
COMP47650 Deep Learning

Final Assignment

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This project tackles one of the most difficult challenges in Machine Learning,
2 Natural Language Processing. Using LSTMs and GRUs a very high level of
3 accuracy was achieved on emotion classification on a Twitter tweets which is
4 known to be a difficult text processing task due to the informal nature of the text.

5 1 Introduction

6 Automatic text classification has become a popular topic in machine learning. The aim of this
7 assignment is to build a Deep Learning model using two different forms of Recurrent Neural
8 Networks, LSTM(Long Short-Term Memory) and GRUs(Gated Recurrent Unit) to build a mode

9 2 Related Work

10 2.1 BERT

11 Transformers are the current state-of-the-art in Natural language processing. The paper "Attention is
12 All you Need" was the paper that they stem from.

13 BERT is the most well known example of a transformer. It was pioneered by a team at Google
14 AI Language. BERT stands for Bidirectional Encoder Representations from Transformers. It is
15 described as being different from recent language representation models as it "is designed to pre-train
16 deep bidirectional representations from unlabeled text by jointly conditioning on both left and right
17 context in all layers." [1] As a result of this, the pre-trained BERT representations can be fine tuned
18 with just one additional output layer to create state-of-the-art models for a large number of tasks, for
19 example question answering and language inference, without substantial task-specific architecture
20 modifications. It uses an encoder and decoder layer to essentially do two jobs: "fill in the blank" and
21 "does sentence B come after sentence A?"

22 BERT-base was trained on 4 cloud TPUs for 4 days and BERT-large was trained on 16 TPUs for 4
23 days. Most people do not have access to this kind of computing power (and cannot afford it either),
24 making the pre-trained paradigm very good at making a very accurate NLP Model available to all.

25 Another advantage of BERT over LSTMs and GRUs is BERT's parallelization capabilities. These
26 models are inherently sequential in nature, precluding parallelization within training examples. This can
27 lead to extremely slow training times when we have a very large dataset with long sequence lengths
28 as memory limits batching across examples.

29 It achieved state-of-the-art results for 11 NLP tasks including a GLUE score of 80.4%.

30 2.2 GPT-2

31 GPT-2 is a 1.5B parameter transformer that achieves state of the art results on 7 of 8 tested language
32 modelling datasets. This language model was trained on a large and diverse data set: 45 million
33 webpages curated/filtered by humans. This create the WebText dataset wwith 80 million rows of data.
34 GPT-2 is a very large transformer model with 48 layers.

35 2.3 XLNet

36 XLNet was created by Carnegie Mellon University in conjunction with the Google Brain team. XLNet
37 outperforms BERT on 20 tasks, often by a large margin[2]. XLNet is a generalized autoregressive
38 pretraining method that utilizes the best auto-regressive language modelling(Transformer-XL) and
39 autoencoding (BERT), while avoiding their limitations. It maximizes the expected log-likelihood of a
40 sequence with respect to all possible permutations of the factorization order.

41 Due to the fact it is an autoregressive language model, XLNet does not rely on data corruption, and
42 thus avoids BERT's limitations due to masking - i.e., pretrain-finetune discrepancy and the assumption
43 that unmasked tokens are independent of each other.

44 2.4 T5

45 T5(Text-to-Text Transfer Transformer) was developed by Google. It utilizes transfer learning, the
46 method of which a model is first "pre-trained on a data rich task before being fine-tuned on a
47 downstream task"[3]. The model is very large with 11B parameters. It has a GLUE score of 89.7 and
48 a superGLUE score of 88.9 which almosts matches human performance (89.8).

49 3 Experimental Setup

50 This section includes a detailed description of the dataset used, the preprocessing applied to the
51 dataset, the proposed algorithm(s), baseline approach and the hyperparameter tuning done.

52 3.1 Dataset Used

53 The dataset used for this assignment was (B.1) Text: "cleaned Emotion Extraction dataset from
54 twitter" taken from Kaggle and is available here.

55 It contains 848,487 rows each containing the emotion of a tweet (disappointed, happy or angry) in the
56 first column, while the second column is cleaned version of the raw tweet data in the third column.

57 The data is balanced with the proportion of disappointed, happy and angry being 34, 33 and 33
58 percent respectively.

59 Since this data is sentence based, this means that the data is sequential. This poses a problem for
60 machine learning algorithms as the right architecture is essential. One of the main difficulty involving
61 sequential data is machine learning algorithms, particularly neural networks, are designed to work
62 with fixed length inputs. Furthermore, the temporal ordering of observations can make it difficult to
63 extract features suitable for use as input to supervised learning models, meaning deep expertise in the
64 domain is typically required.

65 One thing to note about this dataset is how the labels were generated. On the discussion tab of the
66 Kaggle link above, the dataset creator say they used a model based on hashtags and emojis to label
67 the data, meaning it was not done manually. Some tweets, therefore, may be not labelled exactly
68 correct as these emojis or hashtags may have been used sarcastically.

69 3.2 Preprocessing Applied

70 Preprocessing is a very important part of any NLP task. There is a lot of clutter in languages, and this
71 is especially true for twitter data which contains a lot of spelling mistakes and excess punctuation and
72 other inconsistencies.

73 3.2.1 Punctuation and NAs

74 I removed any N/As from the dataset, as well as removing punctuation. While punctuation could be
75 useful for a model like this, again it adds to the vocabulary size which would dramatically increase
76 training time.

77 3.2.2 Number Removal

78 Numbers were removed from the dataset as they are generally not useful for emotion prediction.

79 3.2.3 Lower Casing

80 The text was converted to lowercase to reduce the amount of words in the vocabulary, otherwise,
81 one word would be treated as two in the vocabulary, for example "Dog" and "dog" would be treated
82 as different words due to one having a capital letter at the start. This would lead to a much larger
83 vocabulary that would lead to a significant increase in training time.

84 3.2.4 Stop Words

85 Stop words are common words that appear very often in text data but hold low informational value to
86 a language model. Words like "the", "and" and "it" are some examples. I removed stop words using
87 the CountVectorizer from Scikit-Learn.

88 3.2.5 Stemming

89 I used the nltk package for stemming words. Stemming is the process of reducing inflection in a
90 word into its root form. For example, connected and connection stemmed both become connect. The
91 snowball stemmer from nltk is an effective package for stemming. "The 'english' (snowball) stemmer
92 is better than the original 'porter' stemmer." according to the nltk website linked above.

93 Stemming is an effective way of dealing with sparsity issues. It also helps standardise vocabulary. it
94 reduces down the size of the vocabulary which helps reduce training time.

95 3.2.6 Lemmatizing

96 I used the nltk package for lemmatizing. Lemmatizing is very similar to stemming, but is a more
97 sophisticated approach. It transforms the word to the actual root using a dictionary such as WordNet
98 for mappings. For example, the word "worse" would map to "bad".

99 Lemmatizing, does not appear to improve the accuracy of the model significantly [1], however I have
100 included it anyway as it did not add a significant amount of time to the overall modelling process.

101 3.2.7 Word Encoding

102 Word encoding involves assigning a number to each word. This makes it understandable by the
103 machine learning algorithm as it cannot process strings as it does numbers.

104 3.2.8 Text Enrichment/ Augmentation

105 Word embedding is represents a word in a sentence as a real-valued vector. This vector encodes the
106 meaning of the word such that words that are closer in the vector space are expected to be similar in
107 meaning.

108 I originally used Word2Vec to train my own word embeddings on the data. While the dataset is quite
109 large at approximately 900,00 rows, the model did not perform well using this word embeddings as
110 there was not enough data. For my final model I used pretrained GloVe weights on 2Billion tweets
111 for word embeddings.

112 GloVe stands for "Global Vectors", capturing both global statistics and local statistics of a corpus, in
113 order to make the word vectors. GloVe has an advantage over Word2Vec in that it comes up with
114 global statistics while Word2Vec only comes up with local statistics.

115 3.3 Proposed Algorithms

116 Two algorithms were used during this process: LSTM and GRU

117 3.3.1 LSTM

118 Long short-term memory is a type of Recurrent Neural Network. It is capable of learning order
 119 dependence in sequence prediction problems. LSTMs have a forget gate, an input gate and an output
 120 gate. These gates are fully connected layers with weights that can be trained using back propagation.

121 This enables the LSTM to regulate the flow of information as the gates decide which data is important
 122 and which data should be thrown away. By doing so, it can pass relevant information down the long
 123 chain of sequences to make predictions.

124 The use back propagation via the ADAM algorithm in order to optimize the parameters of these gates.
 125 The ADAM algorithm is an extension to stochastic gradient. The algorithm calculates an exponential
 126 moving average of the gradient and the squared gradient, and the parameters beta1 and beta2 control
 127 the decay rates of these moving averages.

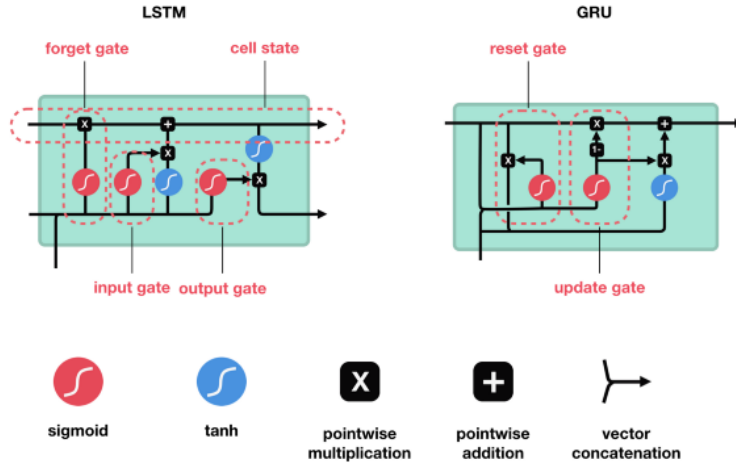


Figure 1: LSTM Architecture

128 3.3.2 GRU

129 GRUs are similar to LSTMs but have different gates. The only have an update and reset gate wheras
 130 the LSTM has the input, output and forget gate. This makes them faster to train.

131 3.4 Baseline Approach

132 For a baseline, I used an LSTM without word embeddings that general defaults for hyperparameters
 133 of 0.01 for the learning rate. As well as it being a single directional model. This had sub-standard
 134 results with accuracy of 70%.

135 3.5 Hyperparameter Tuning

136 Hyperparameter tuning is used to find the optimal model parameters that cannot be found using back-
 137 propagation or similar algorithms. After splitting the data into train, validation and test data, the model
 138 is trained on the train data and predictions are made using the validation dataset. Hyperparameters
 139 that gave the highest accuracy model are then were then chosen for the final model.

140 I used the ray tune package for hyper parameter tuning. It makes use of Async Hyperband Schedul-
 141 ing(ASHA) in order to increase efficiency. Hyperband "focuses on speeding up random search
 142 through adaptive resource allocation and early-stopping"[2].

143 It provides better parallelism and avoids straggler issues during eliminations compared to the vanilla
144 Hyper Band Scheduler.

145 Figure X shows the ASHA Hyperband scheduler in action. We see there some trials are stopped early
146 before 2, 4, 8 and 16(i.e powers of 2) iterations with the reduction factor parameter deciding how
147 often a trial should be checked for early stopping. Some of the more successful trials, usally the ones
148 ending up with 90 % accuracy, are run the full epoch amount. The accuracy plateaus around 8 - 14
149 epochs for the majority of the trials.

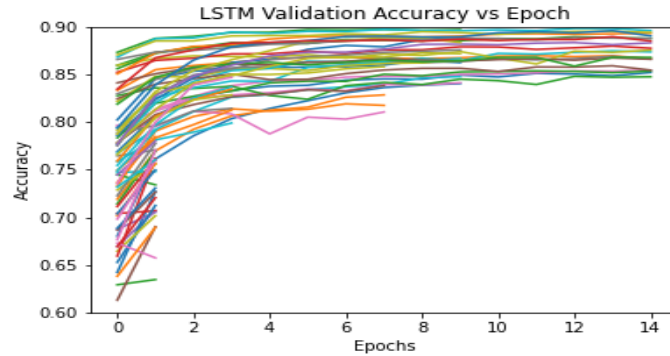


Figure 2: LSTM Hyperparamter Trials

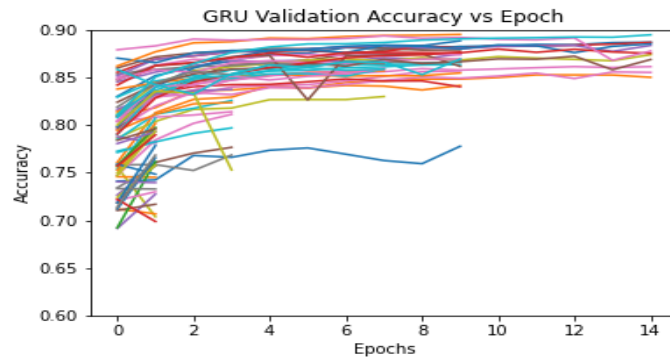


Figure 3: GRU Hyperparamter Trials

150 4 Results

151 Overall we see a very good accuracy on test data from both the LSTM and the GRU, with both having
152 scores of just above 90%.

153 The F1 score on both models is also very high at approx 90%, indicating a good balance of precision
154 and recall.

155 4.1 LSTM

156 We an see that the LSTM is very good at predicting between happy and angry, with almost no
157 mistakes made. It has the greatest difficulty predicting Angry instead of Disappointed. This makes
158 sense as you would imagine people using quite similar words to describe disappointment and anger
159 while using completely different words to describe being happy.

Table 1: LSTM Results

Metric	Train %	Test %
Accuracy	89.67	90.15
F1-Score	89.79	90.20
Precision	90.01	90.52
Recall	89.72	90.19

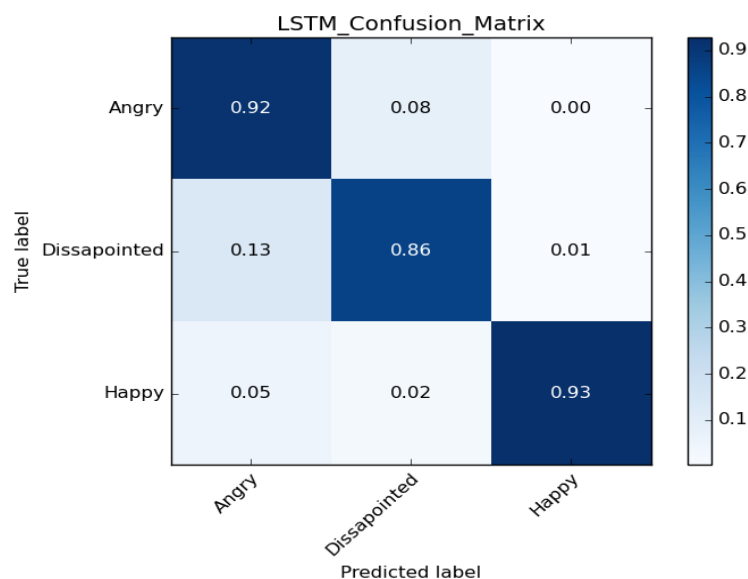


Figure 4: LSTM Confusion

160 4.2 GRU

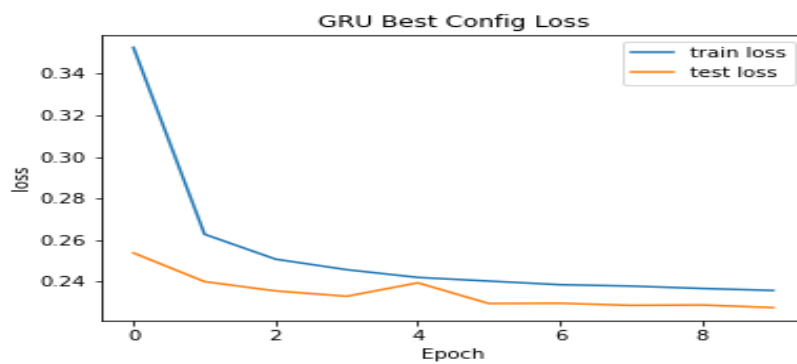


Figure 5: GRU loss per epoch

161 The GRU shows very similar performance to the LSTM, finding it most difficult to tell between a
 162 Dissapointed and Angry tweet.

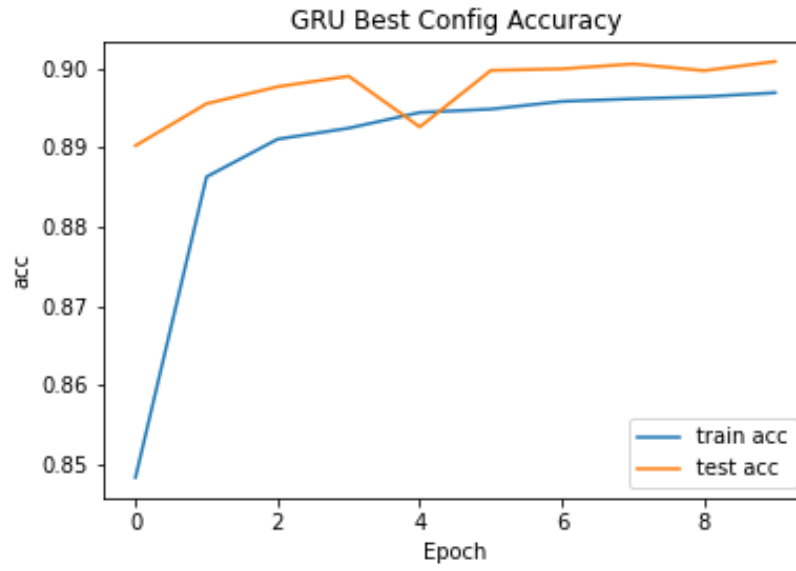


Figure 6: GRU Accuracy per epoch

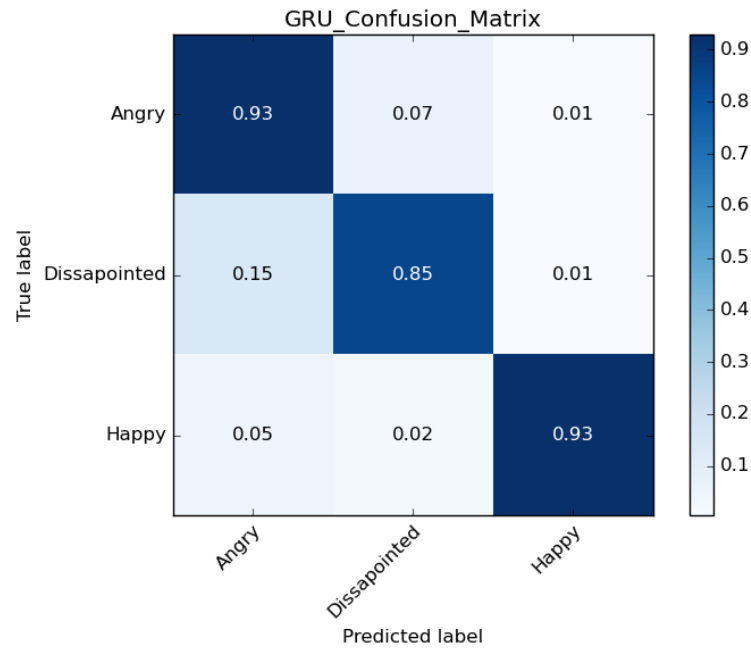


Figure 7: GRU Confusion

Table 2: GRU Results

Metric	Train %	Test %
Accuracy	89.70	90.01
F1-Score	89.81	90.08
Precision	90.09	90.56
Recall	89.75	90.08

5 Conclusion and Future Work

Both LSTMs and GRUs provide very good models for processing sequential data.

The accuracy of the models could be improved by expanding the search space of the hyper-parameters. The word embedding vector size used was 200, and this could be increased to 300 which could potentially lead to an increase in accuracy.

LSTMs and GRUs provide similar performance on this dataset. The GRU's training time is lower and hence a more attractive algorithm to use than the LSTM.

Future work could include building an ensemble of these two models along with additional models such as a convolutional neural network. Ensembles can lead to a reduction in bias and variance while also increasing accuracy.

We could build a model that does not use lowercase words only. People tend to use capital letters when expressing a strong emotion, hence including capital letters may lead to a more accurate model.

We could increase the maximum sequence length to include all 160 words. For my analysis, I used a max sequence length of 40 for the purpose of reducing training time. This obviously leads to the omission of some data which could have otherwise helped to develop a more accurate model. Including punctuation in the dataset should also be considered for a possible increase in accuracy.

Using other model architectures, especially pre-trained models such as BERT could potentially lead to a more accurate model. Many transformers can achieve accuracy near 99% when trained with enough data.

References

- [1] Vaswani, A. & Shazeer, N.M. (2017) Attention is All you Need
- [2] Yang, Z. & Zihang, D. (2020) XLNet: Generalized Autoregressive Pretraining for Language Understanding
- [3] Raffel, C. & Shazeer, N. (2019) Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer
- [4] Camacho-Collados, J. & Pilehvar, M.T. (2018) On the Role of Text Preprocessing in Neural Network Architectures: An Evaluation Study on Text Categorization and Sentiment Analysis. In *Journal of Artificial Intelligence Research (JAIR)* 7, pp. 2. Cambridge, MA: MIT Press.
- [5] Li, L. & Jamieson K. (2018) Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization *Journal of Machine Learning Research* 18 Cornell University

A Appendix