

User Response Classification Challenge

ADRIEN COGNY
Cornell Tech
ac2753@cornell.edu

April 24, 2018

I. INTRODUCTION

Chatbots are used in many applications today: customer support, flight booking, scheduling meeting, ordering food and many more. The application of a chatbot explored in this dataset is for a therapy chatbot. These types of chatbot, while very effective, may require human intervention. Determining when a human should intervene can be quite important, in this case when a person requires help in dealing with a complex situation, and requires tools to identify these situations.

The data set contains 80 examples of responses entered into a therapy chatbot. Each of these responses contains an id as well as an identification. The identification is either "flagged" if the response was flagged for human intervention or "not flagged" if not.

The task at hand was to create an AI agent to classify the user response.

II. TOOLS

The following tools and modules were used to complete this task:

- python 3.6.5 (using conda)
- pandas (0.22.0)
- dy-net (for the RNN) (2.0)
- scikit learn (for the Random Forest) (0.19.1)
- numpy (1.14.2)
- tqdm (for progress bars) (4.22.0)
- csv (For reading the csv embeddings to pandas) (1.0)

- re (regular expressions) (for cleaning the data) (2.2.1)

III. PREPARING DATA

The input data being sentences had to be cleaned up before passing into the models.

The first step was to load the csv file into a pandas dataframe and see what the data looked like. The data was, as mentioned above, a label as well as a sequence of words (not an array implementation yet). Due to the inherent nature of natural language processing both the label and sequence of words had to be converted to something which the machine could understand. That is, the label had to be converted from "flagged" or "not flagged" to 1 or 0 respectively and the sequence had to be converted to a series of word embeddings where each embedding represented a single word.

Natural language contains many words that are very common in sentences. These so-called stop words ("their", "he", "she", etc...) can make classifying a sentence very hard as, when combining word embeddings, they will take over the representation of the sentence simply by sheer number. That is why in the pre-processing of each sentence (before it is given to the model), the stop words were removed from the sentence. By removing the stop words, we do not lose much important information and are able to classify more easily.

The models were created to do the conversion from label to 1 or 0 and from sentence

of words to sequence of embeddings. The embeddings used were the GloVe 6B embeddings which come from wikipedia scraping¹.

IV. MODELS

i. RNN

When looking at sentence classification, one of the first thought was too look at an RNN encoder that would encode the sentence word by word and the computing a probability of being "flagged" or "not flagged". The label with the highest probability would then be applied to the sentence input. This was furhter confirmed by researching text classification and RNNs [Lee and Dernoncourt]

i.1 RNN Description

An RNN was created using dynet, the rnn was an encoder which would take each word in the sentence, find its embedding and compute a hidden state which it would then pass along to the next node in the network. The encoder yielded a final hiddent state on which softmax was applied to find the predictions.

i.2 Tuning Parameters

When the model was created, the different parameters were tuned:

- Embedding Dimmension
- Hidden Dimmension Size
- Number of Epochs Run

The results for all of these tuning experiments are shown in the results section.

ii. Random Forest

After getting results for the RNN encoder and finding the best possible RNN, it was posited (based on research into text classification) that a random forest classifier could be more apt at this task. [Xu, Guo, Ye and Cheng]

¹<https://nlp.stanford.edu/projects/glove/>

ii.1 Random Forest Description

The Random Forest Model was created using scikit learn's random forest classification tool. This model would take the word embeddings for each word in the sentence and combine them in some way (either summation or mean which was a test performed). The resultant embedding would then be passed through the Random Forest model to get a prediction.

ii.2 Tuning Parameters

When the model was created, the different parameters were tuned:

- Number of Estimators
- Sentence to Embedding Methodology

Whereas the number of estimators is a property of the random forest model itself, the sentence to embedding methodology describes how a sentences (or rather sequence of words) is transformed into a single vector wich can be input to the Random forest model.

There were two methods for embedding the sentence. One was to compute the mean of all the word embeddings and the other to compute the sum. The mean would attempt to construct a mean representation of the sentence using all the words in the sentence. The summation would create a sentence which was a sum of its parts.

V. RESULTS

i. RNN

i.1 Hidden dimmension

The hidden dimmension test was done by keeping all parameters of the RNN constant except for the hidden dimmension of the RNN. The following figure shows the results for the hidden dimmension test performed. The test was performed by changing the size of the embedding dimmension from 0 to 9 dimmensions in steps of 1. The training loss, dev set true positive , dev set true negative , dev set false

positive and dev set false negative were computed for each of the models run and a graph was created showing the true positive rate and the false positive rate.

type	test	features	readPos	num_layers	embeddingSize	hiddenDim	train_loss	accuracy	truePos	trailing	falsePos	falseNeg
0	RNN	hidden_dim	hidden_dim	400	1	50	0.261499	0.716667	0	11.2	6.8	0
1	RNN	hidden_dim	hidden_dim	400	1	50	0.5076	0.591667	2.2	11.7	5.5	4.8
2	RNN	hidden_dim	hidden_dim	400	1	50	0.866661	0.504167	2.6	9.5	5.7	6.2
3	RNN	hidden_dim	hidden_dim	400	1	50	0.0041726	0.591667	2.4	11.8	5.3	4.8
4	RNN	hidden_dim	hidden_dim	400	1	50	0.0070414	0.606667	2	9.6	5.1	7.2
5	RNN	hidden_dim	hidden_dim	400	1	50	0.0041676	0.591667	1.9	10.8	4.9	6.4
6	RNN	hidden_dim	hidden_dim	400	1	50	0.0030066	0.5975	1.6	11.3	4.8	6.5
7	RNN	hidden_dim	hidden_dim	400	1	50	0.0041889	0.586667	1.9	11.5	5.8	4.8
8	RNN	hidden_dim	hidden_dim	400	1	50	0.00200365	0.5925	1.8	10.4	6.1	5.6
9	RNN	hidden_dim	hidden_dim	400	1	50	0.0020258	0.491667	2.3	9.5	5.2	7

Figure 1: Shows the table of hidden dimension tests for the RNN

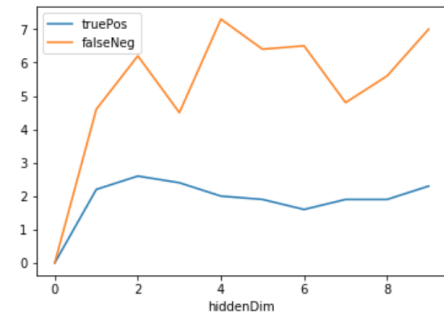


Figure 2: Shows the graph of true positive count vs hidden dimension of RNN

This graph shows that the true positive count is always lower than the false negative count. This means that there are more instances where the model will classify a sentence as "not flagged" when in fact it should be "flagged" than there are instances where the model correctly classifies a sentence as "flagged." While this points towards this particular model (with the hyperparameters described below) not being good, the true positive rate and accuracies are derived (and much more important) metrics to look at.

While this graphs shows the true positive count as well as the false negative counts, a more interesting metric which can be derived from the true positive count and the false negative count is the true positive rate (tpr) which shows how much of the truth the model captures.

The following figure shows the tpr for the RNN model for different hidden dimensions.

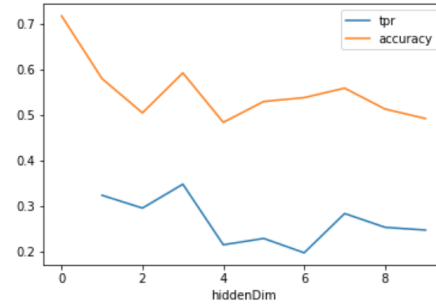


Figure 3: Shows the graph of true positive rate vs hidden dimension of RNN

From this figure, we can clearly see that the tpr is greatest when the hidden dimension is 3, with the other parameters set to: number of epochs ran = 400, number of layers = 1 and embedding size = 50. This model received an accuracy of 0.591667 and a true positive rate of 0.347826.

While the graph of accuracies shows that when the hidden dimension is of size 0, the accuracy jumps to 0.7, this model would not be considered to be a good model as the true positive rate is non-existent because we are not classifying any results as being "flagged", which defeats the whole purpose of the model.

Thus, the hidden dimension will be set to 3 for the other models.

i.2 GloVe embedding size

When the hidden dimension was found, the next hyper-parameter which could be tuned was the size of the GloVe embeddings used to convert the sentences into something the models could understand.

This next figure shows the true positive counts and false negative counts for each embedding dimension available (50, 100, 200 and 300).

type	test	features	readPos	num_layers	embeddingSize	hiddenDim	train_loss	accuracy	truePos	trailing	falsePos	falseNeg
RNN	embedding_dim	embedding_size	embedding_size	400	1	50	0.0011101	0.599667	2	11.8	5.2967	5.1033
RNN	embedding_dim	embedding_size	embedding_size	400	1	100	0.0027265	0.596667	2.5967	9.6	5.0967	6.9033
RNN	embedding_dim	embedding_size	embedding_size	400	1	200	0.00961309	0.591667	2.8967	11.3333	5.0967	4.7033
RNN	embedding_dim	embedding_size	embedding_size	400	1	300	0.00702919	0.595556	2.4	9.7333	5.4	6.4967

Figure 4: Shows the table of embedding dimension tests for the RNN

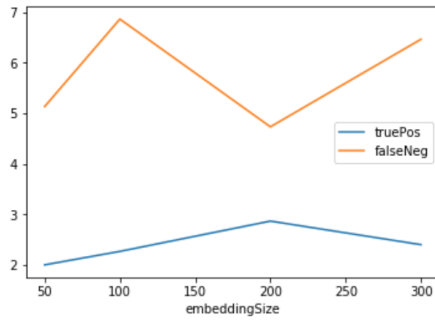


Figure 5: Shows the graph of true positive count and false negative count vs embedding dimensions of RNN

The table and graph show that the true positive count is always smaller than the false negative count. However a dip can be seen in the false negative count at an embedding size of 200. This dip in false negatives is coupled by a rise in the true positive. This indicates that the embedding size of 200 yields the best results. ²

Confirmation of this is seen in the figures below as both the true positive rate and accuracy are highest at an embedding size of 200.

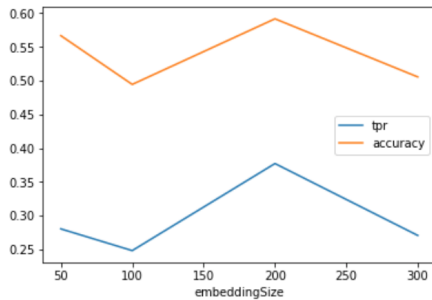
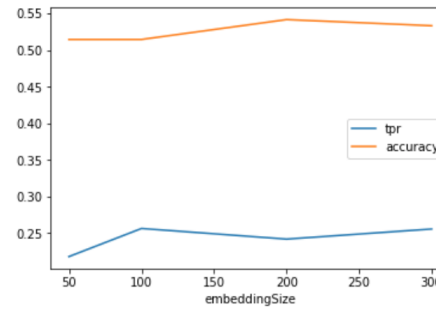
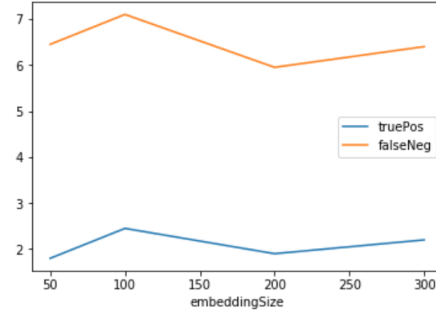


Figure 6: Shows the graph of true positive rate vs embedding dimensions of RNN

However upon further investigation, it seems that this may have been a random occurrence due to the randomization of the data. After running this test several times, the graph did not show any significant improvement. For

²Note that the increase in the true positive count and decrease in false negative counts graphed are not the same magnitude as the samples were drawn randomly from the training set and thus may not always contain the same number of positive and negative examples.

example, the following two figures was another run which illustrates that, depending on the random set of data, the best embeddings change



Therefore, for computation purposes, the embedding size was set to 100. This was thought to be a good compromise between a larger embedding which may provide more information and the speed of loading the embeddings.

i.3 Number Layers

The final hyperparameter to be tuned was the number of layers that the rnn contained. Before this test, the RNNs trained had a single layer. This test was performed twice in order to make sure there was consistency among even more random samples. The number of layers of the rnn were varied from 1 to 25 in increments of 1 for both tests.

The first test's graph shows the number of layers vs true positive count and false negative count is shown below:

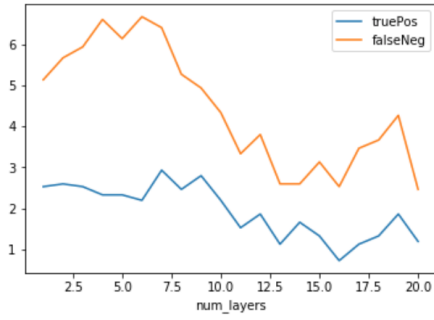


Figure 7: Shows the graph of true positive rate vs number of layers of RNN for the first trial

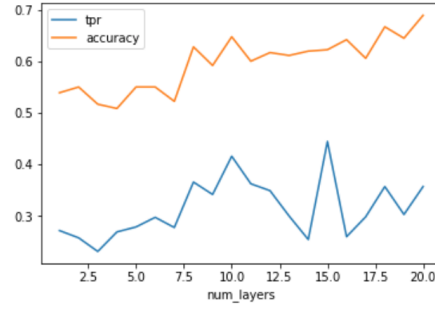


Figure 10: Shows the graph of true positive rate vs number of layers of RNN for the second trial

Both of these tests show that as the number of layers in the rnn increases the number of false negatives also increases. While the exact number of layers cannot be tuned very accurately due to the fairly small dataset, it can be said that around 15 layers seems to be a fairly good point. That is because both the accuracy and the true positive rates are

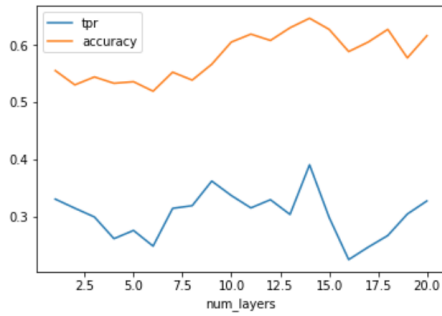


Figure 8: Shows the graph of true positive rate vs number of layers of RNN for the first trial

The second test's graph shows the same statistics as the previous one but for the second trial

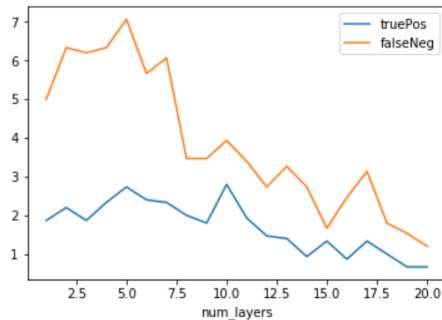


Figure 9: Shows the graph of true positive rate vs number of layers of RNN for the second trial

i.4 RNN Final Model

The final RNN model had the following parameters:

- Hidden Dimmension = 3
- Number of Hidden Layers = 15
- Embedding Size = 100
- Embedding File = 'glove/-glove.6B.100d.txt'
- Number of Epochs trained for = 400

The model was run 15 times to find the best true positive rate (TPR) on the dev set. This model was then saved. The saving procedure for the dynet model was more complex than for the Random Forest Model as dynet is unable to be pickled. Thus, each of the model-wide variables were pickled and the dynet save method was applied to the model. Thus, if the model was loaded from file (which can be done by calling the rnn model constructor with a filename as the filename argument), each piece is loaded separately and then the parameters are loaded through dynet.

The final rnn model got a true positive rate on its dev set of 0.375.

It was saved under the "RNNFinalModel" directory and contains 10 files which are required for it to load.

ii. Random Forest

After doing some more research, it was found that random forest classifiers could be utilized for this task. The RNN having given decent results but not exceptional, this was thought to be worth a try. [Xu, Guo, Ye and Cheng]

ii.1 Summation Methodology

type	test	filename	n_estimators	accuracy	truePos	trueNeg	falsePos	falseNeg	tpr	auc
RandomForest	sum	randomForest_1_estimators	1	0.686967	3.8	12.2	3.86967	4.13333	0.478902	0.489052
RandomForest	sum	randomForest_2_estimators	2	0.672222	5.8	10.3333	2.73333	5.13333	0.520488	0.579668
RandomForest	sum	randomForest_3_estimators	3	0.733333	5.93333	14.0667	3.33333	3.06667	0.532544	0.514563
RandomForest	sum	randomForest_4_estimators	4	0.705056	4.93333	12	2.66667	4.4	0.528571	0.549123
RandomForest	sum	randomForest_5_estimators	5	0.730556	3.4	14.1333	3.66667	2.6	0.586667	0.46789
RandomForest	sum	randomForest_6_estimators	6	0.680556	4.73333	11.6	2.73333	4.93333	0.486655	0.633829
RandomForest	sum	randomForest_7_estimators	7	0.780556	3.06667	15.6667	3.2	2.06667	0.587423	0.489362
RandomForest	sum	randomForest_8_estimators	8	0.718444	3.13333	14.1333	3.86667	2.86667	0.522222	0.447819
RandomForest	sum	randomForest_9_estimators	9	0.736111	3.26667	14.4	3.93333	2.4	0.576471	0.433754
RandomForest	sum	randomForest_10_estimators	10	0.811111	4.33333	15.3333	2.93333	2	0.642111	0.611988
RandomForest	sum	randomForest_11_estimators	11	0.722222	3	14.3333	4.8	2.06667	0.582159	0.364727
RandomForest	sum	randomForest_12_estimators	12	0.741867	3.86667	14.1333	4	2.2	0.625	0.478261
RandomForest	sum	randomForest_13_estimators	13	0.747222	3.06667	14.8667	4.53333	1.53333	0.666667	0.420509
RandomForest	sum	randomForest_14_estimators	14	0.738889	3.86667	14.0667	4.2	2.36667	0.639535	0.466102
RandomForest	sum	randomForest_15_estimators	15	0.761111	3.6	14.8667	4.26667	1.46667	0.715268	0.457827
RandomForest	sum	randomForest_16_estimators	16	0.725	3.4	14	4.13333	2.46667	0.579545	0.451327
RandomForest	sum	randomForest_17_estimators	17	0.772222	3.6	14.9333	3.8	1.86667	0.682544	0.486486
RandomForest	sum	randomForest_18_estimators	18	0.727778	3.33333	14.1333	4.06667	2.46667	0.574713	0.45045
RandomForest	sum	randomForest_19_estimators	19	0.772222	3.8	14.7333	3.26667	2.2	0.633333	0.537736
RandomForest	sum	randomForest_20_estimators	20	0.744444	3.6	14.2667	3.73333	2.4	0.6	0.490909

Figure 11: Shows the table for the experiments where the number of estimators used in the random forest classifier were used. This is for the summation methodology.

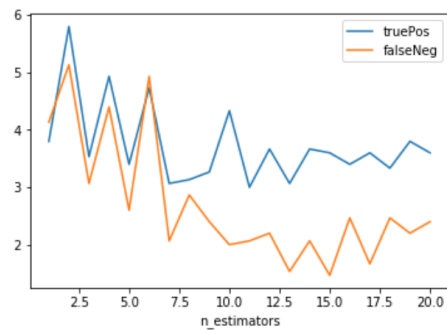


Figure 12: Shows the graph for the experiments where the number of estimators used in the random forest classifier were used. This is for the summation methodology.

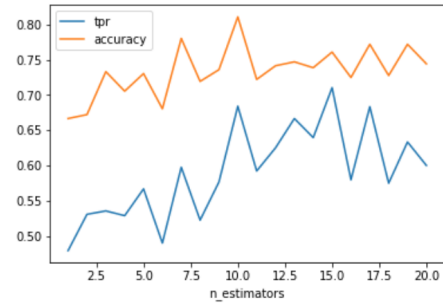


Figure 13: Shows the graph of true positive rate and accuracy vs number of estimators of random forest with summed embeddings as representation for the sentence

ii.2 Mean Methodology

type	test	filename	n_estimators	accuracy	truePos	trueNeg	falsePos	falseNeg	tpr	auc
RandomForest	mean	randomForest_meanNum_1_estimators	1	0.55	2.86667	10.3333	6.2	4.6	0.383029	0.318776
RandomForest	mean	randomForest_meanNum_2_estimators	2	0.586111	5.46667	8.6	3.4	6.53333	0.453006	0.491641
RandomForest	mean	randomForest_meanNum_3_estimators	3	0.555556	2.73333	13	4.66667	3.6	0.421578	0.389389
RandomForest	mean	randomForest_meanNum_4_estimators	4	0.644444	4.06667	11.4	4.06667	4.06667	0.476662	0.5
RandomForest	mean	randomForest_meanNum_5_estimators	5	0.658333	3	12.6	5.2	3	0.5	0.369564
RandomForest	mean	randomForest_meanNum_6_estimators	6	0.661111	4.06667	11.8	3.33333	4.8	0.458647	0.54895
RandomForest	mean	randomForest_meanNum_7_estimators	7	0.675	3.13333	13.0667	4.73333	3.06667	0.503376	0.386305
RandomForest	mean	randomForest_meanNum_8_estimators	8	0.652778	3.06667	12	4	4.33333	0.483333	0.476881
RandomForest	mean	randomForest_meanNum_9_estimators	9	0.683333	2.96667	14.3333	4.86667	2.73333	0.450506	0.286077
RandomForest	mean	randomForest_meanNum_10_estimators	10	0.675	2.53333	13.0667	4.93333	2.86667	0.459136	0.330886
RandomForest	mean	randomForest_meanNum_11_estimators	11	0.705556	3.13333	13.6	4.4	2.66667	0.54023	0.415829
RandomForest	mean	randomForest_meanNum_12_estimators	12	0.628333	2.4	13.4	4.4	3.8	0.387087	0.320541
RandomForest	mean	randomForest_meanNum_13_estimators	13	0.648889	3.0667	13.2667	5.46667	2	0.625253	0.314546
RandomForest	mean	randomForest_meanNum_14_estimators	14	0.55	2.4	13.2	5.06667	3.33333	0.418829	0.211459
RandomForest	mean	randomForest_meanNum_15_estimators	15	0.720556	3.13333	14.4	4.4	2.06667	0.625564	0.415829
RandomForest	mean	randomForest_meanNum_16_estimators	16	0.702778	2.93333	12.9333	4.23333	2.8	0.584158	0.475805
RandomForest	mean	randomForest_meanNum_17_estimators	17	0.686111	2.93333	13.5333	4.26667	3.26667	0.473118	0.425427
RandomForest	mean	randomForest_meanNum_18_estimators	18	0.691667	2.93333	13.6667	4.8	2.6	0.530712	0.37931
RandomForest	mean	randomForest_meanNum_19_estimators	19	0.711111	3.06667	14	4.33333	2.4	0.585876	0.452828
RandomForest	mean	randomForest_meanNum_20_estimators	20	0.711111	3	14.0667	4.46667	2.46667	0.54878	0.451786

Figure 14: Shows the table for the experiments where the number of estimators used in the random forest classifier were used. This is for the mean methodology.

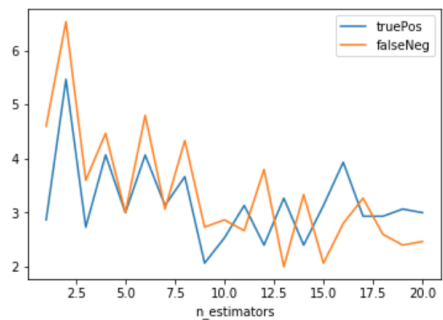


Figure 15: Shows the graph for the experiments where the number of estimators used in the random forest classifier were used. This is for the summation methodology.

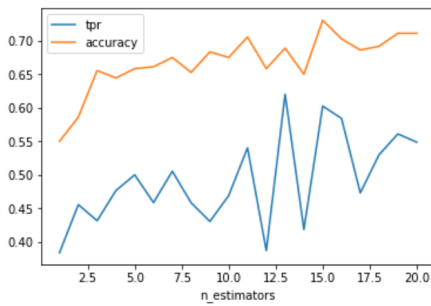


Figure 16: Shows the graph of true positive rate and accuracy vs number of estimators of random forest with summed embeddings as representation for the sentence

ii.3 Mean vs Summation Methodology

From the figures above, it can be seen that, for all numbers of estimators used in the random forest mode the summation accuracy and true positive rates are higher than the mean methodology. This makes sense as the summation method attempts to feed the model the summation of the words rather than the average representation.

ii.4 Random Forest Final Model

The final model for the random forest had the following paramters:

- Number of Estimators = 10
- Embedding Size = 200
- Embedding File = 'glove/-glove.6B.200d.txt'
- Methodology = Summation

In the same way the final model for the rnn was selected, the final model for the random forest was selected by running the model on 15 different training and dev sets, and saving the one which had the highest true postive rate amongst the models with more than 3 flagged truths.

The final Random Forest model got a final true positive rate on its train-test split of 0.8.

It was saved under the "RandomForestFinalModel" directory and is named "bestRandomForest.model"

iii. Final Model

Because the difference between the random forest model and the rnn true positive rates was significant, the Random Forest Model was selected as being the final model. This model is described in the "Random Forest Final Model" section of this report.

VI. USING THE MODEL EXPLANATION OF CODE

i. Loading the Data

Because interacting with the models can be done in a variety of ways, prediction from file, prediction for single input sentence, prediction from array of sentences, loading the data to be utilized (either for training or for prediction) was made as versatile as possible.

i.1 Loading from File

The most common way to load data for training or batch prediction would be loading data from a file. The code provides a "loadData" function which takes a filepath as its parameter. This filepath must lead to a csv file just like the one provided by the challenge. I.E. the data in the csv file should be in the format:

```
response_id, class, response_text, , , , ,
```

note: `class` may or may not be present depending on `if` the data is labeled (train & dev) or not (test)

The function call is as follows:

```
loadedFile =
    loadData('deepnlp/Sheet_1.csv')
```

This function will read in the file line by line, create a pandas dataframe and apply a sentence cleaning method on each sentence.

The function will return a dataframe containing all the information from the csv file as well as a "response text array" column which contains the treated text, ready to be used by the models.

i.2 Loading single sentence

If a single sentence needs to be predicted, the function "convertSentence" will take as its argument a sentence and return the appropriate pandas dataframe which can be loaded into the models for prediction. This function applies the same cleaning function as when loading from file except only on a single sentence.

example:

```
convertedSentence = convertSentence('The  
Sentence to prepare')
```

i.3 Loading array of sentences

If the sentences needing conversion are in an array (not a csv file), then the function "convertSentences" will take the array of sentences and convert it into the pandas data frame after making each sentence go through the cleaning function.

example:

```
convertedSentences =  
    convertSentences(['This is an example  
of multiple sentences.', 'Where each  
sentence is an index in the array.'])
```

ii. Loading Model from File

The models are saved in very specific ways and each type of model needs to be loaded appropriately.

ii.1 RNN Model

Because the RNN model utilizes dynet as its core architecture, the loading procedure requires more work than just reading a pickle file. However, the RNN model class created will take a filepath in its init function which will let the model be loaded from file.

example:

```
curRNN =  
RNNmodel(model_filename='RNNFinalModel/  
bestRNN.model')
```

ii.2 Random Forest Model

The random forest model is saved to a pickle file using the standard pickle library found in python. There are two ways to load the random forest models.

The first is to use pickle directly:

```
model =  
    pickle.load(open('RandomForestFinalModel/  
bestRandomForest.model', 'rb'))
```

The second is to use the function provided in main.py:

```
curRFM = loadRandomForestModel(  
    'RandomForestFinalModel/  
bestRandomForest.model')
```

iii. Saving Models

Both the RNN and Random Forest Models classes have code to save the models. In both cases, simply calling

```
model.saveModel(filepath)
```

will save the model to the filepath.

For the RNN, it is highly recommended to save the model to a subfolder as there are multiple files generated by the save operation.

iv. Training

Both the RNN and Random Forest Models have been configured so that training is done using the same function.

In both cases, the training will occur when

```
model.train(train_data)
```

is called where train data is the data loaded using the previously mentioned ways of loading data.

Note that for the RNN, the train function can take an additional max epochs parameter which specifies how many epochs the model will train for.

v. Predicting

Generating predictions is very similar in both models. After loading the file, sentence or array of sentences as mentioned above, simply call:

```
model.predict(loadedSentences)
```

This function will return an array of predicted value (flagged or not flagged) for each of the texts.

vi. Generating Statistics

To generate statistics on the dev data (labeled testing data that was put aside before training), the function

```
model.computeDevAcc(dev_data)
```

is called.

This function returns the predictions as well as accuracy, true positive count, true negative count, false positive count and false negative count.

The function call takes and optional boolean `printStats` parameter which will dictate whether the stats argument printed to the console when run.

[2] Baoxun Xu, Xiufeng Guo, Yunming Ye, Jiefeng Cheng

<https://pdfs.semanticscholar.org/9b2f/84d85e5b6979bf3>

[3] DyNet Tutorial

http://dynet.readthedocs.io/en/latest/tutorials_notebook.html

[4] Tal Baumel

<https://talbaumel.github.io/blog/rnn%20batch>

[5] Pengfei Liu, Xipeng Qiu, Xuanjing Huang

<https://www.ijcai.org/Proceedings/16/Papers/408.pdf>

The above tutorials [DyNet Tutorial and Tal Baumel] were used as help to construct the dynet model

VII. NOTES

1. While model accuracy is generally very important, in this case the more important metric would be true positive rate as we care more about getting people the help they need (and so getting a high true positive rate is important) and if some people who don't necessarily need help are referred to specialists it is less crucial than if people who need help don't. That is why emphasis was put on true positive rate in the experimentations.

REFERENCES

[1] Ji Young Lee, Franck Dernoncourt

<https://arxiv.org/pdf/1603.03827.pdf>