

Module 9 Assignment ML RNN Neural Nets – NLP with Disaster Tweets

Dan Noel

MSDS 422

May 29, 2022

For me the challenge of this week's assignment involved exploring the data and then developing the best way of cleaning the data such that it would be easy to process via a RNN/transformer encoder model.

Having gained much more knowledge of Kaggle, throughout the last ten weeks I was able to find an approach to cleaning and balancing the data that led to decent results. My first few attempts to vectorize tweets were ineffective at breaking through the median because of my inability to reduce the noise in the data.

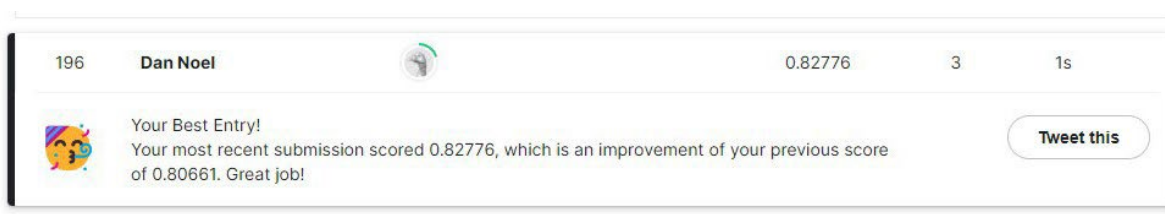
I found the following URL most helpful for explaining some of the various alternative options for building primary tools and getting started with hyperparameters.

[https://www.kaggle.com/code/tuckerarrants/disaster-tweets-eda-glove-rnns-bert/notebook#IV.-More-Complex-Model\(s\)](https://www.kaggle.com/code/tuckerarrants/disaster-tweets-eda-glove-rnns-bert/notebook#IV.-More-Complex-Model(s))

Interestingly, this approach involves a great deal of data cleaning including transforming contractions into whole words and declaring the simple LSTM model.

The important part of this assignment involved an approach could clean data (see Appendix 3). And then the actual Classification algorithm (Appendix 4)

Appendix 1 – Kaggle Metrics



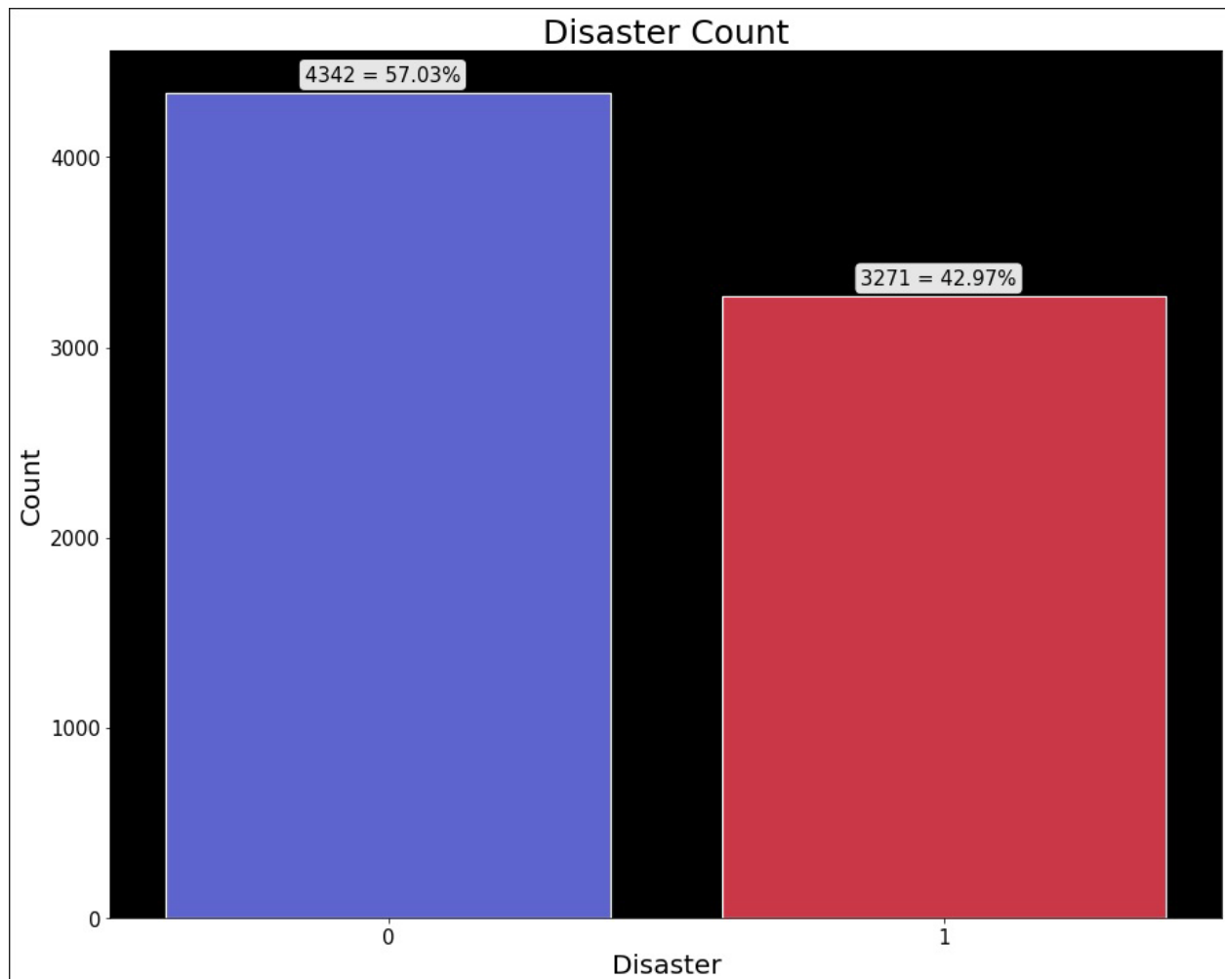
A screenshot of a Kaggle leaderboard entry. At the top, it shows the rank '196', the name 'Dan Noel', a profile picture, the score '0.82776', the number of votes '3', and the time '1s'. Below this, a notification bubble with a party emoji icon says 'Your Best Entry! Your most recent submission scored 0.82776, which is an improvement of your previous score of 0.80661. Great job!'. To the right of the notification is a 'Tweet this' button.

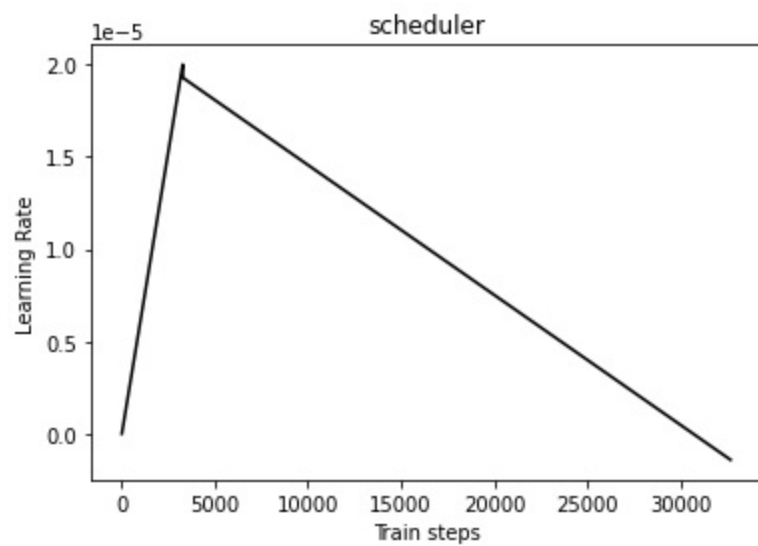
Rank	Name	Score	Votes	Time
196	Dan Noel	0.82776	3	1s

Your Best Entry!
Your most recent submission scored 0.82776, which is an improvement of your previous score of 0.80661. Great job!

[Tweet this](#)

Appendix 2 – EOD - Metrics





Appendix 3 - Cleanup Code

```
1 # %20 is the URL encoding of space, let's replace them with '_'
2 def re_encode_space(input_string):
3     return None if pd.isna(input_string) else input_string.replace('%20', '_')
4
5
6 # Let's try to find hastags
7 import re
8
9 def find_hash_tags(input_string):
10     hash_tags = re.findall(r"#(\w+)", str(input_string))
11     return ','.join(hash_tags)
12
13
14 # Let's turn hashtags to normal words
15 def re_encode_hashtags(input_string):
16     return None if pd.isna(input_string) else input_string.replace('#', '')
17
18
19 # Let's remove URLs from the tweets
20 def remove_links(input_string):
21     res = input_string
22     urls = re.findall(r'(https?://[^\s]+)', res)
23     for link in urls:
24         res = res.strip(link)
25     return res
26
27
28 # Let's remove the state abbreviations
29 def state_renaming(input_string):
30
31     states = {
32         'AK': 'Alaska',
33         'AL': 'Alabama',
34         'AR': 'Arkansas',
35         'AZ': 'Arizona',
36         'CA': 'California',
37         'CO': 'Colorado',
38         'CT': 'Connecticut',
39         'DC': 'District_of_Columbia',
40         'DE': 'Delaware',
41         'FL': 'Florida',
42         'GA': 'Georgia',
43         'HI': 'Hawaii',
44         'IA': 'Iowa',
45         'ID': 'Idaho',
46         'IL': 'Illinois',
47         'IN': 'Indiana',
48         'KS': 'Kansas',
49         'KY': 'Kentucky',
50         'LA': 'Louisiana',
51         'MA': 'Massachusetts',
52         'MD': 'Maryland',
53         'ME': 'Maine',
54         'MI': 'Michigan',
55         'MN': 'Minnesota',
56         'MO': 'Missouri',
57         'MS': 'Mississippi',
```

```
58         'MT': 'Montana',  
59         'NC': 'North_Carolina',
```

```
60     'ND': 'North_Dakota',
61     'NE': 'Nebraska',
62     'NH': 'New_Hampshire',
63     'NJ': 'New_Jersey',
64     'NM': 'New_Mexico',
65     'NV': 'Nevada',
66     'NY': 'New_York',
67     'OH': 'Ohio',
68     'OK': 'Oklahoma',
69     'OR': 'Oregon',
70     'PA': 'Pennsylvania',
71     'RI': 'Rhode_Island',
72     'SC': 'South_Carolina',
73     'SD': 'South_Dakota',
74     'TN': 'Tennessee',
75     'TX': 'Texas',
76     'UT': 'Utah',
77     'VA': 'Virginia',
78     'VT': 'Vermont',
79     'WA': 'Washington',
80     'WI': 'Wisconsin',
81     'WV': 'West_Virginia',
82     'WY': 'Wyoming'
83 }
84
85 result = input_string
86
87 if isinstance(input_string, str):
88     input_candidates = input_string.split(',')
89
90     if len(input_candidates) > 1:
91         for candidate in input_candidates:
92             if candidate in states.keys():
93                 result = states[candidate]
94
95 if input_string in states.keys():
96     result = states[input_string]
97
98 return result
```

Appendix 4 - Bert Classification Model

```

1 from tensorflow.keras.optimizers import Adam
2
3 from sklearn.metrics import roc_auc_score
4 from sklearn.metrics import accuracy_score
5
6 from sklearn.model_selection import StratifiedKFold
7
8 def build_bert_classifier():
9     input_ids = tf.keras.layers.Input(shape=(max_tweet_length,), dtype=tf.int32,
10 name="input_ids")
11     input_mask = tf.keras.layers.Input(shape=(max_tweet_length,), dtype=tf.int32,
12 name="attention_mask")
13     embeddings = bert(input_ids, attention_mask = input_mask)['pooler_output']
14     net = tf.keras.layers.Dropout(0.1)(embeddings)
15     net = tf.keras.layers.Dense(128, activation='relu', name='pre-clf')(net)
16     net = tf.keras.layers.Dropout(0.1)(net)
17     net = tf.keras.layers.Dense(1, activation=None, name='classifier')(net)
18     return tf.keras.Model(inputs=[input_ids, input_mask], outputs=net)
19
20 skf = StratifiedKFold(n_splits=5)
21 train_average_score = 0
22 validation_average_score = 0
23 validation_oof_predictions = np.zeros((len(X['input_ids'].numpy())))
24
25 loss = tf.keras.losses.BinaryCrossentropy(from_logits=True)
26 optimizer = Adam(learning_rate=schedule)
27 epochs = 5
28
29 # It's a good practice to predict the test set on every fold
30 # And average the predictions over the folds
31 averaged_test_predictions = np.zeros((test_array['input_ids'].shape[0]))
32
33 # It's standard practice to use Stratified k-fold cross validation
34 # so we're also using it here
35 for fold_n, (train_idx, test_idx) in enumerate(skf.split(X['input_ids'].numpy(), y)):
36     X_train_ids = X['input_ids'].numpy()[train_idx]
37     X_train_att = X['attention_mask'].numpy()[train_idx]
38     y_train = y[train_idx]
39
40     X_test_ids = X['input_ids'].numpy()[test_idx]
41     X_test_att = X['attention_mask'].numpy()[test_idx]
42     y_test = y[test_idx]
43
44     # Re-build the model at every fold to "reset" it
45     model = build_bert_classifier()
46     model.layers[2].trainable = True
47
48     model.compile(optimizer=optimizer,
49                 loss=loss)
50
51     model.fit(x={'input_ids':X_train_ids,'attention_mask':X_train_att},
52             y=y_train, batch_size=32, epochs=epochs)
53
54     train_predictions =
55     model.predict({'input_ids':X_train_ids,'attention_mask':X_train_att})
56     validation_predictions =
57     model.predict({'input_ids':X_test_ids,'attention_mask':X_test_att})
58
59     train_score = roc_auc_score(y_train, train_predictions)

```



```
56     validation_score = roc_auc_score(y_test, validation_predictions)
57
58     train_average_score += train_score / 5
59     validation_average_score += validation_score / 5
60 validation_oof_predictions[test_idx,] = (validation_predictions >
0.5).astype(int).flatten()
61
62     print(f'Fold: {fold_n}, train auc: {train_score:.3f}, validation auc:
{validation_score:.3f}')
63
64     test_predictions = model.predict({'input_ids':test_array['input_ids'],
65
'attention_mask':test_array['attention_mask']}).flatten()
66     averaged_test_predictions += test_predictions / 5
67
68 print(f'Train average: {train_average_score:.3f}, validation average:
{validation_average_score:.3f}')
69 print(f'OOF Accuracy Score: {accuracy_score(y, validation_oof_predictions)}')
70
71
72
```

Appendix 5 - Tokenization Code

```
1 from nltk.tokenize import TweetTokenizer
2 # The tokenizer is responsible to turn a string of words
3 # into a list of tokens (words) for which we'll get their
4 # vector representation (embeddings)
5 tknznr = TweetTokenizer(
6     preserve_case=False,
7     reduce_len=True,
8     strip_handles=True,
9 )
10
11
12 def tokenize_tweets(tokenizer, input_text):
13     tokens = list(tokenizer.tokenize(input_text))
14     tokens = [re.sub('[^A-Za-z0-9]+', '', i) for i in tokens]
15     return tokens
16
17 original_data['tokens'] = original_data['text']
18 original_data['tokens'] = original_data['tokens'].apply(lambda x:
19     tokenize_tweets(tknznr, x))
20
21 test_data['tokens'] = test_data['text'].apply(lambda x: tokenize_tweets(tknznr, x))
22
23 # We'll pad all embeddings to match the length of the biggest tweet
24 # in order to account for the variability in tweet length
25 # Later on the model is going to mask the padded values, so that
26 # they won't influence the result
27 max_tweet_length = max(original_data['tokens'].apply(lambda x: len(x)).max(),
28     test_data['tokens'].apply(lambda x: len(x)).max())
29
30 X = original_data['text'].tolist()
31 y = np.asarray(original_data['target'].tolist()).astype(np.float32)
32
33 test_array = test_data['text'].tolist()
34
```

