Ryan McCrory
Professor Flanegan
Natural Language Processing
January 20, 2020

# Assignment 2

**Part 1:**

**1.1**    Once I familiarized myself with the starter code, I modified it such that it split the IMDB data into a training set (60% of the data), a validation set (40% of the data), and a test set (100% of the data). I then preprocessed each of the data sets. Then I changed the model from using GRU layers, as it was in the starter code, to using a single SimpleRNN layer. Then I added the validation set to the history = model.fit() line, and began to play around with several hyperparameters. The accuracy of the models with several hyperparameter settings are shown in the tables below. Note that each table below uses the best (highlighted) hyperparameter of each of the tables before it, such that the number of epochs table uses the best optimizer, the training batch size table uses the best number of epochs and the best optimizer, and so on.

### Optimizer

|  | Training | Validation |
|---|---|---|
| Optimizer = "adam" | 0.5101 | 0.4990 |
| Optimizer = "sgd" | 0.5054 | 0.4996 |
| Optimizer = "rmsprop" | 0.5079 | 0.4990 |

### Training batch size

|  | Training | Validation |
|---|---|---|
| batch  = 32 | 0.5717 | 0.5106 |
| batch = 128 | 0.5931 | 0.5219 |
| batch = 512 | 0.5858 | 0.5036 |

### Number of Epochs

|  | Training | Validation |
|---|---|---|
| Epochs  = 1 | 0.5054 | 0.4996 |
| Epochs = 3 | 0.5276 | 0.4996 |
| Epochs = 5 | 0.5708 | 0.5040 |

### Embed size

|  | Training | Validation |
|---|---|---|
| Embed size = 64 | 0.6043 | 0.5108 |
| Embed size = 128 | 0.5931 | 0.5219 |
| Embed size = 256 | 0.5973 | 0.5183 |

**Dropout rate**

|  | Training | Validation |
|---|---|---|
| 0 | 0.5931 | 0.5219 |
| 0.2 | 0.5633 | 0.5151 |
| 0.4 | 0.5527 | 0.5020 |

**1.2**    As shown above, the best performing model had an accuracy of 0.5931 on the training set, and 0.5219 on the validation set. This was with the following hyperparameters:

Optimizer = "sgd",    Number of Epochs = 5,      Training batch size = 128
Embed size = 128,   Dropout rate = 0

Applying these hyperparameters to the test set results in an accuracy of 0.523. This is not an extremely impressive accuracy due to the simplicity of this model, vanishing gradients, and the model forgetting too much during each iteration, but I am sure the GRU model in part 2 will perform better.

**Part 2:**
**2.1**    Because the SimpleRNN model had low accuracy on the test set of only 0.523, I am going to use a gated recurrent unit (GRU) model instead. The results of this model using the same experimental procedure as in part 1 are shown in the tables below

**Optimizer**

|  | Training | Validation |
|---|---|---|
| Optimizer = "adam" | 0.9625 | 0.7592 |
| Optimizer = "sgd" | 0.5043 | 0.4989 |
| Optimizer = "rmsprop" | 0.9378 | 0.7698 |

**Training batch size**

|  | Training | Validation |
|---|---|---|
| batch  = 32 | 0.8761 | 0.7734 |
| batch = 128 | 0.8755 | 0.7765 |
| batch = 512 | 0.8770 | 0.7768 |

**Number of Epochs**

|  | Training | Validation |
|---|---|---|
| Epochs = 1 | 0.6851 | 0.7310 |
| Epochs = 3 | 0.8761 | 0.7734 |
| Epochs = 5 | 0.9378 | 0.7698 |

**Embed size**

|  | Training | Validation |
|---|---|---|
| Embed size = 64 | 0.8639 | 0.7712 |
| Embed size = 128 | 0.8755 | 0.7768 |
| Embed size = 256 | 0.8877 | 0.7761 |

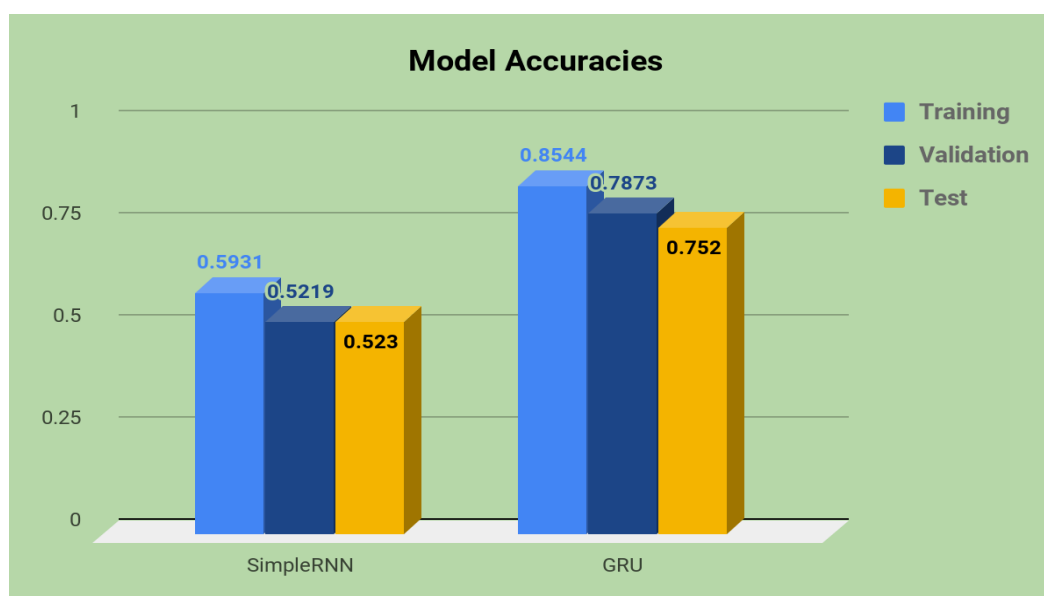**Dropout rate**

|  | Training | Validation |
|---|---|---|
| 0 | 0.8755 | 0.7768 |
| 0.2 | 0.8683 | 0.7764 |
| 0.4 | 0.8544 | 0.7873 |

As shown above, the best performing model had an accuracy of 0.8544 on the training set, and 0.7873 on the validation set. This was with the following hyperparameters:

Optimizer = "rmsprop",     Number of Epochs = 3,     Training batch size = 512
            Embed size = 128,            Dropout rate = 0.4

Applying these hyperparameters to the test set results in an accuracy of 0.752.

**2.2**     We would expect the GRU model to have better accuracy than the simple RNN model on longer sequences. This is because the simple RNN model suffers from vanishing gradients during learning, which makes it hard to propagate error into the distant past. Additionally, state tends to change a lot on each iteration, meaning the model "forgets" too much. We would not expect this to be the case for the GRU model because the update and reset gates allow the model to remember important data and forget less important data. This assumption does hold for my models, as shown below. The GRU model performed better than the Simple RNN model in every data set.

**Part 3:**

**3.1** The embeddings I used in the previous sections were embedding layers which were trained alongside the model whereas the embeddings I used in this section were pretrained from glove (https://nlp.stanford.edu/projects/glove). Using an embedding layer can be beneficial because you can train it with respect to the problem you are trying to solve, whereas using a pretrained embedding may not be good for a given problem, although it does simplify the learning of the model.

**3.2** Below are the charts showing the hyperparameters I tested while optimizing my model using the same experimental procedure as in the previous sections.

**Optimizer**

|  | Training | Validation |
|---|---|---|
| Optimizer = "adam" | 0.9542 | 0.7396 |
| Optimizer = "sgd" | 0.5313 | 0.5056 |
| Optimizer = "rmsprop" | 0.8858 | 0.7608 |

**Training batch size**

|  | Training | Validation |
|---|---|---|
| Batch = 8 | 0.8480 | 0.7654 |
| Batch = 16 | 0.8431 | 0.7658 |
| batch = 32 | 0.8514 | 0.7648 |

**Number of Epochs**

|  | Training | Validation |
|---|---|---|
| Epochs = 3 | 0.8015 | 0.7561 |
| Epochs = 4 | 0.8514 | 0.7648 |
| Epochs = 5 | 0.8858 | 0.7608 |

**Dropout Rate**

|  | Training | Validation |
|---|---|---|
| 0 | 0.8431 | 0.7658 |
| 0.2 | 0.8111 | 0.7673 |
| 0.4 | 0.7629 | 0.7622 |

As shown above, the best performing model had an accuracy of 0.8111 on the training set, and 0.7673 on the validation set. This was with the following hyperparameters:

Optimizer = "rmsprop",    Number of Epochs = 4,    Training batch size = 16
Dropout rate = 0.2

Applying these hyperparameters to the test set results in an accuracy of 0.746. Performance does not change much using pre trained embeddings. My GRU model in which I trained the embedding had an accuracy of 0.752 on the test set, while the GRU model with pre-trained embeddings had an accuracy of 0.746. So, the model in which I trained the embeddings performed 0.006 better than the pre-trained model, which is not

a significant difference. Perhaps a different pre-trained embedding would have made more of a difference for this specific task, as the two I used turned out to be quite similar.

**3.3**     Antonyms may often have similar word embeddings because the models are trained on large amounts of data where antonyms are often used in similar contexts. For example, take the word "good" in the sentence "I can't believe I have a good grade in NLP". If there was another sentence in the training data in which "good" was replaced with "bad", the context of the sentence would be unchanged, although its meaning would be the opposite. This causes antonyms to have similar word embeddings.

**3.4**     Below are the embedding matrices for the words "good" and "bad", which are antonyms.



The embeddings are fairly similar, but not to a great extent. The first entry in the embedding matrix for "good" is -0.40040001 and the first entry in the embedding matrix for "bad" is -0.47384, which are similar. There are several of these similarities throughout the matrices, although not in every position.

Note: I could not get the code to work to test whether or not substituting these antonyms in a sentence changed the sentiment prediction.