

## Natural Language Processing - IMDB Movie Review

	Description	Hyperparameters	Number of Epochs	Training Loss	Training Accuracy	Test Accuracy	Comments
Part 1a	Given model - Word Embedding Layer + Mean Pooling + Fully Connected Layer + Relu + Output Layer	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=500	20	0.025532	0.991	0.9216	The model trained with given hyperparameters returns a rather good result. I trained it with 20 epochs so the accuracy is better than expected. should be around 86-88%
	Custom 1	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=1000	20	0.021467	0.993	0.9237	Describe more about the model/results such as why certain hyperparamters were chosen or the effect it had on the accuracy/training time/overfitting/etc.
	Custom 2	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=15000, HiddenUnits=500	20	0.010972	0.997	0.9244	From the results, we've seen that more hidden units will bring a better performance. Theoretically, more hidden units can also help us capture more types of information, hence also increase the accuracy.
	Custom 3	SGD optimizer with LR=0.1, BatchSize=200, VocabularySize=8000, HiddenUnits=500	20	0.204234	0.921	0.8662	From the results, we've seen that ADAM optimizer works better than SGD if learning rates are properly adjusted.
Part 1b	Given Model - Word Embedding Layer + Mean Pooling + Fully Connected Layer + Relu + Output Layer	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=100000, HiddenUnits=500	100	0.083506	0.97	0.9099	The model trained with given hyperparameters returns a rather good result. I trained it with 100 epochs so the accuracy is better than expected. Should be around ~81-87%
	Custom 1	ADAM optimizer with LR=0.01, BatchSize=200, VocabularySize=100000, HiddenUnits=500	100	0.093555	0.967	0.9084	From the results, we've seen that learning rate 0.001 works a little better than 0.01. Theoretically, the learning rate should be adjusted properly so that it won't be either too slow to converge, or too fast to fail convergence.
	Custom 2	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=100000, HiddenUnits=1000	100	0.056714	0.982	0.9131	From the results, we've seen that more hidden units will bring a better performance. Theoretically, more hidden units can also help us capture more types of information, hence also increase the accuracy.
	Custom 3	SGD optimizer with LR=0.01, BatchSize=200, VocabularySize=100000, HiddenUnits=500	100	0.157885	0.939	0.8877	From the results, we've seen that ADAM optimizer works better than SGD if learning rates are properly adjusted, in the case with GloVe.
Part 2a	Given Model - (write description)	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=500	20	0.082083	0.971	0.8705	The model trained with given hyperparameters returns a rather good result as expected. ~87%
	Custom 1	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=1000	20	0.057841	0.98	0.8627	From the results, we've seen that more hidden units will bring a better performance on training set, but a little worse on testing set. This implies that too many hidden units may cause our model to be a bit overfitting.
	Custom 2	ADAM optimizer with LR=0.01, BatchSize=200, VocabularySize=8000, HiddenUnits=500	20	0.273012	0.888	0.8107	From the results, we've seen that learning rate 0.001 works much better than 0.01 compared to part 1a and 1b. This maybe because for more complicated models, different learning rates are more effective.
	Custom 3	SGD optimizer with LR=0.1, BatchSize=200, VocabularySize=8000, HiddenUnits=500	20	0.267985	0.888	0.783	From the results, we've seen that ADAM optimizer works better than SGD if learning rates are properly adjusted.
Part 2b	Given Model - Word Embedding Layer + Mean Pooling + Fully Connected Layer + Relu + Output Layer	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=100000, HiddenUnits=500, SequenceLength=100	30	0.152165	0.94	0.8751	The model trained with given hyperparameters returns a rather good result, but not as good as expected. The reason may be that more epochs are needed, or the sequence length should increase. ~91%
	Custom 1	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=100000, HiddenUnits=1000, SequenceLength=100	30	0.118245	0.957	0.874	From the results, we've seen that more hidden units will bring a better performance on training set, but a little worse on testing set. This implies that too many hidden units may cause our model to be a bit overfitting.
	Custom 2	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=100000, HiddenUnits=500, SequenceLength=200	30	0.06071	0.978	0.8622	From the results, we've seen that a larger sequence length leads to a better performance on training set, but worse on testing set. This implies that a larger sequence length causes our model to be overfitting, because not all parts in one sentence can reflect its true sentiment, but a very large sequence length will mislead the model to "think" in this way and hence overfits on the training set.
	Custom 3	SGD optimizer with LR=0.01, BatchSize=200, VocabularySize=100000, HiddenUnits=500, SequenceLength=100	30	0.43502	0.798	0.6891	In this case, ADAM optimizer works much better than SGD if learning rates are properly adjusted. This shows the advantage of ADAM over SGD in NLP tasks.

Part 3a	Given Model	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=500, SequenceLength=50, GradientClip=2.0	75	3.916363	0.2571	0.2709	The training and test accuracy keep oscillating around 0.26 and 0.27 in the last 30 epochs, which implies we can use less epochs during training. I think the parameter set is not optimal and can be properly adjusted. For example, we can take a higher sequence length so that more information can be used.	
Part 3b	Generated Review	Temperature=1.0	of family shows and end up dying with nothing . it seemed to be to be funny film that makes it seem like a tv show with a great cast . the actors should not this film . it was a great film , although it was n't based on a true story . i do n't acting was terrible . the plot was weak and the acting was awful . the plot was					
	Generated Review	Temperature=0.75						
	Generated Review	Temperature 0.5						
	Generated Review	Temperature 0.25						
Part 3c The 10 test accuracies are given corresponding to sequence length from 50 to 500.	Given Model	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=500, SequenceLength=100	20	0.064193	0.978	0.7748, 0.8387, 0.867, 0.8832, 0.8928, 0.8973, 0.8996, 0.9021, 0.9036, 0.9037	The model trained with given hyperparameters returns a rather good result as expected.	~91%+
	Custom 1	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=1500, SequenceLength=100	20	0.064773	0.9777	0.7738, 0.8369, 0.8690, 0.8855, 0.8957, 0.9002, 0.904, 0.905, 0.9064, 0.9075	From the results, we've seen that more hidden units make no difference on the training set performance, but slightly increase the test set accuracies. But this is not significant enough to claim that more hidden units are helpful.	
	Custom 2	ADAM optimizer with LR=0.001, BatchSize=200, VocabularySize=8000, HiddenUnits=500, SequenceLength=150	20	0.062714	0.9792	0.7682, 0.8318, 0.8629, 0.8829, 0.8882, 0.8946, 0.8988, 0.9001, 0.9015, 0.902	From the results, we've seen that a larger sequence length leads to a slightly better performance on training set, but slightly worse on testing set. This implies that a larger sequence length may cause our model to be overfitting, but it's not significant enough to ensure.	
	Custom 3	SGD optimizer with LR=0.01, BatchSize=200, VocabularySize=8000, HiddenUnits=500, SequenceLength=100	20	0.316532	0.8633	0.9135, 0.914 0.9027, 0.9072, 0.9105, 0.9126,	This result shows that although optimizer with SGD has the worst training performance among all, the testing accuracies turn out to be the best. We can hence imply that this optimizer can produce good generalization in this case of sentiment analysis.	