

Text Classification of Tweets: Are they about a real disaster or not?

Business Problem

Data has been accumulated from a number of tweets, some of which are about disasters, some of which are not. By creating a model for Natural Language Processing (NLP), we can predict whether or not a given tweet is about a real disaster or not. This can benefit companies who wish to monitor twitter in the event of an emergency.

Data Understanding

Importing necessary packages, libraries, etc.:

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 np.random.seed(42)
4 import nltk
5 nltk.download('punkt')
6 import seaborn as sns
7 import re
8 import matplotlib.pyplot as plt
9 from matplotlib.ticker import MaxNLocator
10 %matplotlib inline
11 from nltk.tokenize import word_tokenize, RegexpTokenizer
12 from sklearn.metrics import f1_score, classification_report, confusion_matrix, ConfusionMat
13 from sklearn.pipeline import Pipeline
14 from sklearn import feature_extraction, linear_model, model_selection, preprocessing
15 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer, TfidfTransfor
16 from nltk.corpus import stopwords
17 from sklearn.model_selection import train_test_split
18 from nltk import FreqDist
19 from sklearn.naive_bayes import MultinomialNB
20 from nltk.stem.snowball import SnowballStemmer
21 from sklearn.model_selection import GridSearchCV
22 from nltk.corpus import stopwords, wordnet
23 from nltk.stem import WordNetLemmatizer
24
25
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   /Users/nicolemichaud/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

Loading the data:

```
In [2]: 1 train_df = pd.read_csv("data/train.csv")
2 test_df = pd.read_csv("data/test.csv")
```

Data Exploration:

Viewing and gaining understanding of the data, its features, number of rows, any missing values, and more so I can preprocess the data accordingly.

In [3]: 1 train_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]: 1 train_df.info()

```
<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 [5]: 1 test_df.info()

```
<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 [6]: 1 *# I won't be working with the 'location' or 'keyword' columns, so I'll just drop them*
2 train_df = train_df.drop(columns = ['location', 'keyword'])
3 test_df = test_df.drop(columns = ['location', 'keyword'])
4 train_df.head()

Out[6]:

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

```
In [7]: 1 train_df['text'].dropna(inplace=True)
        2 train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7613 entries, 0 to 7612
Data columns (total 3 columns):
#   Column   Non-Null Count  Dtype
---  ---
0    id      7613 non-null   int64
1    text     7613 non-null   object
2    target   7613 non-null   int64
dtypes: int64(2), object(1)
memory usage: 178.6+ KB
```

```
In [8]: 1 test_df['text'].dropna(inplace=True)
        2 test_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3263 entries, 0 to 3262
Data columns (total 2 columns):
#   Column   Non-Null Count  Dtype
---  ---
0    id      3263 non-null   int64
1    text     3263 non-null   object
dtypes: int64(1), object(1)
memory usage: 51.1+ KB
```

```
In [9]: 1 # Example of what is NOT a disaster tweet:
        2 train_df[train_df["target"] == 0]["text"].values[6]
```

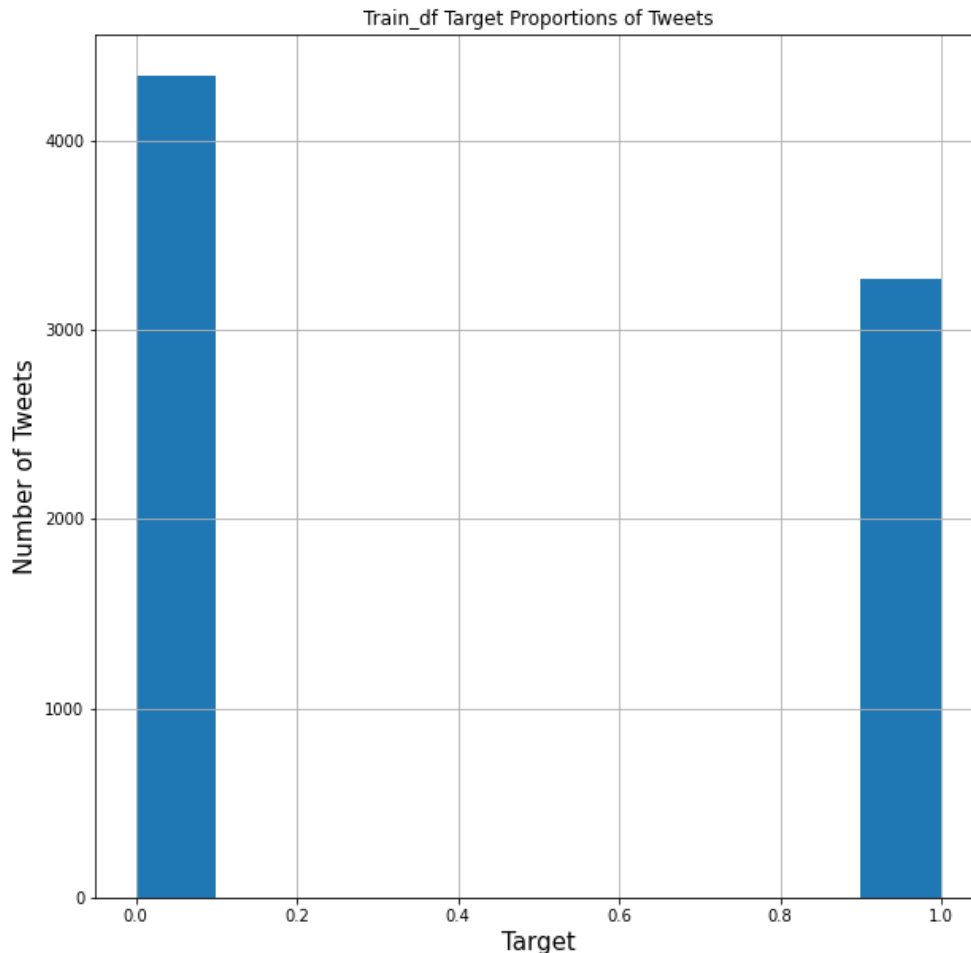
```
Out[9]: 'London is cool ;)'
```

```
In [10]: 1 # Example of what IS a disaster tweet:
         2 train_df[train_df["target"] == 1]["text"].values[20]
```

```
Out[10]: 'Deputies: Man shot before Brighton home set ablaze http://t.co/gWNRhMS08k' (http://t.co/gWNRhMS08k)
```

Visualizing what proportion of the training data are disaster tweets and non-disaster tweets:

```
In [11]: 1 fig, ax = plt.subplots(figsize=(10,10))
2 proportions = train_df['target'].hist()
3 plt.xlabel('Target',fontsize=15)
4 plt.ylabel('Number of Tweets',fontsize=15)
5 plt.title('Train_df Target Proportions of Tweets');
6
7 #fig, ax = plt.subplots(figsize=(10,10))
8 #cm_1 = ConfusionMatrixDisplay(confusion_matrix = cnf_matrix, display_labels = baseline_tree
9 #cm_1.plot(cmap=plt.cm.Greens, ax=ax)
10 #plt.xlabel('Prediction label',fontsize=15)
11 #plt.ylabel('True label',fontsize=15);
```



Calculating the probabilities of disaster and non-disaster tweets in the training data:

```
In [12]: 1 disaster_tweets = train_df[train_df['target']==1]
2
3 other_tweets = train_df[train_df['target']==0]

In [13]: 1 P_disasters = len(disaster_tweets) / (len(disaster_tweets)+len(other_tweets))
2 P_non = len(other_tweets) / (len(other_tweets)+len(disaster_tweets))
3 print(P_disasters)
4 print(P_non)
```

```
0.4296597924602653
0.5703402075397347
```

This tells us that tweets in train_df have a higher probability of not being about a disaster.

Data Preparation

Cleaning text data: Remove urls, tags (contain @), stopwords, punctuation, etc.

```
In [14]: 1 # Creating a function to perform all these cleaning steps at once
2 stopwords_list = stopwords.words('english')
3
4 no_bad_chars = re.compile('[!\"#$%&()*+,-./:;<=>?@[\\]^_`{|}~\\n - ]')
5 no_nums = re.compile('[\\d-]')
6
7 def clean_text(text):
8     text = no_nums.sub('', text)
9     #drop words '&' and 'via'
10    text = re.sub("&", "", text)
11    text = re.sub("via", "", text)
12    text = re.sub("@[A-Za-z0-9]+", "", text) #Remove @ sign
13    text = re.sub(r"(?:\\@|http?:\\/|https?:\\/|www)\\S+", "", text) #Remove http links
14    text = no_bad_chars.sub('', text)
15    text = text.replace("#", "").replace("_", " ") #Remove hashtag sign but keep the text
16    text = text.lower()
17    text = ' '.join(word for word in text.split() if word not in stopwords_list)
18    return text
19
20
21 train_df_cleaned = train_df['text'].apply(clean_text)
22 test_df_cleaned = test_df['text'].apply(clean_text)
23 train_df_cleaned.head(10)
24
25
26
```

```
Out[14]: 0      deeds reason earthquake may allah forgive us
1      forest fire near la ronge sask canada
2      residents asked 'shelter place' notified offic...
3      people receive wildfires evacuation orders cal...
4      got sent photo ruby alaska smoke wildfires pou...
5      rockyfire update california hwy closed directi...
6      flood disaster heavy rain causes flash floodin...
7      i'm top hill see fire woods
8      there's emergency evacuation happening buildin...
9      i'm afraid tornado coming area
Name: text, dtype: object
```

```
In [15]: 1 #Performing a train-test split on the training data to see how our models perform
2 #before applying them to our testing data
3
4 X = train_df.text
5 y = train_df.target
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state = 42)
```

```
In [16]: 1 #Applying the text cleaning function to our data
2
3 X_train_cleaned = X_train.apply(clean_text)
4 X_test_cleaned = X_test.apply(clean_text)
```

Modeling

Building a baseline model

```
In [17]: 1 baseline_model = Pipeline([('vect', CountVectorizer(max_features=None,
2                                                         tokenizer=word_tokenize,
3                                                         stop_words=stopwords_list)),
4                                                         ('clf', MultinomialNB())
5                                                         ])
6 baseline_model.fit(X_train_cleaned, y_train)
7
8
9 y_pred = baseline_model.predict(X_test_cleaned)
10
11 print('Baseline model F1 %s' % f1_score(y_pred, y_test, average="macro"))
12 print(classification_report(y_test, y_pred))
```

/Users/nicolemichaud/opt/anaconda3/envs/learn-env/lib/python3.8/site-packages/sklearn/feature_extraction/text.py:383: UserWarning: Your stop_words may be inconsistent with your preprocessing. Tokenizing the stop words generated tokens [''d'', ''ll'', ''re'', ''s'', ''ve'', 'could', 'might', 'must', 'n't', 'need', 'sha', 'wo', 'would'] not in stop_words.
warnings.warn('Your stop_words may be inconsistent with ')

```
Baseline model F1 0.799114882744776
```

	precision	recall	f1-score	support
0	0.81	0.86	0.84	1091
1	0.80	0.73	0.76	813
accuracy			0.81	1904
macro avg	0.80	0.80	0.80	1904
weighted avg	0.81	0.81	0.80	1904

Stemming the text to see if it improves our model:

```
In [18]: 1 #when stemming/lemmatizing, we are going to do the same to the stopwords,
2 #so we don't want to remove them before this.
3
4 #Creating a function that cleans the text data without removing stopwords:
5
6 def clean_text_nostop(text):
7     text = no_nums.sub('', text)
8     text = re.sub("[A-Za-z0-9]+", "", text)
9     text = re.sub(r"(?:\@|http?\://|https?\://|www)\S+", "", text)
10    text = no_bad_chars.sub('', text)
11    text = text.replace("#", "").replace("_", " ")
12    text = text.lower()
13    return text
14
15 train_df_cleaned_nostop = train_df.copy()
16
17 train_df_cleaned_nostop['text'] = train_df_cleaned_nostop['text'].apply(clean_text_nostop)
18
19 X_nostop = train_df_cleaned_nostop.text
20 y_nostop = train_df_cleaned_nostop.target
21 X_train_nostop, X_test_nostop, y_train_nostop, y_test_nostop = train_test_split(X_nostop, y_nostop,
22
```

```
In [19]: 1 stemmer = SnowballStemmer(language="english")
2 tokenizer=word_tokenize
3
4 def stem_and_tokenize(document):
5     tokens = tokenizer(document)
6     return [stemmer.stem(token) for token in tokens]
7 stemmed_stopwords = [stemmer.stem(word) for word in stopwords_list]
```

```
In [20]: 1 #Stemmed data model
2
3 stem_model = Pipeline([('vect', CountVectorizer(
4     stop_words=stemmed_stopwords,
5     tokenizer=stem_and_tokenize)),
6     ('clf', MultinomialNB()),
7 ])
8 stem_model.fit(X_train_nostop, y_train_nostop)
9
10
11 y_pred_stem= stem_model.predict(X_test_nostop)
12
13 print('Stemmed model F1 %s' % f1_score(y_pred_stem, y_test_nostop, average="macro"))
14 print(classification_report(y_test_nostop, y_pred_stem))
```

/Users/nicolemichaud/opt/anaconda3/envs/learn-env/lib/python3.8/site-packages/sklearn/feature_extraction/text.py:383: UserWarning: Your stop_words may be inconsistent with your preprocessing. Tokenizing the stop words generated tokens ['"', "'d", 'could', 'might', 'must', "n't", 'n eed', 'r', 'sha', 'v', 'wo', 'would'] not in stop_words.

warnings.warn('Your stop_words may be inconsistent with '

```
Stemmed model F1 0.802174272578552
```

	precision	recall	f1-score	support
0	0.82	0.86	0.84	1091
1	0.80	0.74	0.77	813
accuracy			0.81	1904
macro avg	0.81	0.80	0.80	1904
weighted avg	0.81	0.81	0.81	1904

Stemming improved our model.

Conducting GridSearchCV to see if tuning the hyperparameters in our best model will improve it further:

```
In [21]: 1 # First, need to manually tokenize/vectorize data since we won't be using a pipeline for th.
2
3 vectorizer = CountVectorizer()
4 X_train_vectorized = vectorizer.fit_transform(X_train_cleaned)
5 X_test_vectorized = vectorizer.transform(X_test_cleaned)
6
7 cv = CountVectorizer()
8 X_train_vec = cv.fit_transform(X_train_cleaned)
9 X_train_vec = pd.DataFrame.sparse.from_spmatrix(X_train_vec)
10 X_train_vec.columns = sorted(cv.vocabulary_)
11 X_train_vec.set_index(y_train.index, inplace=True)
12
13
14 X_test_vec = cv.transform(X_test_cleaned)
15 X_test_vec = pd.DataFrame.sparse.from_spmatrix(X_test_vec)
16 X_test_vec.columns = sorted(cv.vocabulary_)
17 X_test_vec.set_index(y_test.index, inplace=True)
```

```
In [22]: 1 #GridSearchCV
2 alphas = [0.5, 1.0, 1.5, 2.0, 2.5]
3 p_grid_NB = {'alpha': alphas, 'fit_prior' : [True, False]}
4 NB_cls= MultinomialNB()
5
6 grid = GridSearchCV(estimator = NB_cls, param_grid = p_grid_NB, scoring = 'f1', cv = 3)
7 grid.fit(X_train_vec, y_train)
```

```
Out[22]: GridSearchCV(cv=3, estimator=MultinomialNB(),
                    param_grid={'alpha': [0.5, 1.0, 1.5, 2.0, 2.5],
                                'fit_prior': [True, False]},
                    scoring='f1')
```

```
In [23]: 1 grid.best_params_
```

```
Out[23]: {'alpha': 2.5, 'fit_prior': True}
```

```
In [24]: 1 tuned_model = Pipeline([('vect', CountVectorizer(
2     stop_words=stemmed_stopwords,
3     tokenizer=stem_and_tokenize,
4 )),
5
6     ('clf', MultinomialNB(alpha= 2.0, fit_prior = True)),
7
8 ])
9 tuned_model.fit(X_train_nostop, y_train_nostop)
10
11 y_pred_tuned= tuned_model.predict(X_test_nostop)
12
13 print('Tuned model F1 %s' % f1_score(y_pred_tuned, y_test_nostop, average="macro"))
14
15 print(classification_report(y_test_nostop, y_pred_tuned))
```

/Users/nicolemichaud/opt/anaconda3/envs/learn-env/lib/python3.8/site-packages/sklearn/feature_extraction/text.py:383: UserWarning: Your stop_words may be inconsistent with your preprocessing. Tokenizing the stop words generated tokens ['"', "'d", 'could', 'might', 'must', 'n't', 'n'eed', 'r', 'sha', 'v', 'wo', 'would'] not in stop_words.
warnings.warn('Your stop_words may be inconsistent with '

```
Tuned model F1 0.8022946797889032
              precision    recall  f1-score   support

         0           0.81       0.86       0.84       1091
         1           0.80       0.73       0.77        813

 accuracy          0.81
 macro avg         0.81       0.80       0.80       1904
 weighted avg      0.81       0.81       0.81       1904
```


In [25]:

```

1 sns.set_context("poster")
2 final_clf_report = classification_report(y_test_nostop, y_pred_tuned)
3 # The following code is not mine, I adapted it from a stackoverflow user on how to create a
4 def show_values(pc, fmt="%.2f", **kw):
5     '''
6     Heatmap with text in each cell with matplotlib's pyplot
7     Source: https://stackoverflow.com/a/25074150/395857
8     By HYRY
9     '''
10    pc.update_scalarmappable()
11    ax = pc.axes
12    #ax = pc.axes# FOR LATEST MATPLOTLIB
13    #Use zip BELOW IN PYTHON 3
14    for p, color, value in zip(pc.get_paths(), pc.get_facecolors(), pc.get_array()):
15        x, y = p.vertices[:-2, :].mean(0)
16        if np.all(color[:3] > 0.5):
17            color = (0.0, 0.0, 0.0)
18        else:
19            color = (1.0, 1.0, 1.0)
20        ax.text(x, y, fmt % value, ha="center", va="center", color=color, **kw)
21
22
23 def cm2inch(*tupl):
24     '''
25     Specify figure size in centimeter in matplotlib
26     Source: https://stackoverflow.com/a/22787457/395857
27     By gns-ank
28     '''
29     inch = 2.54
30     if type(tupl[0]) == tuple:
31         return tuple(i/inch for i in tupl[0])
32     else:
33         return tuple(i/inch for i in tupl)
34
35
36 def heatmap(AUC, title, xlabel, ylabel, xticklabels, yticklabels, figure_width=40, figure_h
37     '''
38     Inspired by:
39     - https://stackoverflow.com/a/16124677/395857
40     - https://stackoverflow.com/a/25074150/395857
41     '''
42
43     # Plot it out
44     fig, ax = plt.subplots()
45     #c = ax.pcolor(AUC, edgecolors='k', linestyle= 'dashed', linewidths=0.2, cmap='RdBu', v
46     c = ax.pcolor(AUC, edgecolors='k', linestyle= 'dashed', linewidths=0.3, cmap=cmap, vmir
47
48     # put the major ticks at the middle of each cell
49     ax.set_yticks(np.arange(AUC.shape[0]) + 0.5, minor=False)
50     ax.set_xticks(np.arange(AUC.shape[1]) + 0.5, minor=False)
51
52     # set tick labels
53     #ax.set_xticklabels(np.arange(1,AUC.shape[1]+1), minor=False)
54     ax.set_xticklabels(xticklabels, minor=False)
55     ax.set_yticklabels(yticklabels, minor=False)
56
57     # set title and x/y labels
58     plt.title(title, y=1.25)
59     plt.xlabel(xlabel)
60     plt.ylabel(ylabel)
61
62     # Remove last blank column
63     plt.xlim( 0, AUC.shape[1] )
64
65     # Turn off all the ticks
66     ax = plt.gca()
67     for t in ax.xaxis.get_major_ticks():
68         t.tick1line.set_visible(False)
69         t.tick2line.set_visible(False)
70     for t in ax.yaxis.get_major_ticks():
71         t.tick1line.set_visible(False)
72         t.tick2line.set_visible(False)

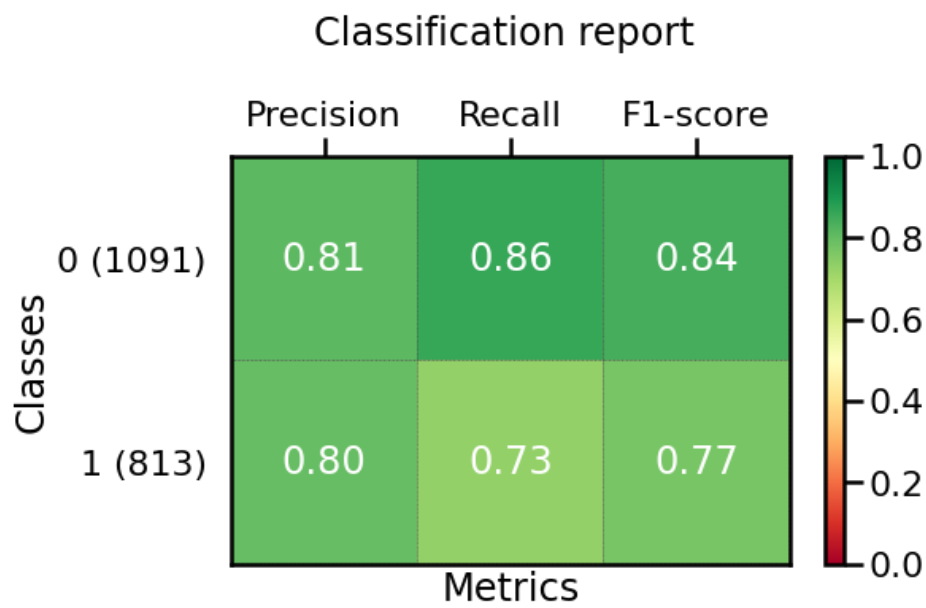
```

```

73
74 # Add color bar
75 plt.colorbar(c)
76
77 # Add text in each cell
78 show_values(c)
79
80 # Proper orientation (origin at the top left instead of bottom left)
81 if correct_orientation:
82     ax.invert_yaxis()
83     ax.xaxis.tick_top()
84
85 # resize
86 fig = plt.gcf()
87 #fig.set_size_inches(cm2inch(40, 20))
88 #fig.set_size_inches(cm2inch(40*4, 20*4))
89 fig.set_size_inches(cm2inch(ffigure_width, figure_height))
90
91
92
93 def plot_classification_report(classification_report, number_of_classes=2, title='Classific
94 '''
95 Plot scikit-learn classification report.
96 Extension based on https://stackoverflow.com/a/31689645/395857
97 '''
98 lines = classification_report.split('\n')
99
100 #drop initial lines
101 lines = lines[2:]
102
103 classes = []
104 plotMat = []
105 support = []
106 class_names = []
107 for line in lines[: number_of_classes]:
108     t = list(filter(None, line.strip().split(' ')))
109     if len(t) < 4: continue
110     classes.append(t[0])
111     v = [float(x) for x in t[1: len(t) - 1]]
112     support.append(int(t[-1]))
113     class_names.append(t[0])
114     plotMat.append(v)
115
116
117 xlabel = 'Metrics'
118 ylabel = 'Classes'
119 xticklabels = ['Precision', 'Recall', 'F1-score']
120 yticklabels = ['{0} ({1})'.format(class_names[idx], sup) for idx, sup in enumerate(sup
121 figure_width = 20
122 figure_height = len(class_names) + 10
123 correct_orientation = True
124 heatmap(np.array(plotMat), title, xlabel, ylabel, xticklabels, yticklabels, figure_widt
125 plt.show()
126
127

```

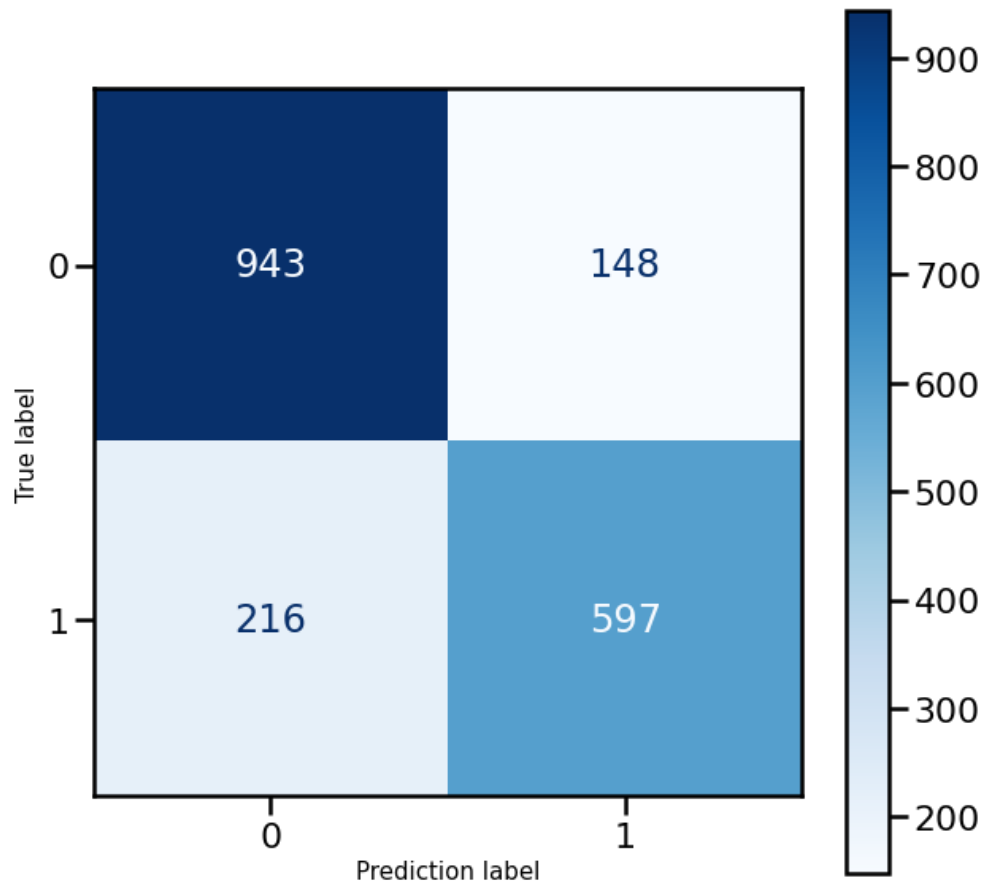
```
In [26]: 1 plot_classification_report(final_clf_report)
```



```
In [27]: 1 cnf_matrix = confusion_matrix(y_test_nostop, y_pred_tuned)
2 print('Confusion Matrix:\n', cnf_matrix)
```

Confusion Matrix:
[[943 148]
 [216 597]]

```
In [28]: 1 #CM Visualization:
2 fig, ax = plt.subplots(figsize=(10,10))
3 cm_1 = ConfusionMatrixDisplay(confusion_matrix = cnf_matrix, display_labels = tuned_model.c
4 cm_1.plot(cmap=plt.cm.Blues, ax=ax)
5 plt.xlabel('Prediction label',fontsize=15)
6 plt.ylabel('True label',fontsize=15);
```



This model gave 943 True Negatives, 597 True Positives, 148 False Positives, and 216 False Negatives.

Both false negatives and false positives are costly in this instance.

Generating Predictions

```
In [29]: 1 sample_submission = pd.read_csv("data/sample_submission.csv")
```

```
In [30]: 1 sample_submission.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3263 entries, 0 to 3262
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  ---
0    id      3263 non-null     int64
1   target  3263 non-null     int64
dtypes: int64(2)
memory usage: 51.1 KB
```

Before creating predictions on the sample submission csv, I want to try it on the test_df to make sure it works:

```
In [31]: 1 #Cleaning the test data:
2 def clean_text_nostop(text):
3     text = no_nums.sub('', text)
4     text = re.sub("@[A-Za-z0-9]+","",text) #Remove @ sign
5     text = re.sub(r"(?:\@|http?\:\/\/|https?\:\/\/|www)\S+", "", text) #Remove http links
6     text = no_bad_chars.sub(' ', text)
7     text = text.replace("#", "").replace("_", " ") #Remove hashtag sign but keep the text
8     text = text.lower()
9     return text
10
11 test_df_cleaned_nostop = test_df.copy()
12
13 test_df_cleaned_nostop['text'] = test_df_cleaned_nostop['text'].apply(clean_text_nostop)
14
```

```
In [32]: 1 test_df_cleaned_nostop.head()
```

Out[32]:

	id	text
0	0	just happened a terrible car crash
1	2	heard about earthquake is different cities s...
2	3	there is a forest fire at spot pond geese are...
3	9	apocalypse lighting spokane wildfires
4	11	typhoon soudelor kills in china and taiwan

```
In [33]: 1 test_df_sample = test_df_cleaned_nostop.copy()
2 test_df_sample['target'] = tuned_model.predict(test_df_sample['text'])
```

```
In [34]: 1 test_df_sample.head()
```

Out[34]:

	id	text	target
0	0	just happened a terrible car crash	1
1	2	heard about earthquake is different cities s...	1
2	3	there is a forest fire at spot pond geese are...	1
3	9	apocalypse lighting spokane wildfires	1
4	11	typhoon soudelor kills in china and taiwan	1

```
In [35]: 1 test_df_sample['target'].value_counts()
```

```
Out[35]: 0    2033
1    1230
Name: target, dtype: int64
```

This appears to have worked, so let's try it on the sample submission:

```
In [36]: 1 sample_submission["target"] = tuned_model.predict(test_df['text'])
```

```
In [37]: 1 sample_submission['target'].value_counts()
```

```
Out[37]: 0    2031
1    1232
Name: target, dtype: int64
```

Conclusion

Recommendations:

- Deploy model

- Continue testing other models to see if performance can be improved

Next Steps:

- Use more tweets data
- Try other types of classifiers/models