 | dings: A thorough evaluation using transfer and aux- | Yuxin Jiang, Linhan Zhang, and Wei Wang. 2022. Im- | 620 |
| 564 | iliary tasks . | proved universal sentence embeddings with prompt- | 621 |
| 565 | Samuel R. Bowman, Gabor Angeli, Christopher Potts, | based contrastive learning and energy-based learning . | 622 |
| 566 | and Christopher D. Manning. 2015. A large anno- | In Findings of the Association for Computational | 623 |
| 567 | tated corpus for learning natural language inference . | <i>Linguistics: EMNLP 2022</i> , pages 3021–3035, Abu | 624 |
| 568 | In <i>Proceedings of the 2015 Conference on Empiri-</i> | Dhabi, United Arab Emirates. Association for Com- | 625 |
| 569 | <i>cal Methods in Natural Language Processing</i> , pages | putational Linguistics. | 626 |
| 570 | 632–642, Lisbon, Portugal. Association for Compu- | Gregory Koch, Richard Zemel, and Ruslan Salakhut- | 627 |
| 571 | tational Linguistics. | dinov. 2015. Siamese neural networks for one-shot | 628 |
| 572 | Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie | image recognition. | 629 |
| 573 | Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind | Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2017. | 630 |
| 574 | Neelakantan, Pranav Shyam, Girish Sastry, Amanda | Adversarial multi-task learning for text classification . | 631 |
| 575 | Askell, Sandhini Agarwal, Ariel Herbert-Voss, | In <i>Proceedings of the 55th Annual Meeting of the</i> | 632 |
| 576 | Gretchen Krueger, Tom Henighan, Rewon Child, | <i>Association for Computational Linguistics (Volume</i> | 633 |
| 577 | Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, | <i>1: Long Papers)</i> , pages 1–10, Vancouver, Canada. | 634 |
| 578 | Clemens Winter, Christopher Hesse, Mark Chen, Eric | Association for Computational Linguistics. | 635 |
| 579 | Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, | Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man- | 636 |
| 580 | Jack Clark, Christopher Berner, Sam McCandlish, | dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, | 637 |
| 581 | Alec Radford, Ilya Sutskever, and Dario Amodei. | Luke Zettlemoyer, and Veselin Stoyanov. 2019. | 638 |
| 582 | 2020. Language models are few-shot learners . | Roberta: A robustly optimized bert pretraining ap- | 639 |
| 583 | Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez- | proach . | 640 |
| 584 | Gazpio, and Lucia Specia. 2017. SemEval-2017 | Minh-Thang Luong, Quoc V. Le, Ilya Sutskever, Oriol | 641 |
| 585 | task 1: Semantic textual similarity multilingual and | Vinyals, and Lukasz Kaiser. 2016. Multi-task se- | 642 |
| 586 | crosslingual focused evaluation . In <i>Proceedings</i> | quence to sequence learning . | 643 |
| 587 | <i>of the 11th International Workshop on Semantic</i> | Marco Marelli, Stefano Menini, Marco Baroni, Luisa | 644 |
| 588 | <i>Evaluation (SemEval-2017)</i> , pages 1–14, Vancouver, | Bentivogli, Raffaella Bernardi, and Roberto Zam- | 645 |
| 589 | Canada. Association for Computational Linguistics. | parelli. 2014. A SICK cure for the evaluation of | 646 |
| 590 | Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, | | |
| 591 | Nicole Limtiaco, Rhomni St. John, Noah Constant, | | |

| | | | |
|-----|--|---|-----|
| 647 | compositional distributional semantic models. In | ding. In <i>Proceedings of the 29th International Con-</i> | 701 |
| 648 | <i>Proceedings of the Ninth International Conference</i> | <i>ference on Computational Linguistics</i> , pages 3898– | 702 |
| 649 | <i>on Language Resources and Evaluation (LREC’14)</i> , | 3907, Gyeongju, Republic of Korea. International | 703 |
| 650 | pages 216–223, Reykjavik, Iceland. European Lan- | Committee on Computational Linguistics. | 704 |
| 651 | guage Resources Association (ELRA). | | |
| 652 | Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, | Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang, | 705 |
| 653 | Sam Gross, Nathan Ng, David Grangier, and Michael | Wei Wu, and Weiran Xu. 2021. <i>Consert: A con-</i> | 706 |
| 654 | Auli. 2019. fairseq: A fast, extensible toolkit for | <i>trastive framework for self-supervised sentence rep-</i> | 707 |
| 655 | sequence modeling. In <i>Proceedings of NAACL-HLT</i> | <i>resentation transfer</i> . | 708 |
| 656 | <i>2019: Demonstrations</i> . | | |
| 657 | Alec Radford, Karthik Narasimhan, Tim Salimans, and | Zhilin Yang, Zihang Dai, Yiming Yang, Jaime | 709 |
| 658 | Ilya Sutskever. 2018. Improving language under- | Carbonell, Ruslan Salakhutdinov, and Quoc V. | 710 |
| 659 | standing by generative pre-training. | Le. 2019. <i>XLnet: Generalized autoregres-</i> | 711 |
| 660 | Alec Radford, Jeff Wu, Rewon Child, David Luan, | <i>sive pretraining for language understand-</i> | 712 |
| 661 | Dario Amodei, and Ilya Sutskever. 2019. Language | <i>ing</i> . Cite arxiv:1906.08237Comment: Pre- | 713 |
| 662 | models are unsupervised multitask learners. | trained models and code are available at | 714 |
| 663 | | https://github.com/zihangdai/xlnet . | 715 |
| 664 | Marek Rei. 2017. <i>Semi-supervised multitask learning</i> | | |
| 665 | <i>for sequence labeling</i> . In <i>Proceedings of the 55th An-</i> | | |
| 666 | <i>nuual Meeting of the Association for Computational</i> | | |
| 667 | <i>Linguistics (Volume 1: Long Papers)</i> , pages 2121– | | |
| 668 | 2130, Vancouver, Canada. Association for Computa- | | |
| 669 | tional Linguistics. | | |
| 670 | Nils Reimers and Iryna Gurevych. 2019. <i>Sentence-bert:</i> | | |
| 671 | <i>Sentence embeddings using siamese bert-networks</i> . | | |
| 672 | Yu Sun, Shuohuan Wang, Shikun Feng, Siyu Ding, | | |
| 673 | Chao Pang, Junyuan Shang, Jiaxiang Liu, Xuyi Chen, | | |
| 674 | Yanbin Zhao, Yuxiang Lu, Weixin Liu, Zhihua Wu, | | |
| 675 | Weibao Gong, Jianzhong Liang, Zhizhou Shang, | | |
| 676 | Peng Sun, Wei Liu, Xuan Ouyang, Dianhai Yu, Hao | | |
| 677 | Tian, Hua Wu, and Haifeng Wang. 2021. <i>Ernie 3.0:</i> | | |
| 678 | <i>Large-scale knowledge enhanced pre-training for lan-</i> | | |
| 679 | <i>guage understanding and generation</i> . | | |
| 680 | Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi | | |
| 681 | Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao | | |
| 682 | Tian, and Hua Wu. 2019a. <i>Ernie: Enhanced repre-</i> | | |
| 683 | <i>sentation through knowledge integration</i> . | | |
| 684 | Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao | | |
| 685 | Tian, Hua Wu, and Haifeng Wang. 2019b. <i>Ernie</i> | | |
| 686 | <i>2.0: A continual pre-training framework for language</i> | | |
| 687 | <i>understanding</i> . | | |
| 688 | Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2019. | | |
| 689 | <i>Representation learning with contrastive predictive</i> | | |
| 690 | <i>coding</i> . | | |
| 691 | Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob | | |
| 692 | Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz | | |
| 693 | Kaiser, and Illia Polosukhin. 2017. <i>Attention is all</i> | | |
| 694 | <i>you need</i> . | | |
| 695 | Adina Williams, Nikita Nangia, and Samuel R. Bow- | | |
| 696 | man. 2018. <i>A broad-coverage challenge corpus for</i> | | |
| 697 | <i>sentence understanding through inference</i> . | | |
| 698 | Xing Wu, Chaochen Gao, Liangjun Zang, Jizhong Han, | | |
| 699 | Zhongyuan Wang, and Songlin Hu. 2022. <i>ESim-</i> | | |
| 700 | <i>CSE: Enhanced sample building method for con-</i> | | |
| | <i>trastive learning of unsupervised sentence embed-</i> | | |