

Domain-level Pairwise Semantic Interaction for Aspect-Based Sentiment Classification

Zhenxin Wu, Jiazheng Gong, Kecen Guo, Guanye Liang,
Qingliang Che, and Bo Liu

2533653096@qq.com, 619790446@qq.com

Abstract. Aspect-based sentiment classification (ABSC) is a very challenging subtask of sentiment analysis (SA) and suffers badly from the class-imbalance. Existing methods only process sentences independently, without considering the domain-level relationship between sentences, and fail to provide effective solutions to the problem of class-imbalance. From an intuitive point of view, sentences in the same domain often have high-level semantic connections. The interaction of their high-level semantic features can force the model to produce better semantic representations, and find the similarities and nuances between sentences better. Driven by this idea, we propose a plug-and-play Pairwise Semantic Interaction (PSI) module, which takes pairwise sentences as input, and obtains interactive information by learning the semantic vectors of the two sentences. Subsequently, different gates are generated to effectively highlight the key semantic features of each sentence. Finally, the adversarial interaction between the vectors is used to make the semantic representation of two sentences more distinguishable. Experimental results on four ABSC datasets show that, in most cases, PSI is superior to many competitive state-of-the-art baselines and can significantly alleviate the problem of class-imbalance.

Keywords: Aspect-based sentiment classification · Pairwise semantic interaction · Class-imbalance

1 Introduction

Aspect-based sentiment classification (ABSC) is a fine-grained sentiment classification subtask of sentiment analysis [10], which aims to identify the sentiment polarity of each aspect in a sentence (positive, negative or neutral). It is widely used in different domains, such as online comments (e.g., movie and restaurant reviews [9]), data mining and e-commerce customer service. For example, sentence 1 in Fig. 1 shows that the customer enjoys the restaurant’s food but thinks the ambience is just not bad. For this sentence, ABSC needs to recognize that the two aspects “ambience” (A1) and “food” (A2) contained in the sentence are “neutral” and “positive”, respectively.

In fact, people’s comments often have obvious emotional preferences, which means that they may suffer the problem of class-imbalance. Since there are

far more comments with “positive” and “negative” in the same domain than those with “neutral”, “neutral” comments are always marginalized and thus misjudged. At present, the commonly used ABSC methods, whether they are traditional methods [20,1] or deep learning models [5,23], none of them has solved the problem of class-imbalance. Moreover, the similarity of semantic contexts between sentences in the same domain has not been fully utilized. In our paper, we define “semantic” as a highly abstract coding vector of sentences extracted by the information extractor, e.g. BERT. If we can make interactive learning of two similar sentences in the same domain, they can learn more domain semantic information from each other and enrich the high-level semantic encoding of sentences. It will also help to find the similarities and nuances between sentences, which can reduce misjudgments due to class-imbalance.

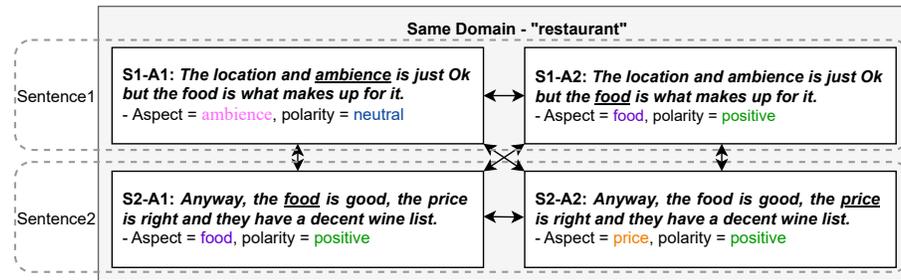


Fig. 1. The two sentences have different aspects, but they all belong to the same domain “restaurant” and have similar semantic context.

For example, making interactive learning of S1-A1 and S1-A2 in Fig. 1, can help to distinguish different sentiment polarities of a sentence which contains different aspects. The interaction between S1-A1 (“neutral”) and S2 (both aspects are “positive”) can make the semantic encoding of “neutral” more discriminative, by comparing it with the strong “positive” sentence S2. At the same time, the interaction between S1-A2, S2-A1, and S2-A2 can also enrich the features of the same sentiment polarity in different aspects or different sentences.

Based on this intuition, we propose a domain-level plug-and-play Pairwise Semantic Interaction (PSI) module for ABSC. For the construction of sentence pairs, it is worth emphasizing that we consider that the sentences in the same dataset belong to the same domain. We do not limit that two sentences must have the same aspect, and encourage richer interactions between sentences (refer to 3.4 for detailed sentence pair construction strategy). For PSI module, firstly, we extract the semantic vectors of the two sentences by semantic extractors (e.g., BERT [5]), respectively. Subsequently, through a gating mechanism, sentences can learn each other’s high-level semantic information adaptively, which enriches the semantic representation of a single sentence. Finally, we additionally use a design similar to the adversarial network [6] to help promote the model to

distinguish the nuances between similar semantic representations. In summary, the contributions of this paper are as follows:

- We introduce a Pairwise Semantic Interaction (PSI) module for the interaction between sentences, help to find the similarities and nuances of sentences, and can significantly reduce the misjudgments due to class-imbalance. Through the interaction between sentences in the same domain, the sentences get better and more discriminative semantic representations.
- The PSI module is plug-and-play and can be easily combined with most mainstream semantic extractors such as BERT.
- The experiments on four prestigious ABSC datasets have justified the efficacy of PSI, achieving or approaching SOTA results.

2 Related work

Existing ABSA researches focus on the use of deep neural networks, such as target dependent LSTM models [18] and Attention-based LSTM [22] for aspect-level sentiment classification. In recent years, the pre-trained language model BERT [5], which has been very successful in many Natural Language Processing tasks, has been applied in ABSA and achieved significant results such as [17,24,11]. However, all of the above studies have ignored semantic relationship between sentences in the same domain. Recently, contrastive learning has achieved great success in both Computer Vision (CV) and Natural Language Processing (NLP). Its main purpose is to make the features of the same category closer to each other, while the distance between the features of different categories is farther. In [26], through an attention interaction, the network can adaptively find delicate clues from two fine-grained images in pairs.

In ABSA, Chen et al. [4] proposed a Cooperative Graph Attention Networks (CoGAN) method for cooperatively learning the aspect-related sentence representation in document level. Tang et al. [19] also use the method of transformer combined with graph, in order to allow dependency graph to guide the representation learning of the transformer encoder. However, their model is based on transformer combined with Graph Networks, which has high computational overhead. By contrast, our proposed PSI is a plug-and-play module, which can achieve decent performance with a little extra overhead.

3 The Proposed Method

In this section, we will describe our PSI module. PSI compares two similar sentences together to find the common semantic representation and semantic differences between them, rather than studying the semantic representation of a single sentence alone.

The module PSI will take two similar sentences as input and go through three carefully designed sub-modules i.e., mutual vector learning, semantic gate generation, and adversarial interaction. The entire structure of PSI is shown in Fig. 2.

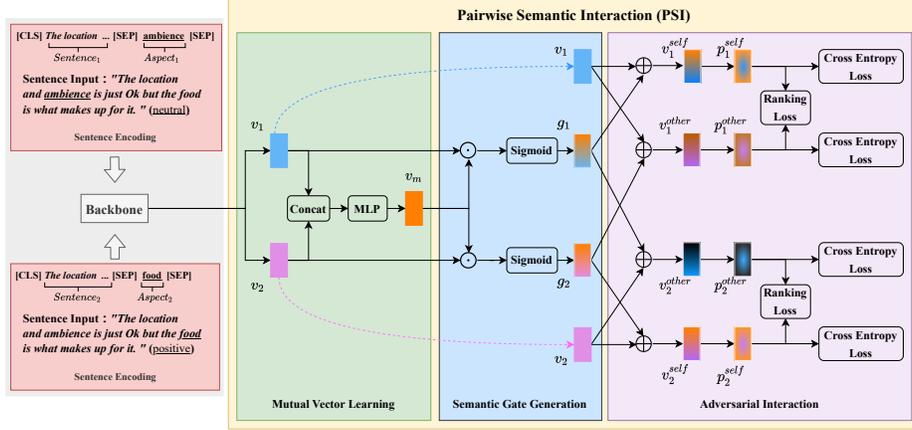


Fig. 2. The structure of PSI. We simply use S1-A1 and S1-A2 in Fig. 1 as an example of pairwise sentences ($Sentence_1, Sentence_2$). In this case, PSI can promote the model to distinguish different sentiment polarities of different aspects in the sentence. It is worth emphasizing that the PSI is a plug-and-play module, i.e., PSI can be combined with most mainstream semantic extraction backbones (e.g., BERT) during the training phase, and flexibly unload it for single-input test sentence.

3.1 Mutual Vector Learning

Before this sub-module, the semantic of these two sentences are extracted by backbone, and two D -dimensional semantic vectors i.e., v_1 and $v_2 \in \mathbb{R}^D$ are generated, respectively. Then, in this sub-module, we learn a mutual vector $v_m \in \mathbb{R}^D$ from individual v_1 and v_2 ,

$$v_m = f_m([v_1, v_2]). \quad (1)$$

where $[]$ is concatenation operation and $f_m(\cdot)$ is a mapping function of $[v_1, v_2]$. Specifically, we use the multi-layer perceptron (MLP) as the mapping function. By summarizing the two feature vectors v_1 and v_2 , the mutual vector v_m is produced accordingly, which contains common high-level semantic information and discriminative semantic clues of two sentences.

3.2 Semantic Gate Generation

After producing the mutual vector v_m , we can use v_m to activate v_1 and v_2 . In order to generate more discriminative information for later comparisons, the dot product of v_m with two feature vectors v_1 and v_2 is carried out according to channels to locate the contrastive information in the two vectors. Then, the gate vectors, i.e., g_1 and g_2 are generated by a sigmoid function,

$$g_i = \text{sigmoid}(v_m \odot v_i), i \in \{1, 2\}. \quad (2)$$

Therefore, g_i can be used as an attention vector to highlight the important semantic representations belonging to individual v_i . For example, the previous example of S1-A1 and S1-A2 in Fig. 1, the gates will help to highlight two different key words “ambience” and “food”, respectively. This helps to distinguish different sentiment polarities belonging to different aspects.

3.3 Adversarial Interaction and Model Training

When comparing two sentences, humans not only focus on the salient parts of one sentence, but also focus on the salient parts of the other one. Based on this, we introduce an adversarial interaction mechanism through residual attention. As shown in Fig. 2, two feature vectors v_1 and v_2 and two gate vectors g_1 and g_2 are combined in pairs, we could then get four attentive semantic vectors with

$$\begin{aligned} v_1^{self} &= v_1 + v_1 \odot g_1, \\ v_2^{self} &= v_2 + v_2 \odot g_2, \\ v_1^{other} &= v_1 + v_1 \odot g_2, \\ v_2^{other} &= v_2 + v_2 \odot g_1. \end{aligned} \quad (3)$$

Intuitively, the semantic vector $v_i (i \in \{1, 2\})$, is guided by the attention of the gate vector $g_j (j \in \{1, 2\})$, strengthens or weakens certain semantic information, and then adds to itself to get the output. While $v_i^{self} \in \mathbb{R}^D$ reinforces the feature region belonging to its own gate vector, and $v_i^{other} \in \mathbb{R}^D$ reinforces the feature region belonging to another gate vector.

Then v_i^j (where $i \in \{1, 2\}$, $j \in \{self, other\}$) are feed into the softmax classifier by $p_i^j = softmax(Wv_i^j + b)$, where $p_i^j \in \mathbb{R}^C$ represents the score vector of prediction, C indicates the number of polarities, and $\{W, b\}$ is the parameter set of the softmax classifier. In order to effectively train the entire PSI module, we define the following loss function as

$$J = J_{ce} + \mu J_{rk}. \quad (4)$$

Among them, J_{ce} is the cross-entropy loss, and J_{rk} is the score ranking regularization loss with a coefficient of μ . Specifically we choose the hinge loss function as the score ranking regularization J_{rk} ,

$$J_{rk} = \sum_{i \in \{1, 2\}} \max(0, p_i^{other}(y_i) - p_i^{self}(y_i) + \varepsilon). \quad (5)$$

where $p_i^j(y_i) \in \mathbb{R}$ represents the score got in the predicted vector p_i^j , and y_i denotes the index of the true polarity of sentence i , and ε is the penalty term. The motivation of this design is that, v_i^{self} is activated by its own gate vector. Hence, compared to v_i^{other} , it should be more discriminative to the corresponding label. That is, the score difference $p_i^{self}(y_i) - p_i^{other}(y_i)$ should be larger than a margin ε , which means that $p_i^{self}(y_i)$ should be larger than $p_i^{other}(y_i)$ and must keep

a distance with $p_i^{other}(y_i)$. At the same time, when cross-entropy loss J_{ce} is optimized, due to having the same label, $p_i^{self}(y_i)$ and $p_i^{other}(y_i)$ will tend to be closer. Therefore, J_{rk} and J_{ce} will be optimized adversarially. As a result, v_i^{other} can learn the semantic information shared by the two sentences, and v_i^{self} can learn its own unique information which will be more discriminative and reduce the noise of sentence pairs.

3.4 Sentence Pair Construction

Next, we'll provide an explanation on how to construct multiple sentence pairs in a batch for end-to-end training. Specifically, we randomly sample N_p polarities in a batch (there are 3 polarities in total, i.e. positive, negative, neutral). For each polarity, we randomly sample N_s training sentences. Consequently, there are $N_p \times N_s$ different sentences in each batch (we set the same sentence to express different aspects, belonging to different sentences). After getting a batch of sentences, we input these sentences into the backbone to generate their respective semantic vectors. For every sentence, we compare its semantic vector with the different sentences in the batch in accordance to Euclidean distance. We do not limit that two different sentences must have the same aspect, and encourage richer interactions between sentences (the following ablation study proves our point). Then, we can construct the inter/intra-pairs in a batch. The inter-pairs are the following sentence pairs which contains two situations. 1) The current sentence and itself (with different aspect and different polarity), e.g., S1-A1 & S1-A2 in Fig. 1; 2) The current sentence and the most similar sentence with different polarities from the current sentence, e.g., S1-A1 & S2(A1/A2). On the contrary, intra-pairs refer to the following sentence pairs. 1) The current sentence and itself (with different aspect and same polarity), e.g., S2-A1 & S2-A2; 2) The current sentence and the most similar sentence with the same polarity from the current sentence, e.g., S1-A2 & S2(A1/A2). This design permits the PSI to learn to distinguish between truly similar and highly overlapping pairs.

3.5 Model Testing

Because PSI is a practical plug-and-play module. In the training phase, the backbone and the PSI module can summarize the comparative clues from sentence pairs, and step by step improve the discriminant capacity of backbone representation for sentences. Therefore, in the testing phase, only the backbone model with updated parameters is used, but not the PSI module, so that the generalization ability of the model can be guaranteed without losing the performance of the model. To be specific, in testing phase, we input a sentence into the backbone, extract its semantic vector $X_* \in \mathbb{R}^D$, and then directly input X_* into the softmax classifier. It is worth emphasizing that the softmax classifier are shared between the training phase and the testing phase. The score vector $P_* \in \mathbb{R}^C$ is applied to label prediction. Thus, our test scheme is the same as a regular backbone, which demonstrates the strong applicability of PSI.

Table 1. Statistics of the datasets.

Polarity	Res14		Lap15		Res16		Lap16	
	train	test	train	test	train	test	train	test
Positive	839	222	765	329	749	204	1084	274
Neutral	500	94	106	79	101	44	188	46
Negative	2179	657	1103	541	1657	611	1637	481
Sum	3518	973	1974	949	2507	859	2909	801

4 Experiments

4.1 Datasets and Metrics

We have carried out experiments on four datasets to verify the performance of our proposed model, PSI. Restaurant14 is from Semeval-2014 Task 4 [16], Laptop15 is from Semeval-2015 Task 12 [15], and the other two datasets (Restaurant16 and Laptop16) are from Semeval-2016 Task 5 [14]. The statistics for these datasets are shown in Table 1. And we use Accuracy (Acc.) and Macro-F1 (F1) as performance metrics.

4.2 Implementation Details

Unless stated otherwise, we implement PSI as follows. For each aspect of a sentence, we concat the corresponding aspect at the end of the sentence. Then we adjust the length of each sentence to 85 (the maximum sentence length after the tokenizer tokenizes is 85). If it is not enough, fill it with zero, and then use it as the input of backbone. Firstly, we extract the semantic vector $v_i \in \mathbb{R}^{786}$ by BERT. Secondly, for all the datasets, we randomly sample 3 polarities, i.e., $N_p = 3$. And for each polarity, we randomly sample 4 sentences to form a batch, i.e., $N_s = 4$. For each sentence, we find its most similar sentence from its own polarity and the rest polarities, according to Euclidean distance between their semantic vectors. As a result, we obtain an intra-pair and an inter-pair for each sentence in the batch. For each pair, we concatenate v_1 and v_2 as input to a two-layer MLP, i.e., FC (1572→512), FC (512→786). Consequently, this operation generates the mutual vector $v_m \in \mathbb{R}^{786}$. All of our models are implemented by Pytorch with a single NVIDIA GTX 2080Ti GPU with 11G Memory. For all datasets, the coefficient μ in Eq. 4 is 1, while the margin ϵ is 0.05 in score ranking regularization. Among them, BERT is optimized by Adam optimizer with $\beta_1 = 0.9$, and the initial learning rate is 0.0001. For our PSI (backbone is BERT or BERT-Large) method, we use another Adam optimizer for training, and the initial learning rate is 0.00002 with $\beta_1 = 0.9$. There are a total of 20 training epochs, and if the loss does not decrease for 5 consecutive epochs, it will invoke early-stop. In addition, we set a fixed seed when training the model to ensure the reproducibility of the results.

Table 2. The comparative results, with data for non-BERT models from [4], kumaGCN from [2], RepWalk from [25], and IMN from [8]. The data for BERT-QA from [17], AC-MIMLLN from [11], and CoGAN from [4]. The experimental configuration for standard BERT and PSI is shown in implementation details. “-” denotes no data available yet. The best result of each dataset is bolded, and the second-best result is underlined.

Models	Res14		Lap15		Res16		Lap16	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
TC-LSTM	0.781	0.675	0.745	0.622	0.813	0.629	0.766	0.578
ATAE-LSTM	0.772	-	0.747	0.637	0.821	0.644	0.781	0.591
RAM	0.802	0.708	0.759	0.639	0.839	0.661	0.802	0.627
IAN	0.793	0.701	0.753	0.625	0.836	0.652	0.794	0.622
Clause-Level ATT	-	-	0.816	0.667	0.841	0.667	0.809	0.634
LSTM+synATT +TarRep	0.806	0.713	0.822	0.649	0.846	0.675	0.813	0.628
kumaGCN	0.814	0.736	-	-	0.894	0.732	-	-
RepWalk	0.838	0.769	-	-	0.896	0.712	-	-
IMN	0.839	0.757	0.831	0.654	0.892	0.71	0.802	0.623
BERT	0.867	0.764	0.818	0.699	0.884	0.755	0.817	0.665
BERT-QA	-	-	0.827	0.595	0.896	0.715	0.812	0.596
AC-MIMLLN	0.893	-	-	-	-	-	-	-
CoGAN	-	-	0.851	0.745	0.920	<u>0.816</u>	<u>0.842</u>	0.707
PSI (BERT)	<u>0.916</u>	<u>0.857</u>	<u>0.860</u>	<u>0.756</u>	0.901	0.788	0.839	<u>0.723</u>
PSI (BERT-Large)	0.924	0.863	0.868	0.760	<u>0.913</u>	0.828	0.87	0.737

4.3 Comparison with SOTA Methods

To fully evaluate the performance of our method, we apply PSI based on BERT or BERT-Large. We compare it with the state-of-the-art (SOTA) baselines including (1) ABSA models without BERT: TC-LSTM [18], ATAE-LSTM [22], RAM [3], IAN [12], Clause-LevelATT [21], LSTM+synATT+TarRep [7], kumaGCN [2], RepWalk [25] and IMN [8]. (2) BERT-based models for ABSA: BERT [5], BERT-QA [17], AC-MIMLLN [11] and CoGAN [4]. Table 2 shows the results of our experiments on four datasets.

From Table 2 we can come to the following conclusion. The performance of PSI (Based on BERT or BERT-Large) on Res14, Lap15 and Lap16 is better than those of all baselines. And in Res16, our PSI module approaches SOTA results. The experiments justify that PSI is a very powerful plug-and-play module, showing the effectiveness of our method.

4.4 Alleviating the Problem of Class-Imbalance

As shown in Table 1, the mainstream datasets have the problem of class-imbalance. For example, in Res14, the “negative” comments is significantly more than the data of other polarities (“negative” accounts for 62%), while the “neutral” comments is far lower for other polarities (“neutral” accounted for 14%). In order to illustrate the advantages of the our method, we compared PSI (Based on BERT)

Table 3. Comparison of accuracy(%) of different polarities between PSI and BERT.

Model	Negative(%)	Neutral(%)	Positive(%)	Overall(%)
BERT	92.5	47.9	86.0	86.7
PSI (BERT)	97.6(+5.1)	60.6(+12.7)	86.9(+0.9)	91.6(+4.9)

with the standard BERT model on the Res14 dataset, for each polarity (positive, negative, and neutral). Table 3 shows the comparison results of the accuracy of different sentiment polarities.

It can be seen that the “positive” and “neutral” accuracy of our model (PSI) is better than that of BERT model. Specifically, PSI significantly improves the performance of “neutral” classification. And the overall accuracy of our model was greatly improved compared with BERT model. It indicates that the sentence pair interaction learning in PSI module can make the semantics between sentences complement each other and effectively learn the nuances between sentences. This can make the semantic representation of different polarities of sentences more distinguishable, and can effectively alleviate the problem of class-imbalance.

4.5 Sample Extraction Strategies

Our proposed method encourages richer interactions between sentences and do not limit that two different sentences must have the same aspect. In order to study the impact of different sample extraction methods (sentiment polarity and aspect) on ABSC, we conducted the following ablation experiments. We use BERT as semantic vector extractor to evaluate different sample extraction methods on Res14.

As mentioned above, our proposed sample extraction method is that Intra/inter pairs are constructed from the same/different sentiment polarities without limiting the range of aspect (**Interacting Polarity, I_P**). And there are three other sample extraction methods, including **1) Interacting Aspect (I_A)**. Intra/inter pairs are constructed from the same/different aspects without limiting the range of sentiment polarity. **2) Interacting Polarity and Limiting Aspect (I_P & L_A)**. Intra/inter pairs are constructed from the same/different polarities by limiting the same aspect. **3) Interacting Aspect and Limiting Polarity (I_A & L_P)**. Intra/inter pairs are constructed from the same/different aspects by limiting the same polarity.

From Table 4, in all the four datasets, the results of the other three ablation experiments are worse than I_P (Ours). For I_P & L_A and I_A & L_P , results demonstrate that we should not limit aspects (in order to better distinguish different aspects of a sentence) and should also allow different polarities to interact in pairs (in order to better distinguish different sentiment polarities). For I_A , due to class-imbalance, if we do not limit the range of sentiment polarity (by constructing the inter-pair of different sentiment polarities), there will be a lot of interactions between sentences belonging to the same majority class

Table 4. Different Sample Extraction Method.

Methods	Res14		Lap15		Res16		Lap16	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
I_P (Ours)	0.916	0.857	0.860	0.756	0.901	0.788	0.839	0.723
I_A	0.914	0.854	0.834	0.699	0.896	0.753	0.830	0.680
I_P & L_A	0.909	0.852	0.840	0.699	0.893	0.787	0.819	0.656
I_A & L_P	0.895	0.826	0.836	0.689	0.873	0.738	0.820	0.641

(“negative”), while the interaction between different sentiment polarities will be insufficient. Finally, for I_P , we explicitly construct inter-pairs (belonging to different polarities) in each batch to ensure that the number of interactions between different polarities is sufficient. In this way, the nuances between different sentiment polarities can be better learned by the model. Therefore, we choose I_P as the sample extraction method for ABSC.

4.6 Ablation Study

In addition, in order to investigate the properties of our proposed PSI, we use Bert as a semantic vector extractor to evaluate its key design on Res14. For fairness, when exploring different strategies for one design, we use the other designs as the base strategies described in the proposed approach and implementation details.

Mutual Vector Generation. To demonstrate the essentiality of x_m in Eq. 1, we investigate different operations to generate it. **1) Individual Operation.** The key of x_m is to learn mutual information from both sentences in the pair. For comparison, we introduce a baseline without it. Specifically, we replace mutual learning in Eq. 1 by individual learning $\tilde{x}_i = f_m(x_i)$, and use \tilde{x}_i to generate the gate vector $g_i = \text{sigmoid}(\tilde{x}_i)$ where $i \in \{1, 2\}$. **2) Elementwise Operations.** We perform a number of widely-used elementwise operations to generate x_m , including Subtract Square: $x_m = (x_1 - x_2)^2$, Sum: $x_m = (x_1 + x_2)$, and Product $x_m = (x_1 * x_2)$. **3) Interactive MLP.** It is the mapping function described in the proposed approach. As shown in Table 5, the Individual operation (i.e., the setting without x_m) performs worst. Hence, it is necessary to learn mutual context by x_m . Based on experimental results, we choose the simple but effective Interactive MLP to generate the mutual vector in our experiments.

Table 5. Different operations of mutual vector.

Mutual Vector	ACC	F1
Individual	0.893	0.824
Sum	0.909	0.850
Product	0.904	0.847
Subtract Square	0.910	0.857
Interactive MLP	0.916	0.857

The Influence of the Number of Interactions. On the basis of sample extraction method I_P , We further studied the influence of the number of sampling polarity N_p and the number of corresponding sentences N_s in each batch. The results are shown in Table 6. It can be seen that PSI is more sensitive to polarity-number than sentence-number. This is because more polarities often lead to richer diversity of sentence pairs. If 3 polarities are selected in a batch, sentence pairs must include the sentence belonging to minority category (neutral). So the sentences belonging to minority polarity (neutral) can learn semantic information from other majority polarity (positive or negative), thereby improving the generalization ability of the model. Therefore, we choose the best setting in our experiments, i.e., $N_p = 3$ and $N_s = 4$.

Table 6. The influence of the num of polarity & sentence in each sampling.

(N_p, N_s)	(2,3)	(2,4)	(2,5)	(3,3)	(3,4)	(3,5)
Acc.	0.892	0.894	0.895	0.914	0.916	0.897
F1	0.819	0.833	0.825	0.848	0.857	0.832

The Influence of Sentence Similarity. Furthermore, We also examine the influence of sentence similarity for intra-pairs and inter-pairs on the Res14 dataset based on I_P . For $12(N_p \times N_s)$ different sentences in each batch, we construct intra/inter-pairs for each sentence, respectively. Therefore, there are 24 sentence pair in a batch. The selection strategies including 1) **Random**. We randomly sampled 24 sentence pairs(intra/inter pairs) in a batch; 2) **Sentence-Distance**. We sampled 24 sentence pairs in accordance to the Euclidean distance between sentences in intra/inter pairs (i.e., Similar (S), Dissimilar (D)). Thus, 8 different Class-Polarity-Sentence settings were generated in Table 7.

Table 7. The influence of sentence similarity.

Pair Construction	Intra	Inter	Acc.	F1
Random	-	-	0.905	0.842
Sentence-Distance	-	D	0.9	0.846
	-	S	0.91	0.847
	D	-	0.882	0.809
	S	-	0.909	0.853
	D	D	0.897	0.826
	S	D	0.91	0.849
	D	S	0.892	0.827
	S	S	0.916	0.857

From Table 7, we can see that most of results of Sentence-Distance are better than Random, which shows that when constructing sentence pairs, the polarity of the sentence and the similarity between them must be considered. Secondly,

in the setting of Sentence-Distance, Whether inter or intra, the results of (S) is higher than (D) in most cases. This proves that similar sentence pairs can increase the training difficulty of the model, which can effectively help model to identify subtle semantic differences.

The Necessity of Ranking Regularization. In this section, we conduct experiments on the necessity of ranking regularization in the sub-module (Adversarial Interaction) of PSI. We compared the standard BERT, PSI without Ranking Regularization J_{rk} (based on BERT) and PSI (based on BERT) in five datasets. It can be seen from the Table 8 that compared with the standard BERT, the performance of the BERT with PSI(with or without J_{rk}) is significantly improved, which proves the effectiveness of the PSI structure. Furthermore, the PSI with J_{rk} is better than PSI without J_{rk} . This demonstrates that, without ranking regularization, the effectiveness of adversarial interaction will be greatly reduced. When without ranking regularization, the model is only optimized by cross-entropy loss J_{ce} , which will make the semantic vectors (v_i^{self}/v_i^{other}) activated by different gates lose their distinction. On the contrary, adding ranking regularization can keep semantic vectors activated by different gates at a distance. This allows the semantic vectors activated by the other’s gate to have the common semantic representation of the two sentences (v_i^{other}), and the semantic vectors activated by the own gate to have unique semantic representation (v_i^{self}), which is more discriminative. In summary, the PSI with ranking regularization helps model to find the similarities and differences between sentence pairs better, and reduce the noise of sentence pairs, which is effective and necessary.

Table 8. The impact of ranking regularization.

Method	Res14		Lap15		Res16		Lap16	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
BERT	0.867	0.764	0.818	0.699	0.884	0.755	0.817	0.665
PSI without J_{rk}	0.907	0.851	0.837	0.719	0.888	0.718	0.833	0.714
PSI with J_{rk}	0.916	0.857	0.860	0.756	0.901	0.788	0.839	0.723

4.7 Visualization Analysis

In order to understand the discriminability of our method, we use UMAP [13] to visualize the polarity separability and compactness in the semantic features extracted from a standard BERT and the PSI (based on BERT) in Res14. In Fig. 3, it is evident that when using our PSI module, the clusters are farther apart and more compact, leading to a more clear distinction of various clusters representing different polarities. This also proves that adding PSI module can promote the model to learn better semantic representation of sentences and make them more distinguishable.

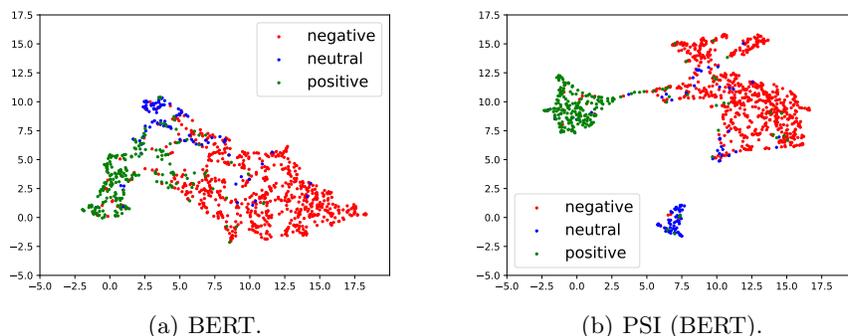


Fig. 3. Discriminability using UMAP to visualize polarity separability and compactness.

5 Conclusion

In this paper, we proposed a domain-level Pairwise Semantic Interaction (PSI) for ABSC. Through the interactions between sentences, PSI can effectively enrich the semantic encoding of sentences and produce better semantic representations. Meanwhile, PSI is plug-and-play module and can further help the model distinguish the nuances between similar sentences and effectively alleviate the problem of class-imbalance. Finally, the empirical results on four prestigious ABSC datasets justified the power of PSI that has achieved SOTA performance in most cases. In future work, we will consider integrating some advanced attention mechanisms into this method.

References

1. A novel lexicalized hmm-based learning framework for web opinion mining. In: ICML. ACM International Conference Proceeding Series, vol. 382, pp. 465–472. ACM (2009), withdrawn.
2. Chen, C., Teng, Z., Zhang, Y.: Inducing target-specific latent structures for aspect sentiment classification. In: EMNLP. pp. 5596–5607. Association for Computational Linguistics (2020)
3. Chen, P., Sun, Z., Bing, L., Yang, W.: Recurrent attention network on memory for aspect sentiment analysis. In: EMNLP. pp. 452–461. Association for Computational Linguistics (2017)
4. Chen, X., Sun, C., Wang, J., Li, S., Si, L., Zhang, M., Zhou, G.: Aspect sentiment classification with document-level sentiment preference modeling. In: ACL. pp. 3667–3677. Association for Computational Linguistics (2020)
5. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. In: NAACL-HLT (1). pp. 4171–4186. Association for Computational Linguistics (2019)

6. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., Bengio, Y.: Generative adversarial nets. In: NIPS. pp. 2672–2680 (2014)
7. He, R., Lee, W.S., Ng, H.T., Dahlmeier, D.: Effective attention modeling for aspect-level sentiment classification. In: COLING. pp. 1121–1131. Association for Computational Linguistics (2018)
8. He, R., Lee, W.S., Ng, H.T., Dahlmeier, D.: An interactive multi-task learning network for end-to-end aspect-based sentiment analysis. In: ACL. pp. 504–515. Association for Computational Linguistics (2019)
9. Kiritchenko, S., Zhu, X., Cherry, C., Mohammad, S.: Nrc-canada-2014: Detecting aspects and sentiment in customer reviews. In: SemEval@COLING. pp. 437–442. The Association for Computer Linguistics (2014)
10. Li, S., Huang, C., Zhou, G., Lee, S.Y.M.: Employing personal/impersonal views in supervised and semi-supervised sentiment classification. In: ACL. pp. 414–423. The Association for Computer Linguistics (2010)
11. Li, Y., Yin, C., Zhong, S., Pan, X.: Multi-instance multi-label learning networks for aspect-category sentiment analysis. In: EMNLP. pp. 3550–3560. Association for Computational Linguistics (2020)
12. Ma, D., Li, S., Zhang, X., Wang, H.: Interactive attention networks for aspect-level sentiment classification. In: IJCAI. pp. 4068–4074. ijcai.org (2017)
13. McInnes, L., Healy, J.: UMAP: uniform manifold approximation and projection for dimension reduction. CoRR [abs/1802.03426](https://arxiv.org/abs/1802.03426) (2018)
14. Pontiki, M., Galanis, D., Papageorgiou, H., Androutsopoulos, I., Manandhar, S., Al-Smadi, M., Al-Ayyoub, M., Zhao, Y., Qin, B., Clercq, O.D., Hoste, V., Apidianaki, M., Tammier, X., Loukachevitch, N.V., Kotelnikov, E.V., Bel, N., Zafra, S.M.J., Eryigit, G.: Semeval-2016 task 5: Aspect based sentiment analysis. In: SemEval@NAACL-HLT. pp. 19–30. The Association for Computer Linguistics (2016)
15. Pontiki, M., Galanis, D., Papageorgiou, H., Manandhar, S., Androutsopoulos, I.: Semeval-2015 task 12: Aspect based sentiment analysis. In: SemEval@NAACL-HLT. pp. 486–495. The Association for Computer Linguistics (2015)
16. Pontiki, M., Galanis, D., Pavlopoulos, J., Papageorgiou, H., Androutsopoulos, I., Manandhar, S.: Semeval-2014 task 4: Aspect based sentiment analysis. In: SemEval@COLING. pp. 27–35. The Association for Computer Linguistics (2014)
17. Sun, C., Huang, L., Qiu, X.: Utilizing BERT for aspect-based sentiment analysis via constructing auxiliary sentence. In: NAACL-HLT (1). pp. 380–385. Association for Computational Linguistics (2019)
18. Tang, D., Qin, B., Feng, X., Liu, T.: Effective lstms for target-dependent sentiment classification. In: COLING. pp. 3298–3307. ACL (2016)
19. Tang, H., Ji, D., Li, C., Zhou, Q.: Dependency graph enhanced dual-transformer structure for aspect-based sentiment classification. In: ACL. pp. 6578–6588. Association for Computational Linguistics (2020)
20. Wagner, J., Arora, P., Cortes, S., Barman, U., Bogdanova, D., Foster, J., Tounsi, L.: DCU: aspect-based polarity classification for semeval task 4. In: SemEval@COLING. pp. 223–229. The Association for Computer Linguistics (2014)
21. Wang, J., Li, J., Li, S., Kang, Y., Zhang, M., Si, L., Zhou, G.: Aspect sentiment classification with both word-level and clause-level attention networks. In: IJCAI. pp. 4439–4445. ijcai.org (2018)
22. Wang, Y., Huang, M., Zhu, X., Zhao, L.: Attention-based LSTM for aspect-level sentiment classification. In: EMNLP. pp. 606–615. The Association for Computational Linguistics (2016)

23. Wu, Z., Ong, D.C.: Context-guided BERT for targeted aspect-based sentiment analysis. In: AACL. pp. 14094–14102. AACL Press (2021)
24. Xu, H., Liu, B., Shu, L., Yu, P.S.: BERT post-training for review reading comprehension and aspect-based sentiment analysis. In: NAACL-HLT (1). pp. 2324–2335. Association for Computational Linguistics (2019)
25. Zheng, Y., Zhang, R., Mensah, S., Mao, Y.: Replicate, walk, and stop on syntax: An effective neural network model for aspect-level sentiment classification. In: AACL. pp. 9685–9692. AACL Press (2020)
26. Zhuang, P., Wang, Y., Qiao, Y.: Learning attentive pairwise interaction for fine-grained classification. In: AACL. pp. 13130–13137. AACL Press (2020)