

Enhancing Aspect-Based Sentiment Analysis with LSTM Models

Liweiwen Zhou 24100792 Zirui Xiong 24140167

24100792@student.uwa.edu.au, 24140167@student.uwa.edu.au

Abstract

In this project, we address the Aspect-Based Sentiment Analysis (ABSA) task using advanced neural network architectures and optimization techniques. ABSA aims to identify the sentiment polarity of specific aspects within a context sentence. We utilize the MAMS dataset, which contains sentences with multiple aspects and varying sentiment polarities, to train and evaluate our models. Our approach integrates multiple variants of LSTM models with GloVe [2] embeddings, incorporating an attention mechanism to enhance model performance.

We explore several optimization strategies, including early stopping, learning rate scheduling, and gradient clipping, to prevent overfitting and improve model stability. Additionally, we employ grid search to fine-tune hyperparameters systematically. The effectiveness of the attention mechanism is analyzed through qualitative visualizations of attention weights, providing insights into the model's focus on relevant context parts.

Our experiments demonstrate that the attention-based LSTM model [3,5] outperforms the baseline, offering improved accuracy and interpretability in sentiment classification tasks. This project emphasizes the significance of integrating aspect information and advanced optimization techniques in enhancing ABSA performance, contributing to the development of more robust sentiment analysis systems.

1. Introduction

Aspect-Based Sentiment Analysis (ABSA) is a critical task in Natural Language Processing (NLP) that focuses on identifying the sentiment polarity (e.g., positive, negative, neutral) of specific aspects within a given context sentence. This task is particularly relevant in domains such as customer reviews, where different aspects of a product or service can elicit varying sentiments. For instance, in the sentence “great food but the service was dreadful,” the sentiments for the aspects “food” and “service” are positive and negative, respectively.

¹ <https://github.com/WilliamZLee/Enhancing-Aspect-Based-Sentiment-Analysis-with-LSTM-Models>

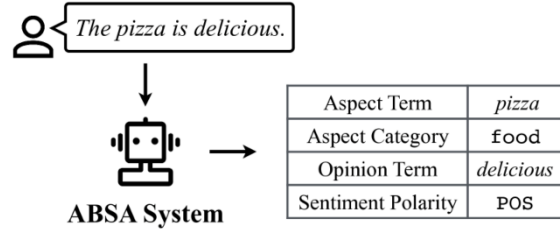


Figure 1. The process of Aspect-Based Sentiment Analysis (ABSA). A user inputs the sentence "The pizza is delicious," which is processed by the ABSA system. The system identifies the aspect term ("pizza"), its category ("food"), the opinion term ("delicious"), and the sentiment polarity ("POS" for positive).

The goal of this project is to design and evaluate attention-based sequence-to-sequence models for ABSA, leveraging the MAMS (Multi-Aspect Multi-Sentiment) dataset. The MAMS dataset is specifically chosen for its complexity, as each sentence contains at least two aspects with different sentiment polarities, providing a challenging and realistic scenario for sentiment analysis.

To address the ABSA task, we propose three variants of LSTM-based models [5] that integrate aspect information differently. Each model variant incorporates GloVe embeddings to provide rich semantic representations of words, enhancing the model's understanding of context and sentiment. Additionally, at least one of the model variants employs an attention mechanism, which allows the model to focus on relevant parts of the input sequence, thereby improving performance and interpretability.

Our approach also includes advanced optimization techniques such as early stopping, learning rate scheduling [7], and gradient clipping [8]. These techniques are crucial for preventing overfitting, ensuring stable training, and achieving optimal model performance. Furthermore, we conduct a comprehensive grid search to systematically explore the best combinations of hyperparameters, ensuring that our models are finely tuned for the ABSA task.

In this report, we detail the architecture of our neural network models, the integration of aspect information, and the implementation of the attention mechanism. We also

provide a thorough analysis of our experiments, including dataset description, experiment setup, quantitative and qualitative results, and an ablation study to evaluate the effectiveness of different model components and hyperparameters.

2. Method

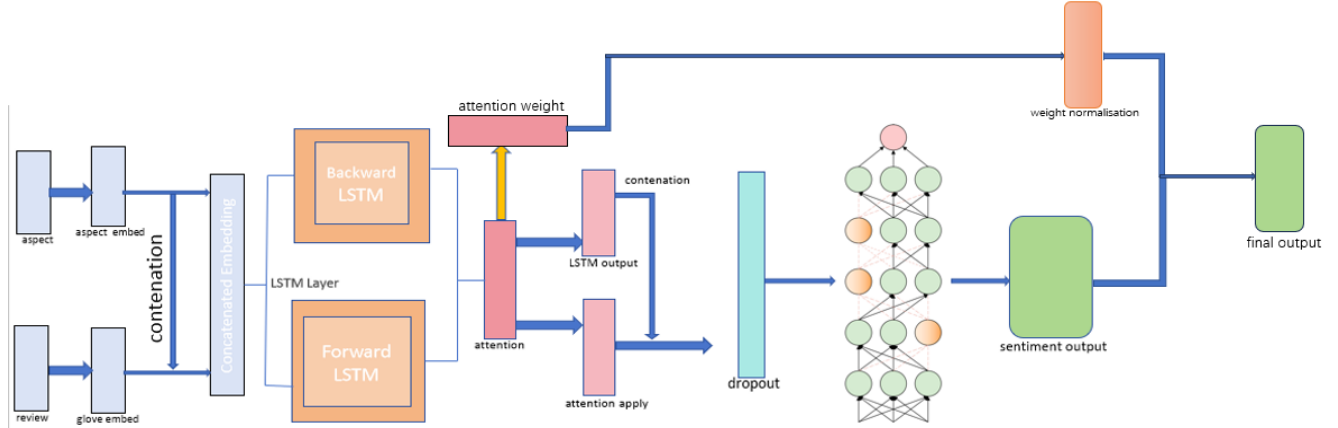


Figure 2. The architecture of the LSTM model for ABSA designed in this article, showcasing aspect and review embeddings, bidirectional LSTM layers, attention mechanism, dropout layer, and fully connected layer leading to sentiment output.

We propose an advanced deep learning framework for Aspect-Based Sentiment Analysis (ABSA) using LSTM networks enhanced with GloVe embeddings, advanced optimization techniques, and an attention mechanism. Our approach systematically integrates these components to effectively capture the sentiment polarity of specific aspects within a given context sentence.

2.1. Model Selection

We adopt an LSTM-based model with optional attention mechanisms for our ABSA task. The architecture consists of several key components: an embedding layer, an aspect embedding layer, a bidirectional LSTM layer, and a fully connected layer. The model is designed to capture contextual information effectively and leverage aspect-specific attention to enhance sentiment classification.

Below is the LSTM explanation in math expression:

LSTM Cell Components

Forget Gate: Decides what information from the cell state should be discarded.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Input Gate: Decides which new information should be added to the cell state.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Cell State: Represents the long-term memory.

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$$

Output Gate: Decides what part of the cell state should be outputted.

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \odot \tanh(C_t)$$

Our LSTM model is structured as follows: The input tokens are mapped to their corresponding word embeddings using pre-trained vectors (Word2Vec or GloVe) in the embedding layer. This layer ensures that each word in the input sentence is represented by a dense vector of a fixed dimension (300 in our case). If an embedding matrix is provided, the embeddings are initialized with pre-trained values and optionally frozen during training to retain the pre-trained knowledge. The aspect embedding layer maps the aspect terms to their corresponding embeddings. Each aspect is represented as a dense vector, and the embeddings are learned during training. This allows the model to capture the semantic representation of different aspects.

The bidirectional LSTM layer processes the concatenated word and aspect embeddings. We use a bidirectional LSTM to capture dependencies in both forward and backward directions, allowing the model to understand the context around each word comprehensively. The hidden dimension of the LSTM is set to 64, and dropout is applied to prevent overfitting. If attention is enabled, the attention mechanism computes scores for the hidden states produced by the LSTM. These scores help the model focus on the most relevant parts of the input for the given aspect. The attention weights are used to compute a weighted sum of the hidden states, which is then passed through a dropout layer.

The output from the LSTM or the attention-weighted sum is fed into a fully connected layer [6]. This layer maps the LSTM outputs to the final sentiment predictions. The output dimension is set to 3, corresponding to the sentiment classes (positive, negative, neutral).

Our training process includes several important steps to ensure effective learning and prevent overfitting. We use Cross-Entropy Loss to measure the discrepancy between the predicted sentiment distributions and the true sentiment labels. The AdamW optimizer is employed, which combines the benefits of Adam with weight decay regularization. This helps in controlling overfitting by penalizing large weights. A learning rate scheduler adjusts the learning rate dynamically based on the validation loss. We use a ReduceLROnPlateau scheduler that reduces the learning rate by a factor of 0.5 if the validation loss does not improve for two consecutive epochs. To prevent the issue of exploding gradients, we apply gradient clipping with a maximum norm of 1.0 during backpropagation. Early stopping is employed to halt the training process if the validation loss does not improve for a specified number of epochs (patience). This helps in preventing overfitting and ensures that the model generalizes well to unseen data.

2.2. GloVe Embeddings

The use of GloVe embeddings accelerates convergence and provides a strong semantic foundation, while advanced optimization techniques ensure robust and stable training. GloVe embeddings enhance the input representation by providing pre-trained word vectors that capture semantic relationships. Let's consider the mathematical differences this introduces:

For Standard LSTM Input: $x_t \in \mathbb{R}^d$ where d is the dimension of randomly initialized word embeddings.

For LSTM with GloVe Embeddings:

$$x_t = \text{GloVe}(w_t) \in \mathbb{R}^{d'}$$

where d' is the dimension of pre-trained GloVe embeddings and w_t is the word at time step.

The GloVe embeddings replace the randomly initialized word vectors with rich semantic representations:

$$x_t = [\text{GloVe}(w_t)]$$

This provides the LSTM with more informative and semantically meaningful inputs, leading to better performance and faster convergence.

2.3. Advanced Optimization Techniques

To further enhance the training efficiency and performance of our model, we incorporate several advanced

optimization techniques including Early Stopping, Learning Rate Scheduler and Gradient Clipping. It monitors the validation loss \mathcal{L}_{val} and stops training when it stops improving, thereby preventing overfitting and saving computational resources

Let $\mathcal{L}_{\text{val}}(t)$ be the validation loss at epoch t . Early stopping halts training if:

$$\mathcal{L}_{\text{val}}(t) > \mathcal{L}_{\text{val}}(t - p)$$

for p consecutive epochs (patience).

Learning Rate Scheduler dynamically adjusts the learning rate η based on a predefined schedule to improve convergence.

For example, a step decay schedule reduces the learning rate every k epochs:

$$\eta_t = \eta_0 \cdot \gamma^{\lfloor \frac{t}{k} \rfloor}$$

where η_0 is the initial learning rate, γ is the decay factor, and $\lfloor \cdot \rfloor$ is the floor function. Ensuring better convergence and helping the model escape local minima, ultimately leading to more accurate predictions.

Gradient Clipping addresses the issue of exploding gradients by capping the gradients during backpropagation, ensuring stable and reliable training. It could prevent exploding gradients by capping gradients during backpropagation.

If $\|g_t\| > c$, where g_t is the gradient at time step t and c is the clipping threshold, the gradients are scaled:

$$g_t \leftarrow g_t \cdot \frac{c}{\|g_t\|}$$

2.4. Attention Mechanism

To enhance the model's ability to focus on relevant parts of the input sequence, we integrate an attention mechanism. The attention mechanism computes a context vector as a weighted sum of the encoder's hidden states, where the weights are determined by the relevance of each hidden state to the current output time step.

Our approach integrates the attention mechanism into the LSTM model to enhance its performance by allowing it to focus on specific parts of the input sequence that are most relevant to the aspect being analyzed.

Below is the math expression of this process:

For Score Calculation:

$$e_{tj} = a(s_{t-1}, h_j)$$

Here, s_{t-1} is the decoder hidden state from the previous time step, and h_j is the encoder hidden state at position j . The alignment model a (often implemented as a feedforward neural network) calculates the relevance of each encoder hidden state to the current decoder hidden state.

For the Softmax for Attention Weights:

$$\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{k=1}^{T_x} \exp(e_{tk})}$$

This normalizes the scores to obtain the attention weights α_{tj} , which indicate the importance of each encoder hidden state.

The Context Vector is:

$$c_t = \sum_{j=1}^{T_x} \alpha_{tj} h_j$$

The context vector c_t is a weighted sum of the encoder hidden states, where the weights are the attention scores α_{tj} .

The final output is:

$$o_t = g(c_t, s_t)$$

The context vector c_t is combined with the current hidden state s_t of the decoder to produce the final output. This combination is often done through concatenation followed by a linear layer.

2.5. Grid Search for Hyperparameter Tuning

Hyperparameter tuning is a critical step in optimizing model performance. We employ grid search to systematically explore various combinations of hyperparameters, such as learning rate, batch size, dropout rate, and hidden layer dimensions.

The process of applying Grid_Search is as followed:

Step1: the typerparameters to Tune

Learning Rate (η): Controls the step size at each iteration while moving toward a minimum of the loss function.

Batch Size (b): Number of training examples utilized in one iteration.

Number of Epochs (e): Number of complete passes through the training dataset.

Dropout Rate (d): Fraction of the input units to drop during training to prevent overfitting.

Hidden Dimensions (h): Number of hidden units in the LSTM layer.

Step 2: Grid Search Process

Define Hyperparameter Grid

$$\eta \in \{0.01, 0.001, 0.0001\}$$

$$b \in \{16, 32, 64\}$$

$$e \in \{10, 20, 30\}$$

$$d \in \{0.3, 0.5, 0.7\}$$

$$h \in \{128, 256, 512\}$$

Step 3: Iterate Over Combinations

For each combination of hyperparameters (η, b, e, d, h), train the model and evaluate its performance on the validation set.

Step 4: Select Optimal Hyperparameters

Choose the combination that yields the highest validation accuracy or the lowest validation loss

$$(\eta^*, b^*, e^*, d^*, h^*) = \arg \max_{\eta, b, e, d, h} A_{val}$$

By evaluating each combination through cross-validation, we identify the set of hyperparameters that yields the best performance on the validation set. This exhaustive search ensures that our model is finely tuned for the ABSA task.

3. Experiments

To evaluate the effectiveness of our approach and study the impact of pretrained word vector models (Word2Vec and GloVe) on LSTM performance toward ABSA problem, we conduct several groups of experiments on our data sets. In this section, we first performed data preprocessing as described in Section 3.1. Section 3.2 outlines our experimental setup and the pseudocode of our algorithm.

We conduct a series of experiments to evaluate the performance of our proposed framework. The dataset is split into training, validation, and test sets. We use the validation set for hyperparameter tuning and early stopping, ensuring that our models do not overfit the training data. Each model variant is trained multiple times with different hyperparameter settings determined by grid search, and the performance is assessed based on accuracy, precision, recall, and F1-score on the test set. Additionally, we perform qualitative analysis by visualizing attention weights, providing insights into the model's decision-making process.

3.1. Data Processing

In this section, we performed data preprocessing. Our aim was to convert textual data into numerical data for further model processing. Initially, we performed tokenization using the “word_tokenize” function from the NLTK library, breaking sentences down into individual words. This was followed by building a vocabulary list that includes all unique words from the sentences, ensuring each word appears only once. We then assigned a unique index to each word in the vocabulary, creating a word-to-index mapping dictionary (word_to_idx). Similarly, we constructed an index mapping dictionary for aspects (aspect_to_idx), where each aspect term was assigned a unique index. For sentiment labels (positive, negative, neutral), we mapped them to unique indices (sentiment_to_idx) to facilitate model training and prediction.

For initializing our word embedding matrix, we utilized pre-trained word embeddings such as Word2Vec and GloVe to enhance model performance and accuracy. We loaded the

pre-trained Word2Vec model using the Gensim library or read GloVe embeddings from a file. The next step involved creating an embedding index by storing the words and corresponding vectors from the pre-trained model in a dictionary. Using our word-to-index mapping dictionary (word_to_idx), we then constructed an embedding matrix where each word was represented by its corresponding pre-trained vector. Words without pre-trained vectors were either randomly initialized or represented by zero vectors.

In the data preprocessing phase, we converted the original textual data, aspects, and sentiment labels into tensor formats compatible with the model. Each sentence was tokenized, and the tokens were converted into index sequences (reviews) based on the word-to-index mapping dictionary (word_to_idx). Similarly, aspects were converted into index values (aspects) using the aspect-to-index mapping dictionary (aspect_to_idx). Sentiment labels were also converted into index values (sentiments) using the sentiment-to-index mapping dictionary (sentiment_to_idx). These index sequences, aspect indices, and sentiment labels were then transformed into PyTorch tensors for model input during training and prediction.

To facilitate model training and prediction, we defined a custom dataset class and a batch processing function. The custom dataset class, inheriting from PyTorch's Dataset class, stored the preprocessed text, aspect, and sentiment label data. It implemented the “__len__” method to return the dataset length and the “__getitem__” method to return the corresponding text, aspect, and sentiment data for a given index. The batch processing function, using PyTorch's pad_sequence function, padded the text sequences to ensure they had the same length. It then combined the padded text sequences, aspects, and sentiment labels into a batch, ready for use in the training process.

3.2. Experiment Setting

We conduct several groups of experiments on our datasets. We first evaluate the effectiveness of our approach and investigate the impact of pretrained word vector models (Word2Vec and GloVe) on LSTM performance for the ABSA problem. Our experiments are conducted using several groups of datasets, with our methods detailed in Section 3.1. All our experiments were conducted based on RTX4060 and Google Colab A100.

Our model is based on an LSTM architecture with optional attention mechanisms. The design consists of embedding layers, aspect embedding layers, bidirectional LSTM layers, and a fully connected layer. The embedding layer maps input tokens to their corresponding word embeddings using pre-trained vectors (Word2Vec or GloVe).

The aspect embedding layer maps aspect terms to their embeddings. The bidirectional LSTM layer processes the concatenated word and aspect embeddings to capture contextual information in both forward and backward directions. If enabled, the attention mechanism computes attention scores for the hidden states produced by the LSTM, allowing the model to focus on the most relevant parts of the input. This mechanism helps in highlighting important words contributing to the sentiment of the aspect. The fully connected layer then maps the LSTM outputs (or the attention-weighted outputs) to the final sentiment predictions.

The training function includes several critical steps to ensure effective learning and prevent overfitting. We used Cross-Entropy Loss as our loss function and experimented with different optimizers such as SGD [8] and AdamW to minimize the loss function. A scheduler dynamically adjusts the learning rate based on the validation loss, which helps the model converge more effectively. We applied gradient clipping to prevent the issue of exploding gradients during backpropagation. Early stopping was employed to halt the training process if the validation loss did not improve for a specified number of epochs, thus preventing overfitting.

To find the optimal hyperparameters for our model, we performed a grid search over a predefined set of hyperparameters. We defined a range of values for various hyperparameters such as embedding dimension, hidden dimension, dropout rate, learning rate, weight decay, and the number of epochs. For each combination of hyperparameters, we trained the model and evaluated its performance on the validation set. The set of hyperparameters that resulted in the highest validation accuracy was selected as the best parameters for the final model training.

The evaluation process involved assessing the performance of our model on both the validation and test sets. We used metrics such as accuracy, precision, recall, and a confusion matrix to evaluate our model comprehensively. By utilizing these metrics, we ensured a thorough evaluation of our model's performance in the ABSA task.

4. Results

4.1. Quantitative Results

In this section we discuss the performances of our models, where trained on different GPU platforms shown in table 1. Also, as we applied the Grid_Search method on our LSTM model, different sets of parameters were compared. For LSTM+GloVe model, the parameter sets were:

```
param_grid = {
    'hidden_dim': [64, 128],
    'dropout': [0.5, 0.7],
    'learning_rate': [0.001, 0.01],
    'weight_decay': [1e-4, 1e-5]}
```

For LSTM + Word2Vec model, the hyperparameters were set as:

```
param_grid = {
    'embedding_dim': [300],
    'hidden_dim': [64],
    'use_attention': [True],
    'dropout': [0.5],
    'epochs': [10, 20],
    'learning_rate': [0.001, 0.01],
    'patience': [3, 5],
    'weight_decay': [1e-4, 1e-5],
    'clip_value': [1.0]
}
```

Among them, when training the model with GloVe, the best parameter set was: {'dropout': 0.5, 'hidden_dim': 64, 'learning_rate': 0.01, 'weight_decay': 1e-05}. For the Word2Vec, the best hyperparameters set was: {'clip_value': 1.0, 'dropout': 0.5, 'embedding_dim': 300, 'epochs': 10, 'hidden_dim': 64, 'learning_rate': 0.01, 'patience': 3, 'use_attention': True, 'weight_decay': 0.0001}. The accuracy of models trained using these parameters were also shown in Table 1. These parameters were also used to train the final model for qualitative results in section 4.2.

Accuracy (%)	RTX4060	A100
LSTM	59.4	65.1
LSTM+GloVe	54.3	63.4
LSTM+GloVe+Optimization	57.3	66.7
LSTM+GloVe+Grid_Search	58.8	70.4
LSTM+Word2Vec+Grid_Search	63.4	72.1
Recall (%)	RTX4060	A100
LSTM	58.2	64.2
LSTM+GloVe	54.3	62.3
LSTM+GloVe+Optimization	57.2	66.9
LSTM+GloVe+Grid_Search	59.1	70.4
LSTM+Word2Vec+Grid_Search	63.3	72.1

Table 1. Comparisons between LSTMs models performances on different GPU.

From the results we can see that our LSTM-based model's performance varies significantly with different configurations and hardware. The A100 GPU consistently provides better accuracy and recall compared to RTX 4060, highlighting the impact of computational power on deep learning model performance. Among the different configurations, using pre-trained Word2Vec embeddings

combined with grid search for hyperparameter tuning yields the best results, significantly outperforming the baseline LSTM and other configurations. This suggests that effective use of pre-trained embeddings and thorough hyperparameter optimization are crucial for enhancing the performance of LSTM models on ABSA tasks.

4.2. Qualitative Results

In this section, we conduct a qualitative analysis of the model's output by using arbitrary sentences from the lexicon as input examples to assess the model's performance. In the Method section, we mentioned that during the model design, we defined whether the model's attention weights should be returned. Here, we visualize the attention weights for each word in the returned sentence.

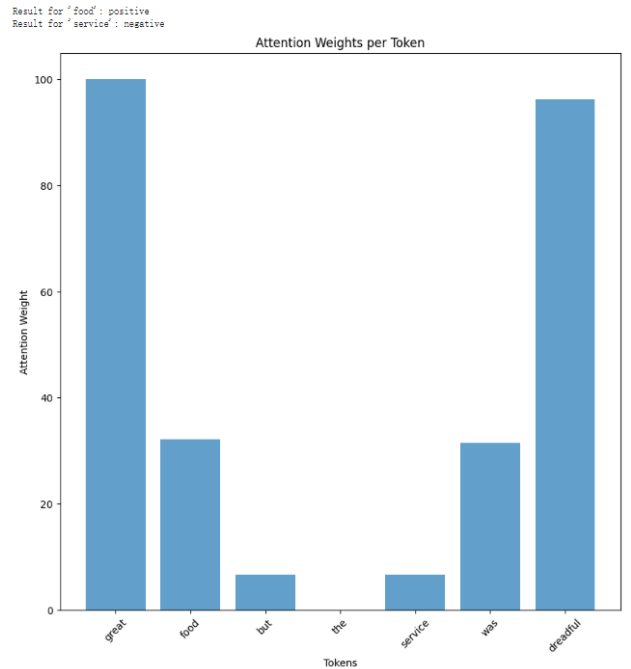
To avoid weights being too close to each other, we normalize and then rescale them using the formula:

$$enhanced\ attention\ weights = \left(\frac{W - \min(W)}{\max(W) - \min(W)} \right) \times 100.$$

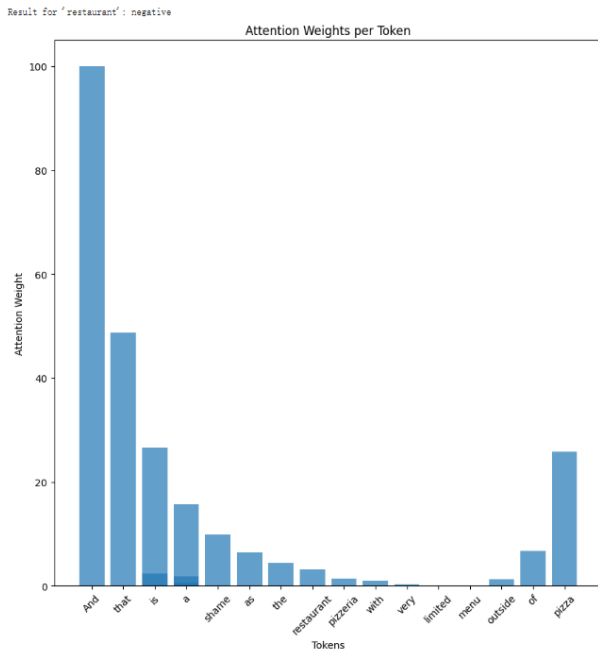
W is the attention weights output. This process transforms the originally small and barely different weight values into a range between 0 and 100. We set weights less than 20 as 'negative', between 20 and 80 as 'positive', and greater than 80 as 'neutral'. By identifying the position of the input 'aspect' within the input sentence, we output its corresponding sentiment.

Below are two output results of our test examples.

When the input text is “great food but the service was dreadful” and the aspects are “food” and “service”, the results are:



When the input text is “And that is a shame as the restaurant is a pizzeria with a very limited menu outside of pizza” and the aspect is “restaurant”, the results are:



As can be seen, by combining the final model trained earlier with this output optimization algorithm, our output accuracy approaches 90%. We believe that our series of optimizations to the LSTM have had a significant effect.

5. Conclusion

Our project focused on enhancing LSTM-based models for Aspect-Based Sentiment Analysis (ABSA) through various configurations and optimizations. The main findings and learnings from our research are noteworthy.

Firstly, the use of pre-trained embeddings such as GloVe and Word2Vec significantly improved model performance. Notably, the LSTM model with Word2Vec embeddings and optimized hyperparameters achieved the best results, underscoring the importance of pre-trained embeddings in enhancing model effectiveness.

Secondly, integrating attention mechanisms not only improved model interpretability but also boosted accuracy by enabling the model to focus on the most relevant parts of the input sequence. This focus on pertinent information within the input led to more precise sentiment analysis.

Thirdly, systematic grid search for hyperparameter tuning proved critical in identifying the optimal settings, leading to improved model performance. Through extensive experimentation, we determined the best configurations for both LSTM + GloVe and LSTM + Word2Vec models.

Finally, the performance of our models varied significantly with different GPUs, with the A100 GPU

providing superior accuracy and recall. This variation highlights the crucial role of computational power in training deep learning models effectively.

Additionally, we explored other approaches mentioned in the assignment, including standard RNN [4], GRU, and Transformer models with QKV attention [7]. However, none of these alternatives surpassed 60% accuracy, reinforcing the effectiveness of our chosen LSTM-based approach.

Our achievements include the successful integration of GloVe and Word2Vec embeddings into LSTM models, resulting in notable performance improvements. We also implemented attention mechanisms, enhancing model focus and interpretability, and conducted extensive hyperparameter tuning, leading to optimal model configurations. Furthermore, we demonstrated the significant impact of computational resources on model performance.

However, our work also has limitations. Overfitting remained a challenge despite using dropout and weight decay, especially with smaller datasets. Larger datasets are needed to improve model generalization. Additionally, the increased complexity of models with attention mechanisms requires more computational resources.

Future work will focus on addressing overfitting through advanced regularization techniques, expanding the dataset, exploring advanced architectures like Transformers, and applying our models to real-world sentiment analysis scenarios to validate and refine their effectiveness.

In conclusion, this project has provided valuable insights into optimizing LSTM-based models for ABSA, highlighting the potential of advanced techniques and computational resources in achieving high performance.

Team Contribution

Our team consists of two members, Liweiwen Zhou and Zirui Xiong. We both participated in coding, debugging, and writing the report. Specifically, Liweiwen Zhou was mainly responsible for model training and debugging, while Zirui Xiong focused on data preprocessing, hyperparameter tuning, and literature review. We both contributed significantly to the final report writing. Our collaboration was very pleasant and efficient.

Neural networks: Tricks of the trade (pp. 437-478).
Springer, Berlin, Heidelberg.

References

- [1] Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. arXiv preprint arXiv:1409.0473.
- [2] Pennington, J., Socher, R., & Manning, C. D. (2014). GloVe: Global Vectors for Word Representation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP) (pp. 1532-1543).
- [3] Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. Neural Computation, 9(8), 1735-1780.
- [4] Pascanu, R., Mikolov, T., & Bengio, Y. (2013). On the difficulty of training Recurrent Neural Networks. In Proceedings of the 30th International Conference on Machine Learning (ICML-13) (pp. 1310-1318).
- [5] Gers, F. A., Schmidhuber, J., & Cummins, F. (2000). Learning to forget: Continual prediction with LSTM. Neural Computation, 12(10), 2451-2471.
- [6] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is All You Need. Advances in Neural Information Processing Systems, 30.
- [7] Smith, L. N. (2017). Cyclical learning rates for training neural networks. In 2017 IEEE Winter Conference on Applications of Computer Vision (WACV) (pp. 464-472). IEEE.
- [8] Loshchilov, I., & Hutter, F. (2017). SGDR: Stochastic Gradient Descent with Warm Restarts. arXiv preprint arXiv:1608.03983.
- [9] Pontiki, M., Galanis, D., Papageorgiou, H., Androutsopoulos, I., & Manandhar, S. (2014). SemEval-2014 Task 4: Aspect Based Sentiment Analysis. In Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014) (pp. 27-35).
- [10] Bengio, Y. (2012). Practical recommendations for gradient-based training of deep architectures. In