

DIWS-LCR-Rot-hop++: A Domain-Independent Word Selector for Cross-Domain Aspect-Based Sentiment Classification

Junhee Lee
Erasmus University Rotterdam
Rotterdam, the Netherlands
ejoone611@gmail.com

Flavius Frasincar
Erasmus University Rotterdam
Rotterdam, the Netherlands
frasincar@ese.eur.nl

Maria Mihaela Truşcă
Buch. Univ. of Econ. Studies
Bucharest, Romania
maria.trusca@csie.ase.ro

ABSTRACT

The Aspect-Based Sentiment Classification (ABSC) models often suffer from a lack of training data in some domains. To exploit the abundant data from another domain, this work extends the original state-of-the-art LCR-Rot-hop++ model that uses a neural network with a rotatory attention mechanism for a cross-domain setting. More specifically, we propose a Domain-Independent Word Selector (DIWS) model that is used in combination with the LCR-Rot-hop++ model (DIWS-LCR-Rot-hop++). DIWS-LCR-Rot-hop++ uses attention weights from the domain classification task to determine whether a word is domain-specific or domain-independent, and discards domain-specific words when training and testing the LCR-Rot-hop++ model for cross-domain ABSC. Overall, our results confirm that DIWS-LCR-Rot-hop++ outperforms the original LCR-Rot-hop++ model under a cross-domain setting in case we impose an optimal domain-dependent attention threshold value for deciding whether a word is domain-specific or domain-independent. For a target domain that is highly similar to the source domain, we find that imposing moderate restrictions on classifying domain-independent words yields the best performance. Differently, a dissimilar target domain requires a strict restriction that classifies a small proportion of words as domain-independent. Also, we observe information loss which deteriorates the performance of DIWS-LCR-Rot-hop++ when we categorize an excessive amount of words as domain-specific and discard them.

CCS Concepts

•Information systems → Sentiment analysis; Information extraction; Web mining;

Keywords

ABSC, cross-domain ABSC, domain-specific word masking, attention threshold, DIWS-LCR-Rot-hop++

1. INTRODUCTION

The evolution of the Web has changed how people think and make decisions. Furthermore, the recent development

of social media and e-commerce platforms has opened the forum for people to exchange their opinions on social events or products and services on the market. As a result, the number of opinionated texts on the Web is rapidly growing, and understanding the voice of the customer is becoming crucial for business success [22]. However, it is difficult to summarize this public opinion because an enormous number of opinionated texts and reviews are available [8]. To solve this issue, the current research has proposed a Natural Language Processing (NLP) task so-called Sentiment Analysis (SA) that identifies the overall sentiment of a given text [19].

One of the branches of SA is Aspect-Based Sentiment Analysis (ABSA) [27]. It consists of two tasks: Aspect Detection (AD), which identifies the aspect in the text [30], and Aspect-Based Sentiment Classification (ABSC) which determines the sentiment about the previously found aspect [2]. Unlike ordinary SA, ABSA can separately detect multiple sentiments in a given text. This feature is especially useful for analyzing customer reviews [12]. On the Web, there exists an excessive amount of positive reviews compared to negative reviews because people with high valuations on a product are more likely to make a purchase and write a review [15]. However, people may partially express negative opinions about a specific aspect of a product even though they are generally satisfied. For instance, one may write a positive review on a new cell phone but complain about its camera quality. ABSA can capture a negative sentiment of a specific aspect in a generally positive review, while ordinary SA ignores it. Figure 1 illustrates an example of ABSA.

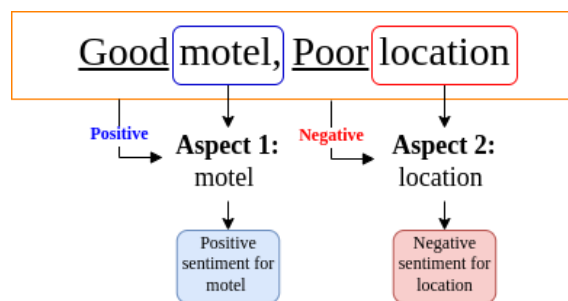


Figure 1: Illustration of ABSA in the hotel domain [23].

This research focuses on ABSC. In practical applications, it is difficult to train ABSC models to a sufficient extent be-

cause of the limited number of input sentence data. This issue is prominent in certain domains while other domains have a sufficient amount of data. To exploit this data imbalance, a number of approaches have been proposed by various studies. For example, [32] fine-tunes the upper layers of LCR-Rot-hop++ to increase cross-domain adaptability. Also, [34] introduces the BERTMasker algorithm that transforms the input sentences into domain-invariant sentences by masking the domain-related words and training the model using domain-agnostic words. Additionally, [16] extends the LCR-Rot-hop++ model with Domain Adversarial Training (DAT) method to construct a cross-domain DAT-LCR-Rot-hop++ model.

Nevertheless, the first fine-tuning approach partially requires the sentiment-labeled target domain data. Similarly, the BERTMasker model performs the best when it utilizes part of the sentiment-labeled target domain data. Also, BERTMasker cannot process the ABSC task. Thus, these models are not suitable for a cross-domain ABSC where sentiment-labeled data are not available in a domain of our interest. To solve this issue, we propose a Domain-Independent Word Selector (DIWS) model and apply it to the LCR-Rot-hop++ model (DIWS-LCR-Rot-hop++). This model utilizes attention weights from a domain classification task to decide whether a word is domain-specific or domain-independent. It only uses domain-independent words to train and test the LCR-Rot-hop++ model. Unlike fine-tuned LCR-Rot-hop++ [32] and BERTMasker [34], we can train the model using only data from other domains with sufficient sentiment-labeled texts. As DAT-LCR-Rot-hop++ shares the same advantage, we use it as a benchmark model to assess the performance of the proposed DIWS-LCR-Rot-hop++ model.

In this paper, we apply the DIWS model to the LCR-Rot-hop++ model yielding the DIWS-LCR-Rot-hop++ model. To this end, we aim to verify whether the combination of DIWS and LCR-Rot-hop++ outperforms the naive cross-domain ABSC performance of LCR-Rot-hop++, where we train and test the model with a completely different domain data without any adjustment to LCR-Rot-hop++. Furthermore, to verify the degree of effectiveness of discarding the domain-specific words, we compare the accuracy level under different strictness levels of deciding whether a word is domain-specific or domain-agnostic. More specifically, we incrementally increase the proportion of discarded words by relaxing the discarding threshold and verifying the trade-off between domain adaptability and information loss. Additionally, we compare the accuracy of the DAT-LCR-Rot-hop++ model under a cross-domain setting to assess the performance of our proposed model.

Differently than in our previous work [17], for which this work is an extension, we provide additional examples to illustrate the addressed problem and its solution, elaborate the hyperparameter optimization process, compare the proposed approach to an additional baseline, and explain our methodology and evaluation in more detail.

This research contributes to the current literature by proposing a method named Domain-Independent Word Selector (DIWS) to better train an existing state-of-the-art ABSC model (LCR-Rot-hop++) for cross-domain sentiment anal-

ysis tasks. Moreover, this proposed model is not only applicable to ABSC but also to other types of SA, which implies its wide applicability in the field of sentiment analysis. To the best of our knowledge, the idea of discarding domain-specific words while training a deep learning method for the ABSC task is new. The Python source code and data are available on the GitHub repository¹.

The structure of the rest of the paper is as follows. In Section 2, we introduce the methodologies for ABSC and cross-domain SA in the current literature in detail and discuss their relevance to our research. In Section 3, we introduce the used datasets and the cleansing process. In Section 4, we elaborate on the theoretical framework, structure, and mathematical formulations of the models that this paper investigates. In Section 5, we give the results and make comparisons between competing methods to answer the research question. Last, Section 6 provides a summary of the findings, discusses the limitations of our research, and proposes future research ideas.

2. RELATED WORK

In this section, we present work related to our research. First, in Subsection 2.1, we present approaches for ABSC. After that, in Subsection 2.2, we describe solutions for cross-domain sentiment classification.

2.1 Aspect-Based Sentiment Classification

There exist two major approaches to ABSC: knowledge-based and machine learning-based. Knowledge-based algorithms, also known as rule-based algorithms use pre-defined sentiment dictionaries such as SenticNet [3] and WordNet [21] to assess the sentiment score of a specific expression corresponding to an aspect in a given sentence. Afterward, this sentiment score is used to define the sentiment polarity of the aspect. Also, [28] argues that a domain-specific ontology can significantly improve the performance of conventional machine learning methods when they are utilized together. An ontology is a set of domain-specific concepts, features, and their semantic inter-relationships [10].

This paper focuses on the second approach which uses machine learning algorithms. In detail, the neural network models and attention models are widely used. For example, [6] introduces Recursive Neural Networks (RecNNs) to this field by proposing an Adaptive Recursive Neural Network (AdaRNN). Also, Recurrent Neural Network (RNN) is one of the popular methods in SA [26; 35]. However, the main drawback of an ordinary RNN methodology is the long-term dependency problem, which refers to the tendency that the prior information to be dissolved when the input sequence is too long [13]. To address this issue, [13] suggests a special type of RNN model called Long Short-Term Memory (LSTM). Unlike traditional RNN models, LSTM employs additional gate nodes to control the information transfer between hidden layers. This structure allows LSTM to efficiently learn long-term relationships of data. Nonetheless, LSTM processes the information sequentially, which leads to a tendency that LSTM output converges to the latest input pattern. To address this concern, [9] suggests a bidirectional

¹<https://github.com/ejoone/DIWS-ABSC>

LSTM (bi-LSTM). It adds a reverse direction LSTM layer to the original LSTM network and uses both forward and backward LSTM layers to obtain a final result.

Furthermore, [39] proposes a Left-Center-Right separated neural network with Rotatory attention (LCR-Rot) that demonstrates high performance when an aspect contains multiple words by capturing the contextual information around the aspect. LCR-Rot is an extension of bi-LSTM and it utilizes three separate bi-LSTM networks, which correspond to the left context, target words, and right context, respectively. Also, rotatory attention helps to better model the relationship between the target words and left/right context, which allows the model to capture the most important words. [39] has confirmed that LCR-Rot outperforms other LSTM-based models. Additionally, LCR-Rot-hop is an extension of LCR-Rot proposed by [33]. It iterates the rotatory attention mechanism multiple times as it better exploits the interactions between the target words and the right/left context. To even better represent the contextual information, [31] proposes LCR-Rot-hop++ which replaces non-contextual word embeddings of LCR-Rot-hop (GloVe) with contextual word embeddings (BERT). Also, it adds an extra attention layer to obtain hierarchical attention. [31] has shown that LCR-Rot-hop++ in combination with a domain ontology (HAABSA++) outperforms other models on ABSC. In this research, we focus on LCR-Rot-hop++ and aim to incorporate the DIWS model inspired by the BERT-Masker network of [34] to extend LCR-Rot-hop++ to the cross-domain setting.

2.2 Cross-Domain Sentiment Classification

Cross-domain sentiment classification aims to solve the insufficient training data problem in one domain by leveraging the data from other domains. It is important to note that cross-domain sentiment classification and multi-domain sentiment classification are different while sharing a similar purpose. Multi-domain sentiment classification aims to train the sentiment classification model using both the source domain (the domain with a sufficient number of training data) and the target domain (the domain that we want to perform sentiment classification for). Hence, it requires the target domain to have training data with sentiment labels although the amount of the data is limited. Instead, cross-domain sentiment classification trains the sentiment classification model by only exploiting the original domain data. Therefore, it can perform sentiment classification tasks on the target domain without any sentiment labels. However, this advantage comes at the cost of accuracy. [34] has verified that its suggested model experiences an accuracy loss of 1.63% when they are applied for cross-domain sentiment analysis, compared to the multi-domain setting. This paper focuses on extending a high-performing ABSC model, i.e., LCR-Rot-hop++ to the cross-domain setting instead of the multi-domain setting.

Unsupervised domain adaptation is one of the approaches to address the training data shortage problem. [40] proposes a representation learning model that selects important domain-independent pivot words. Also, [18] identifies pivots using a hierarchical attention transfer mechanism. Moreover, [11] and [37] extend the domain adaptation to

the multi-domain setting.

Another approach to cross-domain sentiment classification is the shared-private framework [4]. This approach is based on the reasoning that removing the domain-specific tokens would improve the domain-invariance of the input sentence. Hence, the gradient reversal layer is included before the domain classification step and it helps to select the tokens that reduce the performance of the domain classification task and consider corresponding words as domain-agnostic words.

On the other hand, [38] uses an attention mechanism to select the domain-specific information from the shared sentence representation of the input text. This framework is called shared encoder with domain-aware aggregation [34]. To take advantage of the shared-private and shared encoder with domain-aware aggregation paradigms, [34] proposes the BERTMasker model that combines these two frameworks. [34] has demonstrated that the BERTMasker outperforms existing models in both cross-domain and multi-domain settings.

Nevertheless, it is not possible to fully utilize the BERT-Masker network in a cross-domain setting. The BERT-Masker model consists of two parts: shared and private. The shared part masks the domain-specific tokens and uses the unmasked words to train the sentiment classification model. On the other hand, the private part uses the masked tokens to learn the domain-specific sentiment via an attention mechanism using training data of the target domain. The private part may effectively enhance the performance in the multi-domain setting because it is possible to use the labeled target domain data to train the private part of the model. However, in the cross-domain setting, BERTMasker cannot utilize its private part since the target sentiment labels are unavailable.

Several attempts have been made to perform ABSC under the multi-domain and cross-domain settings. For example, [32] applies cross-domain fine-tuning to LCR-Rot-hop++, which is an ABSC model. More specifically, the authors fine-tune the upper layers of LCR-Rot-hop++ because the upper layers contain more domain-specific information while the lower layers represent general language characteristics [32]. However, the fine-tuning procedure requires the training data with sentiment labels from a target domain. Thus, it is a multi-domain ABSC model and we cannot directly compare this model to the cross-domain DIWS-LCR-Rot-hop++ model.

Additionally, [16] suggests the DAT-LCR-Rot-hop++ model that combines Domain Adversarial Training (DAT) [7] with LCR-Rot-hop++ so that it can perform cross-domain ABSC. Unlike [32], it does not require training data from a target domain. It replaces the final Multi-Layer Perceptron (MLP) layer with a domain adversarial component, which consists of two feed-forward MLPs. One is a domain discriminator with a gradient reversal layer. It allocates higher importance to the domain-agnostic words that cannot classify the domain well. The other one is a class discriminator that aims to predict the sentiment label of an aspect in the sentence.

To continue exploring cross-domain aspect-based sentiment classification, this research exploits the attention mechanism

to select the domain-independent words from the input sequence. The DIWS-LCR-Rot-hop++ and DAT-LCR-Rot-hop++ models are based on the same reasoning that paying less or even no attention to the domain-specific words that are crucial for a domain classification task would improve the cross-domain performance of ABSC. However, there are some differences between the two models. First, the DIWS-LCR-Rot-hop++ model sequentially trains the DIWS component and LCR-Rot-hop++ component while DAT-LCR-Rot-hop++ jointly optimize the LCR-Rot-hop++ component and domain class discriminator by letting parameters from LCR-Rot-hop++ affect the discriminator loss. Unlike DAT-LCR-Rot-hop++, in DIWS-LCR-Rot-hop++, DIWS parameters and LCR-Rot-hop++ parameters do not affect each other. Second, DIWS-LCR-Rot-hop++ discretely excludes domain-specific words that pass a certain attention threshold but DAT-LCR-Rot-hop++ allocates less importance to the domain-specific words rather than discarding them.

3. DATA

This research uses review data in five different domains to execute our proposed sentiment analysis. We summarize the used domains, datasets, sample size, and the distribution of sentiments in Table 1.

Table 1: Distribution of sentiment polarities.

Data	Size	Negative		Neutral		Positive	
		Freq.	%	Freq.	%	Freq.	%
Hotel [23]	264	55	21	10	4	199	75
DVD Player [14]	313	172	55	0	0	141	45
Digital Camera [14]	395	74	19	0	0	321	81
MP3 Player [14]	676	262	39	0	0	414	61
Cell Phone [14]	284	70	25	0	0	214	75

Note that the MP3 player review data has the maximum sample size while other domain data such as hotel, DVD player, digital camera, and cell phone contain a relatively small amount of samples. Hence, we use MP3 player review data to train our proposed cross-domain DIWS-LCR-Rot-hop++ model. Such a domain is called the source domain. On the other hand, data in other domains are used to test the performance of trained cross-domain DIWS-LCR-Rot-hop++. Such domains are called target domains. Note that DIWS-LCR-Rot-hop++ requires a pair of one source domain and one target domain to train the model and assess its performance. We fix the MP3 player as a source domain for every pair. Hence, this research examines the performance of DIWS-LCR-Rot-hop++ for the following domain combinations: MP3 Player-Hotel, MP3 Player-DVD Player, MP3 Player-Digital Camera, and MP3 Player-Cell Phone.

In the training process, only the source domain data contains both domain and sentiment labels. On the other hand, target domain data only contain the domain label. In the testing phase, we use DIWS-LCR-Rot-hop++, trained by the source domain (MP3 player), to determine the polarity label of the aspects in the target domain data. Then,

we compare this predicted aspect sentiment to the actual sentiment to measure the test accuracy.

For the robustness of the experiment, we remove some samples in case of the presence of implicit aspects. An implicit aspect refers to a situation where the aspects appear as non-noun words and are implied in the sentence. We removed the samples with such characteristics as the machine learning algorithm cannot process such data [32]. Table 2 shows the results of the cleansing process.

Table 2: Cleansed datasets.

Domain	Removed: implicit aspects (%)
Hotel	22.1
DVD Player	27.4
Digital Camera	19.2
MP3 Player	20.3
Cell Phone	15.9

4. FRAMEWORK

Our proposed DIWS-LCR-Rot-hop++ model uses the DIWS module to identify and discard domain-dependent words from the original input text and process the transformed text using LCR-Rot-hop++. First, in Subsection 4.1, we explain the overall structure of the DIWS-LCR-Rot-hop++. Second, in Subsection 4.2, we elaborate on the DIWS model for the cross-domain sentiment analysis task. Last, in Subsection 4.3, we introduce the LCR-Rot-hop++ model.

4.1 DIWS-LCR-Rot-hop++ Structure

The overall structure of DIWS-LCR-Rot-hop++ is as follows. First, DIWS models the input sentence using pre-trained BERT and obtains corresponding word embeddings. Second, it processes domain classification tasks for both source and target domains via the feed-forward attention layer. In this process, attention weights are computed by a softmax function and we take a linear combination of the attention weights and hidden representation of sentence words as sentence representation. This linear combination is fed into the sigmoid activation function and we obtain the model’s prediction for the domain of the input sentence. Note that the optimal attention weights of each word in an input sentence are determined by gradient descent and back-propagation algorithm. Afterward, we select the domain-independent words by discarding the words that have attention weights higher than a certain attention threshold. We classify such words as domain-specific words. This step is based on the reasoning that the word with high attention weight has a high contribution to the domain classification task, and such words are domain-specific words that specifically appear in a certain domain. For example, the word **sound** specifically appears in the MP3 player domain while the word **nice** may appear in various domains such as MP3 player, hotel, and camera. It is likely that the word **sound** plays an important role in domain classification than **nice** and DIWS masks it to construct domain-agnostic data.

After we identify domain-independent and domain-related words, we move on to the LCR-Rot-hop++ part of the

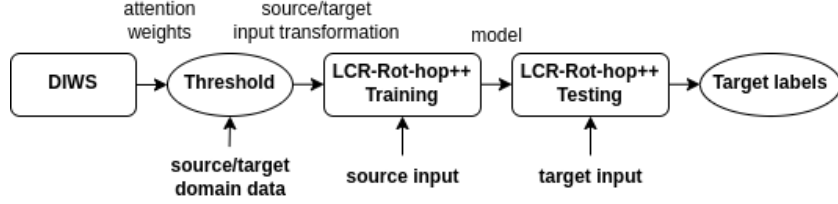


Figure 2: Overall representation of the DIWS-LCR-Rot-hop++ method.

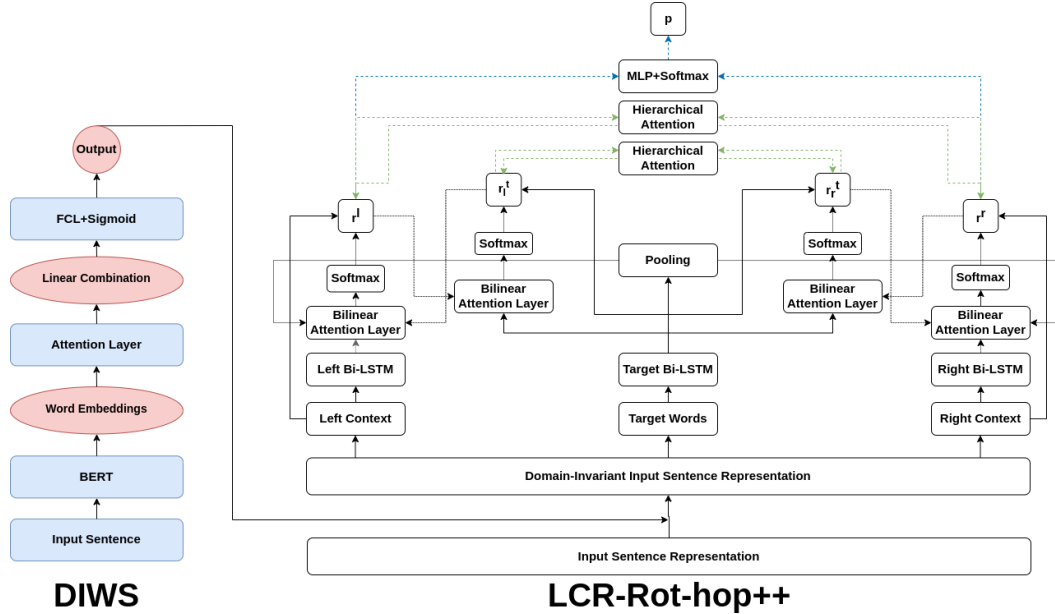


Figure 3: Detailed representation of the DIWS-LCR-Rot-hop++ method.

model. We transform the original input sentence representation into a domain-agnostic input sentence representation by discarding the domain-specific words that we detect using DIWS. Then we split the input representation into three components: left context, target words, and right context. We process them using a rotatory attention mechanism to obtain target-aware left/right context representations and left/right context-aware target representations. Afterward, we process these representations by a hierarchical attention mechanism to update them so that they encode the global information of the input sentence and calculate the final sentiment prediction probabilities. Figure 2 visualizes the method of the model and Figure 3 displays the graphical overview.

4.2 Domain-Independent Word Selector

This section explains the mathematical formulations of Domain-Independent Word Selector (DIWS). An input sentence that we aim to analyze consists of multiple words. Therefore, we denote it as a sequence of words $X = \{x_1, x_2, \dots, x_N\}$, where N corresponds to the number of words in a given input sentence. To let neural network layers recognize and process the input text, we first need to tokenize the sentence as $X = \{[\text{CLS}], x_1, x_2, \dots, x_N, [\text{SEP}]\}$ and

change the character representation of the words to the real-valued vector representation which encodes the meaning of the corresponding word. Here, $[\text{CLS}]$ is a class token that encodes overall information about the corresponding sequence X , and $[\text{SEP}]$ is a separator token that is used to distinguish multiple input sentences. For example, the closer the meaning of the words, the closer the corresponding vectors are in the vector space. This representation allows a computer to execute vector calculations between words. This procedure is called word embedding. It can be performed using language models such as Word2Vec [20], ELMo [24], and BERT [5].

This research uses BERT word embedding because it is a contextual word embedding that considers the context of the word in the sentence when it encodes a word. Therefore, BERT can deal with polysemy, unlike non-contextual word embedding. We use a pre-trained BERT base model with 768 hidden layers and 12 layers to transform the text representation of the input sentence X into a vector representation that the model can process.

Consider a sequence of BERT word embeddings $\{h_1, h_2, \dots, h_N\}$, transformed from input sentence $X = \{x_1, x_2, \dots, x_N\}$. Each word embedding h_i is a $[1 \times d]$ vector, where d is equal to the number of a hidden

layers of the BERT base model. Therefore, each word in our input sentences is transformed into the 768 dimensions numerical vector and fed into the attention layer. The preliminary attention score for the i^{th} word is:

$$\alpha_i = \frac{h_i^T V}{1 \times 1 \quad 1 \times d \quad d \times 1}, \quad (1)$$

where the $V \in \mathbb{R}^d$ is a context vector that is used as a query vector to find a word that is more important and informative for an accurate domain classification. We process these attention scores with the softmax function to obtain corresponding attention weights for every $i = 1, \dots, N$:

$$\alpha_i = \frac{\exp(\alpha_i)}{\sum_{j=1}^N \exp(\alpha_j)}. \quad (2)$$

We use this softmax operation again in LCR-Rot-hop++ and refer to it as the *softmax*(\cdot) function later. Next, we represent the input sentence with length N as a weighted average of word embeddings h_i by their corresponding attention weights a_i where $i = 1, \dots, N$:

$$h = \sum_{i=1}^N \alpha_i \times h_i. \quad (3)$$

This process refers to the linear combination layer in Figure 3. Last, h is fed to the Fully Connected Layer (FCL) in Equation 4 to produce polarity score s and it is fed into a sigmoid function in Equation 5 to produce prediction probability p :

$$s = \frac{h^T W + d}{1 \times 1 \quad 1 \times d \quad d \times 1 \quad 1 \times 1}, \quad (4)$$

$$p = \frac{1}{1 + \exp(-s)}. \quad (5)$$

where $W \in \mathbb{R}^d$ is a weight vector and b is a bias term.

We train the values of the parameters such as weight vector and bias term during the training phase to minimize the loss. The loss function is the binary cross-entropy:

$$L_{DIWS} = - \sum_{j=1}^{M+P} (y_j \log(p^{(j)}) + (1 - y_j) \log(1 - p^{(j)})), \quad (6)$$

where y_j is an actual binary domain label corresponding to j^{th} input sentence and $p^{(j)}$ refers to j^{th} domain prediction probability when there exists M and P (number of) sentences in the source and target domain, respectively.

To verify the domain-classification performance of DIWS, we split the total sample consisting of both source and target domain data into an 80% training sample and 20% testing sample. Additionally, before training the model, we tune the hyperparameters of the DIWS model. Such hyperparameters are learning rate, momentum term, number of epochs, and batch size. As suggested by [33], we use Tree-structured Parzen Estimators (TPE) algorithm [1] for tuning. To optimize the hyperparameters, we split our training sample into 80% pure training samples and 20% validation samples. In this way, we can find the optimal hyperparameters

while maintaining our test sample unseen during the training phase.

The attention weights capture the relative importance of words in an input sentence for predicting the correct domain [29]. Based on this interpretation, we assume that the word with higher attention weights within the sentence is more domain-related because domain-related words play a significant role when determining the domain polarity score compared to articles (a/an/the) or linking verbs (be/is/are) which appear widely regardless of the domain.

In this research, we test different threshold values between domain-specific words with high attention and domain-agnostic words with low attention. Let us define the set of threshold percentiles as $K = \{10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$. For every input sentence,

$$d_i = \begin{cases} \text{domain-specific,} & \text{if } \alpha_i \geq Q_K, \\ \text{domain-independent,} & \text{otherwise,} \end{cases} \quad (7)$$

where d_i is a domain-relatedness label associated with every x_i , and Q_K refers to a K^{th} percentile value of the attention weights for every word in an input sentence. We discard domain-specific words to construct domain-invariant input sentence representations for both source and target domains so that we can perform domain-independent training and domain-independent testing using LCR-Rot-hop++. Last, we fed the transformed input sentences to the LCR-Rot-hop++ model.

4.3 LCR-Rot-hop++

The LCR-Rot-hop++ model uses three bi-LSTM networks and a rotatory, hierarchical attention mechanism to classify the sentiment of a given aspect. This section describes the LCR-Rot-hop++ model and its mathematical formulations.

The LCR-Rot-hop++ model uses a sentence X as its input, where $x_{target} = \{x_1^t, \dots, x_T^t\}$ represents the set of T words describing an aspect target of the sentence X . Then it splits X into three separate components, namely a left context, an aspect target, and a right context. The left context is a set of words that appear before the target x_{target} , and the right context is a set of words that appear after the target x_{target} .

Next, the corresponding word embedding representations of the left context, the target, and the right context are separately processed by bi-LSTM modules. The bi-LSTM model is proposed by [9] to solve the bias issue of the LSTM model. According to [36], LSTM is a biased model in a way that the later input plays a dominant role in determining the final output due to its sequential learning process. This issue is especially significant when we train the model using long input sentences. To solve this issue, LCR-Rot-hop++ employs bi-LSTM. It balances the importance of the beginning input and later input in the learning process by adding a reverse direction LSTM layer. The outputs of the three bi-LSTMs are three hidden states with the dimension of $[2d \times 1]$. It has two times larger dimensions than the dimension of BERT word embedding due to its bidirectional structure. The outputs of these bi-LSTMs are denoted as follows: $[h_1^l, \dots, h_L^l], [h_1^t, \dots, h_T^t], [h_1^r, \dots, h_R^r]$, where L is the number of left context words, T is the number of aspect

target words, and R is the number of right context words. These three hidden states are then processed by the rotatory attention mechanism.

A rotatory attention mechanism aims to distinguish the most important words in the left context, aspect target, and right context for determining the sentiment using a two-step procedure. First, it calculates the target-aware left and right context representations r_l^t and r_r^t by considering the information in the target representation r^t . The initial value of r^t is determined by the pooling operation of the hidden states of the target, which is the output of the target bi-LSTM module:

$$r_{2d \times 1}^t = \text{pooling}([h_1^t, \dots, h_T^t]). \quad (8)$$

r^t is then fed in the bilinear attention layer together with $[h_1^l, \dots, h_L^l]$ and $[h_1^r, \dots, h_R^r]$ separately. We illustrate the mathematical formulation for the left context, but the same logic applies to the right context. The output of the bilinear attention layer is obtained by multiplying the transposed h_i^l , weights (W_c^l) and r^t , and adding the bias term (b_c^l), and in-putting the result to the \tanh activation function for every $i = 1, \dots, L$:

$$f(h_i^l, r^t) = \tanh(h_i^{l\top} \times W_c^l \times r^t + b_c^l). \quad (9)$$

Then we process the obtained score with the softmax function to obtain the attention score α_i^l , and get the target-aware left context representation r^l by computing a weighted average of the left context hidden states in terms of the attention scores:

$$\alpha_i^l = \text{softmax}(f(h_i^l, r^t)), \quad (10)$$

$$r_{2d \times 1}^l = \sum_{i=1}^L \alpha_i^l \times h_i^l. \quad (11)$$

Unlike the first step, the second step uses r^l and r^r to construct the left and right context-aware target representations r_l^t and r_r^t , respectively. The logic is the same as in the first procedure, while we no longer need to use the pooling operation as we already have r^l and r^r from the first step.

We obtain four representations r^l, r^r, r_l^t , and r_r^t as outputs of the rotatory attention mechanism and fed into the hierarchical attention mechanism. It allows these four representations to encode global information around the input sentence, not only the local, left, target, or right contextual information. Two context representations (r^l, r^r) and two target representations (r_l^t, r_r^t) are separately weighted by a hierarchical attention mechanism. We illustrate the mathematical formulation for the context pair, especially the left context, but the same logic applies to the right context and the target pair. First, we compute an attention score at the sentence level:

$$f(r^l) = \tanh(r^{l\top} \times W_h^c + b_h^c), \quad (12)$$

where W_h^c is a weight matrix, and b_h^c is a bias term for the context pair. Then we normalize the attention score with the

softmax function and update the context representation:

$$a_{1 \times 1}^l = \frac{\exp(f(r^l))}{\exp(f(r^l)) + \exp(f(r^r))}, \quad (13)$$

$$r_{2d \times 1}^l = \alpha_{1 \times 1}^l \times r_{2d \times 1}^l. \quad (14)$$

Using the same logic, we update the target pair and obtain the updated representations r^l, r^r, r_l^t , and r_r^t . [33] argues that it is optimal to repeat this procedure three times. We inherit this idea and repeat this mechanism for three hops. The final four representations are concatenated and processed by a Multi-Layer Perceptron (MLP). The mathematical notation \oplus in Equation 15 denotes vector concatenation. Last, we use softmax to calculate the final prediction probability for each sentiment polarity (p), which is a three-dimensional vector as we consider three sentiment polarities, i.e., positive, neutral, and negative:

$$r_{8d \times 1} = r_{2d \times 1}^l \oplus r_{2d \times 1}^r \oplus r_{2d \times 1}^l \oplus r_{2d \times 1}^t, \quad (15)$$

$$p_{3 \times 1} = \text{softmax}(MLP(r)). \quad (16)$$

We calculate the final loss function for LCR-Rot-hop++ sentiment classification by taking the cross-entropy of the predicted sentiment and actual sentiment label of the j^{th} sentence denoted as a_j over M input sentences:

$$L_{sc} = - \sum_{j=1}^M a_j \log(p_{3 \times 1}^{(j)}) + \lambda \|\theta_2\|^2, \quad (17)$$

where $p^{(j)}$ is the prediction probability of the j^{th} instance, $\|\theta_2\|^2$ is a $L2$ -norm regularization term, which determines the penalty of having a certain parameter set, and λ is a weight for this term.

Additionally, before we train the model using source domain data, we tune the hyperparameters of the LCR-Rot-hop++ model. Such hyperparameters are the learning rate, dropout rate, momentum term, and $L2$ -norm regularization term. We use the same TPE algorithm as the hyperparameter tuning process of the DIWS model. After obtaining optimal hyperparameter values, we use full source data to train the model. Note that the training sample of this model purely consists of source domain data and the testing sample purely consists of target domain data. Hence, we split the source domain data into 80% pure training samples and 20% validation samples to perform hyperparameter tuning.

5. EVALUATION

In this section, we present the result of our evaluation. First, in Subsection 5.1, we evaluate the performance of the domain classification. Then, in Subsection 5.2, we present the performance of the aspect-based sentiment classification on the target domain. Next, in Subsection 5.3, we compare the results of our proposed model with the original LCR-Rot-hop++ model and the ones of DAT-LCR-Rot-hop++. Last, in Subsection 5.4, we give insights into the obtained results.

5.1 Domain Classification Performance

Table 3 shows the domain classification accuracy of the DIWS model. The training sample consists of 80% of the randomly mixed source and target domain data, and the testing sample consists of the remaining 20% of the data.

Table 3: Domain classification accuracy of DIWS.

Source-target domain	Domain classification test accuracy
MP3 Player - DVD Player	0.824
MP3 Player - Digital Camera	0.884
MP3 Player - Hotel	0.979
MP3 Player - Cell Phone	0.824
Average	0.878

On average, the DIWS model can well classify the source domain and target domain with an average accuracy of 0.878. It implies the robustness of the attention weights from the DIWS model. In particular, the accuracy for the hotel domain is relatively high compared to the other domains. It signals that the difference between the target domain and the source MP3 player domain is greater for the hotel domain.

5.2 Cross-Domain ABSC Performance

We measure a test accuracy for different values of percentile threshold $K = \{10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$. Here, percentile threshold K implies that there exists $K\%$ of the words with lower attention weight than the corresponding attention threshold. For example, $K = 80$ suggests that we only keep the 80% of the words in the input sentences that are the most domain-independent, and use this transformed input for the LCR-Rot-hop++ model. Figures 4 to 7 display the change in test accuracy level for different target domains as the threshold percentile increases. The lowest accuracy is colored dark gray and the highest accuracy is colored light gray.

Note that $K = 100$ refers to the case that every word is classified as domain-agnostic regardless of their DIWS attention weights. Hence, it represents the original LCR-Rot-hop++ model that is purely trained by MP3 player domain data and tested on the DVD player data. This interpretation of $K = 100$ applies to all target domains.

For the DVD player domain, the DIWS-LCR-Rot-hop++ test accuracy varies from 55% to 71%. The model attains the lowest accuracy when $K = 10$ and attains the highest accuracy when $K = 70$. Also, there exists a general trend that the accuracy increases as we reduce the proportion of domain-specific discarded words, and achieves maximum accuracy when $K = 70$. The accuracy drops at $K = 80$ but bounces again as we reduce the discarded words to the extreme.

For the digital camera domain, the DIWS-LCR-Rot-hop++ test accuracy varies from 66% to 74%. The model attains the lowest accuracy when $K = 40$ and attains the highest accuracy when $K = 60$. Unlike the DVD player domain, DIWS-LCR-Rot-hop++ has a relatively good performance when we discard a large proportion of words. For instance, the accuracy gap between the maximum accuracy at $K = 60$ and accuracy for the low threshold values $K = 10, 20, 30$ is

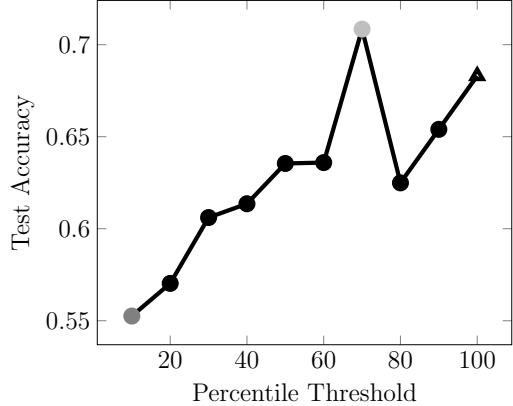


Figure 4: Test accuracy of DVD player domain data under DIWS-LCR-Rot-hop++ model trained with MP3 player domain data.

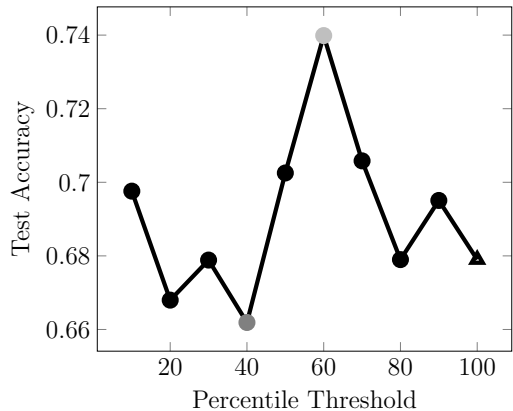


Figure 5: Test accuracy of digital camera domain data under DIWS-LCR-Rot-hop++ model trained with MP3 player domain data.

not as large as for the DVD player domain. Additionally, after having the highest accuracy at $K = 60$, the accuracy diminishes as K increases to the extreme.

For the hotel domain, the DIWS-LCR-Rot-hop++ test accuracy varies from 63% to 73%. The model attains the lowest accuracy at $K = 100$ and attains the highest accuracy at $K = 20$. Unlike the other target domains, DIWS-LCR-Rot-hop++ performs the best for the small K value ($K = 20$), and the accuracy decreases until the model attains the lowest accuracy at $K = 100$, although there are some local peaks at $K = 60$ and $K = 80$.

For the cell phone domain, the test accuracy varies from 64% to 77%. The model attains the minimum accuracy at $K = 80$ and attains the highest accuracy at $K = 90$, without a clear trend. DIWS-LCR-Rot-hop++ has a sudden dip and spike at $K = 80$ and $K = 90$, respectively. After the global maximum accuracy, the accuracy decreases again at $K = 100$.

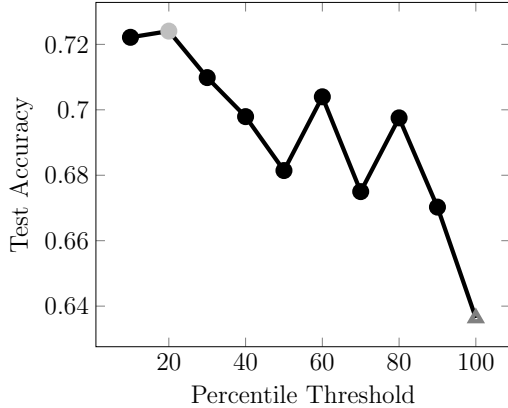


Figure 6: Test accuracy of hotel domain data under DIWS-LCR-Rot-hop++ model trained with MP3 player domain data.

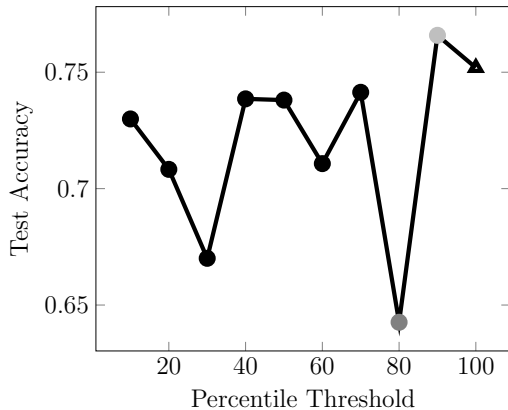


Figure 7: Test accuracy of cell phone domain data under DIWS-LCR-Rot-hop++ model trained with MP3 player domain data.

5.3 Comparison with LCR-Rot-hop++ and DAT-LCR-Rot-hop++

To assess the overall performance of the DIWS-LCR-Rot-hop++, we compare the results to the original LCR-Rot-hop++ model and the DAT-LCR-Rot-hop++ model. For each target domain, we choose the threshold value K that attains maximum accuracy. Table 4 displays the test accuracy of the three models.

On average, DIWS-LCR-Rot-hop++ outperforms LCR-Rot-hop++ and DAT-LCR-Rot-hop++ in three out of four target domains. The hotel domain is the only target domain in which DAT-LCR-Rot-hop++ outperforms DIWS-LCR-Rot-hop++. In comparison to DAT-LCR-Rot-hop++, the absolute performance difference is small for the camera, hotel, and cell phone domains, as it ranges from 0.024 to 0.047. On the other hand, the DVD player domain experiences performance enhancement to a great extent (0.24). Compared to LCR-Rot-hop++, the performance difference is marginal for the DVD player and cell phone domains, while the digital

camera and hotel domains experience considerable performance enhancement. Overall, we conclude that the DIWS-LCR-Rot-hop++ model improves the cross-domain ABSC performance compared to LCR-Rot-hop++ and DAT-LCR-Rot-hop++ for the source and target domains that we use.

5.4 Insights

In general, DIWS-LCR-Rot-hop++ performs better than the random guessing baseline (0.33 if there are three sentiment classes; 0.5 if there are two sentiment classes) in all domains. For the DVD player and cell phone domains, it even outperforms the majority guessing baseline (see Table 1 for the distribution of the sentiment) without having the information about the sentiment distribution. Furthermore, it outperforms the cross-domain performance of the original LCR-Rot-hop++ model ($K = 100$) and the DAT-LCR-Rot-hop++ model if we choose the optimal threshold K for each target domain, besides the hotel domain.

Next, let us investigate the overall pattern of accuracy level as K increases and the reasoning behind the results. First, the results show that the model performs the best when we discard 10% to 40% of the domain-specific words for three out of the four domains (DVD player, camera, and cell phone). If we discard too many words which correspond to the low values of K , the accuracy is below average for most of the data except for the hotel domain. This is due to the excessive information loss. Accordingly, discarding most of the domain-specific words would enhance the domain-invariance of the input texts, but it also makes the remaining sentence useless for ABSC. For example, Figure 8 displays one of the review sentences in the DVD player domain.

When $K = 10$, the only remaining tokens other than the target word `rewind` are `i`, `'`, `ve`, and `had`. Accordingly, they are domain-independent, unlike the words such as `fast`, `forward`, and `smoothly`, but it would be difficult to correctly classify the sentiment with them.

“its `fast-forward` and `rewind` work much more
`smoothly` and consistently than those of other
models `i've had`”

Figure 8: A review sentence in the DVD player domain [14].

On the other hand, the hotel domain obtains its maximum accuracy at $K = 20$, where we discard 80% of the words. This difference is due to the difference in closeness between the source domain and the target domain. Note that the hotel domain is even more distinct from the source MP3 player domain compared to the other target domains such as DVD player, digital camera, and cell phone. The high domain-classification accuracy for the hotel domain (0.979) in Table 3 supports this claim because it would be easier to classify the domain if there exists a large difference between the target and source domains. Thus, for the hotel domain, the accuracy gained from discarding the domain-specific words outweighs the accuracy drop due to the information loss.

There are mixed results about the values of K for which the model performs the worst. For the DVD player domain, the minimum accuracy is measured when $K = 10$, the digital camera domain performs worst when K is around the

Table 4: The accuracy comparison between DIWS-LCR-Rot-hop++ (DIWS++), LCR-Rot-hop++(LCR++), and DAT-LCR-Rot-hop++ (DAT++).

Source-target domain	DIWS++ accuracy	LCR++ accuracy	DAT++ accuracy
MP3 Player - DVD Player	0.708	0.683	0.470
MP3 Player - Digital Camera	0.740	0.679	0.693
MP3 Player - Hotel	0.724	0.636	0.751
MP3 Player - Cell Phone	0.766	0.752	0.742

Note: We train the DAT-LCR-Rot-hop++ with our selection of datasets using its source code as the original paper does not use our datasets.

median ($K = 40$), and the hotel and cell phone domains have their worst accuracy at relatively high values of K ($K = 100$ and $K = 80$, respectively). This observation implies that we should first find an optimal value of K before using the DIWS-LCR-Rot-hop++ model as there is no golden rule that the low values of K always produce worse accuracy and high values of K are superior. This may be due to the difference in distributions of the words between the domains, where some domains may contain more innate domain-independent words such as adverbs and linking words than others.

To conclude, the results show that discarding domain-specific words leads DIWS-LCR-Rot-hop++ to perform better than the original LCR-Rot-hop++ model under the cross-domain ABSC task, while the optimal proportion of remaining words after the dropout depends on the degree of closeness between the source domain and the target domain. In general, discarding an excessive proportion of words even further worsens the performance of DIWS-LCR-Rot-hop++ compared to the original LCR-Rot-hop++ where we do not discard any of the words. Discarding domain-specific words indeed improves the performance of cross-domain aspect-based sentiment analysis when we discard 10% to 40% of the words if the target domain and source domain are not very different. If we drop too many words, the model experiences a performance drop due to the information loss. However, if we recognize that the source and target domains are distinct from each other, we should discard a large percentage of words (80% for the hotel domain) as the accuracy gained from discarding the domain-specific words outweighs the information loss effect. The domain classification accuracy from DIWS would be a good indication of whether a target domain is very different from a source domain or not.

6. CONCLUSION

To apply the state-of-the-art LCR-Rot-hop++ model to the cross-domain setting, this work proposes the DIWS model to select and discard the domain-specific words. The proposed model for cross-domain ABSC is the DIWS-LCR-Rot-hop++ model. It utilizes a domain classification architecture with a feed-forward attention layer to filter out the domain-specific words with attention weights higher than a certain threshold. Then we analyze the performance of our proposed model for ten different threshold values. Based on the experiments on five datasets, we conclude that without any sentiment label of the target domain data, our model effectively enhances the accuracy by discarding the domain-specific words from the source and target domain data.

Furthermore, we have found that there is a danger of information loss and thus we should select the threshold between domain-specific words and domain-agnostic words carefully. In addition, the results imply that the degree of difference between the source domain and target domain affects the performance of the DIWS-LCR-Rot-hop++ model for different threshold values.

Our research has a few limitations. First, due to the lack of computational power of the testing PC environment, we could not apply DIWS-LCR-Rot-hop++ to the large popular datasets in the field of ABSC. For example, such data includes restaurant domain data and laptop domain data from SemEval 2014 [25] in which DAT-LCR-Rot-hop++ demonstrates high performance on the cross-domain ABSC task. Second, the performance gain from DIWS-LCR-Rot-hop++ is not always positive compared to the original LCR-Rot-hop++. If we do not use the optimal threshold value, the accuracy of our model can be even less than the original model. Therefore, we advise users to run the DIWS-LCR-Rot-hop++ using different threshold values and select the optimal one for the final prediction for every source-target domain combination. Third, DIWS-LCR-Rot-hop++ does not consistently outperform DAT-LCR-Rot-hop++. For the hotel domain, the competing model performs better than the proposed model. Last, the DIWS-LCR-Rot-hop++ model sequentially trains the DIWS component and LCR-Rot-hop++ component. This sequential training process may prevent the model to find the global optimal parameter values during optimization.

Several further research directions are available on this topic. First, the sequential training processes can be merged into a simultaneous optimization in which the final loss function is a sum of the DIWS loss function and the LCR-Rot-hop++ loss function. Second, it is possible to directly utilize the optimal domain-classification attention weights by allocating lower importance to words in an input sentence by their attention weights. For example, words with high attention weights are considered less during the LCR-Rot-hop++ training because they are likely to be domain-specific. Last, it is possible to extend the model to the multi-domain setting where the sentiment-labeled target domain data is partially available. In this case, one can exploit the shared-private framework and thus expect even higher performance results.

7. REFERENCES

- [1] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In *24th Inter-*

- national Conference on Neural Information Processing Systems (NIPS 2011)*, pages 2546–2554. Curran Associates, 2011.
- [2] G. Brauwerters and F. Frasincar. A survey on aspect-based sentiment classification. *ACM Computing Surveys*, 55(4):1–37, 2022.
- [3] E. Cambria, Y. Li, F. Z. Xing, S. Poria, and K. Kwok. SenticNet 6: Ensemble application of symbolic and sub-symbolic AI for sentiment analysis. In *29th ACM International Conference on Information and Knowledge Management (CIKM 2020)*, page 105–114. ACM, 2014.
- [4] X. Chen and C. Cardie. Multinomial adversarial networks for multi-domain text classification. In *2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2018)*, pages 1226–1240. ACL, 2018.
- [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019)*. ACL, 2019.
- [6] L. Dong, F. Wei, C. Tan, D. Tang, M. Zhou, and K. Xu. Adaptive recursive neural network for target-dependent Twitter sentiment classification. In *52nd Annual Meeting of the Association for Computational Linguistics (ACL 2014)*, pages 49–54. ACL, 2014.
- [7] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- [8] S. Gao and Y.-K. Ng. Generating extractive sentiment summaries for natural language user queries on products. *SIGAPP Applied Computing Review*, 22(2):5–20, 2022.
- [9] A. Graves and J. Schmidhuber. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, 18(5-6):602–610, 2005.
- [10] T. R. Gruber. A translation approach to portable ontology specifications. *Knowledge Acquisition*, 5(2):199–220, 1993.
- [11] J. Guo, D. J. Shah, and R. Barzilay. Multi-source domain adaptation with mixture of experts. In *2018 Conference on Empirical Methods in Natural Language Processing (EMNLP 2018)*, pages 4694–4703. ACL, 2018.
- [12] Y. Han and Y. Kim. An extracting method of movie genre similarity using aspect-based approach in social media. *SIGAPP Applied Computing Review*, 17(2):36–45, 2017.
- [13] S. Hochreiter and J. Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [14] M. Hu and B. Liu. Mining and summarizing customer reviews. In *10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2004)*, pages 168–177. ACM, 2004.
- [15] N. Hu, J. Zhang, and P. A. Pavlou. Overcoming the J-shaped distribution of product reviews. *Communications of the ACM*, 52(10):144–147, 2009.
- [16] J. Knoester, F. Frasincar, and M. Truşcă. Domain adversarial training for aspect-based sentiment analysis. In *23rd International Conference on Web Information Systems (WISE 2022)*, page 21–37. Springer, 2022.
- [17] J. Lee, F. Frasincar, and M. M. Trusca. A cross-domain aspect-based sentiment classification by masking the domain-specific words. In *38th ACM/SIGAPP Symposium on Applied Computing (SAC 2023)*, page 1595–1602. ACM, 2023.
- [18] Z. Li, Y. Wei, Y. Zhang, and Q. Yang. Hierarchical attention transfer network for cross-domain sentiment classification. In *32nd AAAI Conference on Artificial Intelligence (AAAI 2018)*, pages 5852–5859. AAAI Press, 2018.
- [19] B. Liu. *Sentiment Analysis - Mining Opinions, Sentiments, and Emotions*. Cambridge University Press, second edition, 2020.
- [20] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. In *1st International Conference on Learning Representations (ICLR 2013)*. OpenReview.net, 2013.
- [21] G. A. Miller. WordNet: A lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995.
- [22] D. Naous and C. Legner. A preference-based approach to incorporate the “voice of the customer” in mass-market software product design. *SIGAPP Applied Computing Review*, 20(4):35–49, 2021.
- [23] S. Oepen, M. Kuhlmann, Y. Miyao, D. Zeman, S. Cinková, D. Flickinger, J. Hajič, and Z. Urešová. SemEval 2015 task 18: Broad-coverage semantic dependency parsing. In *9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 915–926. ACL, 2015.
- [24] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. Deep contextualized word representations. In *2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2018)*, pages 2227–2237. ACL, 2018.
- [25] M. Pontiki, D. Galanis, J. Pavlopoulos, H. Papageorgiou, I. Androutsopoulos, and S. Manandhar. SemEval-2014 task 4: Aspect based sentiment analysis. In *8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 27–35. ACL, 2014.
- [26] S. Ruder, P. Ghaffari, and J. G. Breslin. A hierarchical model of reviews for aspect-based sentiment analysis. In *2016 Conference on Empirical Methods in Natural*

- Language Processing (EMNLP 2016)*, pages 999–1005. ACL, 2016.
- [27] K. Schouten and F. Frasincar. Survey on aspect-level sentiment analysis. *IEEE Transactions on Knowledge and Data Engineering*, 28(3):813–830, 2016.
- [28] K. Schouten and F. Frasincar. Ontology-driven sentiment analysis of product and service aspects. In *15th Extended Semantic Web Conference (ESWC 2018)*, pages 608–623. Springer, 2018.
- [29] X. Sun and W. Lu. Understanding attention for text classification. In *58th Annual Meeting of the Association for Computational Linguistics, (ACL 2020)*, pages 3418–3428. ACL, 2020.
- [30] M. M. Truşcă and F. Frasincar. Survey on aspect detection for aspect-based sentiment analysis. *Artificial Intelligence Review*, 56(5):3797–3846, 2023.
- [31] M. M. Truşcă, D. Wassenberg, F. Frasincar, and R. Dekker. A hybrid approach for aspect-based sentiment analysis using deep contextual word embeddings and hierarchical attention. In *20th International Conference of Web Engineering (ICWE 2020)*, pages 365–380. Springer, 2020.
- [32] S. van Berkum, S. van Megen, M. Savelkoul, P. Weterman, and F. Frasincar. Fine-tuning for cross-domain aspect-based sentiment classification. In *20th IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT 2021)*, page 524–531. ACM, 2021.
- [33] O. Wallaart and F. Frasincar. A hybrid approach for aspect-based sentiment analysis using a lexicalized domain ontology and attentional neural models. In *16th Extended Semantic Web Conference (ESWC 2019)*, pages 363–378. Springer, 2019.
- [34] J. Yuan, Y. Zhao, and B. Qin. Learning to share by masking the non-shared for multi-domain sentiment classification. *International Journal of Machine Learning and Cybernetics*, 13(9):2711–2724, 2022.
- [35] M. Zhang, Y. Zhang, and D.-T. Vo. Gated neural networks for targeted sentiment analysis. In *30th AAAI Conference on Artificial Intelligence (AAAI 2016)*, page 3087–3093. AAAI Press, 2016.
- [36] Y. Zhang, J. Wang, and X. Zhang. Conciseness is better: Recurrent attention LSTM model for document-level sentiment analysis. *Neurocomputing*, 462:101–112, 2021.
- [37] H. Zhao, S. Zhang, G. Wu, J. M. F. Moura, J. P. Costeira, and G. J. Gordon. Adversarial multiple source domain adaptation. In *32nd International Conference on Neural Information Processing Systems (NIPS 2018)*, pages 8568–8579. Curran Associates, 2018.
- [38] R. Zheng, J. Chen, and X. Qiu. Same representation, different attentions: Shareable sentence representation learning from multiple tasks. In *27th International Joint Conference on Artificial Intelligence (IJCAI 2018)*, pages 4616–4622. AAAI Press, 2018.
- [39] S. Zheng and R. Xia. Left-Center-Right separated neural network for aspect-based sentiment analysis with rotatory attention. *arXiv preprint arXiv:1802.00892*, 2018.
- [40] Y. Ziser and R. Reichart. Pivot based language modeling for improved neural domain adaptation. In *2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2018)*, pages 1241–1251. ACL, 2018.