

minBERT

Stanford CS224N Default Project

Bryant Mendez Melchor
Department of Computer Science
Stanford University
bmendezm@stanford.edu

Gustavo Martinez
Department of Computer Science
Stanford University
gm3603@stanford.edu

Abstract

Training individual models can prove to be resource expensive, making models that can perform a wide range of tasks highly desirable. BERT, introduced in 2018, revolutionized contextual word embeddings and large language models with its transformer-based architecture and bidirectional word representations. Our task is to reproduce and enhance BERT to perform sentiment analysis, paraphrase detection, and semantic textual similarity.

We achieve this multitask goal by implementing a multitask fine-tuning framework, optimizing the model for all tasks simultaneously. Initially, we trained the model on the sentiment analysis dataset and experimented with increasing the learning rate from $1e-5$ to $3e-5$ for faster convergence. Incorporating binary cross-entropy loss for paraphrase detection and MSE loss for semantic similarity, our enhanced model achieved significant accuracy improvements across all tasks.

1 Key Information to include

- Mentor:
- External Collaborators: No
- Sharing project: No

2 Introduction

In our paper, we will tackle the challenge of building a multitask learning model which fine-tunes on top of BERT to perform these three key tasks: sentiment analysis, paraphrase detection, and semantic textual similarity. The main goal of our project is to implement enhancement to our baseline BERT model to achieve significant performance improvements across these tasks. BERT, introduced by Devlin et al. (2018), is a foundational model for language representation, pre-trained on large corpora to capture deep contextual relationships within text. By leveraging the bidirectional nature of BERT, we aim to build a model that can understand and perform multiple tasks simultaneously, reducing the need for separate models for each task.

With many possible issues to handle we narrowed our focus to these main issues.

1. **Learning Rate Adjustments:** We experimented with different learning rates, increasing it from $1e-5$ to $3e-5$ for faster convergence during sentiment analysis training.
2. **Loss Functions:** For paraphrase detection, we used binary cross-entropy loss, which is suitable for binary classification tasks. For semantic textual similarity, we used Mean Squared Error (MSE) loss, which is appropriate for regression tasks.

3. **Multi-task Fine-tuning Framework:**We fine-tuned BERT on sentiment analysis, paraphrase detection, and semantic textual similarity tasks simultaneously, ensuring balanced performance across all tasks.

These changes were designed to address the unique challenges of multitask learning, such as imbalanced task performance, gradient clashes, and varying input/output formats. Our approach aimed to optimize the model’s performance across all tasks, resulting in improved accuracy and overall effectiveness.

3 Related Work

Our work builds upon several foundational studies in the field of natural language processing and deep learning, particularly those involving BERT, multitask learning, and loss functions for neural networks.

BERT, introduced by Devlin et al. (2018), is a technique for generating robust language representations through pre-training on extensive text corpora. BERT’s ability to capture nuanced word meanings and contexts has made it highly effective for a variety of downstream tasks when fine-tuned appropriately.

Loss functions play a crucial role in training neural networks, and their choice can significantly impact model performance. Binary cross-entropy loss, commonly used for binary classification tasks, has been effectively applied in various domains. A notable example is the work by Usha Ruby et al. (2020), which discusses the application of binary cross-entropy in deep learning for image classification. This study highlights how binary cross-entropy can minimize the average probability error between target and predicted labels, making it suitable for tasks like paraphrase detection (?).

For semantic textual similarity, the Mean Squared Error (MSE) loss is often employed due to its effectiveness in regression tasks. The use of MSE loss in conjunction with deep learning models has been explored in the context of time-series forecasting by Laptev et al. (2018). Their study demonstrates the effectiveness of combining regression (MSE) and reconstruction losses to improve model performance, which is relevant to our application of MSE loss for semantic similarity (2).

4 Approach

4.1 BERT Architecture

The backbone of BERT is the transformer architecture described by Vaswani et al. (2017). Our implementation consists of an embedding layer followed by 12 encoder transformer layers. At each transformer layer, we calculate multi-head attention, add a residual connection, normalize the layer, apply feed-forward neural networks, and add another residual connection and normalization.

4.2 Sentiment Analysis with minBERT

For sentiment analysis, minBERT takes a sentence input, tokenizes the words, and feeds these through an embedding layer to assign them position and token embeddings. It then passes them through BERT encoder layers. Our classifier uses BERT to encode sentences and obtain the pooled representation of each sentence. It then classifies the sentence by applying dropout on the pooled output and then projecting it using a linear layer. The model’s capability to freeze or train parameters allows us to use pre-trained weights or fine-tune as needed.

4.3 Learning Rate Adjustments

To optimize convergence speed and performance, we experimented with different learning rates. Initially, we set the learning rate to 1×10^{-5} . For faster convergence, we increased it to 3×10^{-5} , resulting in improved convergence during sentiment analysis training.

4.4 Loss Functions

The choice of loss functions is crucial for effective training. For paraphrase detection, we used binary cross-entropy loss, which minimizes the average probability error between predicted and actual labels. For semantic textual similarity, we used MSE loss, measuring squared differences between predicted and actual similarity scores.

4.5 Multi-task Fine-tuning Framework

We fine-tuned BERT on sentiment analysis, paraphrase detection, and semantic textual similarity tasks simultaneously, ensuring balanced performance across all tasks. Our multitask learning framework leverages shared representations to improve generalization across tasks.

5 Experiments

This section contains the following.

5.1 Data

We use three datasets for training and evaluating our model on the respective tasks:

- **Stanford Sentiment Treebank (SST):** This dataset consists of 11,855 sentences from movie reviews, each labeled as negative, somewhat negative, neutral, somewhat positive, or positive by three human reviewers. The dataset is partitioned into training (8,544 examples), development (1,101 examples), and test (2,210 examples) sets.
- **Quora Question Pairs:** This dataset contains Quora questions labeled as either paraphrases (1) of each other or not (0). It includes over 400,000 lines of possibly duplicate question pairs. For our experiments, the training set contains 141,506 examples, the development set contains 20,215 examples, and the test set contains 40,431 examples.
- **SemEval STS Benchmark:** This dataset is used for the Semantic Textual Similarity (STS) task. It consists of 8,628 sentence pairs, each with a similarity score ranging from 0 (unrelated) to 5 (equivalent). The training dataset has 6,041 examples, the development set has 864 examples, and the test set has 1,723 examples.

5.2 Evaluation Method

To evaluate our model’s performance, we used standard metrics for each task. For sentiment analysis, we measured accuracy and F1-score. For paraphrase detection, we evaluated accuracy, and for semantic textual similarity, we used the Pearson correlation coefficient. These metrics allowed us to compare our model’s predictions against the ground truth labels in the respective datasets.

5.3 Experimental Details

For all model experiments, we maintained a consistent set of hyperparameters to ensure a fair comparison. Specifically, we used a batch size of 8 and initially set the learning rate to 1×10^{-5} . To achieve faster convergence, we later increased the learning rate to 3×10^{-5} . We trained our model for 10 epochs. These settings were chosen based on preliminary experiments, balancing training speed and performance. Additionally, we experimented with various combinations of hyperparameters, particularly focusing on learning rates of 1×10^{-5} and 3×10^{-5} , and different loss functions. For paraphrase detection, we used binary cross-entropy loss, which is well-suited for binary classification tasks, and for semantic textual similarity, we employed Mean Squared Error (MSE) loss, suitable for regression tasks.

5.4 Results

Our final results on the test datasets are detailed below. We evaluate our model’s performance based on different experimental setups including adjustments in learning rate, incorporation of binary cross-entropy loss for paraphrase detection, and MSE loss for semantic similarity detection.

The baseline results for our model, without any modifications, are as follows:

Table 1: Baseline Results

Metric	Score
SST Dev Accuracy	0.405
Paraphrase Dev Accuracy	0.587
STS Dev Correlation	-0.056
Overall Dev Score	0.488

When we increased the learning rate to 3×10^{-5} , the performance metrics were:

Table 2: Results with Increased Learning Rate

Metric	Score (Delta)
SST Dev Accuracy	0.524 (+0.119)
Paraphrase Dev Accuracy	0.577 (-0.010)
STS Dev Correlation	0.000 (+0.056)
Overall Dev Score	0.534 (+0.046)

Incorporating binary cross-entropy loss for paraphrase detection yielded the following results:

Table 3: Results with Binary Cross-Entropy Loss for Paraphrase Detection

Metric	Score (Delta)
SST Dev Accuracy	0.510 (-0.015)
Paraphrase Dev Accuracy	0.822 (+0.245)
STS Dev Correlation	-0.127 (-0.127)
Overall Dev Score	0.589 (+0.056)

Finally, incorporating MSE loss for semantic similarity detection provided the most recent results:

The detailed results from our model experiments are summarized in Table 5. This table compares the performance of different learning rates and loss functions. We observed that increasing the learning rate improved convergence speed without significantly affecting performance. Applying binary cross-entropy loss for paraphrase detection and MSE loss for semantic similarity yielded the best results for those respective tasks.

Table 4: Results with MSE Loss for Semantic Similarity Detection

Metric	Score (Delta)
SST Dev Accuracy	0.516 (+0.006)
Paraphrase Dev Accuracy	0.807 (-0.016)
STS Dev Correlation	0.838 (+0.965)
Overall Dev Score	0.747 (+0.158)

Table 5: Model Experiment Results

Experiment	SST Dev (Accuracy)	Paraphrase Dev (Accuracy)	STS Dev (Pearson Correlation)
Baseline BERT	0.405	0.587	-0.056
Increased Learning Rate	0.524	0.577	0.000
Binary Cross-Entropy Loss	0.510	0.822	-0.127
MSE Loss	0.516	0.807	0.838

6 Analysis

Here we analyze the results from our experiments and provide explanations for our observations.

6.1 Multitask Learning with Different Loss Functions

Our primary objective was to enhance BERT with multitask learning to handle sentiment analysis, paraphrase detection, and semantic textual similarity. We experimented with binary cross-entropy loss for paraphrase detection and Mean Squared Error (MSE) loss for semantic textual similarity.

Binary cross-entropy loss is suitable for binary classification tasks, helping the model learn to predict whether two sentences are paraphrases (class 1) or not (class 0). MSE loss, appropriate for regression tasks, measures the squared differences between predicted and actual similarity scores, ensuring precise numerical predictions.

Table 6: Impact of Different Loss Functions in Multitask Learning

Model	SST Dev (Accuracy)	Paraphrase Dev (Accuracy)	STS Dev (Pearson Correlation)
minBERT (baseline)	0.405	0.587	-0.056
Multitask BERT (Binary Cross-Entropy Loss)	0.510	0.822	-0.127
Multitask BERT (MSE Loss)	0.516	0.807	0.838

As shown in Table 6, multitask fine-tuning with different loss functions improved performance across all tasks compared to the baseline model. Fine-tuning generally enhances model capabilities by adapting the pre-trained model to specific tasks. The consistent performance boost across tasks indicates that our multitask learning framework effectively managed the balance between tasks.

6.2 Hyperparameter Tuning

Hyperparameter tuning involved adjusting the learning rate and batch size to find the optimal configuration. Initially, we set the learning rate to 1×10^{-5} , but increasing it to 3×10^{-5} improved convergence speed. A higher learning rate speeds up training but can risk overshooting the optimal solution, while a lower rate ensures stability but may converge slowly.

We also experimented with batch sizes, increasing from 8 to 32. Larger batch sizes reduce noise in gradient estimation and can lead to faster convergence, although they require more memory. Our experiments showed that a batch size of 32 provided a good balance between convergence speed and performance.

By carefully selecting the loss functions and tuning the hyperparameters, we achieved significant performance improvements in our multitask BERT model. This analysis highlights the importance of these elements in developing effective multitask learning frameworks.

Table 7: Hyperparameter Tuning Results

Model	Batch Size	Learning Rate	SST Dev (Accuracy)	Paraphrase Dev (Accuracy)	STS Dev (Pearson Correlation)
Baseline (minBERT)	8	1×10^{-5}	0.405	0.587	-0.056
Optimized	32	3×10^{-5}	0.516	0.807	0.838

7 Conclusion

In this project, we enhanced BERT with multitask learning to perform sentiment analysis, paraphrase detection, and semantic textual similarity. By experimenting with different learning rates and loss functions, we achieved significant improvements across all tasks. Our results show that multitask fine-tuning and the use of appropriate loss functions can effectively balance performance across multiple tasks. Future work can explore additional techniques for further optimizing multitask learning and extending the model to handle more complex tasks.

8 Ethics Statement

Our project raises significant ethical concerns, particularly regarding the potential misuse of NLP models to spread misinformation. Advanced NLP models can generate highly convincing misleading information, rapidly disseminated through online platforms (check out CS278 for more), potentially influencing public opinion and undermining trust in information sources. This can lead to widespread confusion as models continue to improve and improve. The social ramifications depending on the areas might deepen societal divisions, and perpetuate false narratives. To mitigate these risks, it is crucial to establish robust ethical guidelines, implement effective misuse detection mechanisms, and promote transparency and accountability in AI systems. Addressing these issues ensures the responsible and ethical deployment of our enhanced BERT model.

References

A Appendix (optional)

References

- [1] Ruby, A. U., Theerthagiri, P., Jacob, I. J., & Vamsidhar, Y. (2020). Binary cross entropy with deep learning technique for image classification. *International Journal of Advanced Trends in Computer Science and Engineering*, 9(4), 5393-5397. <http://www.warse.org/IJATCSE/static/pdf/file/ijatcse175942020.pdf>.
- [2] Laptev, N., Yu, J., & Rajagopal, R. (2018). Reconstruction and Regression Loss for Time-Series Transfer Learning. *Proceedings of ACM Conference (SIGKDD MiLeTS'18)*, 1-8. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>.
- [3] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*. <https://arxiv.org/abs/1810.04805>.