

```
In [1]: 1 #import important libraries
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import seaborn as sns
6
7 # import libraries for EDA and preprocessing
8 from datetime import datetime
9 import nltk
10 nltk.download('stopwords')
11 from nltk.corpus import stopwords
12 nltk.download('wordnet')
13 from nltk.stem import WordNetLemmatizer
14 from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
15 %matplotlib inline
16
17 # Train test split
18 from sklearn.model_selection import train_test_split
19
20 # Text pre-processing
21 import tensorflow as tf
22 from sklearn import feature_extraction, linear_model, model_selection, preprocessing
23 from tensorflow.keras.preprocessing.text import Tokenizer
24 from tensorflow.keras.preprocessing.sequence import pad_sequences
25 from tensorflow.keras.callbacks import EarlyStopping
26 from keras.callbacks import ModelCheckpoint
27
28 # Modeling
29 from tensorflow.keras.models import Sequential
30 from tensorflow.keras.layers import LSTM, GRU, Dense, Embedding, Dropout, GlobalAveragePooling1D, Flatten, SpatialD
31
32 # Evaluating
33 from sklearn.metrics import roc_curve, auc, roc_auc_score
34 from sklearn.metrics import classification_report, confusion_matrix
35
36 import os
37 os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'
```

```
[nltk_data] Downloading package stopwords to /usr/share/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /usr/share/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

Step 1: Brief description of the problem and data

1.1 Problem

Twitter is one of social media that has become an important communication channel in different situations, for instance, in times of emergency. The smartphones enable people to announce an emergency they see in real-time. Because of that, there is an challenge that how to recognize whether a tweet text is talking about a real disaster or uses those keywords as a metaphor, which can lead to huge mislabeling of tweets. Hence, this project aims on using Natural Language Processing (NLP) and classification models to distinguish between real and fake disaster tweets. NLP is the ability of a computer program to understand human language as it is spoken and written, referred to as natural language, uses artificial intelligence to take real-world input, process it, and make sense of it in a way a computer can understand.

So, to do this, first, we will inspect, visualize, clean and vectorize the data by Count Vectorizer, then split train data into train_df (85%) and valid_df (15%) and train three models:

- (1) Long Short Term Memory (LSTM)
- (2) Bidirectional Long Short Term Memory (Bi-LSTM)
- (3) Gated Recurrent Unit (GRU)

Then, I will compare these three deep learning models by validation accuracy score and tune hyperparameter (dropout) to get the best model and use this best model for predicting test data and print out the submission file.

Reference Source:

- (1) <https://www.kaggle.com/code/philculliton/nlp-getting-started-tutorial/notebook> (<https://www.kaggle.com/code/philculliton/nlp-getting-started-tutorial/notebook>)
- (2) <https://medium.com/mlearning-ai/the-classification-of-text-messages-using-lstm-bi-lstm-and-gru-f79b207f90ad> (<https://medium.com/mlearning-ai/the-classification-of-text-messages-using-lstm-bi-lstm-and-gru-f79b207f90ad>)

1.2 Data

In this project, I use data from Kaggle, were downloaded from the link:

<https://www.kaggle.com/competitions/nlp-getting-started/data> (<https://www.kaggle.com/competitions/nlp-getting-started/data>)

There are two data from this resource, included train and test data. Train data has 7613 observations and 5 columns included : id, keyword, location, text and target. Test data has 3263 observations and 4 columns included: id, keyword, location and text.

```
In [2]: 1 # read train data
2 df = pd.read_csv('../input/nlp-getting-started/train.csv')
3 #df = pd.read_csv('train.csv')
4
5 # take a look at some rows of train data
6 df.head()
7
```

```
Out[2]:
```

	id	keyword	location	text	target
0	1	NaN	NaN	Our Deeds are the Reason of this #earthquake M...	1
1	4	NaN	NaN	Forest fire near La Ronge Sask. Canada	1
2	5	NaN	NaN	All residents asked to 'shelter in place' are ...	1
3	6	NaN	NaN	13,000 people receive #wildfires evacuation or...	1
4	7	NaN	NaN	Just got sent this photo from Ruby #Alaska as ...	1

```
In [3]: 1 # read test data
2 test = pd.read_csv('../input/nlp-getting-started/test.csv')
3 #test = pd.read_csv('test.csv')
4
5 # take a look at some rows of test data
6 test.head()
```

```
Out[3]:
```

	id	keyword	location	text
0	0	NaN	NaN	Just happened a terrible car crash
1	2	NaN	NaN	Heard about #earthquake is different cities, s...
2	3	NaN	NaN	there is a forest fire at spot pond, geese are...
3	9	NaN	NaN	Apocalypse lighting. #Spokane #wildfires
4	11	NaN	NaN	Typhoon Soudelor kills 28 in China and Taiwan

```
In [4]: 1 # view some rows of the sample submission
2 sample_submission = pd.read_csv('../input/nlp-getting-started/sample_submission.csv')
3 #sample_submission = pd.read_csv('sample_submission.csv')
4 sample_submission.head()
5
```

```
Out[4]:
```

	id	target
0	0	0
1	2	0
2	3	0
3	9	0
4	11	0

```
In [5]: 1 # the shape of train data
2 df.shape
3
```

```
Out[5]: (7613, 5)
```

```
In [6]: 1 # the shape of test data
2 test.shape
3
```

```
Out[6]: (3263, 4)
```

Step 2: Exploratory Data Analysis (EDA) - Inspect, Visualize, and Clean the Data

3.1 Inspect the data

```
In [7]: 1 # look at the first example of a not disaster tweet
2 fake_disaster_tweet = df[df["target"] == 0]
3 fake_disaster_tweet["text"].values[0]
4
```

```
Out[7]: "What's up man?"
```

```
In [8]: 1 # look at the first example of a disaster tweet
2 real_disaster_tweet = df[df["target"] == 1]
3 real_disaster_tweet["text"].values[0]
4
```

Out[8]: 'Our Deeds are the Reason of this #earthquake May ALLAH Forgive us all'

```
In [9]: 1 # get a quick description of the data
2 df.describe()
3
```

Out[9]:

	id	target
count	7613.000000	7613.000000
mean	5441.934848	0.42966
std	3137.116090	0.49506
min	1.000000	0.00000
25%	2734.000000	0.00000
50%	5408.000000	0.00000
75%	8146.000000	1.00000
max	10873.000000	1.00000

```
In [10]: 1 # check null values in data
2 df.isnull().sum()
3
```

Out[10]:

id	0
keyword	61
location	2533
text	0
target	0

dtype: int64

```
In [11]: 1 # check for duplicate articles
2 df.duplicated(keep=False).sum()
3
```

Out[11]: 0

```
In [12]: 1 # the structure of data also tells us the number of rows (observations) and columns (variables)
2 df.info()
3
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7613 entries, 0 to 7612
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0    id          7613 non-null   int64
1    keyword     7552 non-null   object
2    location    5080 non-null   object
3    text        7613 non-null   object
4    target      7613 non-null   int64
dtypes: int64(2), object(3)
memory usage: 297.5+ KB
```

```
In [13]: 1 # get the label of data
2 df['target'].unique()
3
```

Out[13]: array([1, 0])

```
In [14]: 1 # the structure of data also tells us the number of rows (observations) and columns (variables)
2 test.info()
3
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3263 entries, 0 to 3262
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  -
0    id          3263 non-null   int64
1    keyword     3237 non-null   object
2    location    2158 non-null   object
3    text        3263 non-null   object
dtypes: int64(1), object(3)
memory usage: 102.1+ KB
```

```
In [15]: 1 # check null values in test data
          2 test.isnull().sum()
          3
```

```
Out[15]: id          0
          keyword     26
          location    1105
          text        0
          dtype: int64
```

```
In [16]: 1 # check for duplicate observations in test data
          2 test.duplicated(keep=False).sum()
          3
```

```
Out[16]: 0
```

From the output above, we can summarize that:

- There are 7613 observations and 5 columns in train data.
- There is no missing values in "id", "text" and "target" column in train data.
- There is no missing values in "id" and "text" column in test data.
- There is no duplicated observations in both train and test data.
- There are 2 targets: 0 (fake disaster tweet) and 1 (real disaster tweet).

3.2 Visualize the data

Next, let's calculate and visualize the count and the proportion of each target.

```
In [17]: 1 # calculate the count of each target
          2 df['target'].value_counts()
          3
```

```
Out[17]: 0    4342
          1    3271
          Name: target, dtype: int64
```

```
In [18]: 1 # calculate the proportion of each label
          2 df['target'].value_counts()/len(df)*100
          3
```

```
Out[18]: 0    57.034021
          1    42.965979
          Name: target, dtype: float64
```

```
In [19]: 1 # plot the count of each label
2 fig, ax = plt.subplots(figsize=(6,6))
3 sns.countplot(data=df, y='target', ax=ax).set(title='\nFigure 1. The Count of Each Target\n')
4
5 # plot the proportion of each category
6 labels = df['target'].unique().tolist()
7 counts = df['target'].value_counts()
8 sizes = [counts[v] for v in labels]
9 fig1, ax1 = plt.subplots()
10 ax1.pie(sizes, labels=labels, autopct='%0.2f%%')
11 ax1.axis('equal')
12 plt.title("\nFigure 2. The Proportion of Each Target\n")
13 plt.tight_layout()
14 plt.show()
15
```

Figure 1. The Count of Each Target

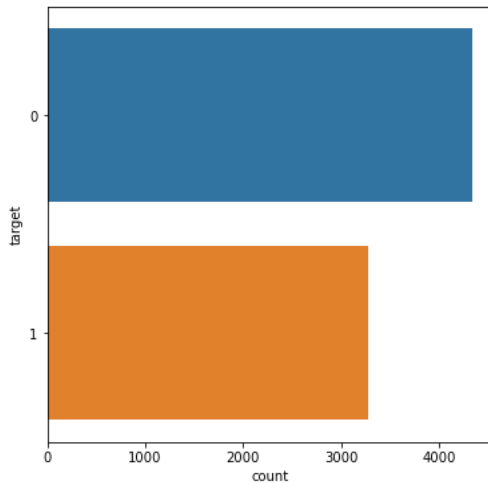


Figure 2. The Proportion of Each Target

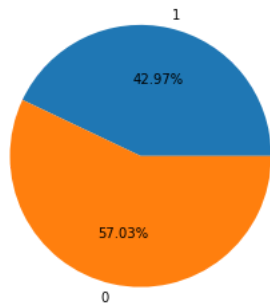


Figure 1 shows the count of each label and figure 2 shows the proportions of each target. Looking at these two figures, we can see that in overall, the number of observations for "1" target is larger than the number of observations for "0" target. So, to solve this problem, I will balance the data by downsampling the fake disaster tweets before build models since if one or two categories was severely underrepresented or, in contrast, overrepresentative in the train data, then it may cause our model to be biased and/or perform poorly on some or all of the test data.

3.3 Clean the data/ Data Preprocessing

3.3.1 Clean the data

To clean the data for training models, some works has to be done such as:

- balance the data by downsampling the fake disaster tweets
- drop unused columns in train data included: id, keyword and location.

To preprocess our text simply means to bring our text into a form that is predictable and analyzable for our task. So, what I am going to do is:

- (1) lowercasing all our text data
- (2) remove punctuation
- (3) remove stop words: stop words are a set of commonly used words in a language. Examples of stop words in English are "a", "the", "is", "are" and etc. The intuition behind using stop words is that, by removing low information words from text, we can focus on the important words instead.

(4) lemmatization: lemmatization usually refers to doing things properly with the use of a vocabulary and morphological analysis of words, normally aiming to remove inflectional endings only and to return the base or dictionary form of a word, which is known as the lemma. For example, runs, running, ran are all forms of the word run, therefore run is the lemma of all these words.

Since I'm planning to redo these cleaning steps for a test data without target as well, thus for convenience, I will create a clean_text function for this data and reuse it for cleaning untargeted test data later.

```
In [20]: 1 # downsample the fake disaster tweet
2 fake_disaster = fake_disaster_tweet.sample(n = len(real_disaster_tweet), random_state = 44)
3 df_news = pd.concat([fake_disaster, real_disaster_tweet], axis=0).reset_index(drop=True)
4 df_news["target"].value_counts()
5
```

```
Out[20]: 0    3271
1    3271
Name: target, dtype: int64
```

```
In [21]: 1 # drop id, keyword and location columns
2 df_news = df_news.loc[:, ["text", "target"]]
3
4 # view some sample rows of df_news
5 df_news.sample(10)
```

```
Out[21]:
```

	text	target
284	Detonation fashionable mountaineering electron...	0
6291	#Earthquake #Sismo M 1.9 - 5km S of Volcano Ha...	1
6190	Trauma injuries involving kids and sport usual...	1
1	@DyannBridges @yeshayad Check out this #rockin...	0
2135	@Eric_Tsunami worry about yourself	0
966	collapsed the moment i got home last night lol	0
3522	#anthrax #bioterrorism CDC To Carry Out Extens...	1
2138	I liked a @YouTube video http://t.co/FNpDJwVw1...	0
1479	someone's gonna get screamed at for getting th...	0
6490	Related News: \n\nPlane Wreckage Found Is Part...	1

```
In [22]: 1 def clean_text(data, text):
2     # lowercasing all text data
3     data[text] = data[text].str.lower()
4     # remove punctuation
5     data[text] = data[text].str.replace('[^\w\s]', '', regex=True)
6     # remove stop words
7     stop_words = stopwords.words('english')
8     data[text] = data[text].apply(lambda x: ' '.join([word for word in x.split() if word not in (stop_words)]))
9     # lemmatization
10    lemmatizer = WordNetLemmatizer()
11    data[text] = data[text].apply(lambda x: ' '.join([lemmatizer.lemmatize(word) for word in x.split()]))
12    return
13
```

```
In [23]: 1 # clean news data
2 clean_text(df_news, "text")
3
4 # view text in a row after cleaning all text data
5 df_news["text"][1]
6
```

```
Out[23]: 'dyannbridges yeshayad check rockin preview claytonbryant danger zone coming soon httpstcoipgmf4ttdx artistsunited'
```

```
In [24]: 1 # calculate the count of word per observation
2 df_news["Word_Count"] = df_news["text"].apply(lambda x: len(x.split()))
3
```

```
In [25]: 1 # view some first rows of news data
2 df_news.head()
3
```

```
Out[25]:
```

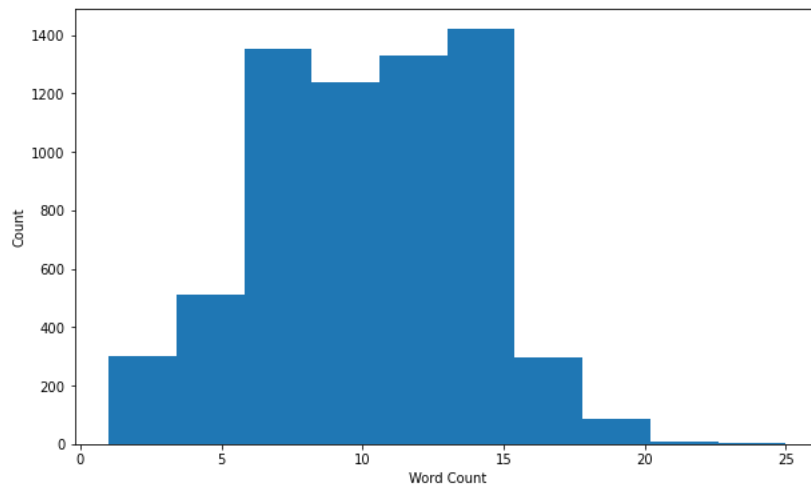
	text	target	Word_Count
0	bcfticketlady mr_aamir_javid see inundated a...	0	12
1	dyannbridges yeshayad check rockin preview cla...	0	12
2	hot funtenna hijacking computer send data soun...	0	14
3	nasasolarsystem jupiter great red spot violent...	0	12
4	learn gained access secret top earner amp used...	0	14

```
In [26]: 1 # The average count of word per observation
2 print("The average count of word per observation", round(np.mean(df_news.Word_Count)))
3
4 # The maximum count of word per observation
5 print("The maximum count of word per observation", round(np.max(df_news.Word_Count)))
6
7 # The minimum count of word per observation
8 print("The minimum count of word per observation", round(np.min(df_news.Word_Count)))
9
```

The average count of word per observation 10
The maximum count of word per observation 25
The minimum count of word per observation 1

```
In [27]: 1 # plot the count of word per observation
2 fig, ax = plt.subplots(figsize=(10,6))
3 df_news['Word_Count'].plot(kind='hist')
4 plt.xlabel("Word Count")
5 plt.xticks(rotation=360)
6 plt.ylabel("Count")
7 plt.title("Figure 3. The count of words per observation\n")
8 plt.show()
9
```

Figure 3. The count of words per observation

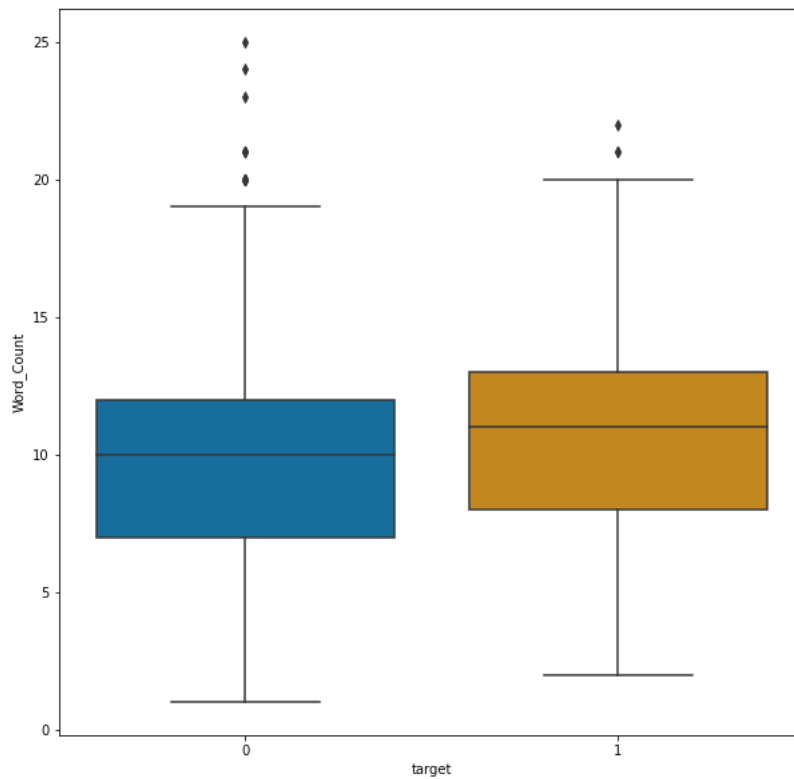


```

In [28]: 1 # visualize the count of words per label
2 fig, ax = plt.subplots(figsize=(10, 10))
3 sns.boxplot(data = df_news, x = 'target', y = 'Word_Count', palette = 'colorblind')
4           ).set(title = 'Figure 4. The count of words per target\n')
5 plt.show()
6

```

Figure 4. The count of words per target

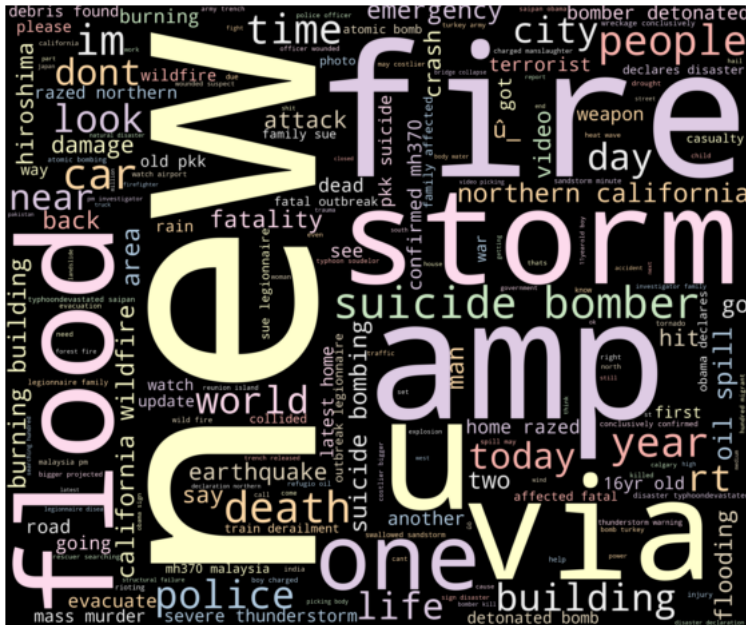


Looking at figure 4, we observe that the mean and the variability of word count of fake and real disaster tweets are not different much.


```

1 # Generate a word cloud image for real disaster tweets
2 real_disaster_text = ' '.join(df_news[df_news["target"] == 1]["text"])
3 real_disaster_text_cloud = WordCloud(width = 3000,
4                                     height = 2500,
5                                     stopwords=STOPWORDS,
6                                     background_color = "black",
7                                     colormap='Pastell').generate(real_disaster_text)
8 plt.figure(figsize=(10,10))
9 plt.imshow(real_disaster_text_cloud, interpolation='bilinear')
10 plt.axis('off') # turn off axis
11 plt.show()

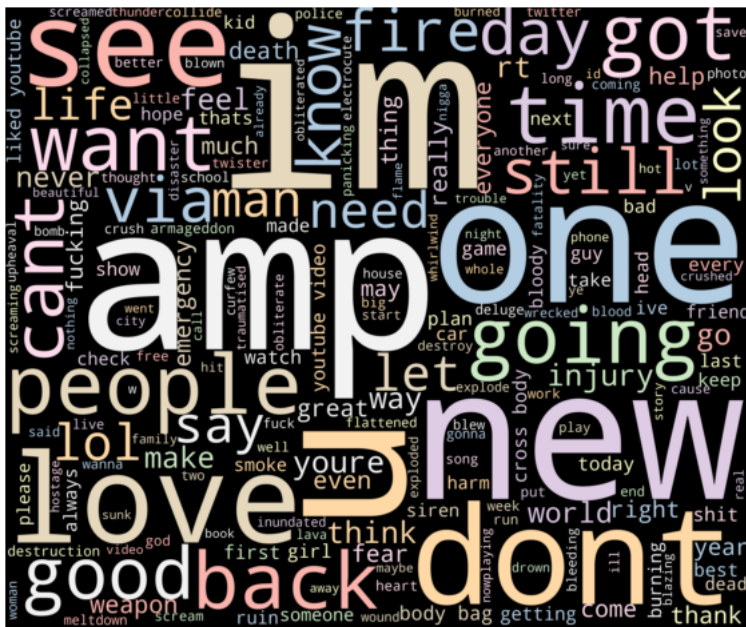
```



```

1 # Generate a word cloud image for fake disaster tweets
2 fake_disaster_text = ' '.join(df_news[df_news["target"] == 0]["text"])
3 fake_disaster_text_cloud = WordCloud(width=3000,
4                                       height=2500,
5                                       stopwords=STOPWORDS,
6                                       background_color="black",
7                                       colormap='Pastell').generate(fake_disaster_text)
8 plt.figure(figsize=(10,10))
9 plt.imshow(fake_disaster_text_cloud, interpolation='bilinear')
10 plt.axis('off') # turn off axis
11 plt.show()
12

```



3.3.2 Data Preprocessing

Text data requires a special approach to machine learning. This is because text data can have hundreds of thousands of dimensions (words and phrases) but tends to be very sparse. Machines, unlike humans, cannot understand the raw text. Machines can only see numbers. Particularly, statistical techniques such as machine learning can only deal with numbers. Therefore, we need to convert the text data into numerical representation, so the model will understand it.

Different approaches exist to convert text into the corresponding numerical form. In this case I will use Count Vectorizer to transform a given text into a vector on the basis of the frequency (count) of each word that occurs in the entire text. I choose this method because:

1. It is one of the simplest ways of doing text vectorization.
2. It creates a document term matrix, which is a set of dummy variables that indicates if a particular word appears in the document.
3. Count vectorizer will fit and learn the word vocabulary and try to create a document term matrix in which the individual cells denote the frequency of that word in a particular document, which is also known as term frequency, and the columns are dedicated to each word in the corpus.

```
In [31]: 1 # drop word_count column
2 df_news = df_news.drop(columns='Word_Count')
3
```

```
In [32]: 1 X = df_news["text"]
2 Y = df_news["target"]
```

```
In [33]: 1 count_vectorizer = feature_extraction.text.CountVectorizer()
2
3 ## let's get counts for the first 5 tweets in the data
4 example_train_vectors = count_vectorizer.fit_transform(X[0:5])
```

```
In [34]: 1 ## we use .todense() here because these vectors are "sparse" (only non-zero elements are kept to save space)
2 print(example_train_vectors[0].todense().shape)
3 print(example_train_vectors[0].todense())

(1, 64)
[[0 0 0 1 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 1
  0 1 0 0 1 0 0 0 0 0 0 0 1 0 1 0 0 0 0 0 0 0 0 0 1 0 1 0 0]]
```

This example shows us that:

- There are 42 unique words (tokens) in the first five tweets.
- The first tweet contains only some of those unique tokens - all of the non-zero counts above are the tokens that do exist in the first tweet.

```
In [35]: 1 # create vector for train set
2 train_vectors = count_vectorizer.fit_transform(X)
3
```

```
In [36]: 1 # create a list to store validation accuracy score
2 valid_auc_score = []
```

Split data

After cleaning and vectorizing data by CountVectorizer, to prepare for building and training models, I'll split 15% of the data into validation set. Noted that, I'll use sklearn train_test_split to split the data, with default shuffle = True and stratify=target, means this method will split our data into random train and test subsets and have the same proportion of target in df_news.

```
In [37]: 1 # shuffle and split the data into train and test set
2 x_train, x_valid, y_train, y_valid = train_test_split(train_vectors, Y, test_size=0.15,
3                                                     shuffle=True,
4                                                     random_state = 42,
5                                                     stratify=df_news.target)
6
7 ## get shape of train and validation dataset after splitting
8 print(x_train.shape, y_train.shape)
9 print(x_valid.shape, y_valid.shape)
10

(5560, 18980) (5560,)
(982, 18980) (982,)
```

```
In [38]: 1 # view train data
2 print('Training set:')
3 x_train = x_train.toarray()
4 x_train
5
```

Training set:

```
Out[38]: array([[0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 ...,
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0]])
```

```
In [39]: 1 # view validation data
2 print('Validation set:')
3 x_valid = x_valid.toarray()
4 x_valid
5
```

Validation set:

```
Out[39]: array([[0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 ...,
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0],
 [0, 0, 0, ..., 0, 0, 0]])
```

```
In [40]: 1 # check target value count in train and validation set
2 print(y_train.value_counts())
3 print(y_valid.value_counts())
4
```

```
1    2780
0    2780
Name: target, dtype: int64
0     491
1     491
Name: target, dtype: int64
```

Step 4: Building and training models

4.1 Long Short Term Memory (LSTM)

Long Short Term Memory (LSTM) was designed to overcome the problems of simple Recurrent Neural Network (RNN) by allowing the network to store data in a sort of memory that it can access at a later times. The key of the LSTM model is the cell state. The cell state is updated twice with few computations that resulting stabilize gradients. It has also a hidden state that acts like a short term memory.

In LSTM there are Forget Gate, Input Gate and Output Gate.

- (1) The first step is to decide what information we're going to throw away from the cell state. This decision is made by a sigmoid layer called the "Forget Gate" layer.
- (2) The second step is to decide what new information that we're going to store in the cell state. This has two parts. First, a sigmoid layer called the "Input Gate" layer decides which values we'll update. Next, a tanh layer which creates a vector of new candidate values that could be added to the state.
- (3) Finally, we need to decide what we are going to output. This output will be based on our cell state, but will be a filtered version. First, we run a sigmoid layer which decides what parts of the cell state we're going to output. Then, we put the cell state through tanh (to push the values to be between -1 and 1) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided.

We use the binary_crossentropy as a loss function because the output of the model is binary and for the optimizer, we use adam which makes use of momentum to avoid local minima.

- epoch : number of times the learning algorithm will work through the entire training data.
- callbacks : to pass the early stopping parameter. EarlyStopping(monitor='val_loss', patience=2) was used to define that we want to monitor the validation loss and if the validation loss is not improved after 2 epochs, then the model training will stop. This technique helps to avoid overfitting problem.
- verbose : 2 , it will show us loss and accuracy on each epoch.

```
In [41]: 1 # Define the LSTM model architecture
2
3 # Define parameter
4 n_lstm = 200
5 embedding_dim = 128
6 max_len = train_vectors.shape[1]
7 drop_lstm = 0.2
8 vocab_size = len(set(" ".join(X).split()))
9 print(vocab_size)
10
```

19017

```
In [42]: 1 # Define LSTM Model
2 model1 = Sequential()
3 model1.add(Embedding(vocab_size, embedding_dim, input_length=max_len))
4 model1.add(SpatialDropout1D(drop_lstm))
5 model1.add(LSTM(n_lstm, return_sequences=False))
6 model1.add(Dropout(drop_lstm))
7 model1.add(Dense(1, activation='sigmoid'))
8
9 # summary model1
10 model1.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
embedding (Embedding)	(None, 18980, 128)	2434176

spatial_dropout1d (SpatialDr	(None, 18980, 128)	0

lstm (LSTM)	(None, 200)	263200

dropout (Dropout)	(None, 200)	0

dense (Dense)	(None, 1)	201
=====		
Total params: 2,697,577		
Trainable params: 2,697,577		
Non-trainable params: 0		

```
In [43]: 1 # compile the model
2 model1.compile(loss = 'binary_crossentropy',
3               optimizer = 'adam',
4               metrics = ['accuracy', tf.keras.metrics.AUC()])
```

```
In [44]: 1 num_epochs = 10
2 early_stop = EarlyStopping(monitor='val_loss', patience=2)
3 mp = ModelCheckpoint(filepath='model1_cp', monitor='val_loss', save_best_only=True)
4 history = model1.fit(x_train,
5                     y_train,
6                     epochs=num_epochs,
7                     validation_data=(x_valid, y_valid),
8                     callbacks=[early_stop, mp],
9                     verbose=2)
```

Epoch 1/10

174/174 - 200s - loss: 0.6943 - accuracy: 0.4950 - auc: 0.4947 - val_loss: 0.6940 - val_accuracy: 0.5000 - val_auc: 0.5000

Epoch 2/10

174/174 - 195s - loss: 0.6936 - accuracy: 0.5018 - auc: 0.5029 - val_loss: 0.6934 - val_accuracy: 0.5000 - val_auc: 0.5000

Epoch 3/10

174/174 - 195s - loss: 0.6936 - accuracy: 0.4960 - auc: 0.4931 - val_loss: 0.6934 - val_accuracy: 0.5000 - val_auc: 0.5000

Epoch 4/10

174/174 - 195s - loss: 0.6936 - accuracy: 0.4980 - auc: 0.4946 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000

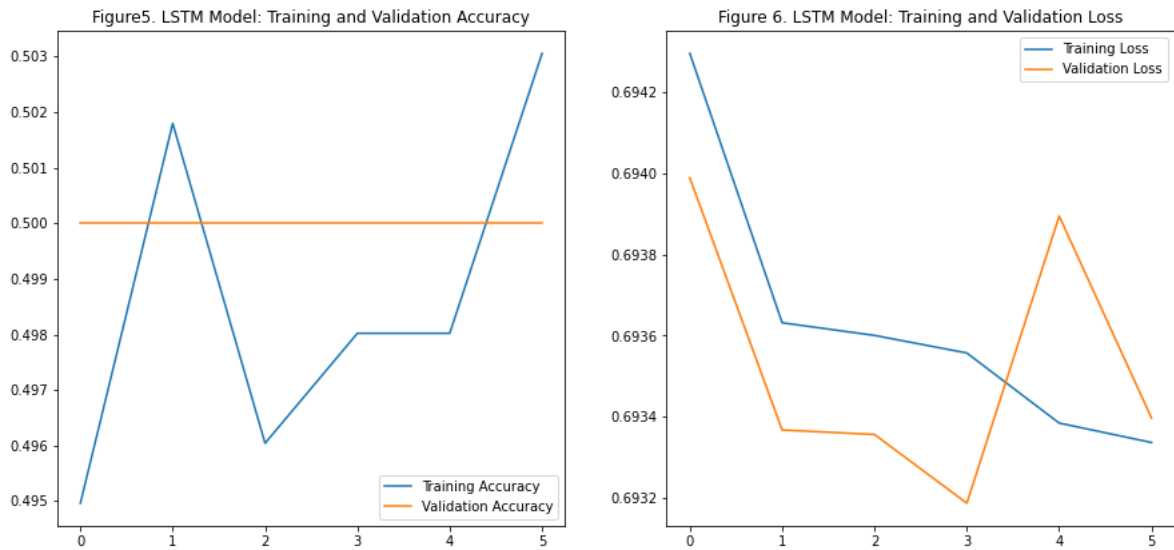
Epoch 5/10

174/174 - 195s - loss: 0.6934 - accuracy: 0.4980 - auc: 0.5000 - val_loss: 0.6939 - val_accuracy: 0.5000 - val_auc: 0.5010

Epoch 6/10

174/174 - 195s - loss: 0.6933 - accuracy: 0.5031 - auc: 0.5051 - val_loss: 0.6934 - val_accuracy: 0.5000 - val_auc: 0.5000

```
In [45]: 1 # plot the graph of accuracy
2 acc = history.history['accuracy']
3 val_acc = history.history['val_accuracy']
4 loss = history.history['loss']
5 val_loss = history.history['val_loss']
6
7 epochs_range = range(6)
8
9 plt.figure(figsize=(15, 15))
10 plt.subplot(2, 2, 1)
11 plt.plot(epochs_range, acc, label='Training Accuracy')
12 plt.plot(epochs_range, val_acc, label='Validation Accuracy')
13 plt.legend(loc='lower right')
14 plt.title('Figure5. LSTM Model: Training and Validation Accuracy')
15
16 # plot the graph of loss
17 plt.subplot(2, 2, 2)
18 plt.plot(epochs_range, loss, label='Training Loss')
19 plt.plot(epochs_range, val_loss, label='Validation Loss')
20 plt.legend(loc='upper right')
21 plt.title('Figure 6. LSTM Model: Training and Validation Loss')
22 plt.show()
```



4.2 Bidirectional Long Short Term Memory (Bi-LSTM)

A Bidirectional LSTM, or BiLSTM, is a sequence processing model that consists of two LSTMs: one taking the input in a forward direction, and the other in a backwards direction. BiLSTMs effectively increase the amount of information available to the network, improving the context available to the algorithm (e.g. knowing what words immediately follow and precede a word in a sentence). Unlike standard LSTM, the input flows of Bi-LSTM in both directions, and it's capable of utilizing information from both sides. It's also a powerful tool for modeling the sequential dependencies between words and phrases in both directions of the sequence.

BiLSTM adds one more LSTM layer, which reverses the direction of information flow. Briefly, it means that the input sequence flows backward in the additional LSTM layer. Then we combine the outputs from both LSTM layers in several ways, such as average, sum, multiplication, or concatenation.

```
In [46]: 1 # define Bi_LSTM model
2 model2 = Sequential()
3 model2.add(Embedding(vocab_size,
4                     embedding_dim,
5                     input_length = max_len))
6 model2.add(Bidirectional(LSTM(n_lstm,
7                             return_sequences = False)))
8 model2.add(Dropout(drop_lstm))
9 model2.add(Dense(1, activation='sigmoid'))
10
11 # summary model2
12 model2.summary()
```

Model: "sequential_1"

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 18980, 128)	2434176
bidirectional (Bidirectional)	(None, 400)	526400
dropout_1 (Dropout)	(None, 400)	0
dense_1 (Dense)	(None, 1)	401
Total params: 2,960,977		
Trainable params: 2,960,977		
Non-trainable params: 0		

```
In [47]: 1 # compile model2
2 model2.compile(loss = 'binary_crossentropy',
3               optimizer = 'adam',
4               metrics=['accuracy', tf.keras.metrics.AUC()])
```

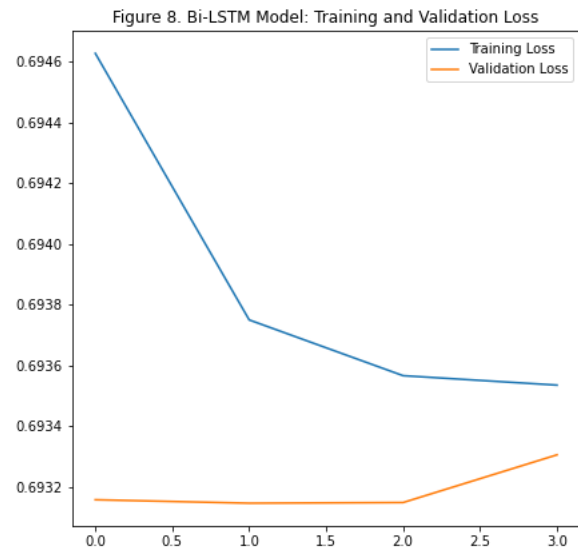
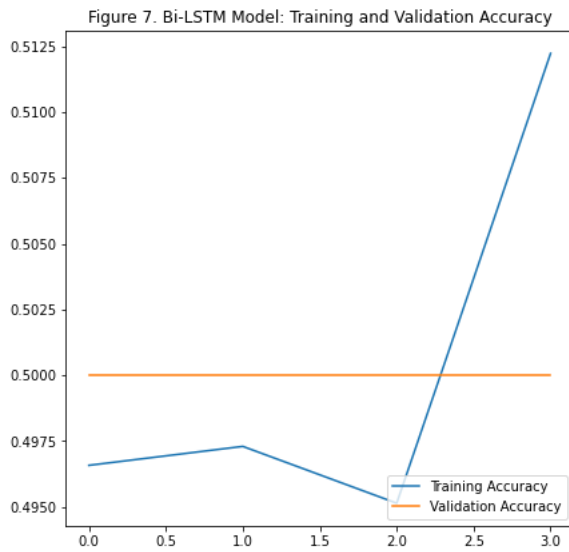
```
In [48]: 1 # train model2
2 num_epochs = 10
3 early_stop = EarlyStopping(monitor = 'val_loss',
4                             patience = 2)
5 mp = ModelCheckpoint(filepath='model2_cp', monitor='val_loss', save_best_only=True)
6 history2 = model2.fit(x_train,
7                       y_train,
8                       epochs = num_epochs,
9                       validation_data = (x_valid, y_valid),
10                      callbacks = [early_stop, mp],
11                      verbose = 2)
```

Epoch 1/10
174/174 - 389s - loss: 0.6946 - accuracy: 0.4966 - auc_1: 0.4957 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc_1: 0.5000
Epoch 2/10
174/174 - 385s - loss: 0.6937 - accuracy: 0.4973 - auc_1: 0.4892 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc_1: 0.5000
Epoch 3/10
174/174 - 385s - loss: 0.6936 - accuracy: 0.4951 - auc_1: 0.4930 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc_1: 0.5000
Epoch 4/10
174/174 - 385s - loss: 0.6935 - accuracy: 0.5122 - auc_1: 0.5093 - val_loss: 0.6933 - val_accuracy: 0.5000 - val_auc_1: 0.5000

```

In [49]: 1 # plot the graph of accuracy
2 acc2 = history2.history['accuracy']
3 val_acc2 = history2.history['val_accuracy']
4 loss2 = history2.history['loss']
5 val_loss2 = history2.history['val_loss']
6
7 epochs_range = range(4)
8
9 plt.figure(figsize=(15, 15))
10 plt.subplot(2, 2, 1)
11 plt.plot(epochs_range, acc2, label='Training Accuracy')
12 plt.plot(epochs_range, val_acc2, label='Validation Accuracy')
13 plt.legend(loc='lower right')
14 plt.title('Figure 7. Bi-LSTM Model: Training and Validation Accuracy')
15
16 # plot the graph of loss
17 plt.subplot(2, 2, 2)
18 plt.plot(epochs_range, loss2, label='Training Loss')
19 plt.plot(epochs_range, val_loss2, label='Validation Loss')
20 plt.legend(loc='upper right')
21 plt.title('Figure 8. Bi-LSTM Model: Training and Validation Loss')
22 plt.show()

```



4.3 Gated Recurrent Unit (GRU)

A Gated Recurrent Unit, or GRU, is a type of recurrent neural network. It is similar to an LSTM, but only has two gates — a reset gate and an update gate and notably lacks an output gate. Fewer parameters means GRUs are generally easier/faster to train than their LSTM counterparts.

```
In [50]: 1 # define GRU model
2 model3 = Sequential()
3 model3.add(Embedding(vocab_size,
4                     embedding_dim,
5                     input_length = max_len))
6 model3.add(SpatialDropout1D(0.2))
7 model3.add(GRU(128, return_sequences = False))
8 model3.add(Dropout(0.2))
9 model3.add(Dense(1, activation = 'sigmoid'))
10
11 # summary model3
12 model3.summary()
13
```

Model: "sequential_2"

Layer (type)	Output Shape	Param #
embedding_2 (Embedding)	(None, 18980, 128)	2434176
spatial_dropout1d_1 (Spatial	(None, 18980, 128)	0
gru (GRU)	(None, 128)	99072
dropout_2 (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 1)	129
Total params: 2,533,377		
Trainable params: 2,533,377		
Non-trainable params: 0		

```
In [51]: 1 # compile model3
2 model3.compile(loss = 'binary_crossentropy',
3               optimizer = 'adam',
4               metrics=['accuracy', tf.keras.metrics.AUC()])
5
```

```
In [52]: 1 # train model3
2 num_epochs = 10
3 early_stop = EarlyStopping(monitor='val_loss', patience=2)
4 mp = ModelCheckpoint(filepath='model3_cp', monitor='val_loss', save_best_only=True)
5 history3 = model3.fit(x_train,
6                      y_train,
7                      epochs=num_epochs,
8                      validation_data=(x_valid, y_valid),
9                      callbacks=[early_stop, mp],
10                      verbose=2)
```

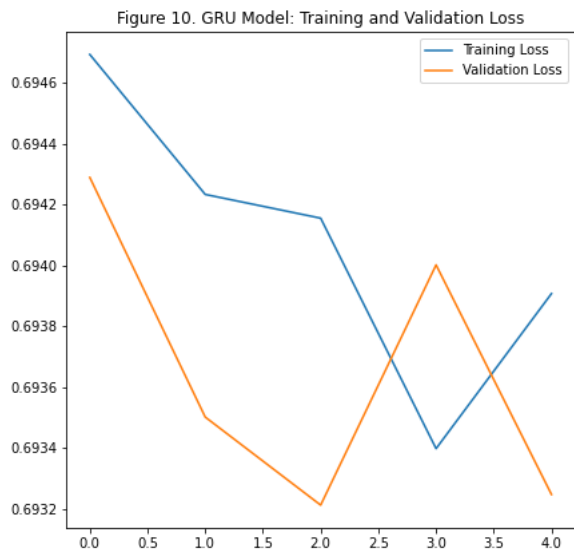
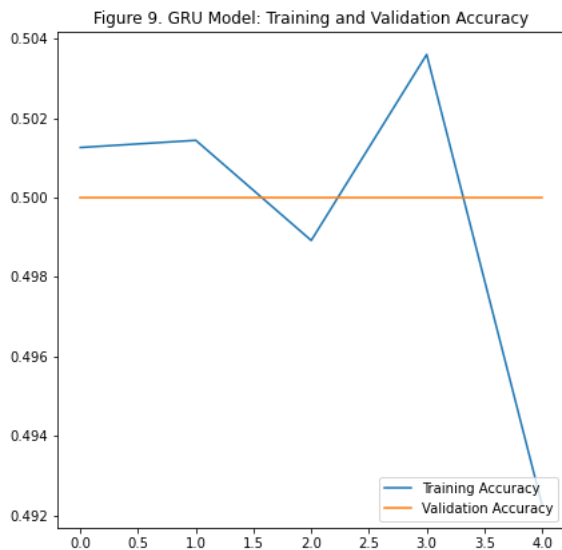
```
Epoch 1/10
174/174 - 134s - loss: 0.6947 - accuracy: 0.5013 - auc_2: 0.5011 - val_loss: 0.6943 - val_accuracy: 0.5000 - val_auc_
2: 0.5010
Epoch 2/10
174/174 - 132s - loss: 0.6942 - accuracy: 0.5014 - auc_2: 0.4944 - val_loss: 0.6935 - val_accuracy: 0.5000 - val_auc_
2: 0.5000
Epoch 3/10
174/174 - 132s - loss: 0.6942 - accuracy: 0.4989 - auc_2: 0.4953 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc_
2: 0.5000
Epoch 4/10
174/174 - 132s - loss: 0.6934 - accuracy: 0.5036 - auc_2: 0.5048 - val_loss: 0.6940 - val_accuracy: 0.5000 - val_auc_
2: 0.5000
Epoch 5/10
174/174 - 132s - loss: 0.6939 - accuracy: 0.4923 - auc_2: 0.4925 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc_
2: 0.5000
```



```

In [54]: 1 # plot the graph of accuracy
2 acc3 = history3.history['accuracy']
3 val_acc3 = history3.history['val_accuracy']
4 loss3 = history3.history['loss']
5 val_loss3 = history3.history['val_loss']
6
7 epochs_range = range(5)
8
9 plt.figure(figsize=(15, 15))
10 plt.subplot(2, 2, 1)
11 plt.plot(epochs_range, acc3, label='Training Accuracy')
12 plt.plot(epochs_range, val_acc3, label='Validation Accuracy')
13 plt.legend(loc='lower right')
14 plt.title('Figure 9. GRU Model: Training and Validation Accuracy')
15
16 # plot the graph of loss
17 plt.subplot(2, 2, 2)
18 plt.plot(epochs_range, loss3, label='Training Loss')
19 plt.plot(epochs_range, val_loss3, label='Validation Loss')
20 plt.legend(loc='upper right')
21 plt.title('Figure 10. GRU Model: Training and Validation Loss')
22 plt.show()
23

```



Step 4: Results and Analysis

4.1 Results

Model 1

```

In [57]: 1 score = model1.evaluate(x_valid, y_valid)

31/31 [=====] - 12s 382ms/step - loss: 0.6934 - accuracy: 0.5000 - auc: 0.5000

```

```

In [58]: 1 print('LSTM model loss:', score[0])
2 print('LSTM model accuracy:', score[1])

```

```

LSTM model loss: 0.6933974623680115
LSTM model accuracy: 0.5

```

```

In [59]: 1 # add validation accuracy score into list
2 v_auc_score1 = history.history["val_auc"]
3 valid_auc_score.append(v_auc_score1)
4
5 # best validation accuracy result
6 best_val_auc1 = max(v_auc_score1)
7 print("LSTM Best Validation AUC: ", best_val_auc1)
8

```

```

LSTM Best Validation AUC: 0.5010183453559875

```

```
In [61]: 1 # make predictions on the validation dataset
2 #load_model1 = keras.models.load_model('modell_cp')
3 y_pred1 = model1.predict(x_valid)
4 y_pred1 = np.where(y_pred1>0.5, 1, 0)
5
6 # print out classification report
7 print(classification_report(y_valid, y_pred1))
8
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	491
1	0.50	1.00	0.67	491
accuracy			0.50	982
macro avg	0.25	0.50	0.33	982
weighted avg	0.25	0.50	0.33	982

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

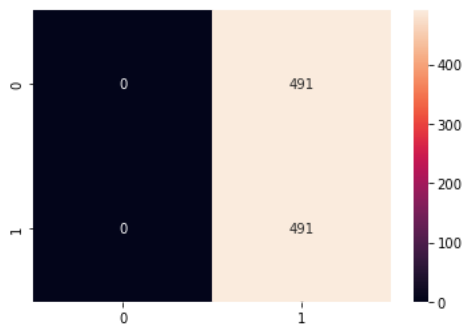
_warn_prf(average, modifier, msg_start, len(result))

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

```
In [62]: 1 # print out the confusion matrix
2 cm1 = confusion_matrix(y_valid, y_pred1)
3 sns.heatmap(cm1, annot=True, fmt=".0f")
4
```

Out[62]: <AxesSubplot:>



Model 2

```
In [63]: 1 score2 = model2.evaluate(x_valid, y_valid)
2 print('Bi-LSTM model loss:', score2[0])
3 print('Bi-LSTM model accuracy:', score2[1])
4
```

31/31 [=====] - 24s 776ms/step - loss: 0.6933 - accuracy: 0.5000 - auc_1: 0.5000
Bi-LSTM model loss: 0.6933056712150574
Bi-LSTM model accuracy: 0.5

```
In [69]: 1 # add validation accuracy score into list
2 v_auc_score2 = history2.history["val_auc_1"]
3 #v_auc_score2 = history2.history["val_auc"]
4 valid_auc_score.append(v_auc_score2)
5
6 # best validation accuracy result
7 best_val_auc2 = max(v_auc_score2)
8 print("Bi-LSTM Best Validation AUC: ", best_val_auc2)
9
```

Bi_LSTM Best Validation AUC: 0.5

```
In [65]: 1 # make predictions on the validation dataset
2 #load_model2 = keras.models.load_model('model2_cp')
3 y_pred2 = model2.predict(x_valid)
4 y_pred2 = np.where(y_pred2>0.5, 1, 0)
5
6 # print out classification report
7 print(classification_report(y_valid, y_pred2))
8
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	491
1	0.50	1.00	0.67	491
accuracy			0.50	982
macro avg	0.25	0.50	0.33	982
weighted avg	0.25	0.50	0.33	982

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

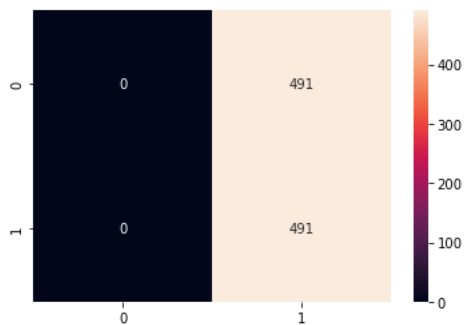
_warn_prf(average, modifier, msg_start, len(result))

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

```
In [80]: 1 # print out the confusion matrix
2 cm2 = confusion_matrix(y_valid, y_pred2)
3 sns.heatmap(cm2, annot=True, fmt=".0f")
4
```

Out[80]: <AxesSubplot:>



Model 3

```
In [67]: 1 score3 = model3.evaluate(x_valid, y_valid)
2 print('GRU model loss:', score3[0])
3 print('GRU model accuracy:', score3[1])
4
```

31/31 [=====] - 9s 275ms/step - loss: 0.6932 - accuracy: 0.5000 - auc_2: 0.5000
GRU model loss: 0.69324791431427
GRU model accuracy: 0.5

```
In [70]: 1 # add validation accuracy score into list
2 v_auc_score3 = history3.history["val_auc_2"]
3 valid_auc_score.append(v_auc_score3)
4
5 # best validation accuracy result
6 best_val_auc3 = max(v_auc_score3)
7 print("GRU Best Validation AUC: ", best_val_auc3)
8
```

GRU Best Validation AUC: 0.5010183453559875

```
In [71]: 1 # make predictions on the validation dataset
2 #load_model3 = keras.models.load_model('model3_cp')
3 y_pred3 = model3.predict(x_valid)
4 y_pred3 = np.where(y_pred3>0.5, 1, 0)
5
6 # print out classification report
7 print(classification_report(y_valid, y_pred3))
8
```

```

              precision    recall  f1-score   support

     0       0.50      1.00      0.67         491
     1       0.00      0.00      0.00         491

 accuracy          0.50         0.50         0.50         982
 macro avg          0.25         0.50         0.33         982
 weighted avg          0.25         0.50         0.33         982

```

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

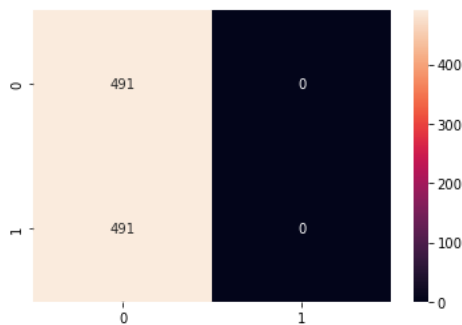
_warn_prf(average, modifier, msg_start, len(result))

/opt/conda/lib/python3.7/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

```
In [72]: 1 # print out the confusion matrix
2 cm3 = confusion_matrix(y_valid, y_pred3)
3 sns.heatmap(cm3, annot=True, fmt=".0f")
4
```

Out[72]: <AxesSubplot:>



4.2 Comparing the three different models

```
In [75]: 1 # create compare dataframe to compare three models
2 models = {'Model': ["LSTM", "Bi_LSTM", "GRU"],
3            'Accuracy': [score[1], score2[1], score3[1]],
4            'Loss': [score[0], score2[0], score3[0]],
5            'Best Validation AUC': {best_val_auc1, best_val_auc2, best_val_auc3}}
6 compare_data = pd.DataFrame(models)
7 compare_data = compare_data.sort_values(by='Best Validation AUC', ascending = False, ignore_index=True)
8 print("Compare three deep learning models: ")
9 display(compare_data)
10
```

Compare three deep learning models:

	Model	Accuracy	Loss	Best Validation AUC
0	LSTM	0.5	0.693397	0.501018
1	GRU	0.5	0.693248	0.501018
2	Bi_LSTM	0.5	0.693306	0.500000

We observe that LSTM and GRU models are better than Bi-LSTM model with higher best validation AUC.

4.3 Run Dropout Tuning

Dropout is a technique where randomly selected neurons are ignored during training. This means that their contribution to the activation of downstream neurons is temporally removed on the forward pass, and any weight updates are not applied to the neuron on the backward pass.

When we have training data, if we try to train your model too much, it might overfit, and dropout regularization is one technique used to tackle overfitting problems in deep learning.

In this project, we will use LSTM model and try 3 different Dropout: [0.1, 0.2, 0.3].

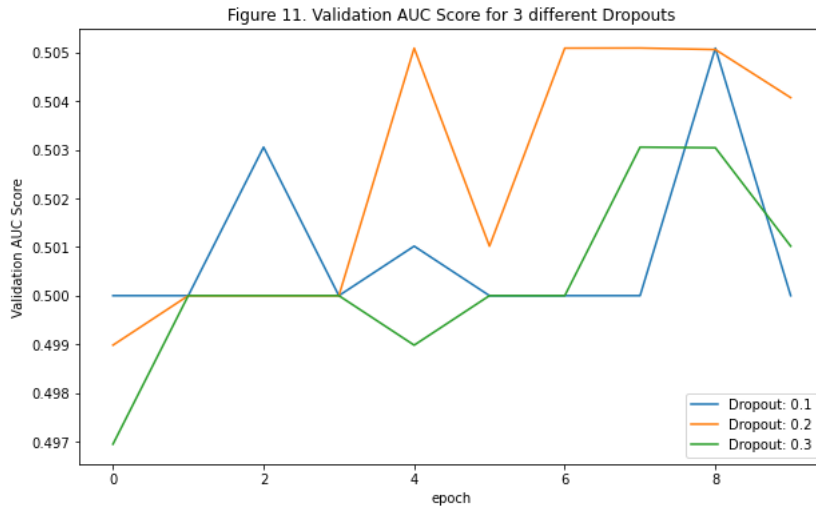
Because we are creating many models in a loop, this global state will consume an increasing amount of memory over time, so we should clear it. Calling `clear_session()` releases the global state: this helps avoid clutter from old models and layers, especially when memory is limited.

```
In [79]: 1 # create a list to store the result
2 dropout_val_auc = []
3
4 for dropout in [0.1, 0.2, 0.3]:
5     # clear session:\
6     tf.keras.backend.clear_session()
7
8     # define new model
9     new_model = Sequential()
10    new_model.add(Embedding(vocab_size,
11                            embedding_dim,
12                            input_length = max_len))
13    new_model.add(LSTM(n_lstm, return_sequences=False))
14    new_model.add(Dropout(dropout))
15    new_model.add(Dense(1, activation='sigmoid'))
16
17    # compile new model
18    new_model.compile(loss = 'binary_crossentropy',
19                     optimizer = 'adam',
20                     metrics=['accuracy', tf.keras.metrics.AUC()])
21
22    # train new model
23    num_epochs = 10
24    new_history = new_model.fit(x_train,
25                                y_train,
26                                epochs=num_epochs,
27                                validation_data=(x_valid, y_valid),
28                                verbose=2)
29
30    # best result
31    print("Best Validation AUC for Dropout: ", dropout, "is: ", max(new_history.history["val_auc"]))
32
33    # add validation AUC score into list
34    dropout_val_auc.append(new_history.history["val_auc"])
```

Epoch 1/10
174/174 - 198s - loss: 0.6943 - accuracy: 0.4924 - auc: 0.4891 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 2/10
174/174 - 195s - loss: 0.6935 - accuracy: 0.4926 - auc: 0.4934 - val_loss: 0.6934 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 3/10
174/174 - 195s - loss: 0.6936 - accuracy: 0.4946 - auc: 0.4933 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5031
Epoch 4/10
174/174 - 195s - loss: 0.6934 - accuracy: 0.4878 - auc: 0.4917 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 5/10
174/174 - 196s - loss: 0.6934 - accuracy: 0.4944 - auc: 0.4932 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5010
Epoch 6/10
174/174 - 195s - loss: 0.6935 - accuracy: 0.4844 - auc: 0.4811 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 7/10
174/174 - 195s - loss: 0.6935 - accuracy: 0.4942 - auc: 0.4960 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 8/10
174/174 - 195s - loss: 0.6933 - accuracy: 0.4995 - auc: 0.4952 - val_loss: 0.6931 - val_accuracy: 0.5010 - val_auc: 0.5000
Epoch 9/10
174/174 - 195s - loss: 0.6935 - accuracy: 0.4919 - auc: 0.4895 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5051
Epoch 10/10
174/174 - 196s - loss: 0.6934 - accuracy: 0.4991 - auc: 0.5011 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000
Best Validation AUC for Dropout: 0.1 is: 0.505091667175293
Epoch 1/10
174/174 - 197s - loss: 0.6942 - accuracy: 0.4957 - auc: 0.4947 - val_loss: 0.6938 - val_accuracy: 0.5000 - val_auc: 0.4990
Epoch 2/10
174/174 - 195s - loss: 0.6939 - accuracy: 0.4984 - auc: 0.4885 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 3/10
174/174 - 195s - loss: 0.6934 - accuracy: 0.4982 - auc: 0.4971 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 4/10
174/174 - 194s - loss: 0.6933 - accuracy: 0.5020 - auc: 0.5016 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 5/10
174/174 - 195s - loss: 0.6933 - accuracy: 0.5061 - auc: 0.5017 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5051
Epoch 6/10
174/174 - 194s - loss: 0.6933 - accuracy: 0.5038 - auc: 0.4958 - val_loss: 0.6931 - val_accuracy: 0.5051 - val_auc: 0.5010
Epoch 7/10
174/174 - 195s - loss: 0.6944 - accuracy: 0.5068 - auc: 0.5024 - val_loss: 0.6931 - val_accuracy: 0.5031 - val_auc: 0.5051
Epoch 8/10
174/174 - 195s - loss: 0.6935 - accuracy: 0.4955 - auc: 0.4908 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5051
Epoch 9/10
174/174 - 194s - loss: 0.6934 - accuracy: 0.5005 - auc: 0.4928 - val_loss: 0.6931 - val_accuracy: 0.5031 - val_auc: 0.5051
Epoch 10/10
174/174 - 195s - loss: 0.6934 - accuracy: 0.4964 - auc: 0.4902 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5041
Best Validation AUC for Dropout: 0.2 is: 0.5050937533378601
Epoch 1/10
174/174 - 197s - loss: 0.6945 - accuracy: 0.4885 - auc: 0.4878 - val_loss: 0.6933 - val_accuracy: 0.5000 - val_auc: 0.4969
Epoch 2/10
174/174 - 195s - loss: 0.6938 - accuracy: 0.4982 - auc: 0.4909 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 3/10
174/174 - 195s - loss: 0.6937 - accuracy: 0.4935 - auc: 0.4935 - val_loss: 0.6933 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 4/10
174/174 - 194s - loss: 0.6934 - accuracy: 0.4975 - auc: 0.4980 - val_loss: 0.6934 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 5/10
174/174 - 195s - loss: 0.6938 - accuracy: 0.4993 - auc: 0.4886 - val_loss: 0.6931 - val_accuracy: 0.5010 - val_auc: 0.4990
Epoch 6/10
174/174 - 194s - loss: 0.6935 - accuracy: 0.4883 - auc: 0.4887 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 7/10
174/174 - 195s - loss: 0.6935 - accuracy: 0.5059 - auc: 0.5018 - val_loss: 0.6934 - val_accuracy: 0.5000 - val_auc: 0.5000
Epoch 8/10
174/174 - 194s - loss: 0.6936 - accuracy: 0.5002 - auc: 0.4962 - val_loss: 0.6931 - val_accuracy: 0.5000 - val_auc: 0.5031

Epoch 9/10
 174/174 - 195s - loss: 0.6934 - accuracy: 0.5016 - auc: 0.5023 - val_loss: 0.6931 - val_accuracy: 0.5020 - val_auc: 0.5030
 Epoch 10/10
 174/174 - 195s - loss: 0.6934 - accuracy: 0.4885 - auc: 0.4947 - val_loss: 0.6932 - val_accuracy: 0.5000 - val_auc: 0.5010
 Best Validation AUC for Dropout: 0.3 is: 0.5030549764633179

```
In [81]: 1 # plot validation AUC score for three different droupouts
2 plt.figure(figsize=(10, 6))
3 plt.xlabel("epoch")
4 plt.ylabel("Validation AUC Score")
5 epoch_range = list(range(10))
6 labels = ["Dropout: 0.1", "Dropout: 0.2", "Dropout: 0.3"]
7 for i in range(len(dropout_val_auc)):
8     plt.plot(epoch_range, dropout_val_auc[i], label=f'{labels[i]}')
9 plt.legend(loc='lower right')
10 plt.title("Figure 11. Validation AUC Score for 3 different Dropouts")
11 plt.show()
12
```



Looking at the result above, we can conclude that in this case, the best dropout is 0.2

4.3 Use the best model to predict test data without target

Looking at the compare dataframe above, we can see that GRU is the best model because it has the highest validation AUC score and least loss compare to LSTM and Bi-LSTM models and the best dropout is 0.2. That is the model 3 (GRU) was built above, so now I will:

- clean text in test data
- vectorizing text
- use model 3 for predicting test data
- create submission file

```
In [82]: 1 # clean test data
2 clean_text(test, "text")
3
4 # view text in a row after cleaning all text data
5 test["text"][0]
6
```

Out[82]: 'happened terrible car crash'

```
In [95]: 1 # create vector for test set
2 test_vectors = count_vectorizer.fit_transform(test["text"])
3
```

```
In [85]: 1 # load model3
2 best_model = tf.keras.models.load_model('model3_cp')
```

```
In [97]: 1 # predict test data
2 test_vectors = test_vectors.toarray()
3 y_pred = best_model.predict(test_vectors)
4 test_predictions = np.where(y_pred>0.5, 1, 0)
```



```
In [98]: 1 # create dataframe of result
2 submission = pd.DataFrame()
3 submission['id'] = test['id']
4 submission['target'] = test_predictions
5 submission.head()
6
```

```
Out[98]:
```

	id	target
0	0	1
1	2	1
2	3	1
3	9	1
4	11	1

```
In [99]: 1 # view test prediction counts
2 submission['target'].value_counts()
3
```

```
Out[99]: 1    3249
0         14
Name: target, dtype: int64
```

```
In [100]: 1 # plot the count of each label
2 fig, ax = plt.subplots(figsize=(6,6))
3 sns.countplot(data=submission, y='target', ax=ax).set(title='\nFigure 5. The Count of Each Target\n')
4
5 # plot the proportion of each label
6 labels = submission['target'].unique().tolist()
7 counts = submission['target'].value_counts()
8 sizes = [counts[v] for v in labels]
9 fig1, ax1 = plt.subplots()
10 ax1.pie(sizes, labels=labels, autopct='%0.2f%%')
11 ax1.axis('equal')
12 plt.title("\nFigure 12. The Proportion of Each Target\n")
13 plt.tight_layout()
14 plt.show()
15
```

Figure 5. The Count of Each Target

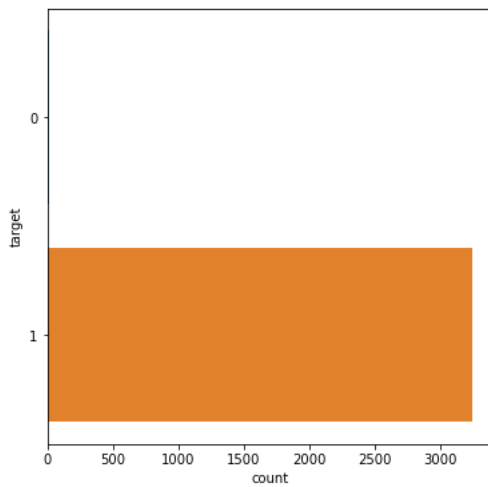
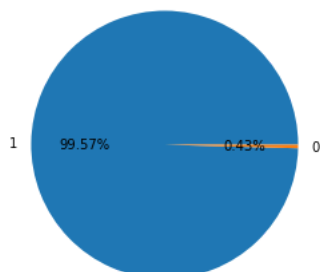


Figure 12. The Proportion of Each Target



```
In [101]: 1 # convert to csv to submit to competition
          2 submission.to_csv('submission.csv', index=False)
          3
```

Step 5: Conclusion

In this project, there are 5 parts:

- (1) Brief description of the problem and data.
- (2) EDA - Inspect, Visualize, and Clean the data.
- (3) Building and training models:

- LSTM
- Bi-LSTM
- GRU

(4) Results and Analysis.

(5) Conclusion.

The goal of this project is to detect fake and real disaster tweets. By comparing three different deep learning models including: LSTM, Bi-LSTM and GRU, we can conclude that in this case, the GRU model is the model that has the best performance with the highest validation AUC value of 0.501018 and the loss value = 0.693248. I know this result was not good, however, because of the limitation of data and the running time was too costly, the models just were trained on limited approach. I think there are many other ways can improve this kind of project such as: building more deep learning models by tuning hyperparameters to get optimal results, or we can run models with more epoch, or use other type of Word Embeddings such as: Tokenization or Bag-of-Words.

Because of the curiosity, I would like to print out the predictions of all three models and let's see how these model's performance are.

```
In [ ]: 1 # use three models predict test set and print out the submission
        2
        3 for m in ["model1", "model2", "model3"]:
        4     model = tf.keras.models.load_model(f'{m}_cp')
        5
        6     # predict test data
        7     y_pred = model.predict(test_vectors)
        8     test_predictions = np.where(y_pred > 0.5, 1, 0)
        9
        10    # create dataframe of result
        11    submission = pd.DataFrame()
        12    submission['id'] = test['id']
        13    submission['target'] = test_predictions
        14    submission.to_csv(f'{m}_submission.csv', index=False)
```



Submissions

Submission and Description		Public Score ⓘ
<div> <div>All</div> <div>Successful</div> <div>Errors</div> </div> <div>Recent ▾</div>		
<div>✓</div> model3_submission.csv Complete · 2m ago	0.42844	
<div>✓</div> model2_submission.csv Complete · 4m ago	0.57155	
<div>✓</div> model1_submission.csv Complete · 5m ago	0.42844	

It's interesting that model 2 (Bi-LSTM) has better score.