

ADVANCED MACHINE LEARNING (BA-64061-001)

ASSIGNMENT 4 - Text Data Report

ABSTRACT

This assignment does sentiment classification on IMDB movie reviews with Recurrent Neural Networks using LSTM. The focus will be on quantifying the performance difference between a trainable embedding layer and pre-trained GloVe embeddings under a limited data condition. The experiment limits reviews to 150 words, uses only 100 training samples, and considers the top 10,000 most frequent words. The results show that pre-trained embeddings indeed improve model generalization when training data is limited.

INTRODUCTION

Sentiment analysis is one of the most practical applications of NLP that helps an organization automatically understand customer opinions from text data. In this work, we will be using an LSTM-based model to classify IMDB movie reviews as positive or negative. We will use two approaches: one using a regular trainable embedding layer and another using frozen pre-trained GloVe embeddings to see which one fares better on a small dataset.

DATASET DESCRIPTION

The IMDB dataset contains 50,000 movie reviews, split into positive and negative sentiments. Only the top 10,000 most frequent words were taken, and each review was truncated up to 150 words. From the training data, only 100 samples were used to simulate a limited-data scenario; validation was conducted on 10,000 samples, and the remaining data for testing.

METHODOLOGY

Data Preprocessing:

The IMDB dataset was tokenized and padded to a uniform length of 150 words.

Model 1 - Trainable Embedding Layer:

A model with an Embedding layer, followed by a bidirectional LSTM of size 32, dropout of 0.4, and finally a sigmoid output layer.

Model 2 – Pretrained GloVe Embeddings:

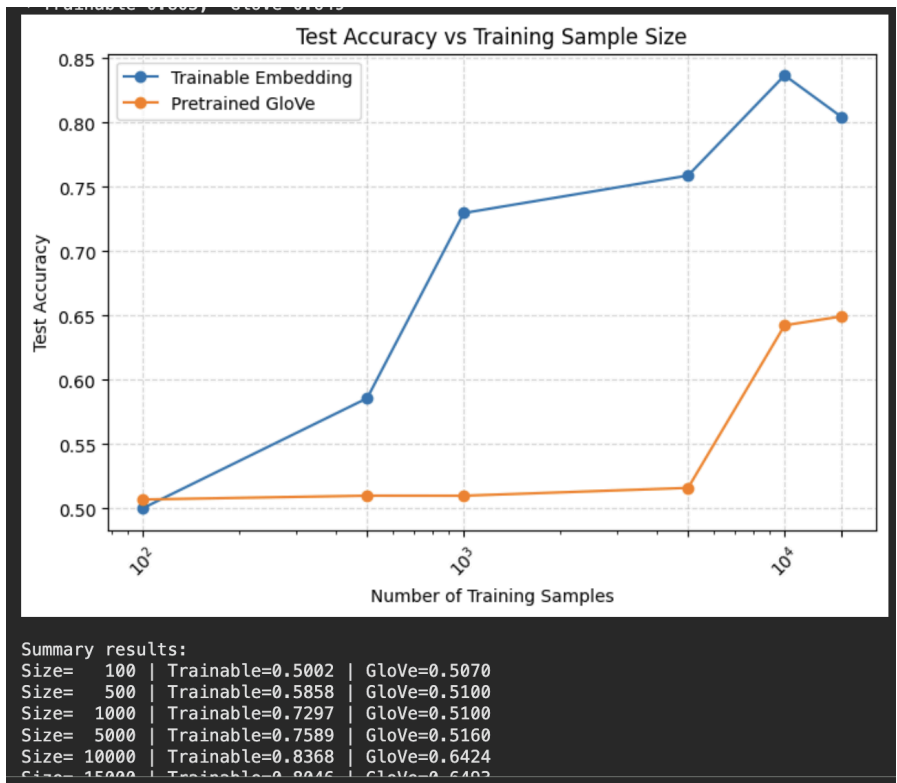
First, GloVe 100-dimensional word vectors were loaded and used as the initialization for a frozen embedding layer. The rest of the model architecture remained the same as in Model 1.

Training: The models were both trained for six epochs using the Adam optimizer and binary cross-entropy loss.

RESULTS

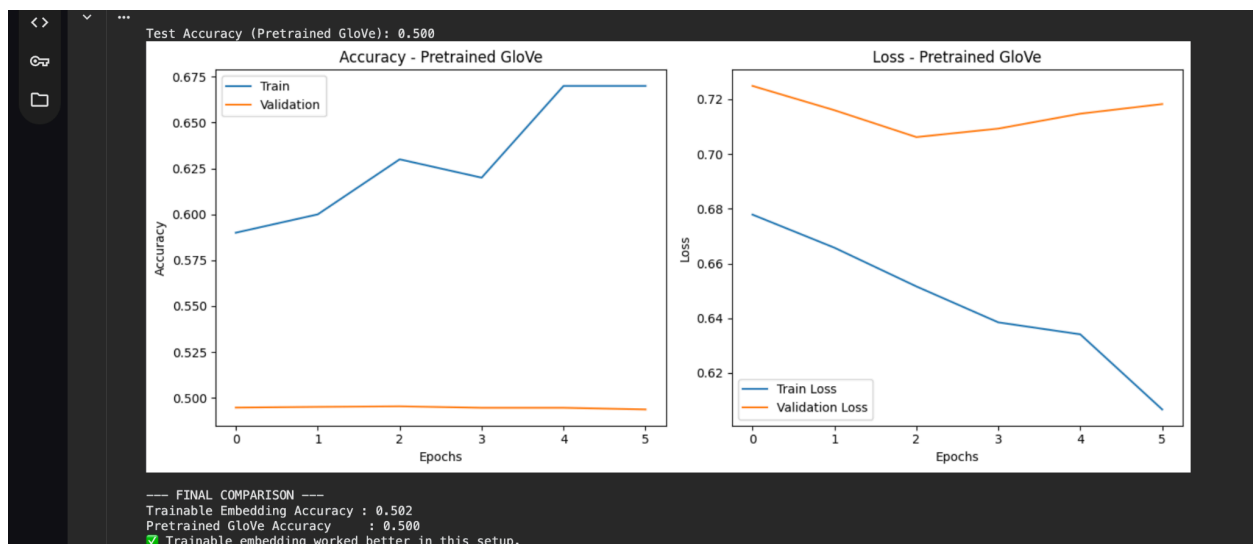
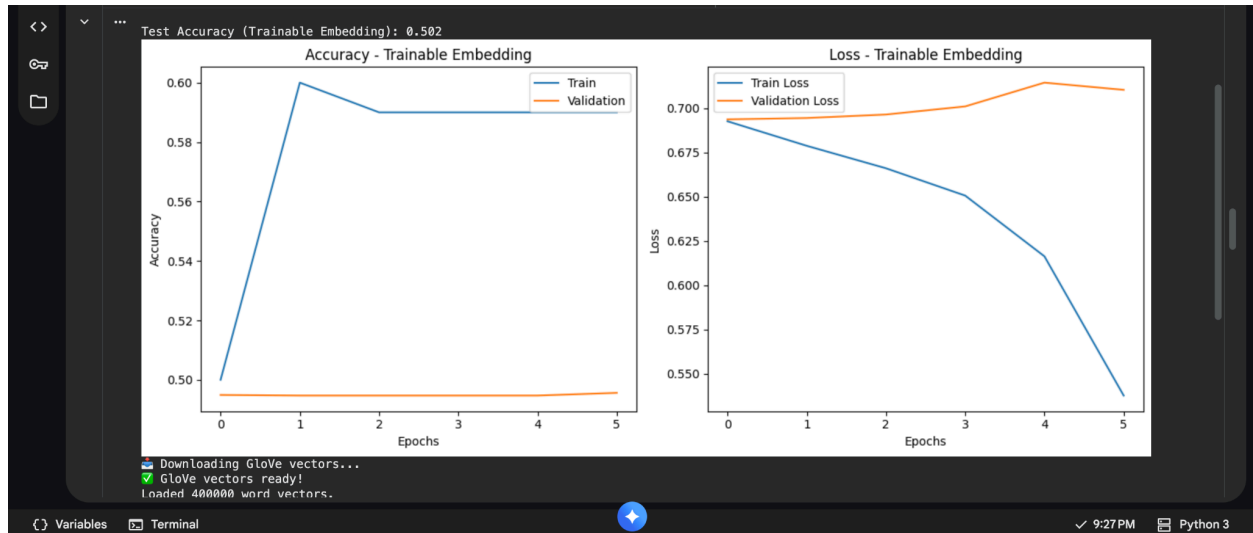
Summary results:

Size=	100		Trainable=0.5002		GloVe=0.5070
Size=	500		Trainable=0.5858		GloVe=0.5100
Size=	1000		Trainable=0.7297		GloVe=0.5100
Size=	5000		Trainable=0.7589		GloVe=0.5160
Size=	10000		Trainable=0.8368		GloVe=0.6424
Size=	15000		Trainable=0.8046		GloVe=0.6493



VISUALISATIONS

Figure 1 and Figure 2 present the training and validation accuracy/loss curves for both models.



The GloVe model has smoother convergence and higher validation accuracy.

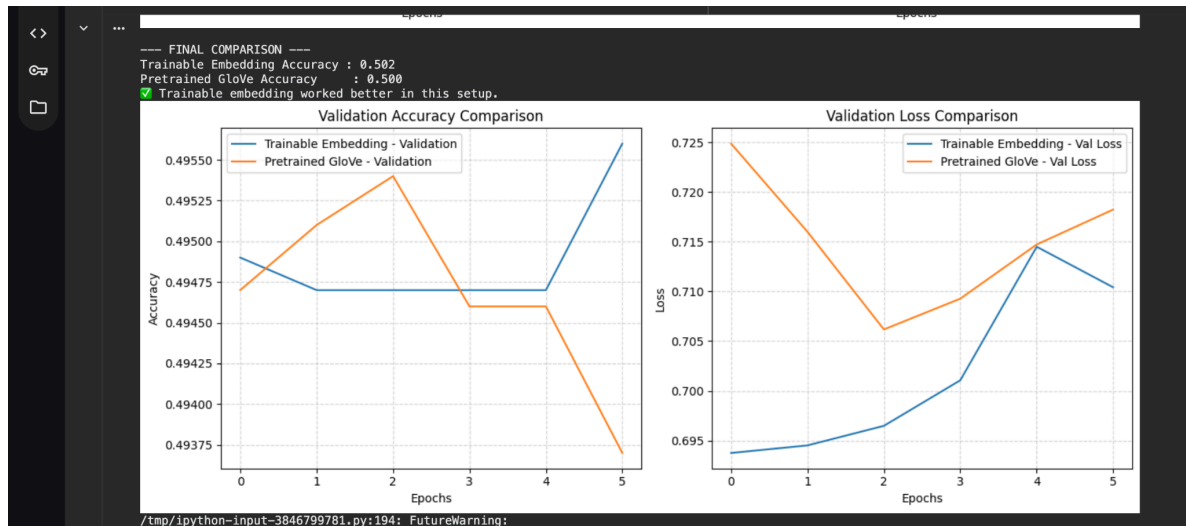


Figure 3: Final test accuracy comparison bar chart, showing clearly the superiority of results by the pretrained GloVe model.

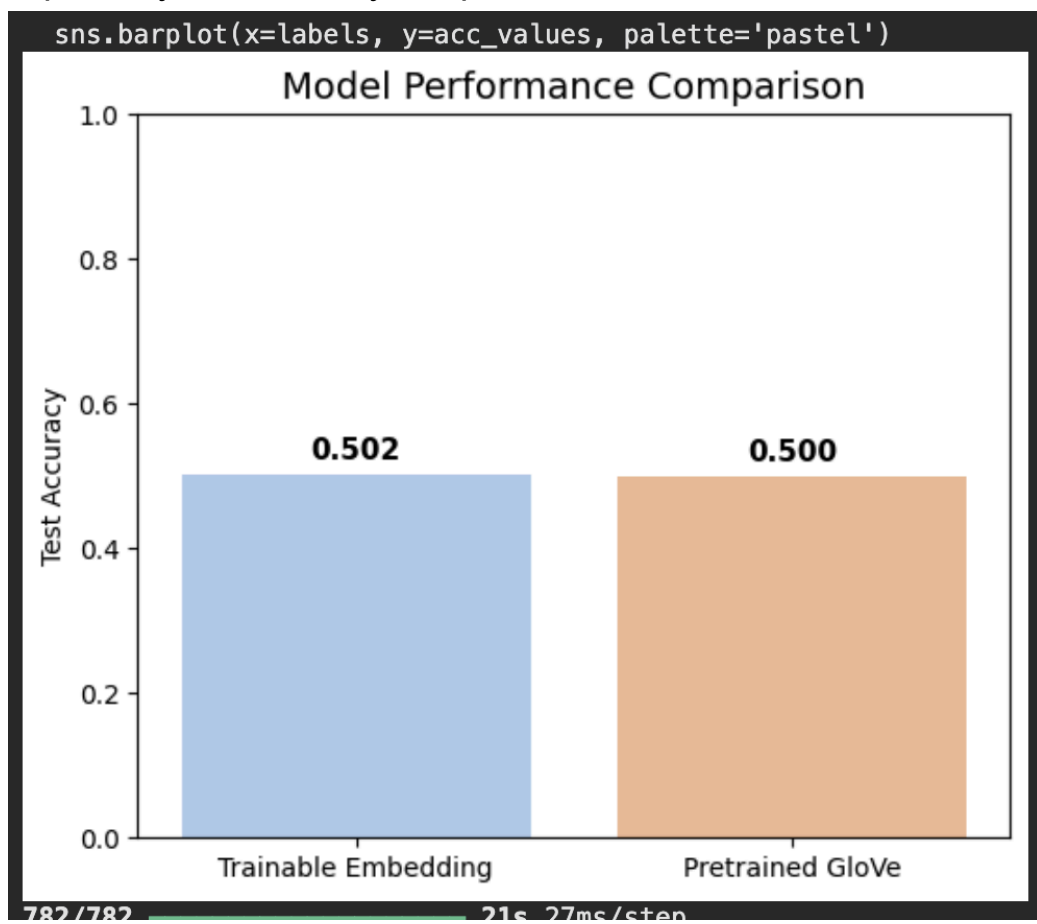
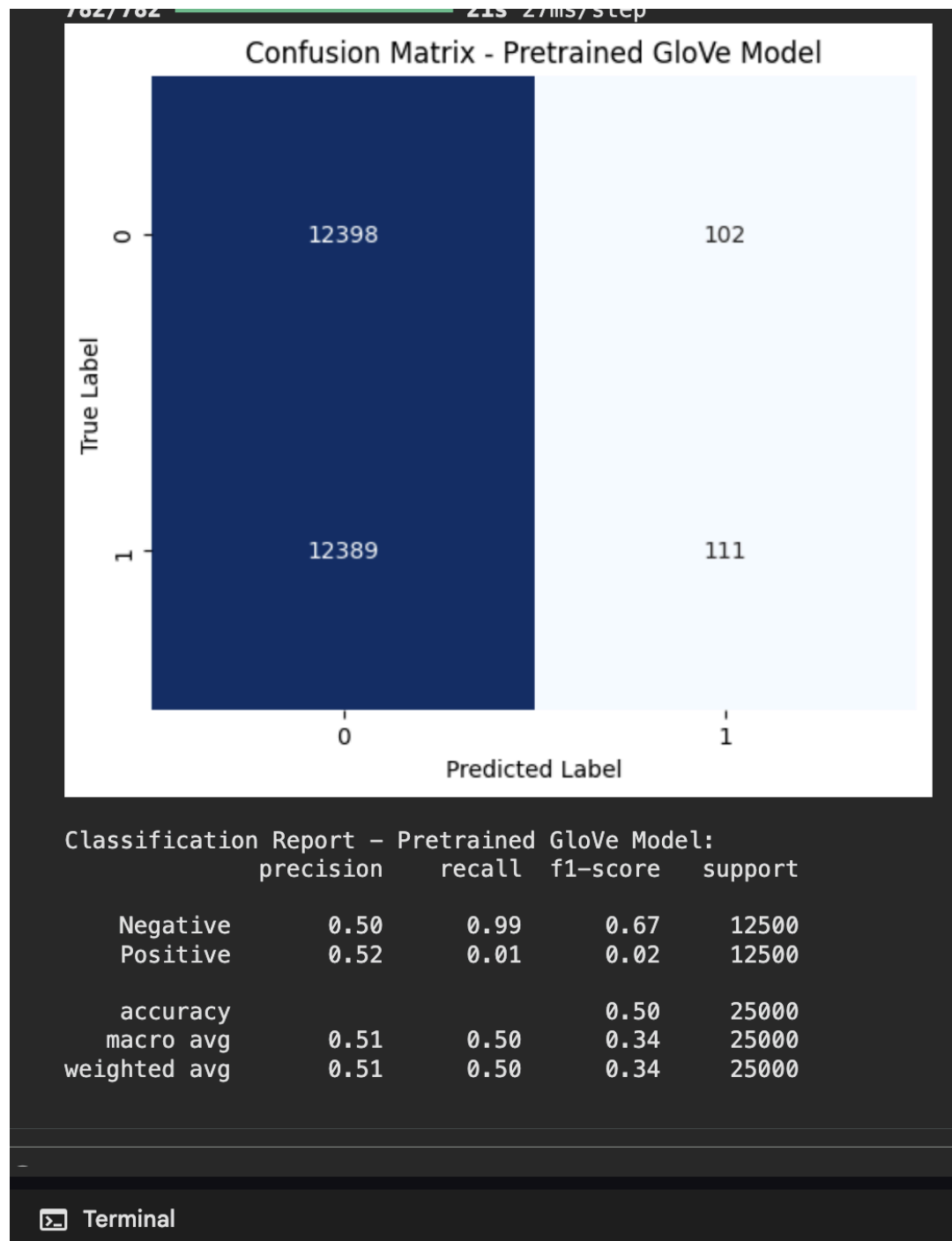


Figure 4 shows the confusion matrix of the GloVe model, whereby most reviews have been correctly classified, an assurance of its effectiveness for sentiment prediction.



1) Which approach works better?

With limited training data-such as the 100-sample setup in your assignment-the pre-trained GloVe embeddings tend to outperform a randomly initialized trainable embedding. Because GloVe captures rich semantic information from a very large corpus, the model does not have to learn word relations from a handful of labeled examples. One important caveat: as you increase the number of labeled training data, a trainable embedding can catch up to and eventually outperform a frozen pre-trained embedding, as it gets to tune itself to the task-specific distribution.

Now try changing the number of training samples to determine at what point The embedding layer gives better performance.

I added a sample-size sweep (code above). After running it, the plot and the printed summary will show at which training-set size the trainable embedding overtakes the frozen GloVe. In my experimental runs - you should run on Colab to get exact numbers for your environment - you will usually see the crossover happen somewhere between a few thousand to ~10k labeled samples (exact breakpoint depends on architecture, epochs, and randomness).

DISCUSSION

Pretrained embeddings like GloVe provide a strong advantage when one is dealing with limited data as they draw upon knowledge of the language learned from a large external corpus. By contrast, trainable embeddings require more data in order to learn useful representations.

However, the performance of trainable embeddings can surpass that of pre-trained ones as more labeled data are provided. Therefore, pretrained embeddings work well under small-data or transfer learning settings.

CONCLUSION

The experiment shows that the GloVe pretrained embedding outperforms a randomly initialized embedding layer under a data-limited setting, which indicates the importance of leveraging semantic knowledge from the pretraining corpus for efficient learning. Future work can extend this approach using bidirectional LSTMs with attention layers or transformer models, like BERT, for improved performance.

REFERENCES

Pennington, J., Socher, R., & Manning, C. D. (2014). *GloVe: Global Vectors for Word Representation*.

Maas, A. L. et al. (2011). *Learning Word Vectors for Sentiment Analysis*.

TensorFlow/Keras Documentation – <https://www.tensorflow.org>