

CSCI 5922 -Neural Networks and Deep Learning

Lab - 3 Solutions

Gowri Shankar Raju Kurapati
Student ID 110568555

October 10, 2022

1 Impact of RNN Architecture

In this section, we will explore the impact of various RNN architectures by analyzing spam/ ham classification of SMS messages. We have considered the dataset from the UCI repository(LINK) which consists of nearly 5.5k examples with varying lengths. We split the dataset into 70/30 train/test and create three models which are Vanilla RNN, LSTM and GRU. All the sentences are tokenized and the maximum length is set to 150. Standard Tokenizers are used for this experiment to vectorize the words in the message and to maintain the consistency of sentence lengths, padding is also applied to give a max length of 150 for batch size.

For all three models, the input dimension is the size of the vocabulary which is around 8k and the dense embedding is set to the size of 8. All the cell states have a dimension of 32. The output of the last cell state is passed through a fully connected layer and sigmoid layer to predict if it is spam(target label is 1) or not (target label is 0). The loss function used for back-propagating the gradients is Binary Cross Entropy.

| Model / Metric | RNN | LSTM | GRU |
|----------------|------|-------------|-------------|
| Precision | 0.94 | 0.97 | 0.98 |
| Recall | 0.92 | 0.96 | 0.95 |

The above table gives the precision and recall scores for the different models used. We can see a gradual increase in precision from RNNs to GRUs with a significant increase from RNN to LSTM/GRU. This is expected as the LSTM is more robust to take the information needed with selective read, write and forget gates. Also, the recall increases which means that the LSTM is better at classifying the spam messages correctly among all the ones that actually needed to be predicted as spam.

To explore how better the model is at classifying varying-length sentences, we have further divided the test set into three equal parts i.e short, medium, and long inputs. Below is the histogram for the length of the message and the corresponding count.

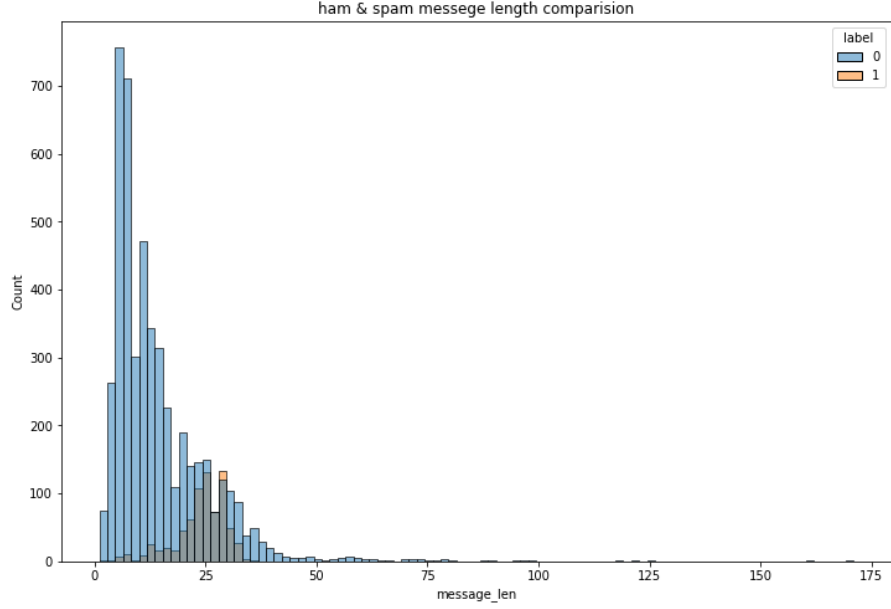


Figure 1: Histogram for message length and count in Dataset

As you can see from figure 1, the trisecting points to equally divide the test data are at message lengths of 8 and 18. All the messages of length less than or equal to 8 fall under the short category, between 8 and 18 fall under medium, and greater than 18 fall in the long input category. We have evaluated the models in these three categories and the metrics are tabulated below.

Short Input Lengths

| Model / Metric | RNN | LSTM | GRU |
|----------------|------|------|-----|
| Precision | 0.98 | 1 | 1 |
| Recall | 0.25 | 0.5 | 0.5 |

Medium Input Lengths

| Model / Metric | RNN | LSTM | GRU |
|----------------|------|------|------|
| Precision | 0.99 | 1 | 1 |
| Recall | 0.75 | 0.88 | 0.88 |

Large Input Lengths

| Model / Metric | RNN | LSTM | GRU |
|----------------|------|------|------|
| Precision | 0.98 | 0.99 | 0.99 |
| Recall | 0.92 | 0.98 | 0.99 |

From *figure - 1* One important thing to observe is that most of the spam messages are of large input size. This holds true even for the training split as the splitting is done in a way to preserve class distribution. That is why the recall scores for large input sizes are very similar to the scores on the whole dataset.

The number of actual spam messages is 8 and 16 among the small and medium input sizes, out of which 6 and 4 are miss-classified by the RNN and LSTM/GRU, and 4 and 14 are mis-classified by LSTM/GRU respectively. That is why the recall scores are 0.25, 0.75 for small input size, and 0.5 and 0.88 for large input size. Since the data points are very scarce for spam messages, the inferences made here may not hold strongly but we can see that LSTM and GRU are classifying the spam messages

better compared to RNNs.

Also, we can see that the recall score gradually increased from RNN to LSTM and GRU, which states that LSTM is better at handling long-length sequences compared to vanilla RNNs. This observation can be directly tied back to the vanishing gradients problem in the RNNs which makes it the model hard to remember the information in the initial inputs.

There is not much difference in performance between LSTM and GRU models for any category of the input-sized dataset. Due to the skewness and limited size of the dataset, a clear comparison between the performances of LSTM and GRU cannot be made.

2 Impact of Pretrained Word Embedding

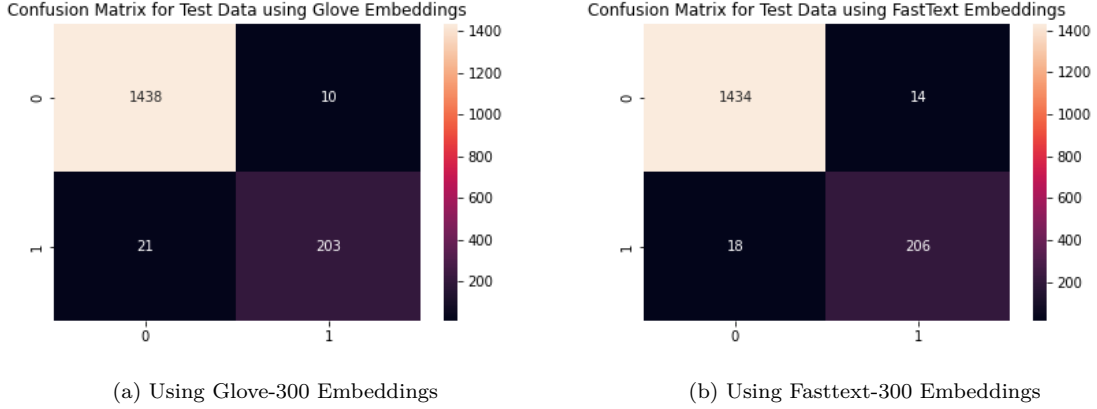
To understand the impact of the pre-trained word embeddings, we have trained two GRU models by using two different type of embeddings which are Glove embeddings(300 Dimension) and FastText(300 Dimension) Wiki embeddings. The GRU model consists of the embedding layer followed by the GRU Unit and later the output of the last layer of GRU is connected to fully connected network which has 2 hidden layers of 250 and 120 nodes with ReLu activation. The output node has the sigmoid activation to predict the class. We have used early stopping as regularizer and Adam optimizer with learning rate of 0.001, which weren't present in the models used for first part of the question. Also, the same loss as in first question which is Binary Cross Entropy is used.

We have also used to 100 dimensional Glove Embeddings to make the comparison of score on the impact of the size of the embeddings. The Precision and Recall for the models are given below.

| Model / Metric | GRU - Glove(100D) | GRU- Glove(300D) | FastText (300D) |
|----------------|-------------------|------------------|-----------------|
| Precision | 0.968 | 0.98 | 0.95 |
| Recall | 0.88 | 0.91 | 0.94 |

As we can point out, there is a increase in both precision and recall score when we move from 100 dimensional glove embeddings to 300 dimensional glove embeddings but the increase is really not significant. Since the the dataset is sparse which consists of around just 8k words, 100 dimensional dense word embedding is able to capture the information needed to make the classification. Maybe if the dataset had many words and with closely related semantics, 300 dimensional embeddings would have a made difference when compared with 100 dimensional counter part. Glove-300 embeddings seem to perform better when competed against the FastText-300 Embeddings in precision scores. Glove-300 is slightly better at classifying the true spams better from all of the predicted spams.

Figure 2: Confusion Matrix for Test Data



From the Precision & Recall Table and also from the confusion matrices (*figure 2*), though the precision is slightly better for Glove -300 compared to Fasttext-300, Fasttext is unconditionally better at recall by classifying the true spams better from all of the predicted spams. We usually wouldn't want to miss on ham messages and we would be okay to let a few spam messages in, but we would never want to miss out on a ham message. That is when precision is important compared to recall and we would go with Glove embeddings in that case as its precision is higher. But if the use case is other way around and we are stringent on having to remove the spam messages at the cost of leaving a few ham messages behind, the out emphasis would be on recall and choose FastText model in that scenario. But ideally, we would want both the precision and recall to be good enough. Though pretrained word embeddings aren't making much difference when compared to the tokenized input representation for our dataset, we can expect them to work better when the vocabulary of the dataset is large and includes words used under different semantics.

3 CODE