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

*Project by Nicole Michaud, 12/30/2023*

## 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, confusi
         13  from sklearn.pipeline import Pipeline
         14  from sklearn import feature_extraction, linear_model, model_selectio
         15  from sklearn.feature_extraction.text import TfidfVectorizer, CountVe
         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]:
```python
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]:
```python
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]:
```python
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]:
```python
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
```

Dropping the 'keyword' and 'location' columns, as I won't be working with them. This project focuses on the text of the tweet. See miscellaneous notebook for investigation into the 'keyword' feature.

In [6]:
```
1  train_df = train_df.drop(columns = ['location', 'keyword'])
2  test_df = test_df.drop(columns = ['location', 'keyword'])
3  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 |

There doesn't appear to be any null values in the text column, but just to be sure I will drop null values in both the training and testing data.

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
```

Previewing a random tweet from both datasets to get an idea of how they might look before cleaning:

In [9]:
```python
1  # Example of what is NOT a disaster tweet:
2  train_df[train_df["target"] == 0]["text"].values[6]
```
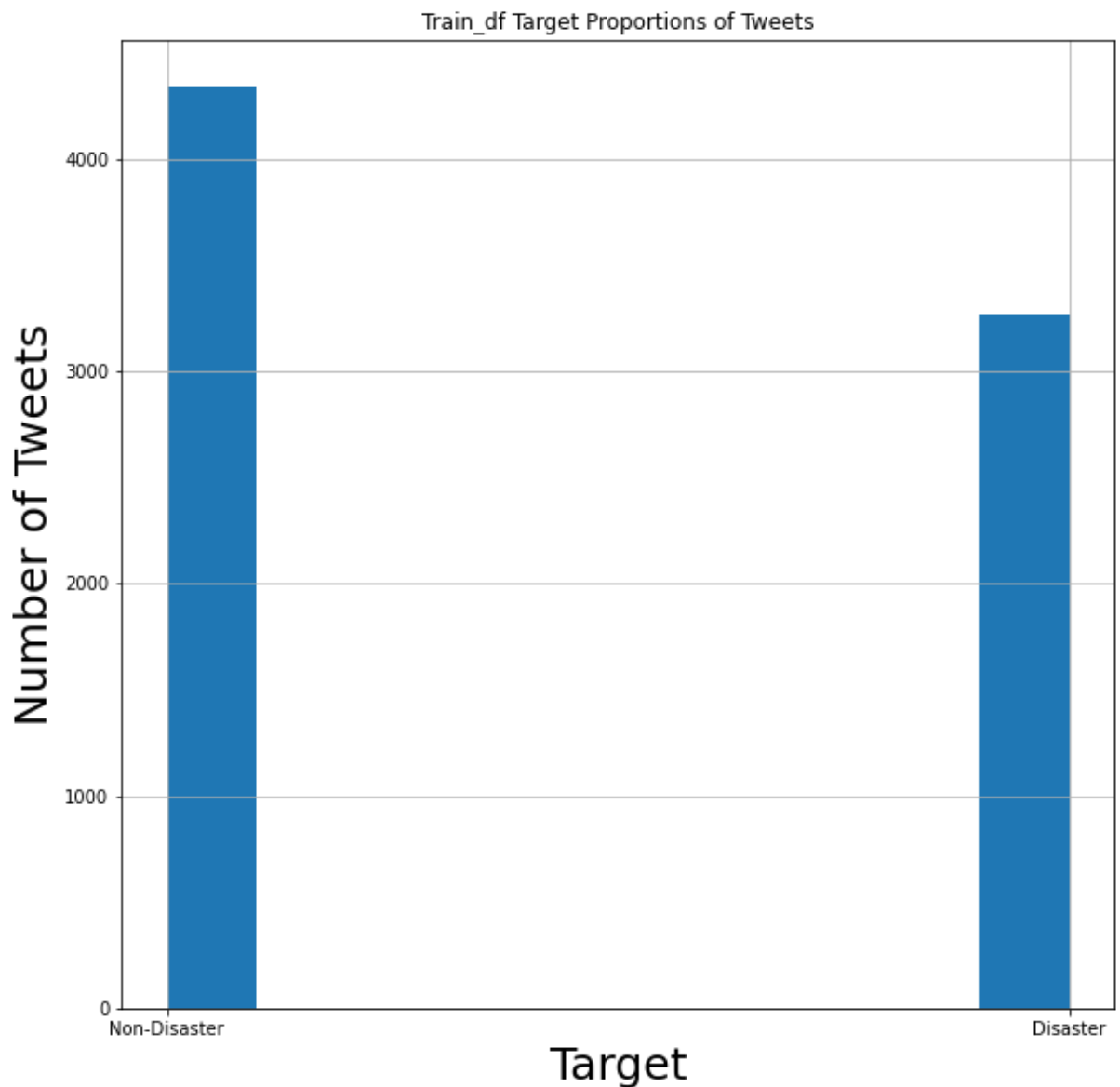
Out[9]: 'London is cool ;)'

In [10]:
```python
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/gWNRhMS
O8k' (http://t.co/gWNRhMSO8k')

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

In [11]:

```python
#from matplotlib import ticker
fig, ax = plt.subplots(figsize=(10,10))
proportions = train_df['target'].hist()

plt.locator_params(axis='x', nbins=2)

labels = [item.get_text() for item in ax.get_xticklabels()]
labels[1] = 'Non-Disaster'
labels[2] = 'Disaster'

ax.set_xticklabels(labels)

plt.xlabel('Target',fontsize=25)
plt.ylabel('Number of Tweets',fontsize=25)
plt.title('Train_df Target Proportions of Tweets');
```

```
<ipython-input-11-b887d99680f8>:11: UserWarning: FixedFormatter should
only be used together with FixedLocator
  ax.set_xticklabels(labels)
```

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_
            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 '&amp' and 'via', see miscellaneous notebook for exp
           10      text = re.sub("&amp", "", 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) #
           14      text = no_bad_chars.sub(' ', text)
           15      text = text.replace("#", "").replace("_", " ") #Remove hashtag s
           16      text = text.lower()
           17      text = ' '.join(word for word in text.split() if word not in stop
           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
```

Performing a train-test split on the ***training data*** to see how our models perform before applying them to our testing data and then applying the cleaning function to the split data:

```
In [15]:   1  X = train_df.text
           2  y = train_df.target
           3  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=(
```

```
In [16]:   1  X_train_cleaned = X_train.apply(clean_text)
           2  X_test_cleaned = X_test.apply(clean_text)
```

# Modeling

## Building a baseline model

In [17]:
```python
1  # Baseline Multinomial Naive Bayes model vectorized with CountVector.
2  # no tuned parameters
3  baseline_model = Pipeline([('vect', CountVectorizer(max_features=Non
4                                            tokenizer=word_to
5                                            stop_words=stopwo
6                         ('clf', MultinomialNB())
7                 ])
8  baseline_model.fit(X_train_cleaned, y_train)
9
10
11 y_pred = baseline_model.predict(X_test_cleaned)
12
13 print('Baseline model F1 %s' % f1_score(y_pred, y_test, average="mac
14 print(classification_report(y_test, y_pred))
```

```
/Users/nicolemichaud/opt/anaconda3/envs/learn-env/lib/python3.8/site-pa
ckages/sklearn/feature_extraction/text.py:383: UserWarning: Your stop_w
ords may be inconsistent with your preprocessing. Tokenizing the stop w
ords generated tokens ["'d", "'ll", "'re", "'s", "'ve", 'could', 'migh
t', '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:

First, need to re-create the text cleaning function to not remove the stopwords, because we will want to have the stopwords stemmed as well before removal for consistency.

In [18]:
```python
#Creating a function that cleans the text data without removing stop

def clean_text_nostop(text):
    text = no_nums.sub('', text)
    text = re.sub("@[A-Za-z0-9]+","",text)
    text = re.sub(r"(?:\@|http?\://|https?\://|www)\S+", "", text)
    text = no_bad_chars.sub(' ', text)
    text = text.replace("#", "").replace("_", " ")
    text = text.lower()
    return text

train_df_cleaned_nostop = train_df.copy()

train_df_cleaned_nostop['text'] = train_df_cleaned_nostop['text'].ap

X_nostop = train_df_cleaned_nostop.text
y_nostop = train_df_cleaned_nostop.target
X_train_nostop, X_test_nostop, y_train_nostop, y_test_nostop = train
```

In [19]:
```python
#Initializing the stemmer and creating a list of stemmed stopwords
stemmer = SnowballStemmer(language="english")
tokenizer=word_tokenize

def stem_and_tokenize(document):
    tokens = tokenizer(document)
    return [stemmer.stem(token) for token in tokens]
stemmed_stopwords = [stemmer.stem(word) for word in stopwords_list]
```

In [20]:

```python
#Stemmed data model

stem_model = Pipeline([('vect', CountVectorizer(
                            stop_words=stemmed_stopwords,
                            tokenizer=stem_and_tokenize)),
                ('clf', MultinomialNB()),
                ])
stem_model.fit(X_train_nostop, y_train_nostop)


y_pred_stem= stem_model.predict(X_test_nostop)

print('Stemmed model F1 %s' % f1_score(y_pred_stem, y_test_nostop, a
print(classification_report(y_test_nostop, y_pred_stem))
```

```
/Users/nicolemichaud/opt/anaconda3/envs/learn-env/lib/python3.8/site-pa
ckages/sklearn/feature_extraction/text.py:383: UserWarning: Your stop_w
ords may be inconsistent with your preprocessing. Tokenizing the stop w
ords generated tokens ["'", "'d", 'could', 'might', 'must', "n't", 'nee
d', '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]:
```python
# First, need to manually tokenize/vectorize data since we won't be

vectorizer = CountVectorizer()
X_train_vectorized = vectorizer.fit_transform(X_train_cleaned)
X_test_vectorized = vectorizer.transform(X_test_cleaned)

cv = CountVectorizer()
X_train_vec = cv.fit_transform(X_train_cleaned)
X_train_vec  = pd.DataFrame.sparse.from_spmatrix(X_train_vec)
X_train_vec.columns = sorted(cv.vocabulary_)
X_train_vec.set_index(y_train.index, inplace=True)


X_test_vec = cv.transform(X_test_cleaned)
X_test_vec  = pd.DataFrame.sparse.from_spmatrix(X_test_vec)
X_test_vec.columns = sorted(cv.vocabulary_)
X_test_vec.set_index(y_test.index, inplace=True)
```

Multinomial Naive Bayes doesn't have many tunable parameters. See miscellaneous notebook for explanation of the parameters and parameter values used here.

In [22]:
```python
#GridSearchCV
alphas = [0.5, 1.0, 1.5, 2.0, 2.5]
p_grid_NB = {'alpha': alphas, 'fit_prior' : [True, False]}
NB_cls= MultinomialNB()

grid = GridSearchCV(estimator = NB_cls, param_grid = p_grid_NB, scor
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]:
```python
grid.best_params_
```

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

Note: The GridSearch Determined the best value of alpha to be 2.5 and of fit_prior to be True. I ran the model below with alpha=2.5 and it did not perform as well as alpha=2.0. I'm not sure the reason for this, but this is why the value 2.0 is used instead.

In [24]:
```python
tuned_model = Pipeline([('vect', CountVectorizer(
    stop_words=stemmed_stopwords,
    tokenizer=stem_and_tokenize,
)),

                        ('clf', MultinomialNB(alpha= 2.0,fit_prior =
            ])

tuned_model.fit(X_train_nostop, y_train_nostop)

y_pred_tuned= tuned_model.predict(X_test_nostop)

print('Tuned model F1 %s' % f1_score(y_pred_tuned, y_test_nostop, av

print(classification_report(y_test_nostop, y_pred_tuned))
```

```
/Users/nicolemichaud/opt/anaconda3/envs/learn-env/lib/python3.8/site-pa
ckages/sklearn/feature_extraction/text.py:383: UserWarning: Your stop_w
ords may be inconsistent with your preprocessing. Tokenizing the stop w
ords generated tokens ["'", "'d", 'could', 'might', 'must', "n't", 'nee
d', '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      1904
   macro avg       0.81      0.80      0.80      1904
weighted avg       0.81      0.81      0.81      1904
```

Creating a visualization of the classification report:

In [25]:

```python
# The following code is not mine, I adapted it from a stackoverflow

sns.set_context("poster")
final_clf_report = classification_report(y_test_nostop, y_pred_tuned
def show_values(pc, fmt="%.2f", **kw):
    '''
    Heatmap with text in each cell with matplotlib's pyplot
    Source: https://stackoverflow.com/a/25074150/395857
    By HYRY
    '''
    pc.update_scalarmappable()
    ax = pc.axes
    #ax = pc.axes# FOR LATEST MATPLOTLIB
    #Use zip BELOW IN PYTHON 3
    for p, color, value in zip(pc.get_paths(), pc.get_facecolors(),
        x, y = p.vertices[:-2, :].mean(0)
        if np.all(color[:3] > 0.5):
            color = (0.0, 0.0, 0.0)
        else:
            color = (1.0, 1.0, 1.0)
        ax.text(x, y, fmt % value, ha="center", va="center", color=c


def cm2inch(*tupl):
    '''
    Specify figure size in centimeter in matplotlib
    Source: https://stackoverflow.com/a/22787457/395857
    By gns-ank
    '''
    inch = 2.54
    if type(tupl[0]) == tuple:
        return tuple(i/inch for i in tupl[0])
    else:
        return tuple(i/inch for i in tupl)


def heatmap(AUC, title, xlabel, ylabel, xticklabels, yticklabels, fi
    '''
    Inspired by:
    - https://stackoverflow.com/a/16124677/395857
    - https://stackoverflow.com/a/25074150/395857
    '''

    # Plot it out
    fig, ax = plt.subplots()
    #c = ax.pcolor(AUC, edgecolors='k', linestyle= 'dashed', linewid
    c = ax.pcolor(AUC, edgecolors='k', linestyle= 'dashed', linewidt

    # put the major ticks at the middle of each cell
    ax.set_yticks(np.arange(AUC.shape[0]) + 0.5, minor=False)
    ax.set_xticks(np.arange(AUC.shape[1]) + 0.5, minor=False)

    # set tick labels
    #ax.set_xticklabels(np.arange(1,AUC.shape[1]+1), minor=False)
    ax.set_xticklabels(xticklabels, minor=False)
    ax.set_yticklabels(yticklabels, minor=False)
```

```python
58        # set title and x/y labels
59        plt.title(title, y=1.25)
60        plt.xlabel(xlabel)
61        plt.ylabel(ylabel)
62
63        # Remove last blank column
64        plt.xlim( (0, AUC.shape[1]) )
65
66        # Turn off all the ticks
67        ax = plt.gca()
68        for t in ax.xaxis.get_major_ticks():
69            t.tick1line.set_visible(False)
70            t.tick2line.set_visible(False)
71        for t in ax.yaxis.get_major_ticks():
72            t.tick1line.set_visible(False)
73            t.tick2line.set_visible(False)
74
75        # Add color bar
76        plt.colorbar(c)
77
78        # Add text in each cell
79        show_values(c)
80
81        # Proper orientation (origin at the top left instead of bottom l
82        if correct_orientation:
83            ax.invert_yaxis()
84            ax.xaxis.tick_top()
85
86        # resize
87        fig = plt.gcf()
88        #fig.set_size_inches(cm2inch(40, 20))
89        #fig.set_size_inches(cm2inch(40*4, 20*4))
90        fig.set_size_inches(cm2inch(figure_width, figure_height))
91
92
93
94    def plot_classification_report(classification_report, number_of_clas
95        '''
96        Plot scikit-learn classification report.
97        Extension based on https://stackoverflow.com/a/31689645/395857
98        '''
99        lines = classification_report.split('\n')
100
101        #drop initial lines
102        lines = lines[2:]
103
104        classes = []
105        plotMat = []
106        support = []
107        class_names = []
108        for line in lines[: number_of_classes]:
109            t = list(filter(None, line.strip().split('  ')))
110            if len(t) < 4: continue
111            classes.append(t[0])
112            v = [float(x) for x in t[1: len(t) - 1]]
113            support.append(int(t[-1]))
114            class_names.append(t[0])
```
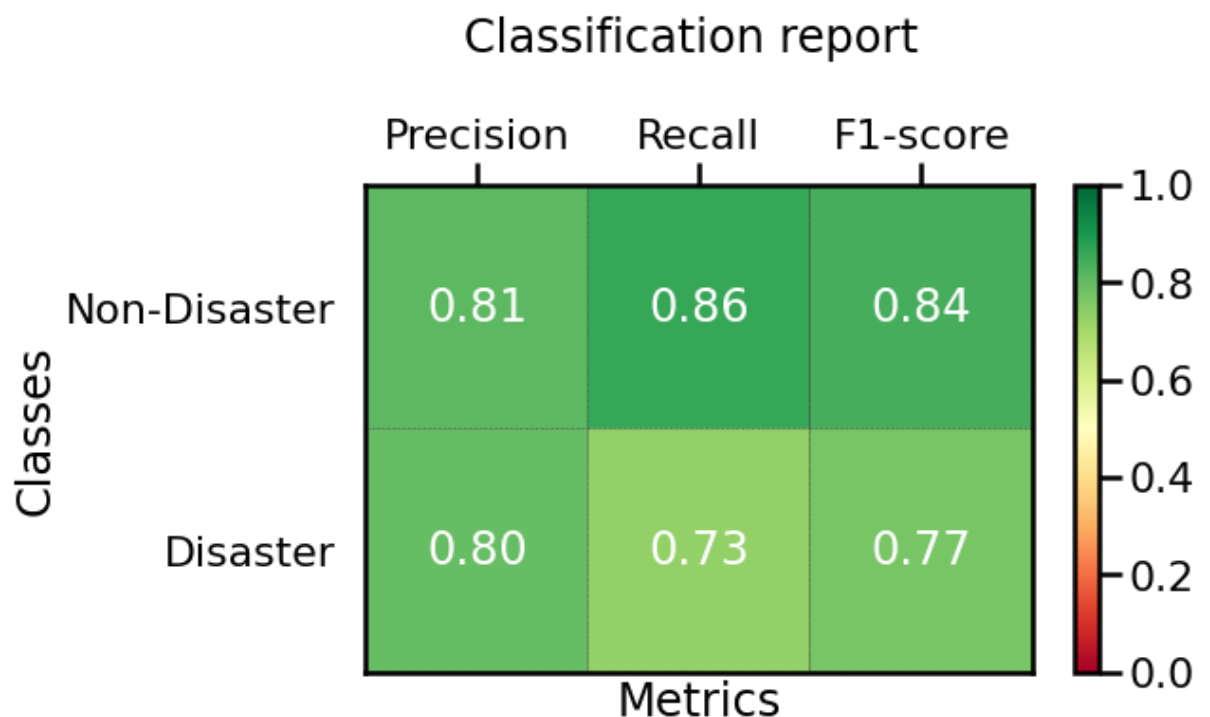
```
115              plotMat.append(v)
116
117
118      xlabel = 'Metrics'
119      ylabel = 'Classes'
120      xticklabels = ['Precision', 'Recall', 'F1-score']
121      yticklabels = ['Non-Disaster', 'Disaster']
122      #'{0} ({1})'.format(class_names[idx], sup)  for idx, sup  in enum
123      figure_width = 20
124      figure_height = len(class_names) + 10
125      correct_orientation = True
126      heatmap(np.array(plotMat), title, xlabel, ylabel, xticklabels, y
127      plt.show()
128
129
```

In [26]:    1  plot_classification_report(final_clf_report)
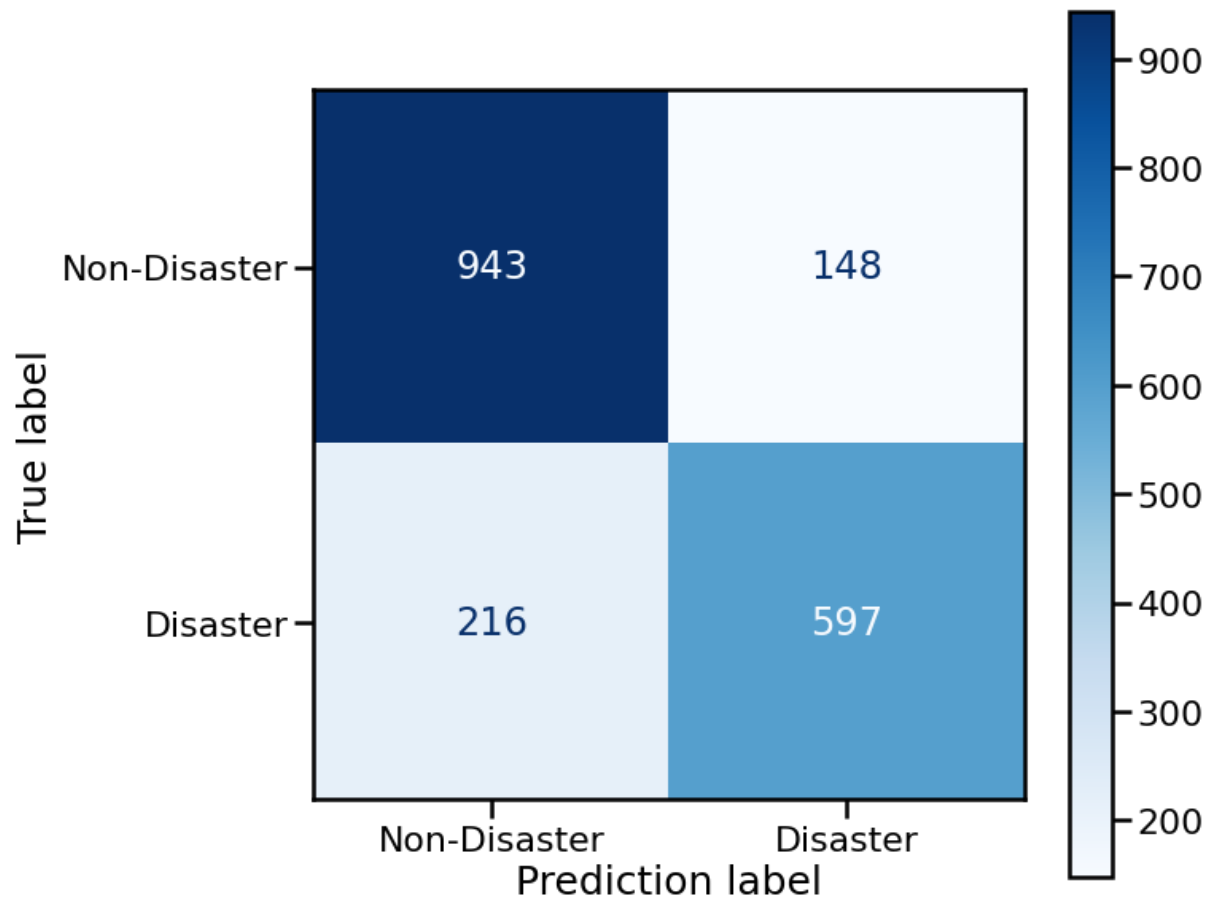
## Classification report



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]]
```

Creating a visualization of this confusion matrix:

In [28]:
```python
fig, ax = plt.subplots(figsize=(10,10))
cm_1 = ConfusionMatrixDisplay(confusion_matrix = cnf_matrix, display
cm_1.plot(cmap=plt.cm.Blues, ax=ax)
plt.xlabel('Prediction label',fontsize=25)
plt.ylabel('True label',fontsize=25);
```



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]:
```python
sample_submission = pd.read_csv("data/sample_submission.csv")
```

In [30]:
```python
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]:
```python
 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) #
 6      text = no_bad_chars.sub(' ', text)
 7      text = text.replace("#", "").replace("_", " ") #Remove hashtag s
 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
14
```

In [32]:
```python
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]:
```python
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:

- Model should be deployed to monitor twitter for disaster tweets

## Next Steps:

- Continue testing other models to see if performance can be improved
  - Use more tweets data
  - Use word embeddings
  - Flag top words/phrases of disaster tweets

# Contact me:

- LinkedIn: https://www.linkedin.com/in/nicole-michaud2/ (https://www.linkedin.com/in/nicole-michaud2/)
- Email: michaud.nicole00@gmail.com (mailto:michaud.nicole00@gmail.com)
- Blog: https://medium.com/@nicolemichaud03 (https://medium.com/@nicolemichaud03)