

NLP Disaster Tweets Kaggle Mini-Project

1. Introduction

Tweeter is a powerful monbile communication with a wide audiences. This is very useful in time of emergency. However, it can be challenging for machine learning model to differentiate a text message to be classified as a "disaster" as some word can be use interchangably in a non-disaster situation. The goal of this project is to build a RNN machine learning model to predict if a tweet is about a real disaster.

```
In [71]: # Load the Libraries
import os
import re
import string
import nltk
import string
import unicodedata
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import time
from collections import Counter
from wordcloud import WordCloud
from nltk.corpus import stopwords
nltk.download('stopwords')
from collections import Counter
import tensorflow as tf
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, SimpleRNN, Dense, Dropout, LSTM, GRU
from tensorflow.keras.callbacks import EarlyStopping
from sklearn.model_selection import train_test_split
from sklearn.metrics import roc_curve, auc, confusion_matrix, classification_report
from sklearn.utils import class_weight
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] /home/dataengineer/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

1.1 Data Source

- The of Source of the dataset are taken from Kaggle: <https://www.kaggle.com/c/nlp-getting-started/overview>
- There are 2 main data file used in this project
 - `train.csv` - the training set
 - `test.csv` - the test set

```
In [2]: #TRAIN_DIR = '/kaggle/input/nlp-getting-started/train.csv'
#TEST_DIR = '/kaggle/input/nlp-getting-started/test.csv'

TRAIN_DIR = '/mnt/d/Data/nlp-getting-started/train.csv'
TEST_DIR = '/mnt/d/Data/nlp-getting-started/test.csv'

trained_labeled_df = pd.read_csv(TRAIN_DIR)
test_df = pd.read_csv(TEST_DIR)
```

1.2 Training Dataset train.csv

- 7613 rows
- 5 columns:
 - `id` : unique identifier in int64 type
 - `keyword` : keyword from the tweet in object type
 - `location` : location info in object type
 - `text` : tweet text in object in object type
 - `target` : 0 = Non-Disaster, 1 = Disaster in int64 type
- There are no nulls for `id`, `text`, `target` columns
- There are nulls in `keyword` and `location` columns which might impact the cleaning process and training of model
 - keyword: 61
 - location: 2533
- The `text` column consist of abbrivation, symbols like #, =>, etc, so will need to be treated during the cleaning process

```
In [3]: trained_labeled_df.head()
```

```
Out[3]:
```

	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 [4]: print('Train DataFrame Shape:', trained_labeled_df.shape)
```

Train DataFrame Shape: (7613, 5)

```
In [5]: print('\nSummary of Train DataFrame:\n',trained_labeled_df.info())

# Shows the total null count for every column in the DataFrame
print('\nNulls in Train DataFrame:\n',trained_labeled_df.isnull().sum())
```

```
<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
```

Summary of Train DataFrame:
None

Nulls in Train DataFrame:

id	0
keyword	61
location	2533
text	0
target	0

dtype: int64

1.3 Test Dataset test.csv

- 3263 rows
- 4 columns:
 - `id` : unique identifier in int64 type
 - `keyword` : keyword from the tweet in object type
 - `location` : location info in object type
 - `text` : tweet text in object in object type
- Similar dataframe structure to `train.csv` except `target` column
- There are nulls in `keyword` and `location` columns which might impact the cleaning process and training of model
- The `text` column consist of abbrivation, symbols like #, =>, etc, so will need to be treated during the cleaning process

```
In [7]: test_df.head()
```

```
Out[7]:
```

	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 [9]: print('Test DataFrame shape :', test_df.shape)
```

```
Test DataFrame shape : (3263, 4)
```

```
In [10]: print('\nSummary of Test DataFrame:\n',test_df.info())

# Shows the total null count for every column in the DataFrame
print('\nNulls in Test DataFrame:\n',test_df.isnull().sum())
```

```
<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
```

```
Summary of Test DataFrame:
```

```
None
```

```
Nulls in Test DataFrame:
```

```
id          0
keyword     26
location    1105
text        0
dtype: int64
```

2. Exploratory Data Analysis (EDA)

2.1 Target Distribution of the Train Dataframe

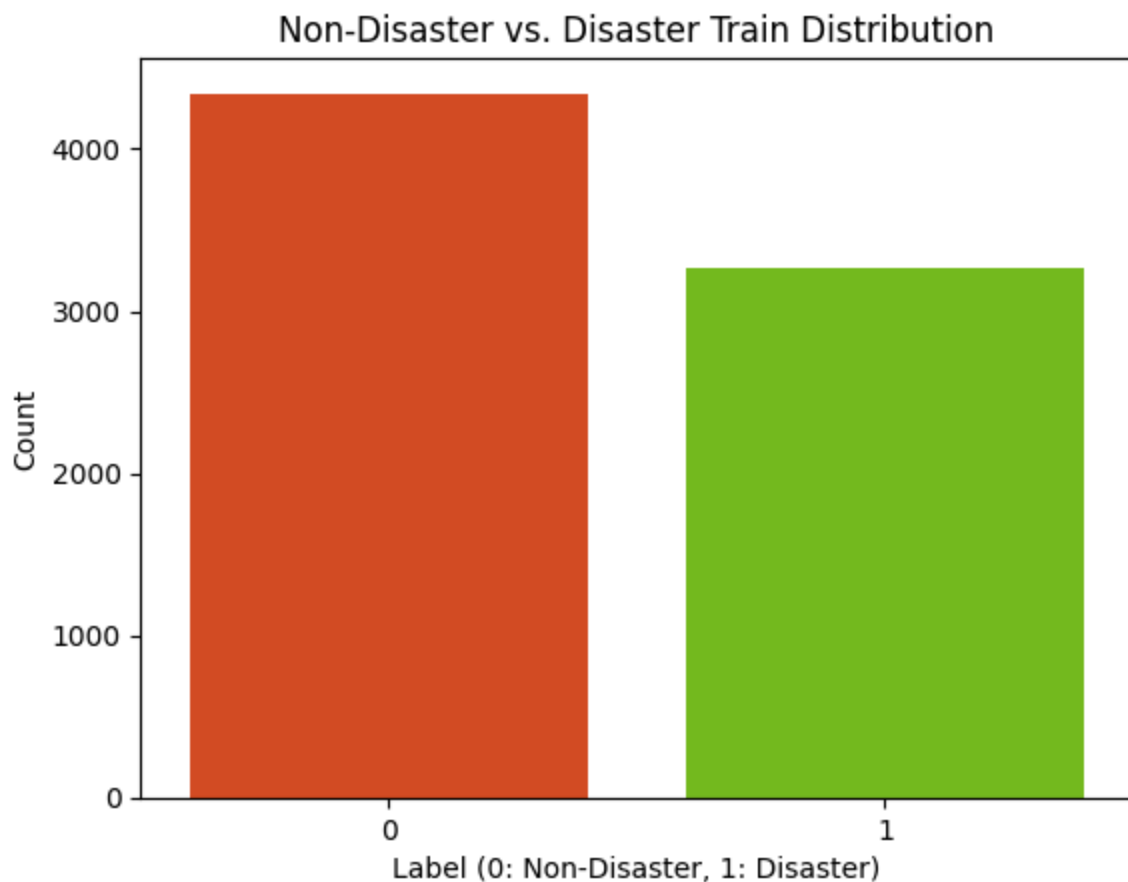
Target	%
0 = Non-Disaster	57%
1 = Disaster	42%

The slight imbalance would need to be take into consideration during the calculation for bias.

```
In [12]: print('\nTarget distribution:\n', trained_labeled_df['target'].value_counts(normali
```

```
Target distribution:
target
0    0.57034
1    0.42966
Name: proportion, dtype: float64
```

```
In [25]: sns.countplot(
    x=trained_labeled_df['target'],
    palette=["#f33d05", "#7bd307"],
    hue=trained_labeled_df['target'],
    legend=False
)
plt.title('Non-Disaster vs. Disaster Train Distribution', fontsize=12)
plt.xlabel('Label (0: Non-Disaster, 1: Disaster)')
plt.ylabel('Count')
plt.show()
```



2.2 Word Length Per Tweet Distribution

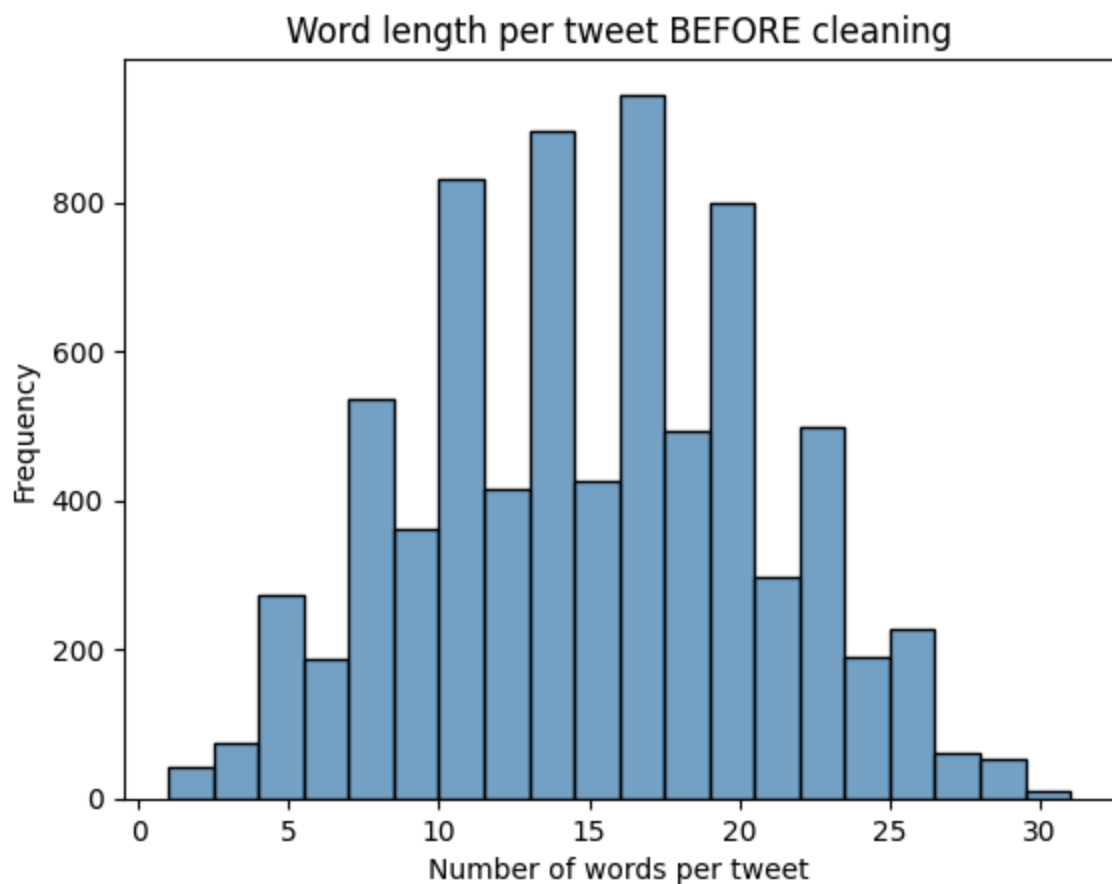
- Histogram of training dataset text shows a normal distribution, with most tweets having between 5 and 25 words

- The peak is at 17 words per tweet, with an occurrence of about 10,000 times
- For an RNN, every input sequence must have the same length
- The distribution will change after removing irrelevant words, and we will then decide whether padding or truncation is needed

```
In [27]: def plot_word_length_distribution(df, text_col='text', bins=20, title='Word length
df['len_words'] = df[text_col].astype(str).str.split().apply(len)
sns.histplot(data=df, x='len_words', bins=bins, color='steelblue')
plt.title(title)
plt.xlabel('Number of words per tweet')
plt.ylabel('Frequency')

plt.show()

plot_word_length_distribution(trained_labeled_df, title='Word length per tweet BEFO
```



2.2 Keyword Frequencies

List and plot the top 10 keywords and their frequencies for both disaster and non-disaster tweets.

```
In [28]: # Find and plot the top most frequent words
def word_count(df, text_column):
    all_words = []
```

```

# Lowercase
text_data = df[text_column].astype(str).str.lower()
for text in text_data:
    # use regex to find each word
    words = re.findall(r'\b\w+\b', text)
    cleaned_words = [re.sub(r'^20', '', word) for word in words]
    all_words.extend(cleaned_words)

# Count frequencies of each word to a dictionary
word_counts = Counter(all_words)
return word_counts

# print(word_counts)

def plot_top_words(word_counts, n=50, plot_size=(15, 6), title=''):

    # List top N to a new dataframe
    top_n_words = word_counts.most_common(n)
    freq_df = pd.DataFrame(top_n_words, columns=['Word', 'Count'])

    # Plot with Barplot
    #plt.figure(figsize=(15, 6))
    plt.figure(figsize=plot_size)
    sns.barplot(
        x='Count',
        y='Word',
        data=freq_df,
        palette='viridis',
        hue='Word',
        legend=False
    )

    if not title:
        title = f'Top {n} Most Frequent Words'

    plt.title(title, fontsize=16)
    plt.xlabel('Frequency (Count)', fontsize=12)
    plt.ylabel('Word', fontsize=12)
    plt.xticks(rotation=75, ha='right', fontsize=10)
    plt.tight_layout()
    plt.show()

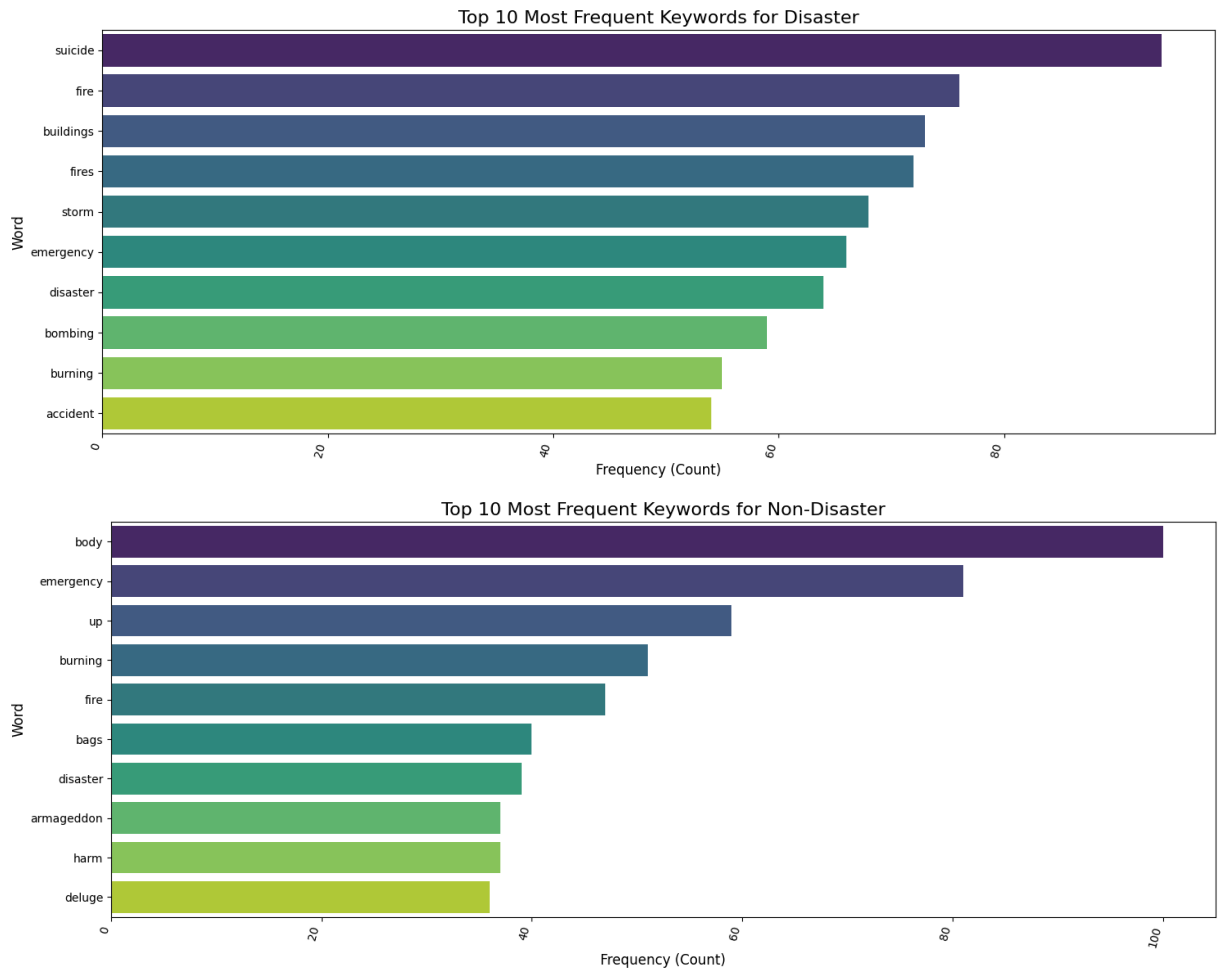
```

```

In [29]: top_n = 10
plot_title = f'Top {top_n} Most Frequent Keywords for Disaster'
disaster_keyword_count_results = word_count(trained_labeled_df[trained_labeled_df.t
plot_top_words(disaster_keyword_count_results, n=top_n, title=plot_title)

plot_title = f'Top {top_n} Most Frequent Keywords for Non-Disaster'
non_disaster_keyword_count_results = word_count(trained_labeled_df[trained_labeled_
plot_top_words(non_disaster_keyword_count_results, n=top_n, title=plot_title)

```

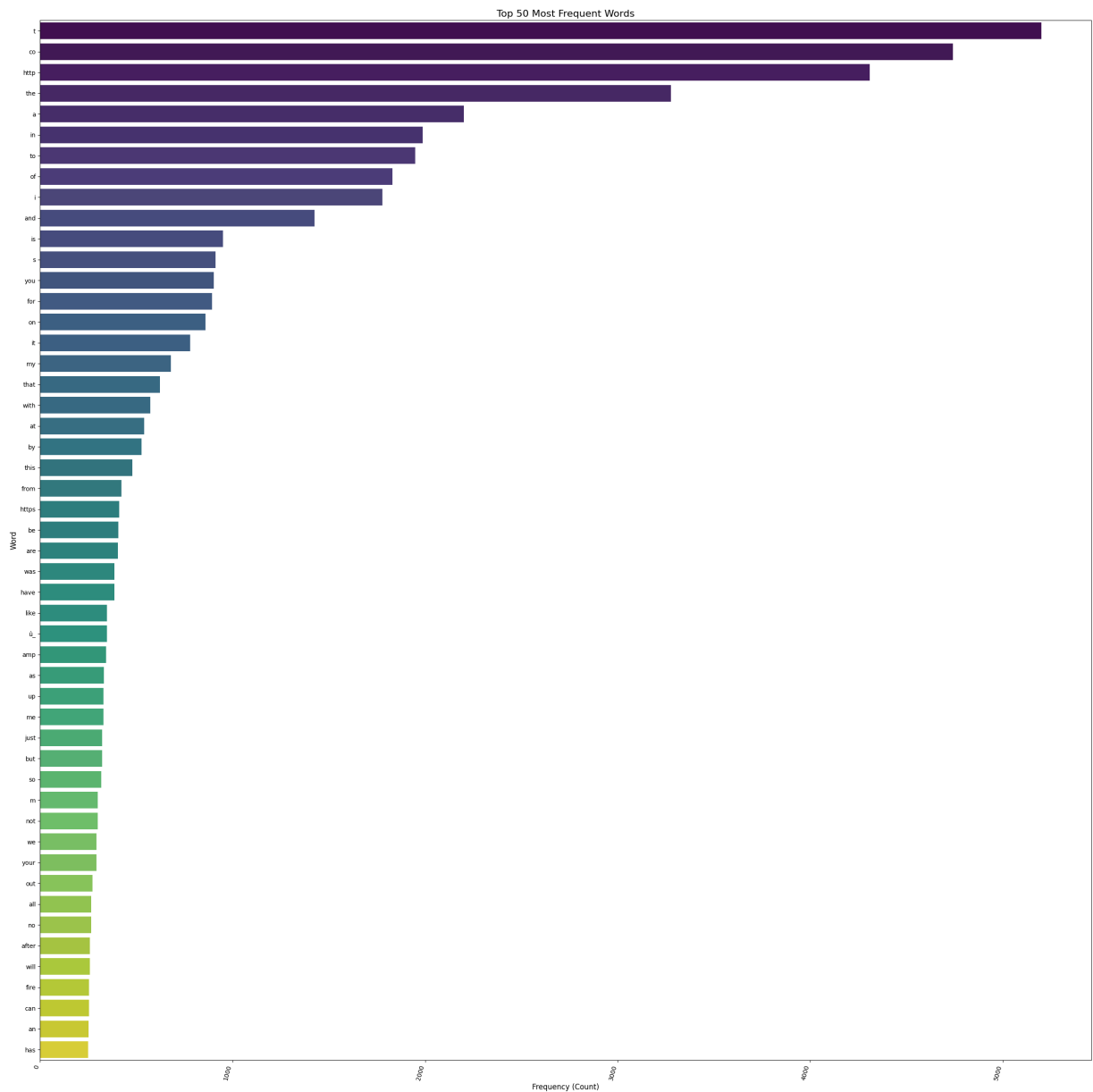


2.3 Visual Inspection of the Text

2.3.1 Word in Text Frequencies

- List and plot the top 50 words and their frequencies in tweets text.
- Usually high amount of irrelevant words such as "t", "co", "http", etc. should be remove in the cleaning process

```
In [30]: text_count_results = word_count(trained_labeled_df, 'text')
plot_top_words(text_count_results, n=50, plot_size=(25, 25))
```

2.3.2 Word Cloud for Train Dataset

- A word cloud was used to visually inspect the top 100 words in the training dataset

```
In [31]: # Word Cloud - Cleaned Text
def word_cloud_plot(df_text, plot_title):

    # Combine all tweets into one string
    combined_text = " ".join(df_text.astype(str))

    # Generate word cloud
    wc_data = WordCloud(width=800, height=400, background_color='white', max_words=

    if not plot_title:
        plot_title = f'Word Cloud'

    plt.figure(figsize=(10,5))
    plt.imshow(wc_data, interpolation='bilinear')
```

```
plt.axis('off')
plt.title(plot_title)
plt.show()
```

```
plot_title = f'Word Cloud Training Dataset Text Before Cleaning'
word_cloud_plot(trained_labeled_df['text'], plot_title)
```



2.3.3 Word Cloud for Test Dataset

- A word cloud was also used to visually inspect the top 100 words in the test dataset.
- Based on observation, it can be assumed that similar words found in the training dataset are also in the test dataset.
- Therefore, the same cleaning logic can be applied to both datasets.

```
plot_title = f'Word Cloud Test Dataset Text Before Cleaning'
word_cloud_plot(test_df['text'], plot_title)
```



```

])

url_re = re.compile(r'https?://\S+|www\.\S+')
mention_re = re.compile(r'@\w+')
html_re = re.compile(r'<.*?>')
emoji_re = re.compile("[\U00010000-\U0010ffff]", flags=re.UNICODE)

def clean_text(s, remove_stopwords=True):
    if not isinstance(s, str):
        return ''
    s = s.lower()
    s = url_re.sub(' ', s)
    s = mention_re.sub(' ', s)
    s = html_re.sub(' ', s)
    s = emoji_re.sub(' ', s)
    s = s.replace('\n', ' ')
    s = s.translate(str.maketrans('', '', string.punctuation))
    s = re.sub('\s+', ' ', s).strip()

    if remove_stopwords:
        words = [w for w in s.split() if w not in CUSTOM_STOPWORDS]
        s = ' '.join(words)
    return s

# Clean train dataset
trained_labeled_df['text_clean'] = trained_labeled_df['text'].apply(remove_accents)
trained_labeled_df['text_clean'] = trained_labeled_df['text_clean'].apply(clean_text)

# Clean test dataset
test_df['text_clean'] = test_df['text'].apply(remove_accents)
test_df['text_clean'] = test_df['text_clean'].apply(clean_text)

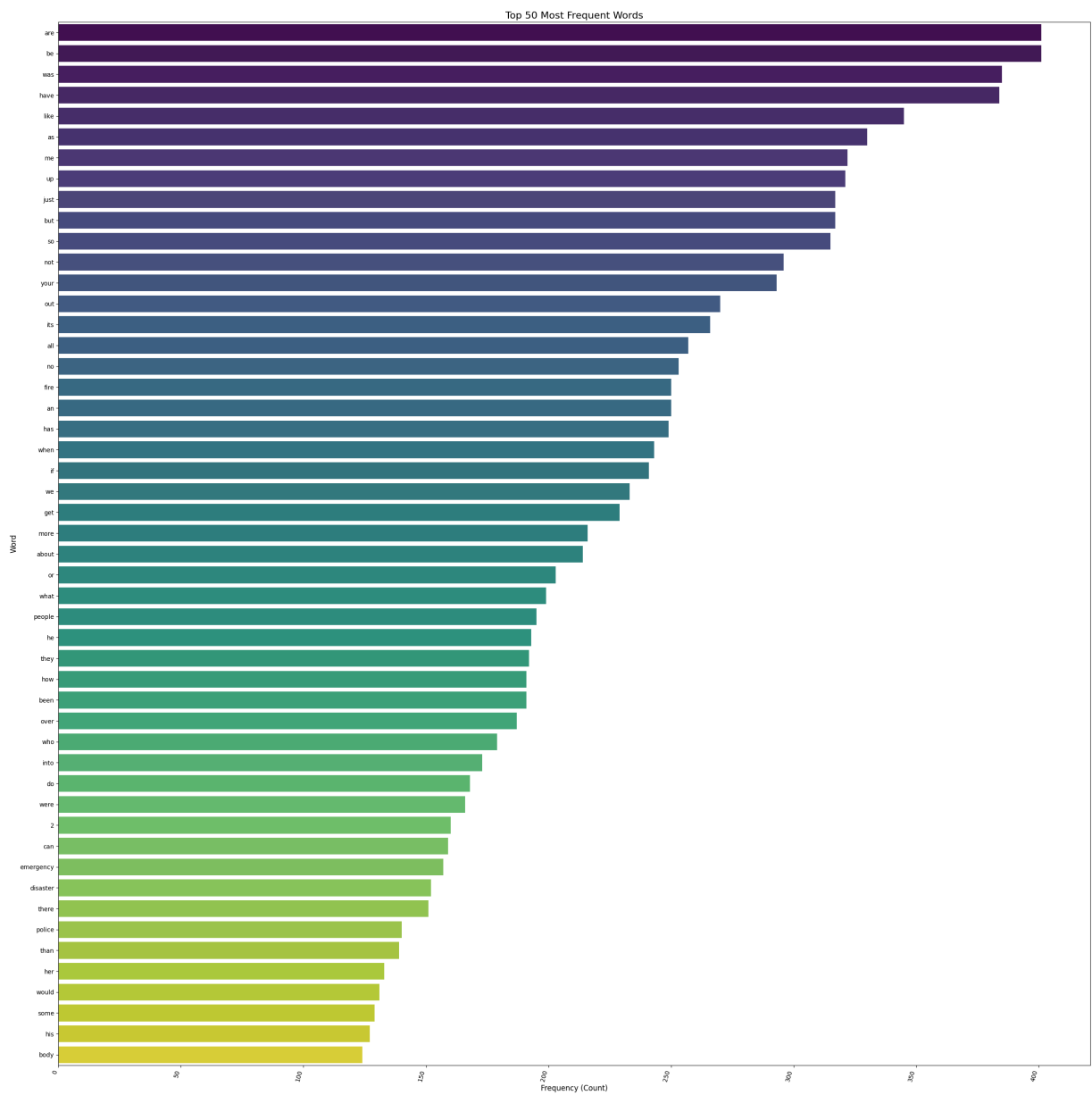
```

2.4.2 Post Data Cleaning Inspection

```

In [35]: text_clean_count_results = word_count(trained_labeled_df, 'text_clean')
plot_top_words(text_clean_count_results, n=50, plot_size=(25, 25))

```



```
In [34]: plot_title = f'Word Cloud Text After Cleaning'
word_cloud_plot(trained_labeled_df['text_clean'], plot_title)
```


- The data are split into a training set (80%) and a validation set (20%) to test the model
- Tested vocabulary sizes between 600 and 10,000 and found that a size of 600 gave the best accuracy and performance
- Words are then tokenized to:
 - Limit the vocabulary to the 600 most common words
 - Convert text into numerical sequences
 - Standardize word length to 40 by either padding or truncation
- A stop rule is used to terminate training early if the model's performance is not improving

```
In [46]: # Prep the data
# feature variable: input cleaned text
X = trained_labeled_df['text_clean']
# target variable: labels to predict
y = trained_labeled_df['target']

# Split the data into 80% training set and 20% validation set.
X_train, X_val, y_train, y_val = train_test_split(X, y, stratify=y, test_size=0.2,

# Limit number of words to use in the vocabulary
max_words = 600
# Standardize word length to 40 by either padding or truncation
max_len = 40

# Tokenization
# set tokenization parameters and initialize it
tokenizer = Tokenizer(num_words=max_words)
tokenizer.fit_on_texts(X)

# convert the text into numbers and ensure all text per tweet are the same length
X_train_tokenized = pad_sequences(tokenizer.texts_to_sequences(X_train), maxlen=max
X_val_tokenized = pad_sequences(tokenizer.texts_to_sequences(X_val), maxlen=max_len
X_test_tokenized = pad_sequences(tokenizer.texts_to_sequences(test_df['text_clean']

# EarlyStopping stops the training early if the model isn't improving for 3 epochs
# the best version was saved before it stopped
es = EarlyStopping(monitor='val_loss', patience=3, restore_best_weights=True)
```

3. Model Architecture

3.1 Model Implementation

Comparing three RNN-family architectures, each using the same input embedding layer and a single dense output unit:

- SimpleRNN
- LSTM
- GRU

Same settings across three models:

- Embedding layer: converts tokenized text into dense vectors with input_dim = max_words and output_dim = 32
- Dropout: 0.3 to reduce overfitting
- Dense output: 1 unit of sigmoid activation to predict either disaster or non-disaster
- Batch size: 64

```
In [55]: # model_rnn
model_rnn = Sequential([
    Embedding(input_dim=max_words, output_dim=32),
    SimpleRNN(32), #32
    Dropout(0.3), #0.3
    Dense(1, activation='sigmoid')
])

model_rnn.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model_rnn.summary()

#history_model_rnn = model_rnn.fit(X_train_tokenized, y_train, validation_data=(X_val, y_val))
history_model_rnn = model_rnn.fit(X_train_tokenized, y_train, validation_data=(X_val, y_val))
```

Model: "sequential_8"

Layer (type)	Output Shape	Param #
embedding_8 (Embedding)	?	0 (unbuilt)
simple_rnn_5 (SimpleRNN)	?	0 (unbuilt)
dropout_8 (Dropout)	?	0
dense_8 (Dense)	?	0 (unbuilt)

Total params: 0 (0.00 B)

Trainable params: 0 (0.00 B)

Non-trainable params: 0 (0.00 B)

Epoch 1/10
 96/96 ————— 3s 16ms/step - accuracy: 0.6232 - loss: 0.6554 - val_accuracy: 0.7124 - val_loss: 0.5851
 Epoch 2/10
 96/96 ————— 1s 13ms/step - accuracy: 0.7716 - loss: 0.4984 - val_accuracy: 0.7800 - val_loss: 0.4836
 Epoch 3/10
 96/96 ————— 1s 15ms/step - accuracy: 0.8089 - loss: 0.4321 - val_accuracy: 0.7794 - val_loss: 0.4874
 Epoch 4/10
 96/96 ————— 2s 12ms/step - accuracy: 0.8305 - loss: 0.3903 - val_accuracy: 0.7695 - val_loss: 0.5050
 Epoch 5/10
 96/96 ————— 3s 28ms/step - accuracy: 0.8471 - loss: 0.3553 - val_accuracy: 0.7735 - val_loss: 0.5216
 Epoch 6/10
 96/96 ————— 6s 32ms/step - accuracy: 0.8677 - loss: 0.3251 - val_accuracy: 0.7643 - val_loss: 0.5620
 Epoch 7/10
 96/96 ————— 2s 23ms/step - accuracy: 0.8811 - loss: 0.2994 - val_accuracy: 0.7420 - val_loss: 0.6154
 Epoch 8/10
 96/96 ————— 2s 16ms/step - accuracy: 0.8869 - loss: 0.2771 - val_accuracy: 0.7452 - val_loss: 0.6499
 Epoch 9/10
 96/96 ————— 1s 15ms/step - accuracy: 0.8957 - loss: 0.2583 - val_accuracy: 0.7459 - val_loss: 0.7016
 Epoch 10/10
 96/96 ————— 1s 11ms/step - accuracy: 0.9053 - loss: 0.2419 - val_accuracy: 0.7452 - val_loss: 0.7336

```
In [56]: model_lstm = Sequential([
    Embedding(input_dim=max_words, output_dim=32),
    LSTM(32),
    Dropout(0.3),
    Dense(1, activation='sigmoid')
])

model_lstm.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model_lstm.summary()

#history_model_lstm = model_lstm.fit(X_train_tokenized, y_train, validation_data=(X_
history_model_lstm = model_lstm.fit(X_train_tokenized, y_train, validation_data=(X_
```

Model: "sequential_9"


Layer (type)	Output Shape	Param #
embedding_9 (Embedding)	?	0 (unbuilt)
lstm_2 (LSTM)	?	0 (unbuilt)
dropout_9 (Dropout)	?	0
dense_9 (Dense)	?	0 (unbuilt)

Total params: 0 (0.00 B)


Trainable params: 0 (0.00 B)

Non-trainable params: 0 (0.00 B)


Epoch 1/10

96/96  3s 22ms/step - accuracy: 0.6090 - loss: 0.6584 - val_accuracy: 0.7374 - val_loss: 0.5832


Epoch 2/10

96/96  2s 19ms/step - accuracy: 0.7739 - loss: 0.4936 - val_accuracy: 0.7820 - val_loss: 0.4789


Epoch 3/10

96/96  3s 20ms/step - accuracy: 0.8030 - loss: 0.4396 - val_accuracy: 0.7840 - val_loss: 0.4735


Epoch 4/10

96/96  2s 21ms/step - accuracy: 0.8146 - loss: 0.4241 - val_accuracy: 0.7873 - val_loss: 0.4760


Epoch 5/10

96/96  2s 19ms/step - accuracy: 0.8166 - loss: 0.4166 - val_accuracy: 0.7722 - val_loss: 0.4861


Epoch 6/10

96/96  2s 19ms/step - accuracy: 0.8205 - loss: 0.4142 - val_accuracy: 0.7853 - val_loss: 0.4794


Epoch 7/10

96/96  2s 21ms/step - accuracy: 0.8197 - loss: 0.4082 - val_accuracy: 0.7768 - val_loss: 0.4813


Epoch 8/10

96/96  2s 21ms/step - accuracy: 0.8236 - loss: 0.4020 - val_accuracy: 0.7774 - val_loss: 0.4836

Epoch 9/10

96/96  2s 24ms/step - accuracy: 0.8271 - loss: 0.3968 - val_accuracy: 0.7827 - val_loss: 0.4892

Epoch 10/10

96/96  2s 22ms/step - accuracy: 0.8299 - loss: 0.3931 - val_accuracy: 0.7781 - val_loss: 0.4876

```
In [57]: model_gru = Sequential([
            Embedding(input_dim=max_words, output_dim=32),
            GRU(32),
            Dropout(0.3),
            Dense(1, activation='sigmoid')
        ])

model_gru.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model_gru.summary()

#history_model_gru = model_gru.fit(X_train_tokenized, y_train, validation_data=(X_val_tokenized, y_val))
history_model_gru = model_gru.fit(X_train_tokenized, y_train, validation_data=(X_val_tokenized, y_val))
```

Model: "sequential_10"


Layer (type)	Output Shape	Param #
embedding_10 (Embedding)	?	0 (unbuilt)
gru_1 (GRU)	?	0 (unbuilt)
dropout_10 (Dropout)	?	0
dense_10 (Dense)	?	0 (unbuilt)

Total params: 0 (0.00 B)


Trainable params: 0 (0.00 B)

Non-trainable params: 0 (0.00 B)


Epoch 1/10

96/96  5s 31ms/step - accuracy: 0.6223 - loss: 0.6485 - val_accuracy: 0.7551 - val_loss: 0.5695


Epoch 2/10

96/96  3s 25ms/step - accuracy: 0.7796 - loss: 0.4839 - val_accuracy: 0.7794 - val_loss: 0.4824


Epoch 3/10

96/96  2s 25ms/step - accuracy: 0.8048 - loss: 0.4359 - val_accuracy: 0.7827 - val_loss: 0.4818


Epoch 4/10

96/96  2s 24ms/step - accuracy: 0.8092 - loss: 0.4264 - val_accuracy: 0.7754 - val_loss: 0.4860


Epoch 5/10

96/96  2s 25ms/step - accuracy: 0.8126 - loss: 0.4171 - val_accuracy: 0.7695 - val_loss: 0.4938


Epoch 6/10

96/96  2s 24ms/step - accuracy: 0.8199 - loss: 0.4085 - val_accuracy: 0.7715 - val_loss: 0.4924


Epoch 7/10

96/96  2s 23ms/step - accuracy: 0.8215 - loss: 0.4020 - val_accuracy: 0.7649 - val_loss: 0.5117


Epoch 8/10

96/96  2s 24ms/step - accuracy: 0.8325 - loss: 0.3923 - val_accuracy: 0.7722 - val_loss: 0.5122

Epoch 9/10

96/96  3s 27ms/step - accuracy: 0.8340 - loss: 0.3846 - val_accuracy: 0.7768 - val_loss: 0.5018

Epoch 10/10

96/96  5s 25ms/step - accuracy: 0.8384 - loss: 0.3783 - val_accuracy: 0.7708 - val_loss: 0.5254

3.2 Model Comparison Summary

Accuracy and Loss

- The SimpleRNN model performs the worst, as it starts overfitting quickly.
- The LSTM and GRU models handle the data better and do not overfit as much. Their performance is very close.

F1 Scores

Model	Validation F1 Score
SimpleRNN	0.692
LSTM	0.722
GRU	0.731

Overall

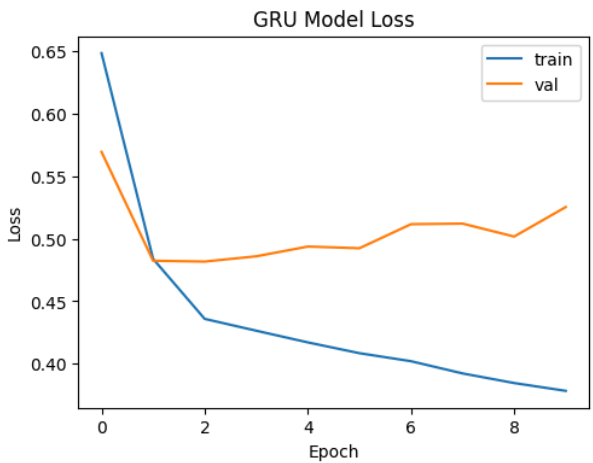
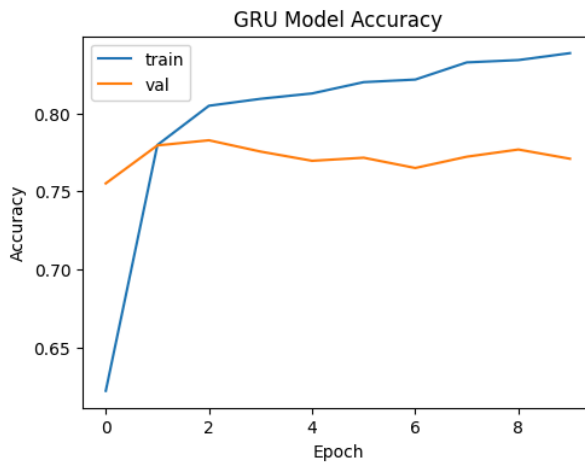
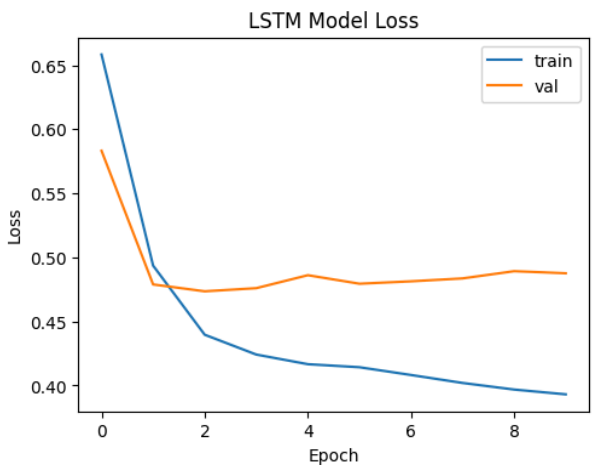
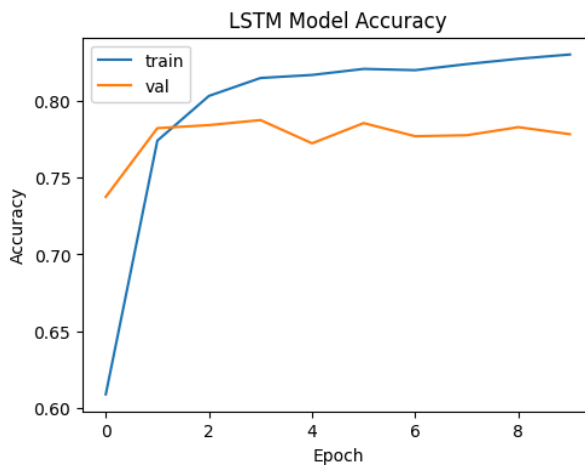
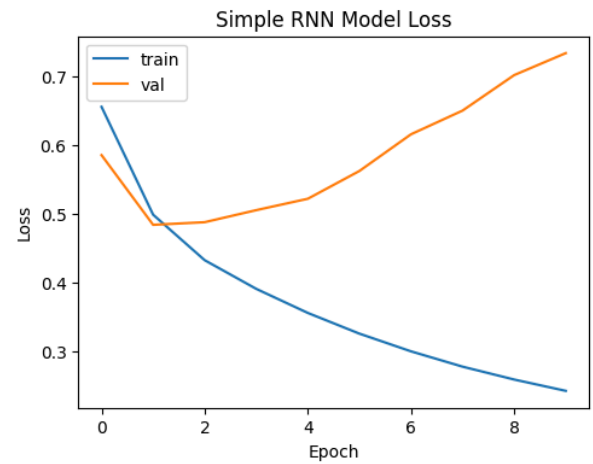
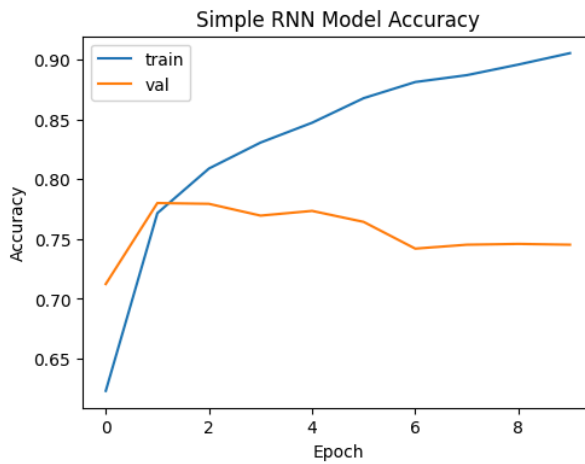
- The GRU model is the best choice for this project because it had the highest F1 score.
- It can remember important information over a long series of data points better than a SimpleRNN, similar to an LSTM.
- The GRU is more lightweight compared to an LSTM because it has fewer gates and parameters. This allows it to train faster and use less memory.

```
In [58]: # Plot training history
def plot_history(history, title):
    plt.figure(figsize=(12,4))
    plt.subplot(1,2,1)
    plt.plot(history.history['accuracy'], label='train')
    plt.plot(history.history['val_accuracy'], label='val')
    plt.title(f'{title} Accuracy')
    plt.xlabel('Epoch')
    plt.ylabel('Accuracy')
    plt.legend()

    plt.subplot(1,2,2)
    plt.plot(history.history['loss'], label='train')
    plt.plot(history.history['val_loss'], label='val')
    plt.title(f'{title} Loss')
    plt.xlabel('Epoch')
    plt.ylabel('Loss')
    plt.legend()

    plt.show()

plot_history(history_model_rnn, 'Simple RNN Model')
plot_history(history_model_lstm, 'LSTM Model')
plot_history(history_model_gru, 'GRU Model')
```



```
In [59]: val_preds_model_rnn_1 = (model_rnn.predict(X_val_tokenized) > 0.5).astype(int).ravel()
print('SimpleRNN Validation F1:', f1_score(y_val, val_preds_model_rnn_1))

val_preds_model_lstm = (model_lstm.predict(X_val_tokenized) > 0.5).astype(int).ravel()
print('LSTM Validation F1:', f1_score(y_val, val_preds_model_lstm))

val_preds_model_gru = (model_gru.predict(X_val_tokenized) > 0.5).astype(int).ravel()
print('GRU Validation F1:', f1_score(y_val, val_preds_model_gru))
```

48/48 ————— 0s 8ms/step
SimpleRNN Validation F1: 0.692063492063492
48/48 ————— 1s 11ms/step
LSTM Validation F1: 0.722495894909688
48/48 ————— 1s 9ms/step
GRU Validation F1: 0.730917501927525

3.3 Hyperparameter Tuning

Since the GRU model achieved the best score among the tested architectures, hyperparameter tuning was performed on this model using the following key parameters:

- Batch sizes: [16, 32, 64]
- Dropout rates: [0.5, 0.3, 0.05]

Hyperparameter optimization results (sorted by highest validation accuracy):

Batch Size	Dropout Rate	Best Val Accuracy	Time Taken
32	0.05	0.7892	45.31
64	0.30	0.7886	23.18
64	0.05	0.7879	26.35
16	0.50	0.7859	80.60
16	0.05	0.7846	67.53
32	0.50	0.7846	53.09
16	0.30	0.7840	84.70
32	0.30	0.7807	42.20
64	0.50	0.7800	27.37

From the results, most hyperparameter combinations achieved similar validation scores around 0.78.

However, the best hyperparameters for the GRU model, balancing accuracy and processing time, are:

- Batch size: 64
- Dropout rate: 0.3

```
In [ ]: # hyperparameters to tune
batch_sizes = [16, 32, 64]
dropout_rates = [0.5, 0.3, 0.05]

results = []

for batch in batch_sizes:
    for dropout in dropout_rates:
```

```

print(f"\nRunning GRU with batch_size={batch}, dropout={dropout}...")

# define model
model_gru = Sequential([
    Embedding(input_dim=max_words, output_dim=32),
    GRU(32),
    Dropout(dropout),
    Dense(1, activation='sigmoid')
])

model_gru.compile(optimizer='adam', loss='binary_crossentropy', metrics=['a

# track time
start_time = time.time()

history = model_gru.fit(
    X_train_tokenized, y_train,
    validation_data=(X_val_tokenized, y_val),
    epochs=10,
    batch_size=batch,
    verbose=0 # silent training
)

end_time = time.time()
elapsed_time = end_time - start_time

# get best val accuracy
best_val_acc = max(history.history['val_accuracy'])

# save results
results.append({
    'Batch Size': batch,
    'Dropout Rate': dropout,
    'Best Val Accuracy': round(best_val_acc, 4),
    'Time Taken (s)': round(elapsed_time, 2)
})

# tabulate
df_results = pd.DataFrame(results)
df_results = df_results.sort_values(by='Best Val Accuracy', ascending=False).reset_
df_results['Rank'] = df_results.index + 1

print("\nHyperparameter Tuning Results:")
print(df_results)

```

Running GRU with batch_size=16, dropout=0.5...

Running GRU with batch_size=16, dropout=0.3...

Running GRU with batch_size=16, dropout=0.05...

Running GRU with batch_size=32, dropout=0.5...

Running GRU with batch_size=32, dropout=0.3...

Running GRU with batch_size=32, dropout=0.05...

Running GRU with batch_size=64, dropout=0.5...

Running GRU with batch_size=64, dropout=0.3...

Running GRU with batch_size=64, dropout=0.05...

Hyperparameter Tuning Results:

	Batch Size	Dropout Rate	Best Val Accuracy	Time Taken (s)	Rank
0	32	0.05	0.7892	45.31	1
1	64	0.30	0.7886	23.18	2
2	64	0.05	0.7879	26.35	3
3	16	0.50	0.7859	80.60	4
4	16	0.05	0.7846	67.53	5
5	32	0.50	0.7846	53.09	6
6	16	0.30	0.7840	84.70	7
7	32	0.30	0.7807	42.20	8
8	64	0.50	0.7800	27.37	9

4. Predict Test Dataset

4.1 Train with Best Hyperparameters

Proceed to train GRU Model with the best hyperparameters combination:

- Batch size: 64
- Dropout rate: 0.3

Note: This were the same hyperparameters used in the initial model comparison.

4.2 Evaluation Training Process

- Training accuracy improved around 81% by epoch 3
- Validation accuracy stabilized around 78% as the model is learning without overfitting
- The F1 score on the validation set is 0.64, as F1 balances precision and recall with imbalanced classes (disasters and non-disaster) Target distribution: target 0 0.57034 1 0.42966 Name: proportion, dtype: float64

```
In [68]: # Compute class weights based on the training labels
class_weights_values = class_weight.compute_class_weight(
```



```

class_weight='balanced',
classes=np.unique(y_train),
y=y_train
)
class_weights_dict = dict(enumerate(class_weights_values))

# Build the GRU model
model_gru = Sequential([
    Embedding(input_dim=max_words, output_dim=32),
    GRU(32),
    Dropout(0.3),
    Dense(1, activation='sigmoid')
])

model_gru.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model_gru.summary()

# Train with class weights
history_model_gru = model_gru.fit(
    X_train_tokenized, y_train,
    validation_data=(X_val_tokenized, y_val),
    epochs=10,
    batch_size=64,
    class_weight=class_weights_dict, # <-- added class weights
    callbacks=[es]
)

# Plot training history
plot_history(history_model_gru, 'GRU Model')

# Validation predictions and F1 score
val_preds_model_gru = (model_gru.predict(X_val_tokenized) > 0.5).astype(int).ravel()
print('GRU Validation F1:', f1_score(y_val, val_preds_model_gru))

```

Model: "sequential_25"

Layer (type)	Output Shape	Param #
embedding_25 (Embedding)	?	0 (unbuilt)
gru_16 (GRU)	?	0 (unbuilt)
dropout_25 (Dropout)	?	0
dense_25 (Dense)	?	0 (unbuilt)

Total params: 0 (0.00 B)

Trainable params: 0 (0.00 B)

Non-trainable params: 0 (0.00 B)

Epoch 1/10

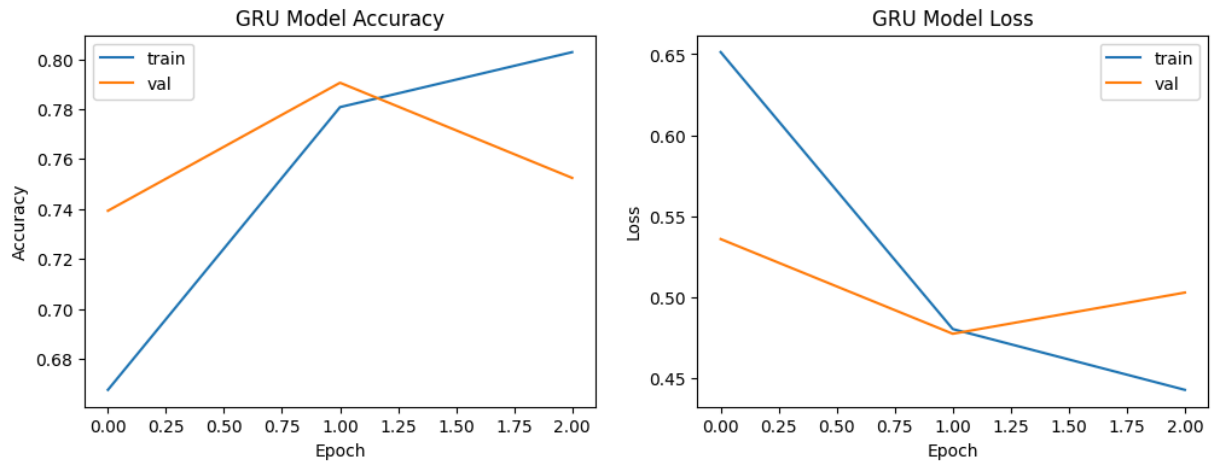
96/96 ————— **5s** 31ms/step - accuracy: 0.6677 - loss: 0.6513 - val_accuracy: 0.7393 - val_loss: 0.5358

Epoch 2/10

96/96 ————— **5s** 26ms/step - accuracy: 0.7808 - loss: 0.4802 - val_accuracy: 0.7905 - val_loss: 0.4773

Epoch 3/10

96/96 ————— **3s** 36ms/step - accuracy: 0.8028 - loss: 0.4427 - val_accuracy: 0.7525 - val_loss: 0.5028



48/48 ————— **1s** 9ms/step

GRU Validation F1: 0.6821457165732586

4.3 ROC Curve

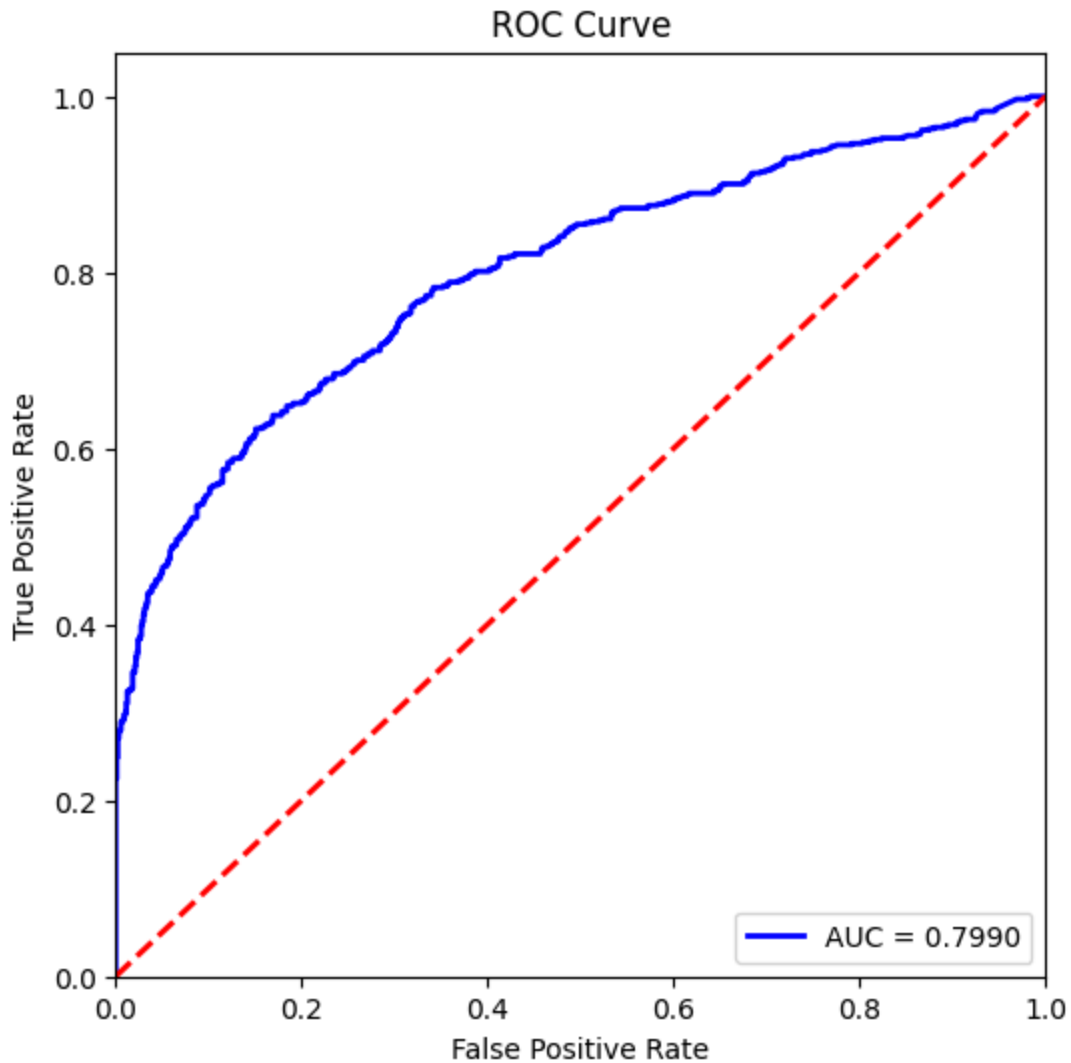
The ROC curve shows GRU model has ability to distinguish between the classes. The curve is close to the top-left corner, which means the model has a high True Positive Rate and a low False Positive Rate. The Area Under the Curve (AUC) is 0.80 close to 1.

```
In [73]: # Get predictions as probabilities
prediction_probabilities = model_gru.predict(X_val_tokenized).ravel()

# ROC curve
false_positive_rates, true_positive_rates, thresholds = roc_curve(y_val, prediction_probabilities)
roc_auc = auc(false_positive_rates, true_positive_rates)

# Plot ROC
plt.figure(figsize=(6, 6))
plt.plot(false_positive_rates, true_positive_rates, color="blue", lw=2, label=f"AUC {roc_auc}")
plt.plot([0, 1], [0, 1], color="red", lw=2, linestyle="--")
plt.title("ROC Curve")
plt.legend(loc="lower right")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.show()
```

48/48 ————— **0s** 5ms/step



4.4 Predict Test Dataset

The trained GRU model was used to predict on the test dataset. The results were submitted to Kaggle and achieved a score of 0.77106.

```
In [75]: test_preds_rnn = (model_gru.predict(X_test_tokenized) > 0.5).astype(int).ravel()

# Save to submission.csv
submission = pd.DataFrame({
    "id": test_df['id'],
    "target": test_preds_rnn
})
submission.to_csv("submission.csv", index=False)

print("submission.csv saved with", len(submission), "entries")
print(submission.head())
```

102/102 1s 5ms/step

submission.csv saved with 3263 entries

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

Natural Language Processing with Disaster Tweets

Predict which Tweets are about real disasters and which ones are not



[Overview](#) [Data](#) [Code](#) [Models](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Team](#) [Submissions](#)

Leaderboard

[Raw Data](#) [Refresh](#)

YOUR RECENT SUBMISSION

✓

submission.csv
Submitted by Peculiar Data - Submitted 4 minutes ago

Score: 0.77106

[Jump to your leaderboard position](#)

This leaderboard is calculated with all of the test data.

#	Team	Members	Score	Entries	Last	Join
433	Peculiar Data		0.77106	1	4m	

5. Conclusion

5.1 Model Performance Analysis

- The GRU is the best performing RNN-family model and fulfills the requirement.
- SimpleRNN works well but may underperform compared to LSTM and GRU.

5.2 Troubleshooting

- I spent a lot of time cleaning the text to see if removing more irrelevant words would improve the results. In fact, the F1 score got worse. So, I went back and only removed some of the stopwords. Removing too much data could lose important context or meaning to the text.

5.3 Insights from Hyperparameter Tuning

- The GRU model consistently scored around 0.78 regardless of the hyperparameters.

- The top-performing model used a batch size of 32 and a dropout rate of 0.05, achieving a validation accuracy of 0.7892.
- A model with a batch size of 64 and a dropout rate of 0.30 had a validation accuracy of 0.7886, which is close to the top score, but it trained in about half the time. This is because a larger batch size processes more samples at once, making each training step faster.
- The model with the worst performance had a batch size of 16 and a dropout rate of 0.30 with the longest training time at 84.70 seconds.
- Dropout rate did not significantly affect training time, but it did impact accuracy. Lower dropout rates, like 0.05, consistently produced better results than higher rates. This is because a low dropout rate ignores fewer neurons, which lets the model learn most of its connections and not oversimplify the data.

5.4 Takeaways

- By adjusting the batch size and dropout rate to achieve a good balance between accuracy and training time.
- Cleaning text data is important for RNN models because it gets rid of noise. However, removing too much or too little data can hurt the model's performance.

5.5 Future improvements

- To improve the model's performance in the future, I will use pre-trained GloVe embeddings with the GRU model. GloVe has already learned the meaning and relationships between millions of words from massive text sources like Wikipedia and Common Crawl. This will give the model a much better starting point than training word representations from scratch.