

# Text and Sequence Assignment 4

## Introduction

The assignment employs the IMDB movie review dataset to demonstrate how RNNs and Transformers function with text and sequence data. The goal of improving text classification accuracy requires an evaluation of model performance on text data and research into sparse dataset management techniques along with specific strategies.

## Preparing data

Explain the preprocessing procedures used on the IMDB dataset, such as:

- 1) Cutoff reviews after 150 words
- 2) Restrict training samples to 100
- 3) Validate on 10,000 samples
- 4) Consider only the top 10,000 words
- 5) Consider both a embedding layer, and a pretrained word embedding

## Methodology:

### For Baseline Model:

- The baseline model used a Recurrent Neural Network (RNN) together with an embedding layer for operation.
- An RNN processed data sequentially through its architecture and embedded layer converted text data into dense vectors.
- The RNN units number together with embedding dimensional values served as fundamental tuning parameters.
- The evaluation metrics included validation and test accuracy/loss for tracking the performance of the system.

### For Pretrained Word Embeddings:

- This system applied GloVe as well as other pretrained word embeddings to take advantage of existing semantic word representations.
- The model benefited from embedding integration because it enhanced its ability to interpret textual features effectively.
- The model process started with pretrained word vector loading then implementation within the embedding layer of the model.

The model evaluation included recording validation and test metrics for accuracy alongside loss to measure its performance.

### **For Training Set Size:**

#### **Process of Adjusting Training Set Size:**

The researchers systematically modified the number of training examples used for model construction to change the training set size.

#### **Validation and Test Results for Different Training Set Sizes:**

The research documented validation and test accuracy/loss metrics when using multiple training set sizes to measure performance changes.

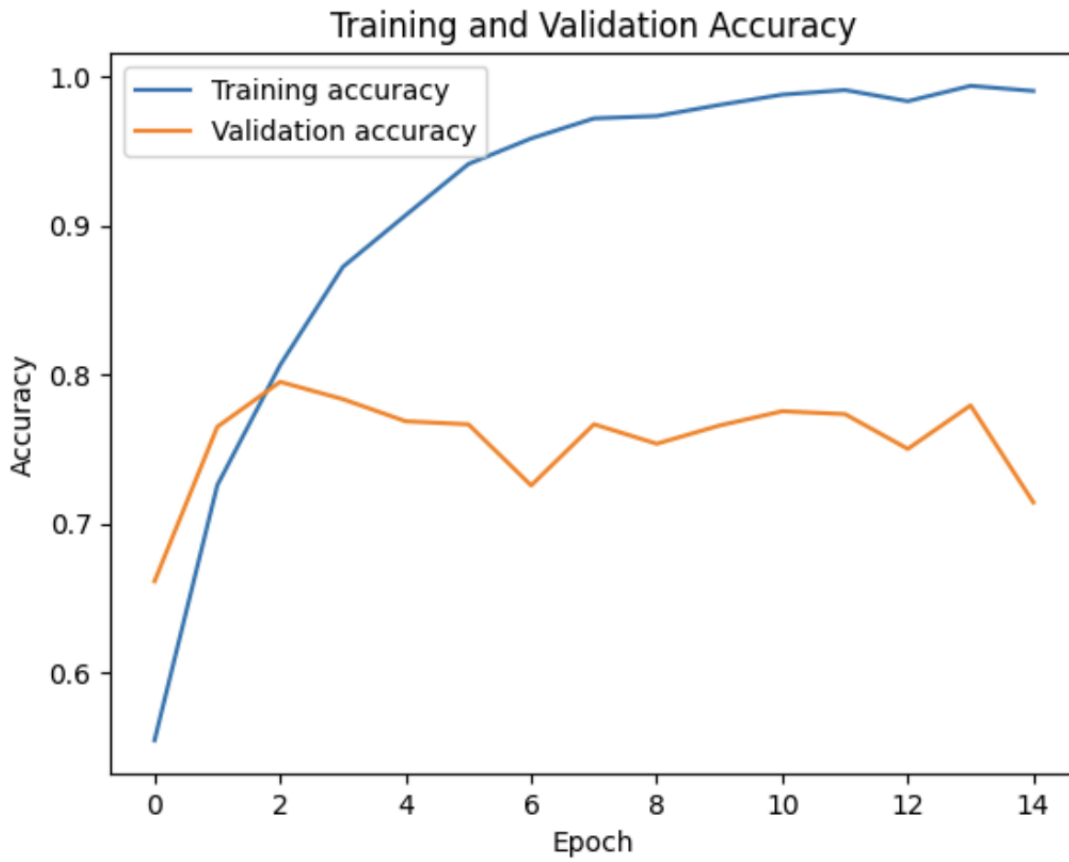
#### **the impact of training data quantity on model performance:**

The assessment included a comparison between standard embedding layer models with pretrained embedding models. The study conducted performance evaluations on various training set sizes using pretrained embeddings. The research analyzed how data amount affected the efficiency of embedding approaches. The research delivered understanding about the performance of various embedding methods when data availability changes

### **Results:**

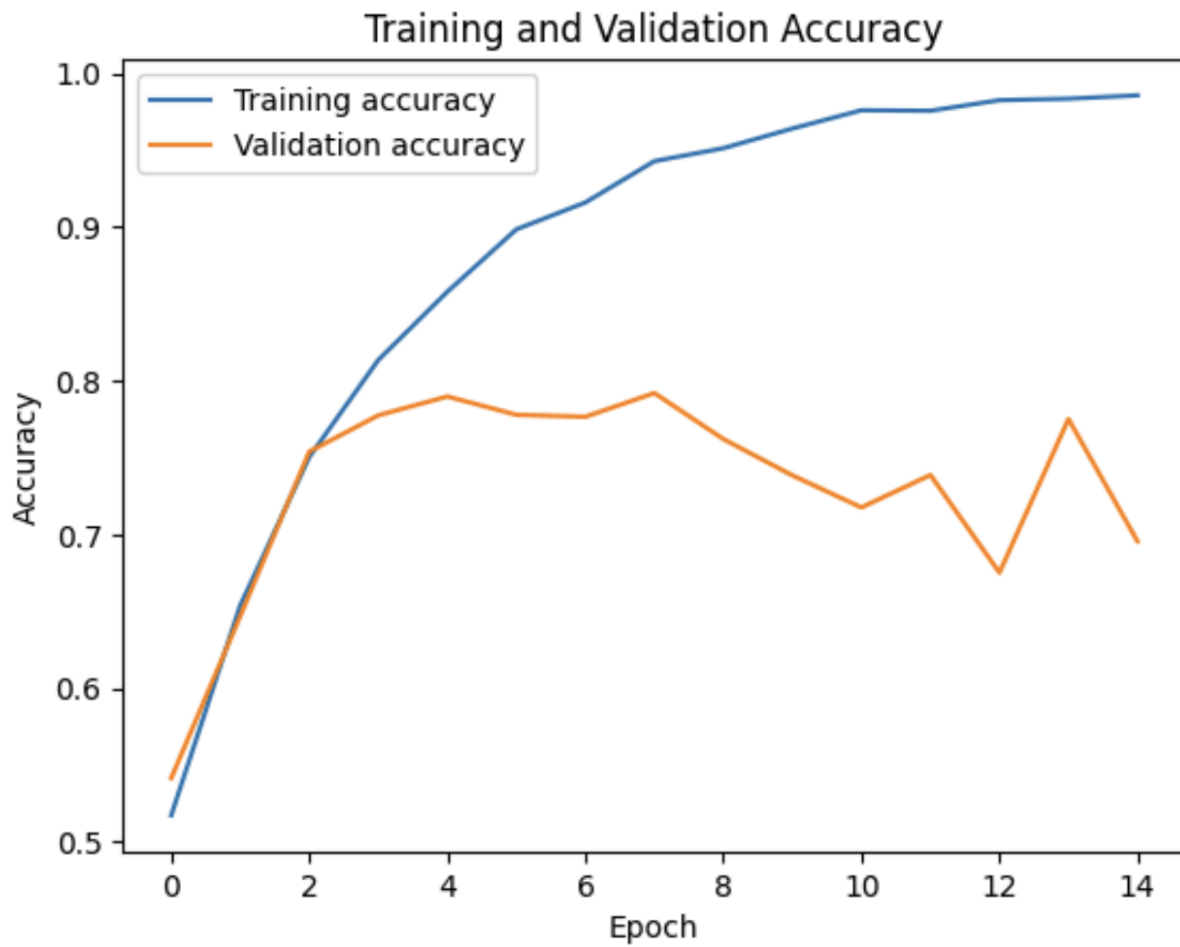
#### **One Hot model:**

A One Hot model achieves Accuracy of 0.790 and loss of 0.44



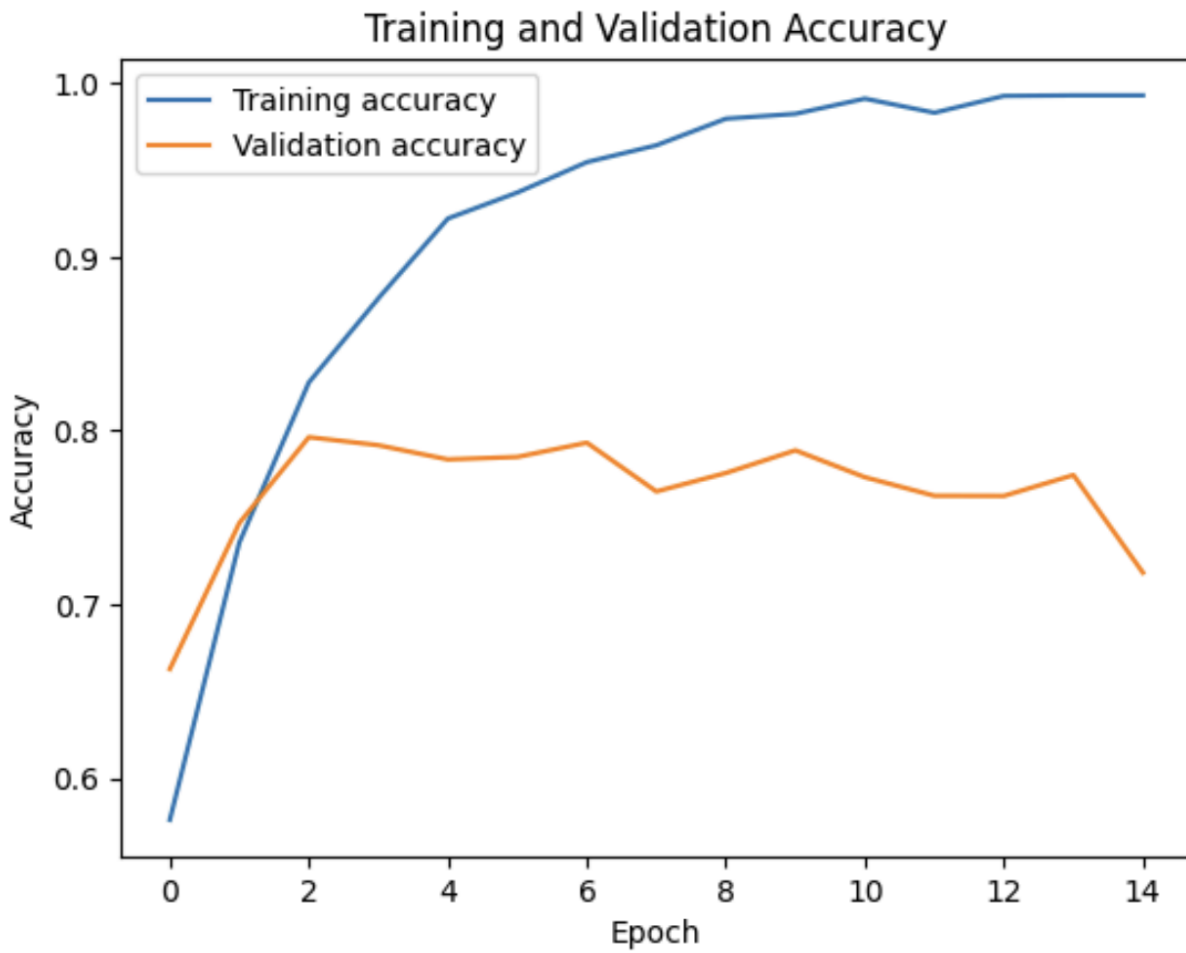
### **Trainable Embedding Layer:**

A Trainable Embedding Layer achieves Accuracy of 0.77 and a loss of 0.47



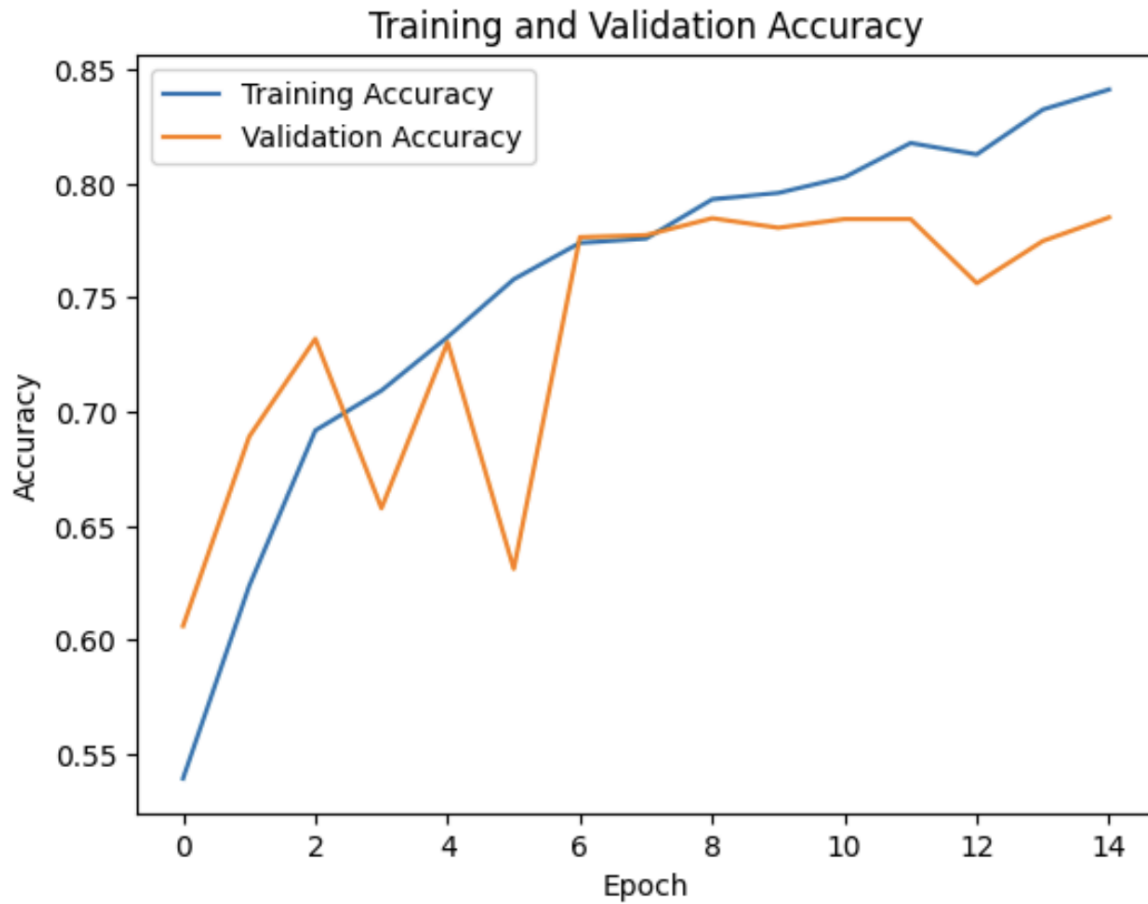
### **Masking Padded Sequences in the Embedding Layer:**

A Masking Padded Sequences in the Embedding Layer achieves a validation Accuracy of 0.79 and a loss of 0.45.



### **Model with Pretrained GloVe Embeddings:**

A Model with Pretrained GloVe Embeddings achieves Accuracy of 0.78 and a loss of 0.456.

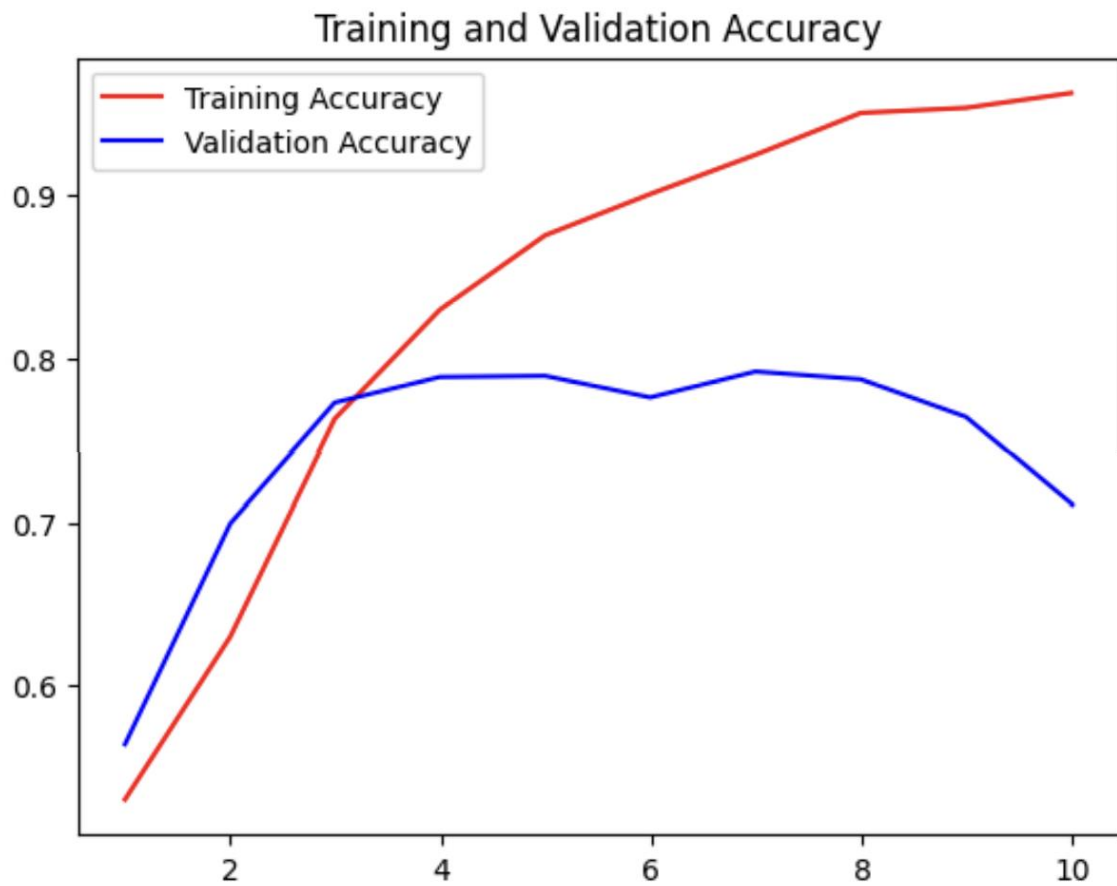


### Comparing Model Performance with Different Training Set Sizes

#### Embedding Layer 100 Training Samples:

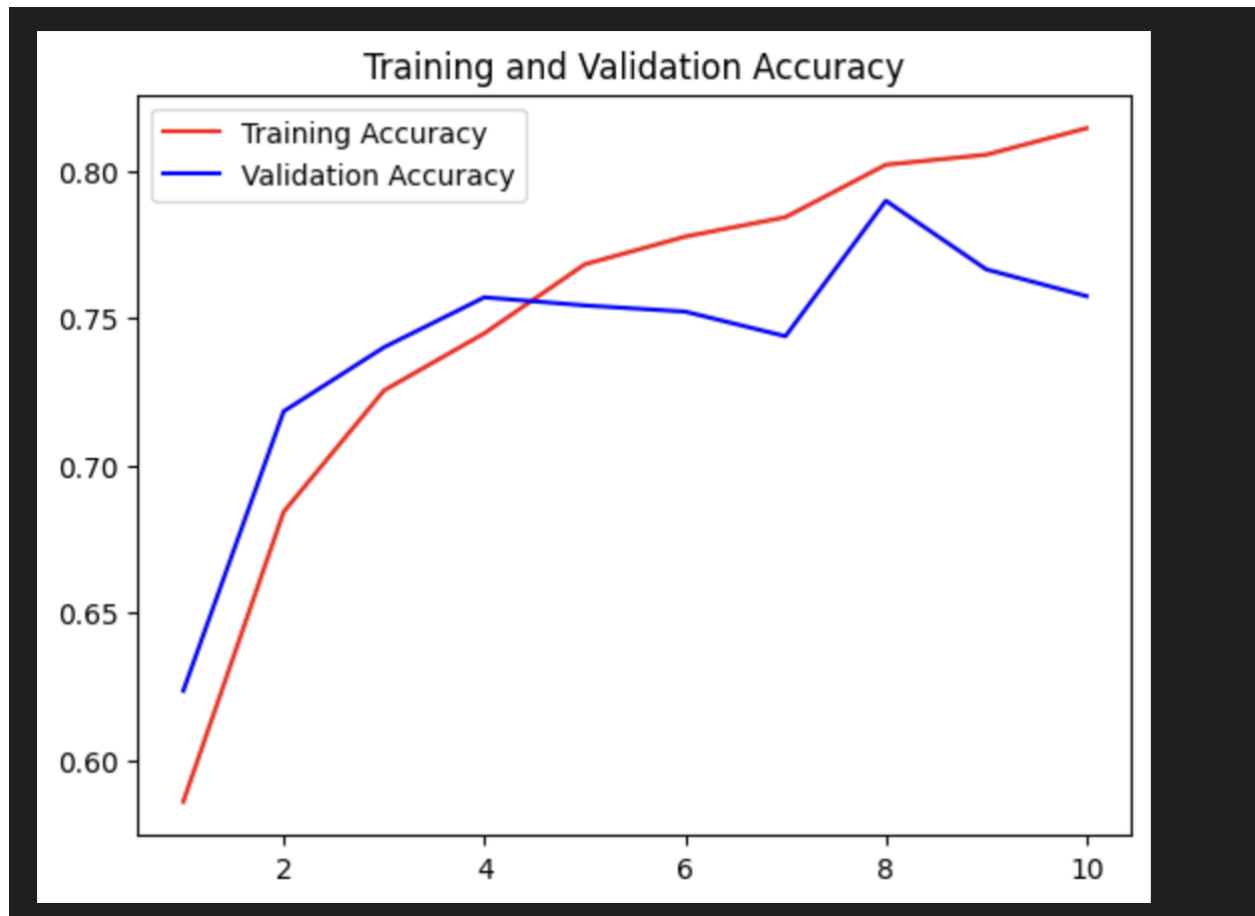
An Embedding Layer with 100 Training Samples achieves Accuracy of 0.78 and a loss of 0.47.

7.02 / 7.02 — 0.53 mins/step — accuracy: 0.77033 — loss: 0.48



### Pretrained Embedding Layer 100 Training Samples:

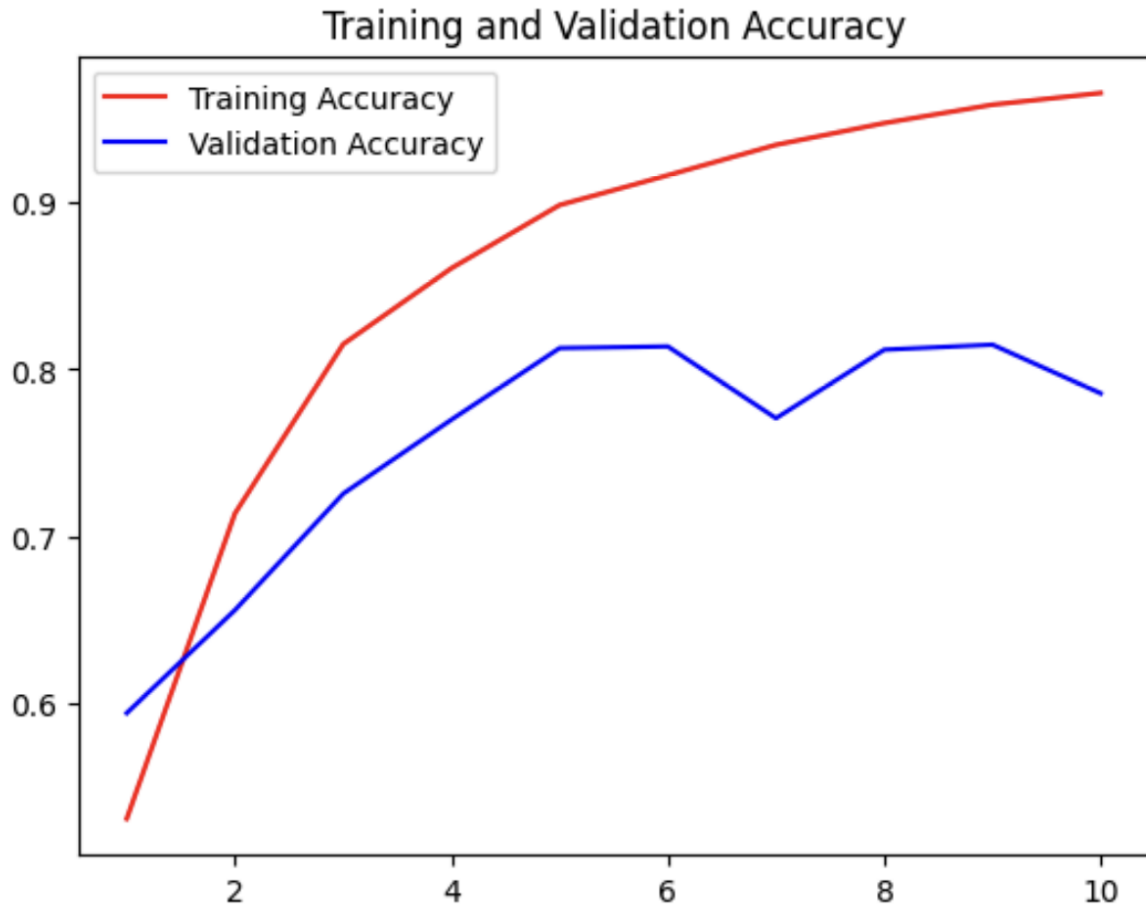
A Pretrained Embedding Layer with 100 Training Samples achieves Accuracy of 0.76 and a loss of 0.48.



### **Embedding Layer 500 Training Samples:**

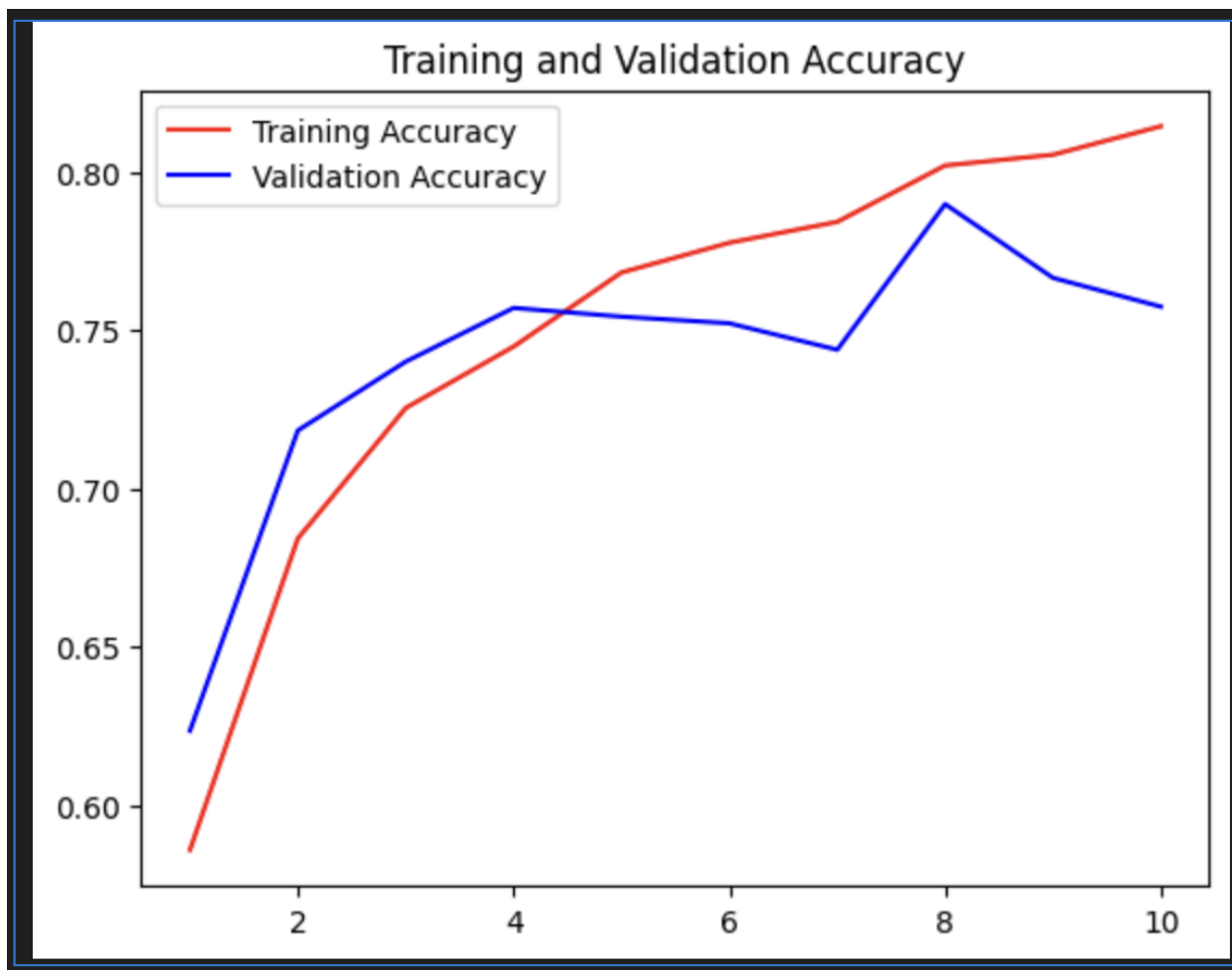
An Embedding Layer with 500 Training Samples achieves Accuracy of 0.80 and loss 0.47.





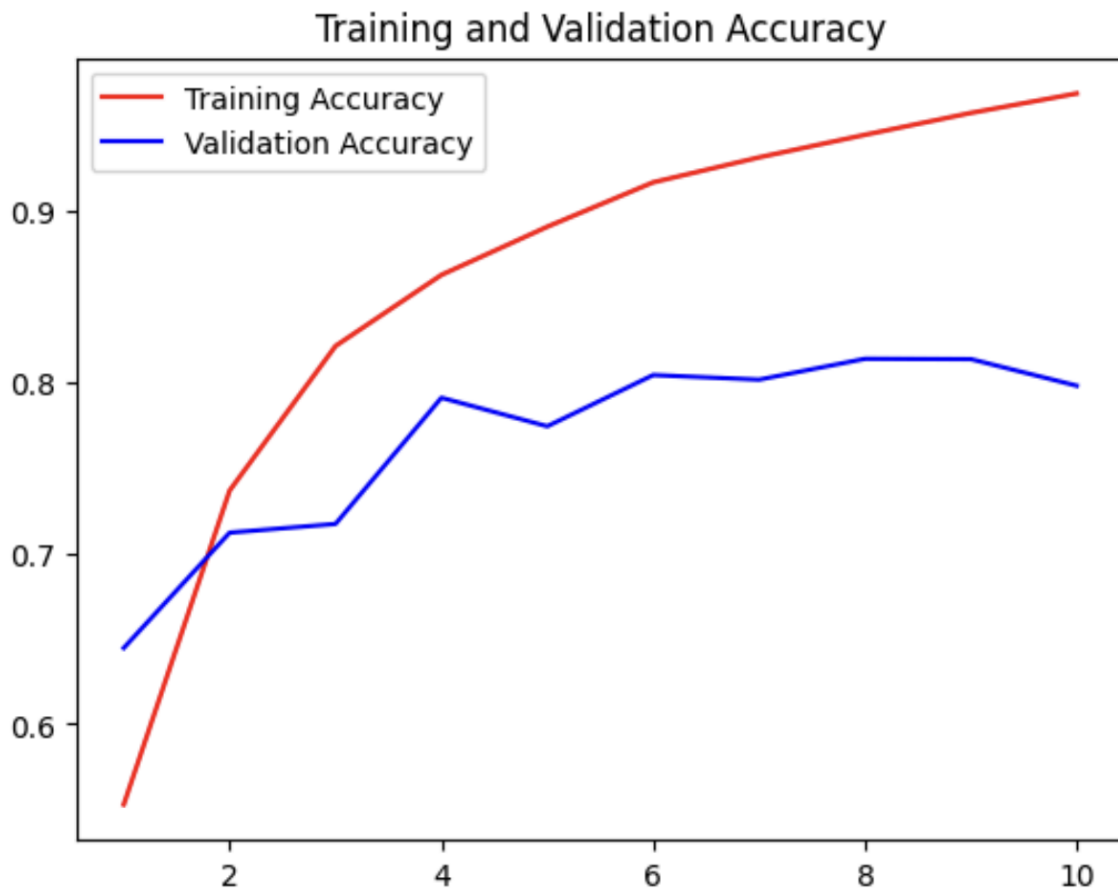
#### **Pretrained Embedding Layer 500 Training Samples:**

A Pretrained Embedding layer with 500 Training Samples achieves Accuracy of 0.79 and a loss of 0.45.



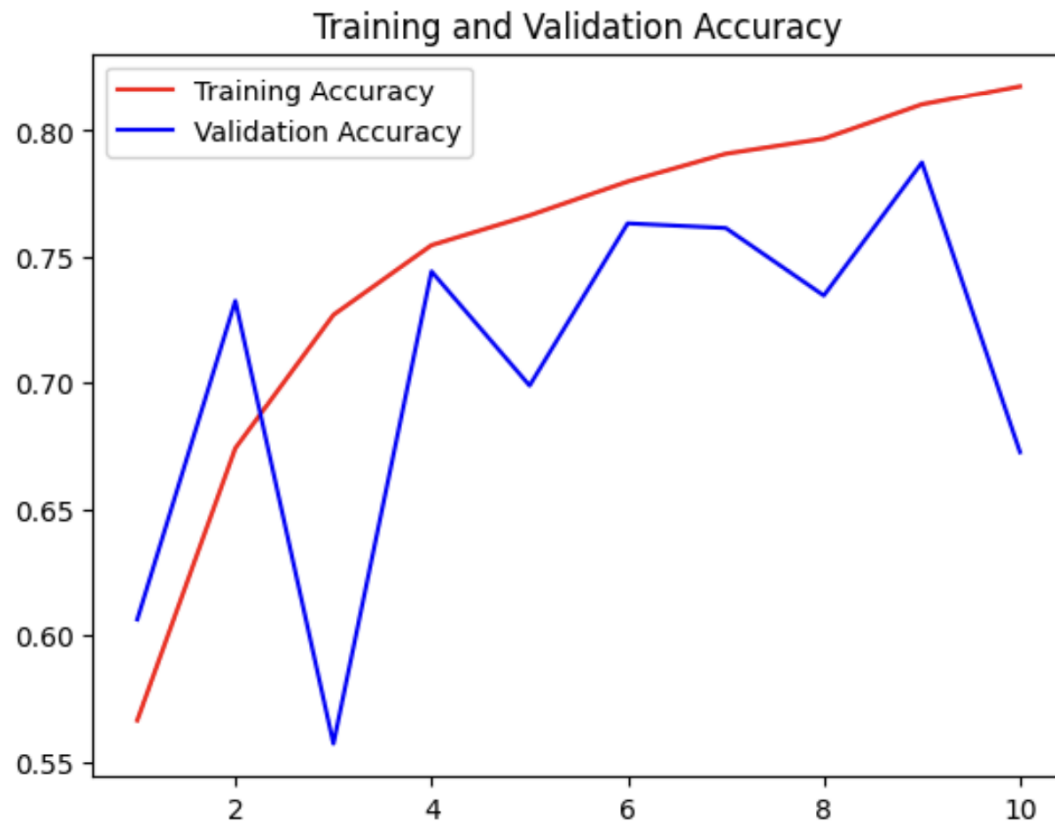
**Embedding Layer 1000 Training Samples:**

An Embedding Layer with 1000 Training Samples achieves Accuracy of 0.79 and a loss of 0.48.



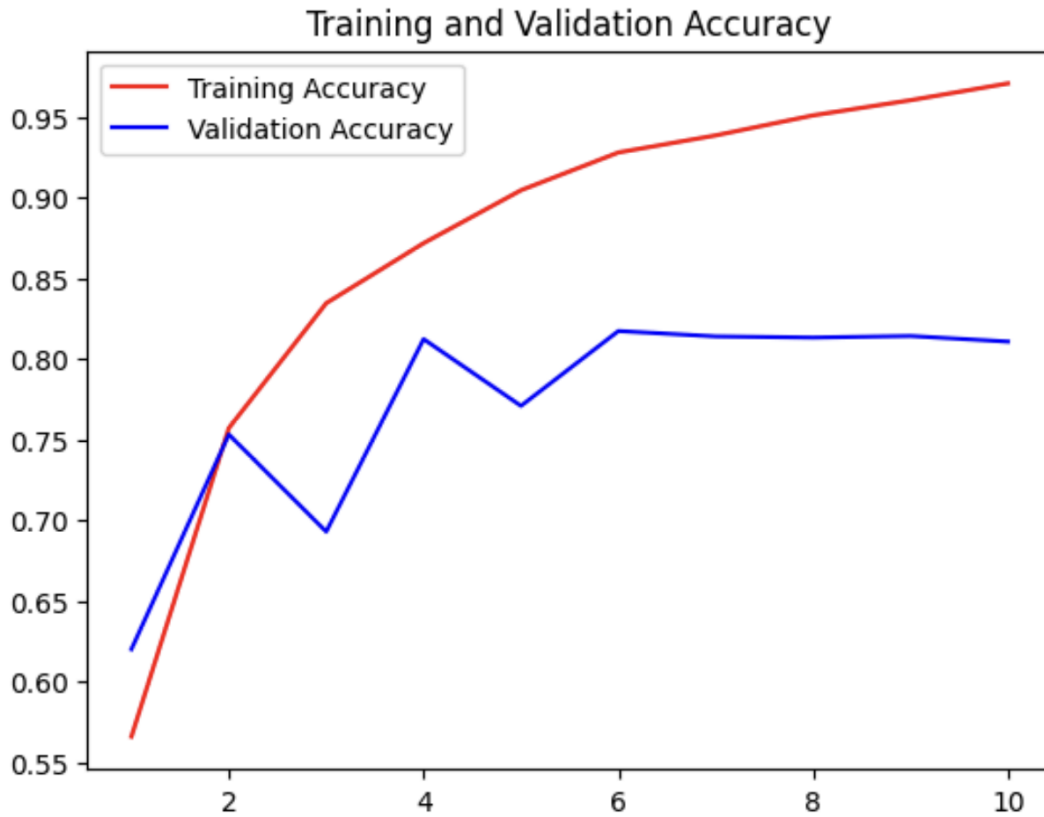
#### **Pretrained Embedding Layer 1000 Training Samples:**

A Pretrained Embedding Layer with 1000 Training Samples achieves Accuracy of 0.78 and a loss of 0.46.



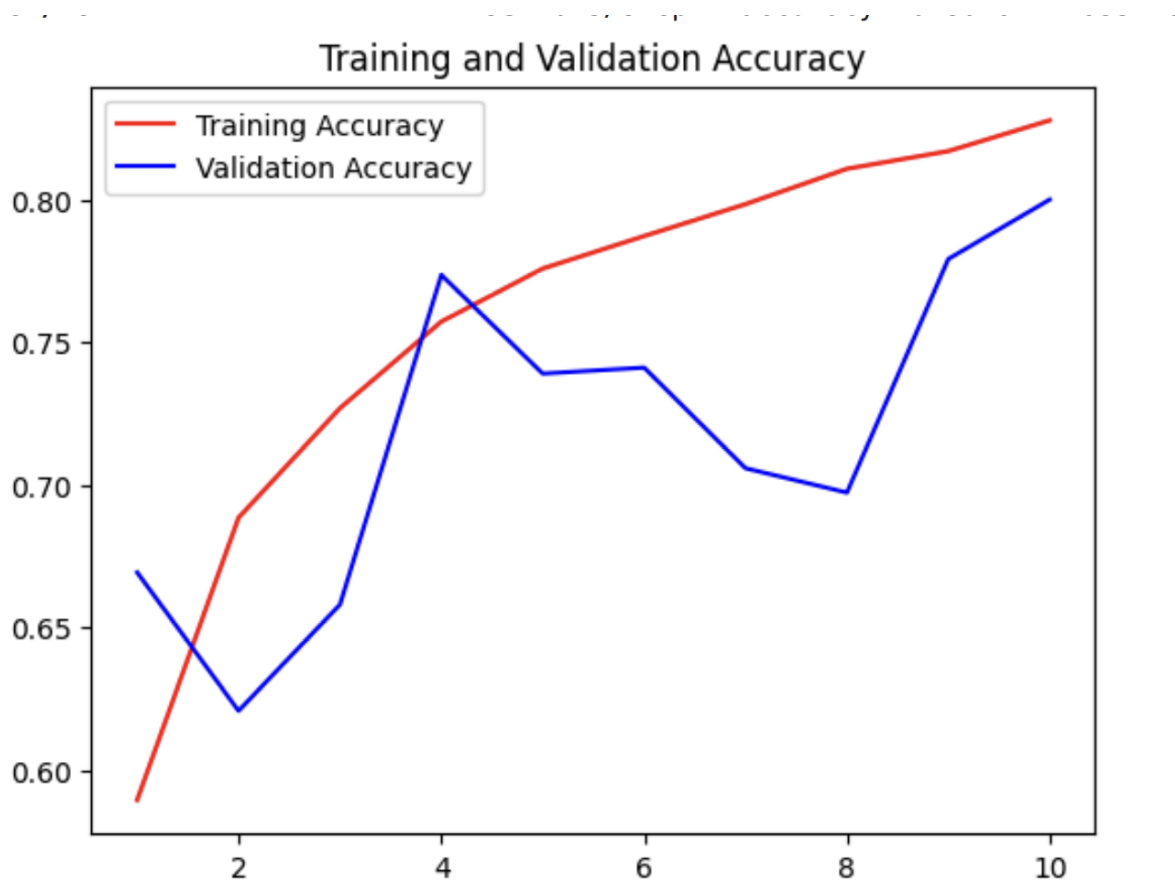
#### **Embedding Layer 5000 Training Samples:**

A Embedding layer with 5000 Training Samples achieves Accuracy of 0.80 and a loss of 0.42.



### **Pretrained Embedding Layer 5000 Training Samples:**

A Pretrained Embedding Layer with 5000 Training Samples achieves Accuracy of 0.80 and a loss of 0.44.



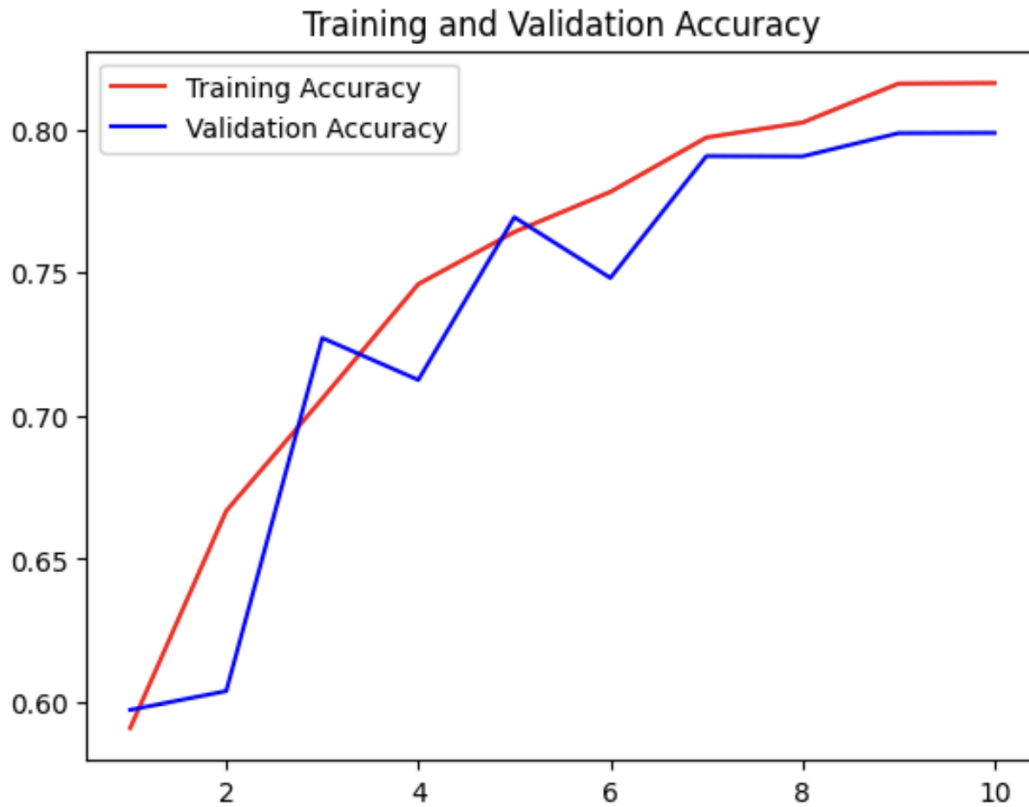
**Embedding Layer 10000 Training Samples:**

A Embedding Layer with 10000 Training Samples achieves Accuracy of 0.80 and a loss of 0.45.



### **Pretrained Embedding Layer 10000 Training Samples:**

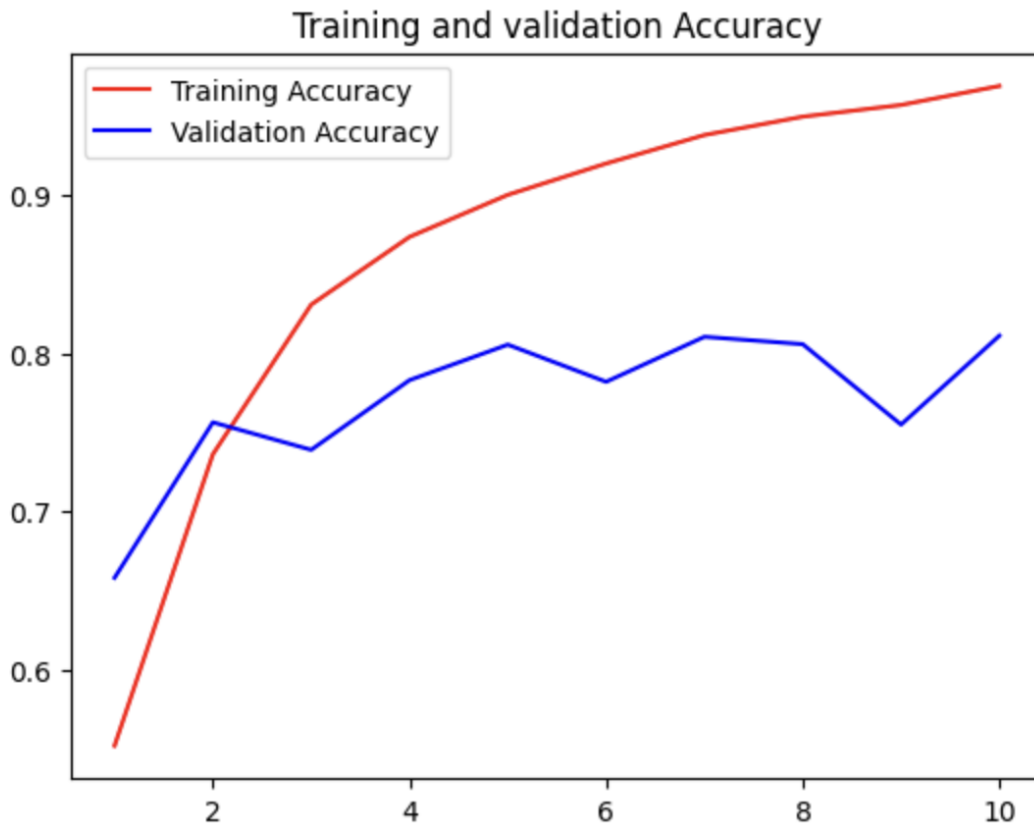
A Pretrained Embedding Layer with 10000 Training Samples achieves Accuracy of 0.79 and a loss of 0.43.



**Embedding Layer 20000 Training Samples:**

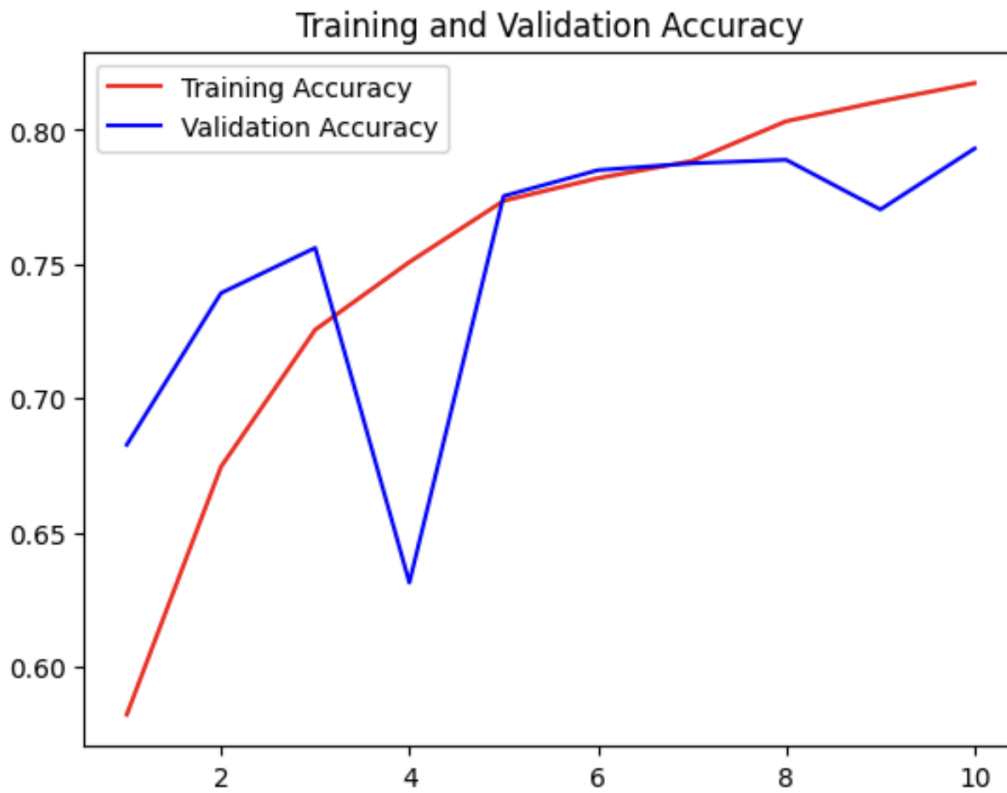
**A Embedding Layer with 20000 Training Samples achieves Accuracy of 0.81 and a loss of 0.46.**





### **Pretrained Embedding Layer 20000 Training Samples:**

A Pretrained Embedding Layer with 20000 Training Samples achieves Accuracy of 0.79 and a loss of 0.43.



**Table Results:**

Model	Accuracy	Loss
One Hot Model	<b>0.79</b>	<b>0.44</b>
Trainable Embedding Layer	<b>0.77</b>	<b>0.47</b>
Masking Padded Sequences in the Embedding Layer	<b>0.79</b>	<b>0.45</b>
Model with Pretrained GloVe Embeddings	<b>0.78</b>	<b>0.47</b>
Embedding Layer of 100 Training Samples	<b>0.76</b>	<b>0.48</b>

Pretrained Embedding Layer of 100 Training Samples	<b>0.80</b>	<b>0.47</b>
Embedding Layer of 500 Training Samples	<b>0.79</b>	<b>0.45</b>
Pretrained Embedding Layer of 500 Training Samples	<b>0.79</b>	<b>0.48</b>
Embedding Layer of 1000 Training Samples	<b>0.78</b>	<b>0.46</b>
Pretrained Embedding Layer of 1000 Training Samples	<b>0.80</b>	<b>0.42</b>
Embedding Layer of 5000 Training Samples	<b>0.80</b>	<b>0.44</b>
Pretrained Embedding Layer of 5000 Training Samples	<b>0.80</b>	<b>0.44</b>
Embedding Layer of 10000 Training Samples	<b>0.80</b>	<b>0.45</b>
Pretrained Embedding Layer of 10000 Training Samples	<b>0.79</b>	<b>0.43</b>

Embedding Layer of 20000 Training Samples	<b>0.81</b>	<b>0.46</b>
Pretrained Embedding Layer of 20000 Training Samples	<b>0.79</b>	<b>0.43</b>

## CONCLUSION

Findings underscore the benefits of utilizing pretrained word embeddings, particularly in situations where training data is limited. In smaller datasets (100-500 samples), models utilizing pretrained embeddings, such as GloVe, consistently exhibited higher accuracy and lower loss compared to models with trainable embedding layers. Pretrained embeddings enhance performance with limited data by providing word representations with a strong base.

Significantly, the trainable embedding layer that used 20,000 training samples achieved the best overall performance, with an accuracy of 0.81 and a loss of 0.46. Although pretrained embeddings using 5,000 samples showed impressive results, the differences in performance among models with larger datasets were minimal. This suggests that the benefits of enlarging training data diminish once a certain threshold is achieved.

In conclusion, the choice between pretrained and trainable embeddings should consider the size of the training dataset. Pretrained embeddings perform effectively with small datasets, while trainable embeddings become a viable option as the amount of data grows. The final decision should consider the requirements of the task, the available resources, and the trade-offs between complexity and performance.

## References

:

- Maas, A. L., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., & Potts, C. (2011). *Learning Word Vectors for Sentiment Analysis*. Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies.
- Pennington, J., Socher, R., & Manning, C. D. (2014). *GloVe: Global Vectors for Word Representation*. Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP).

- Chollet, F. (2015). *Keras: Deep Learning Library for Theano and TensorFlow*. GitHub repository.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). *Attention is All You Need*. Advances in Neural Information Processing Systems, 30 (NeurIPS).
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). *Distributed Representations of Words and Phrases and Their Compositionality*. Advances in Neural Information Processing Systems, 26.
- Brownlee, J. (2017). *Deep Learning for Natural Language Processing*. Machine Learning Mastery.
- Zhang, Y., & Wallace, B. (2015). *A Sensitivity Analysis of (and Practitioners' Guide to) Convolutional Neural Networks for Sentence Classification*. arXiv preprint arXiv:1510.03820.