

HateBot Machine Learning

1. Introduction

Hate speech is something that is a big issue nowadays especially with social media like Twitter. Negativity that comes out of it can take a toll on a person's mental health. To minimize the damage of that, the HateBot can return if a tweet is hate speech or not by being given a Tweet. This way we can ensure that the world becomes if not a better place then at least a place where we can influence the amount of negativity we get.

Note: This is the second notebook of this project. The EDA is conducted in a separate Jupyter notebook that can be found in the submission.

2. Setup

In this part of the notebook, I am going to do the importing of libraries/modules and the dataset and setup the global settings needed for the figures.

2.1 Imports

```
In [1]: # Imports
import pandas as pd

from collections import Counter

import matplotlib
import matplotlib.pyplot as plt

import scattertext as st
from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator

from sklearn.svm import SVC
from sklearn.model_selection import train_test_split, cross_val_score
import sklearn.metrics as metrics
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.feature_extraction.text import HashingVectorizer, CountVectorizer

print('pandas version:', pd.__version__)
print('matplotlib version:', matplotlib.__version__)
print('scattertext version:', st.__version__)
```

```
%matplotlib inline
```

```
pandas version: 1.1.3  
matplotlib version: 3.3.2  
scattertext version: 0.1.2
```

2.2 Importing the dataset that will be used

The dataset is stored in a CSV format (comma-separated values file) in the file `HateBotDataset.csv`. This dataset is already cleaned in the EDA phase.

```
In [2]: df = pd.read_csv("CleanedHateBotDataset.csv")
```

2.2.1 Check the data in the dataset

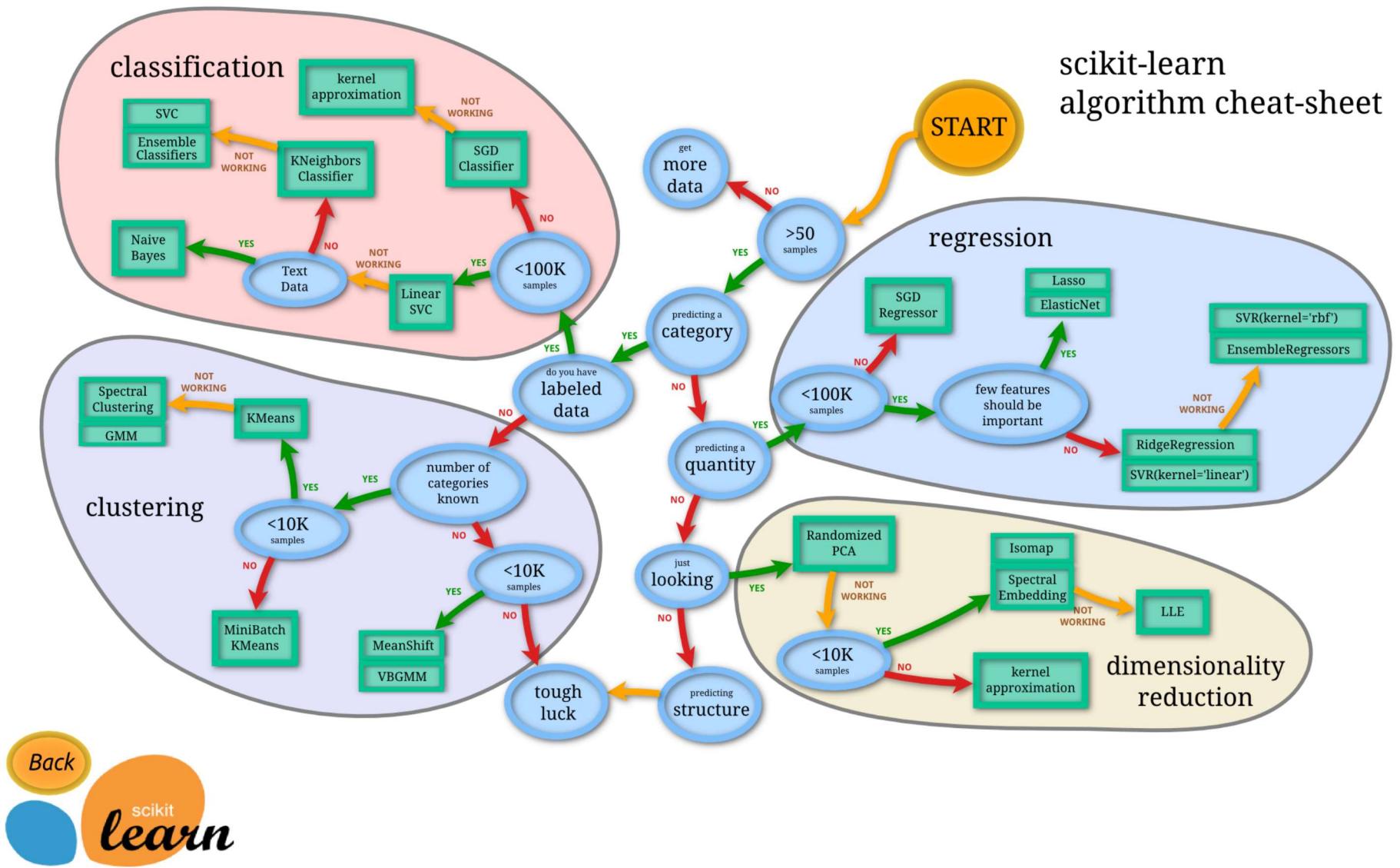
To be sure, we want to see the first five entries of the dataset to make sure that we have loaded the correct file.

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

```
Out[3]:   count  neither  class          tweet  marked  is_hate_speech
0       3        3     2  RT As a woman you shouldnt complain about cl...      0    False
1       3        0     1  RT boy dats coldtyga dwn bad for cuffin dat ...      3     True
2       3        0     1  RT Dawg RT You ever fuck a bitch and she st...      3     True
3       3        1     1           RT she look like a tranny      2     True
4       6        0     1  RT The shit you hear about me might be true ...      6     True
```

2.3 Selection of an algorithm for the Machine Learning

scikit-learn algorithm cheat-sheet



The algorithms that will be used for this challenge are Linear Support Vector Classification (Linear SVC) and Naive Bayes. The following algorithms were selected because we have our sample size (less than 100k) and because we are using text data (because the Tweets are text). Naive bayes does quite well when the training data doesn't contain all possibilities so it can be very good with low amounts of data. Decision trees work better with lots of data compared to Naive Bayes but our sample size is not big enough.

2.4 Data description

From the EDA phase, it has been decided that we are going to train via the `tweet` column (data type: `string`) and the `is_hate_speech` (data type: `Boolean`) column.

3. Machine learning using a Linear SVC

3.1 Training using a SVC with a CountVectorizer

The CountVectorizer converts a collection of text documents to a matrix of token counts. This is needed because machine learning algorithms cannot run on raw text data.

3.1.1 Get all the words that are in all the tweets

```
In [4]: d = Counter(" ".join(df.tweet).split(" ")).items()
most_used_words = dict(sorted(d, key=lambda item: item[1], reverse=True))
words = {key for key, val in most_used_words.items()}
```

3.1.2 Create and fit the CountVectorizer

Here we tokenize the words and build the vocabulary.

```
In [5]: countVectorizer = CountVectorizer()
countVectorizer.fit(words)
```

```
Out[5]: CountVectorizer()
```

3.1.3 Encode all the words from the Tweets

```
In [6]: countVector = countVectorizer.transform(words)
print(countVector.shape)

(34011, 28861)
```

3.1.4 Transform every tweet into a pandas Series

We add a new column called `vectorised_words` which has for every row in the dataset the tweet itself but after it was transformed via our count vectorizer.

```
In [7]: df['vectorised_words'] = countVectorizer.transform(df.tweet.values)
```

```
In [8]: type(df['vectorised_words'])
```

```
Out[8]: pandas.core.series.Series
```

3.2 Training the SVC

3.2.1 Create the needed variables

```
In [9]: svc_1 = SVC(kernel='linear')
```

```
In [10]: X = countVectorizer.transform(df.tweet.values)
y = df['is_hate_speech'].astype(int)
```

```
In [11]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state=46)
```

3.2.2 Fit the SVC

```
In [12]: svc_1.fit(X_train, y_train)
```

```
Out[12]: SVC(kernel='linear')
```

3.2.3 Evaluate the predictions

```
In [13]: y_pred = svc_1.predict(X_test)
print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.85	0.84	0.84	1239
1	0.97	0.97	0.97	6196
accuracy			0.95	7435
macro avg	0.91	0.90	0.91	7435
weighted avg	0.95	0.95	0.95	7435

```
In [14]: print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
print("Precision:",metrics.precision_score(y_test, y_pred))
print("Recall:",metrics.recall_score(y_test, y_pred))
```

```
Accuracy: 0.947679892400807
```

```
Precision: 0.9679290894439968
```

```
Recall: 0.9693350548741123
```

We can see that this model is quite succesfull - it's accuracy is almost at 95% which in our case is more than sufficient.

3.2.4 Hyperparameter tuning

```
In [16]: from sklearn.model_selection import GridSearchCV
```

```
In [17]: param_grid = {'C': [0.1, 1, 10, 100], 'gamma': [1, 0.1, 0.01, 0.001], 'kernel': ['rbf', 'poly', 'sigmoid']}
grid = GridSearchCV(SVC(), param_grid, refit=True, verbose=2)
```

```
In [18]: grid.fit(X_train, y_train)
print(grid.best_estimator_)
```

```
Fitting 5 folds for each of 48 candidates, totalling 240 fits
[CV] C=0.1, gamma=1, kernel=rbf .....
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[CV] ..... C=0.1, gamma=1, kernel=rbf, total= 1.1min
[CV] C=0.1, gamma=1, kernel=rbf .....
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 1.1min remaining: 0.0s
[CV] ..... C=0.1, gamma=1, kernel=rbf, total= 58.9s
[CV] C=0.1, gamma=1, kernel=rbf .....
[CV] ..... C=0.1, gamma=1, kernel=rbf, total= 1.0min
[CV] C=0.1, gamma=1, kernel=rbf .....
[CV] ..... C=0.1, gamma=1, kernel=rbf, total= 1.1min
[CV] C=0.1, gamma=1, kernel=rbf .....
[CV] ..... C=0.1, gamma=1, kernel=rbf, total= 1.1min
[CV] C=0.1, gamma=1, kernel=poly .....
[CV] ..... C=0.1, gamma=1, kernel=poly, total= 29.2s
[CV] C=0.1, gamma=1, kernel=poly .....
[CV] ..... C=0.1, gamma=1, kernel=poly, total= 31.6s
[CV] C=0.1, gamma=1, kernel=poly .....
[CV] ..... C=0.1, gamma=1, kernel=poly, total= 30.3s
[CV] C=0.1, gamma=1, kernel=poly .....
[CV] ..... C=0.1, gamma=1, kernel=poly, total= 30.8s
[CV] C=0.1, gamma=1, kernel=poly .....
[CV] ..... C=0.1, gamma=1, kernel=poly, total= 30.4s
[CV] C=0.1, gamma=1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=1, kernel=sigmoid, total= 8.5s
[CV] C=0.1, gamma=1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=1, kernel=sigmoid, total= 8.3s
[CV] C=0.1, gamma=1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=1, kernel=sigmoid, total= 7.9s
[CV] C=0.1, gamma=1, kernel=sigmoid .....
```

```
[CV] ..... C=0.1, gamma=1, kernel=sigmoid, total= 7.9s
[CV] C=0.1, gamma=1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=1, kernel=sigmoid, total= 8.2s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.1, kernel=rbf, total= 12.5s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.1, kernel=rbf, total= 12.6s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.1, kernel=rbf, total= 13.1s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.1, kernel=rbf, total= 12.3s
[CV] C=0.1, gamma=0.1, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.1, kernel=rbf, total= 12.5s
[CV] C=0.1, gamma=0.1, kernel=poly .....
[CV] ..... C=0.1, gamma=0.1, kernel=poly, total= 27.7s
[CV] C=0.1, gamma=0.1, kernel=poly .....
[CV] ..... C=0.1, gamma=0.1, kernel=poly, total= 27.7s
[CV] C=0.1, gamma=0.1, kernel=poly .....
[CV] ..... C=0.1, gamma=0.1, kernel=poly, total= 29.6s
[CV] C=0.1, gamma=0.1, kernel=poly .....
[CV] ..... C=0.1, gamma=0.1, kernel=poly, total= 28.4s
[CV] C=0.1, gamma=0.1, kernel=poly .....
[CV] ..... C=0.1, gamma=0.1, kernel=poly, total= 27.6s
[CV] C=0.1, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=0.1, kernel=sigmoid, total= 9.0s
[CV] C=0.1, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=0.1, kernel=sigmoid, total= 9.2s
[CV] C=0.1, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=0.1, kernel=sigmoid, total= 9.0s
[CV] C=0.1, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=0.1, kernel=sigmoid, total= 9.1s
[CV] C=0.1, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=0.1, gamma=0.1, kernel=sigmoid, total= 9.2s
[CV] C=0.1, gamma=0.01, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.01, kernel=rbf, total= 9.7s
[CV] C=0.1, gamma=0.01, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.01, kernel=rbf, total= 10.0s
[CV] C=0.1, gamma=0.01, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.01, kernel=rbf, total= 9.8s
[CV] C=0.1, gamma=0.01, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.01, kernel=rbf, total= 9.6s
[CV] C=0.1, gamma=0.01, kernel=rbf .....
[CV] ..... C=0.1, gamma=0.01, kernel=rbf, total= 10.3s
[CV] C=0.1, gamma=0.01, kernel=poly .....
[CV] ..... C=0.1, gamma=0.01, kernel=poly, total= 9.3s
[CV] C=0.1, gamma=0.01, kernel=poly .....
[CV] ..... C=0.1, gamma=0.01, kernel=poly, total= 8.7s
```

```
[CV] C=0.1, gamma=0.01, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.01, kernel=poly, total= 9.0s  
[CV] C=0.1, gamma=0.01, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.01, kernel=poly, total= 9.1s  
[CV] C=0.1, gamma=0.01, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.01, kernel=poly, total= 9.1s  
[CV] C=0.1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.01, kernel=sigmoid, total= 9.9s  
[CV] C=0.1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.01, kernel=sigmoid, total= 10.1s  
[CV] C=0.1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.01, kernel=sigmoid, total= 10.0s  
[CV] C=0.1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.01, kernel=sigmoid, total= 10.0s  
[CV] C=0.1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.01, kernel=sigmoid, total= 10.1s  
[CV] C=0.1, gamma=0.001, kernel=rbf .....  
[CV] ..... C=0.1, gamma=0.001, kernel=rbf, total= 11.4s  
[CV] C=0.1, gamma=0.001, kernel=rbf .....  
[CV] ..... C=0.1, gamma=0.001, kernel=rbf, total= 11.6s  
[CV] C=0.1, gamma=0.001, kernel=rbf .....  
[CV] ..... C=0.1, gamma=0.001, kernel=rbf, total= 10.9s  
[CV] C=0.1, gamma=0.001, kernel=rbf .....  
[CV] ..... C=0.1, gamma=0.001, kernel=rbf, total= 11.0s  
[CV] C=0.1, gamma=0.001, kernel=rbf .....  
[CV] ..... C=0.1, gamma=0.001, kernel=rbf, total= 10.9s  
[CV] C=0.1, gamma=0.001, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.001, kernel=poly, total= 6.8s  
[CV] C=0.1, gamma=0.001, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.001, kernel=poly, total= 6.8s  
[CV] C=0.1, gamma=0.001, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.001, kernel=poly, total= 6.8s  
[CV] C=0.1, gamma=0.001, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.001, kernel=poly, total= 6.9s  
[CV] C=0.1, gamma=0.001, kernel=poly .....  
[CV] ..... C=0.1, gamma=0.001, kernel=poly, total= 6.8s  
[CV] C=0.1, gamma=0.001, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.001, kernel=sigmoid, total= 10.8s  
[CV] C=0.1, gamma=0.001, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.001, kernel=sigmoid, total= 10.8s  
[CV] C=0.1, gamma=0.001, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.001, kernel=sigmoid, total= 10.8s  
[CV] C=0.1, gamma=0.001, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.001, kernel=sigmoid, total= 11.3s  
[CV] C=0.1, gamma=0.001, kernel=sigmoid .....  
[CV] ..... C=0.1, gamma=0.001, kernel=sigmoid, total= 10.9s  
[CV] C=1, gamma=1, kernel=rbf .....
```

```
[CV] ..... C=1, gamma=1, kernel=rbf, total= 1.4min
[CV] C=1, gamma=1, kernel=rbf .....
[CV] ..... C=1, gamma=1, kernel=rbf, total= 1.4min
[CV] C=1, gamma=1, kernel=rbf .....
[CV] ..... C=1, gamma=1, kernel=rbf, total= 1.5min
[CV] C=1, gamma=1, kernel=rbf .....
[CV] ..... C=1, gamma=1, kernel=rbf, total= 1.4min
[CV] C=1, gamma=1, kernel=rbf .....
[CV] ..... C=1, gamma=1, kernel=rbf, total= 1.3min
[CV] C=1, gamma=1, kernel=poly .....
[CV] ..... C=1, gamma=1, kernel=poly, total= 23.2s
[CV] C=1, gamma=1, kernel=poly .....
[CV] ..... C=1, gamma=1, kernel=poly, total= 24.1s
[CV] C=1, gamma=1, kernel=poly .....
[CV] ..... C=1, gamma=1, kernel=poly, total= 24.4s
[CV] C=1, gamma=1, kernel=poly .....
[CV] ..... C=1, gamma=1, kernel=poly, total= 24.7s
[CV] C=1, gamma=1, kernel=poly .....
[CV] ..... C=1, gamma=1, kernel=poly, total= 25.1s
[CV] C=1, gamma=1, kernel=sigmoid .....
[CV] ..... C=1, gamma=1, kernel=sigmoid, total= 8.0s
[CV] C=1, gamma=1, kernel=sigmoid .....
[CV] ..... C=1, gamma=1, kernel=sigmoid, total= 8.0s
[CV] C=1, gamma=1, kernel=sigmoid .....
[CV] ..... C=1, gamma=1, kernel=sigmoid, total= 7.8s
[CV] C=1, gamma=1, kernel=sigmoid .....
[CV] ..... C=1, gamma=1, kernel=sigmoid, total= 7.5s
[CV] C=1, gamma=1, kernel=sigmoid .....
[CV] ..... C=1, gamma=1, kernel=sigmoid, total= 7.4s
[CV] C=1, gamma=0.1, kernel=rbf .....
[CV] ..... C=1, gamma=0.1, kernel=rbf, total= 20.7s
[CV] C=1, gamma=0.1, kernel=rbf .....
[CV] ..... C=1, gamma=0.1, kernel=rbf, total= 21.1s
[CV] C=1, gamma=0.1, kernel=rbf .....
[CV] ..... C=1, gamma=0.1, kernel=rbf, total= 21.4s
[CV] C=1, gamma=0.1, kernel=rbf .....
[CV] ..... C=1, gamma=0.1, kernel=rbf, total= 21.9s
[CV] C=1, gamma=0.1, kernel=rbf .....
[CV] ..... C=1, gamma=0.1, kernel=rbf, total= 20.9s
[CV] C=1, gamma=0.1, kernel=poly .....
[CV] ..... C=1, gamma=0.1, kernel=poly, total= 34.3s
[CV] C=1, gamma=0.1, kernel=poly .....
[CV] ..... C=1, gamma=0.1, kernel=poly, total= 35.0s
[CV] C=1, gamma=0.1, kernel=poly .....
[CV] ..... C=1, gamma=0.1, kernel=poly, total= 35.3s
[CV] C=1, gamma=0.1, kernel=poly .....
[CV] ..... C=1, gamma=0.1, kernel=poly, total= 34.0s
```

```
[CV] C=1, gamma=0.1, kernel=poly .....  
[CV] ..... C=1, gamma=0.1, kernel=poly, total= 34.8s  
[CV] C=1, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.1, kernel=sigmoid, total= 4.6s  
[CV] C=1, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.1, kernel=sigmoid, total= 4.7s  
[CV] C=1, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.1, kernel=sigmoid, total= 4.8s  
[CV] C=1, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.1, kernel=sigmoid, total= 5.1s  
[CV] C=1, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.1, kernel=sigmoid, total= 4.7s  
[CV] C=1, gamma=0.01, kernel=rbf .....  
[CV] ..... C=1, gamma=0.01, kernel=rbf, total= 8.3s  
[CV] C=1, gamma=0.01, kernel=rbf .....  
[CV] ..... C=1, gamma=0.01, kernel=rbf, total= 8.5s  
[CV] C=1, gamma=0.01, kernel=rbf .....  
[CV] ..... C=1, gamma=0.01, kernel=rbf, total= 8.7s  
[CV] C=1, gamma=0.01, kernel=rbf .....  
[CV] ..... C=1, gamma=0.01, kernel=rbf, total= 8.8s  
[CV] C=1, gamma=0.01, kernel=rbf .....  
[CV] ..... C=1, gamma=0.01, kernel=rbf, total= 8.9s  
[CV] C=1, gamma=0.01, kernel=poly .....  
[CV] ..... C=1, gamma=0.01, kernel=poly, total= 14.5s  
[CV] C=1, gamma=0.01, kernel=poly .....  
[CV] ..... C=1, gamma=0.01, kernel=poly, total= 14.5s  
[CV] C=1, gamma=0.01, kernel=poly .....  
[CV] ..... C=1, gamma=0.01, kernel=poly, total= 14.9s  
[CV] C=1, gamma=0.01, kernel=poly .....  
[CV] ..... C=1, gamma=0.01, kernel=poly, total= 14.2s  
[CV] C=1, gamma=0.01, kernel=poly .....  
[CV] ..... C=1, gamma=0.01, kernel=poly, total= 14.6s  
[CV] C=1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.01, kernel=sigmoid, total= 8.7s  
[CV] C=1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.01, kernel=sigmoid, total= 8.7s  
[CV] C=1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.01, kernel=sigmoid, total= 8.8s  
[CV] C=1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.01, kernel=sigmoid, total= 8.8s  
[CV] C=1, gamma=0.01, kernel=sigmoid .....  
[CV] ..... C=1, gamma=0.01, kernel=sigmoid, total= 8.7s  
[CV] C=1, gamma=0.001, kernel=rbf .....  
[CV] ..... C=1, gamma=0.001, kernel=rbf, total= 9.6s  
[CV] C=1, gamma=0.001, kernel=rbf .....  
[CV] ..... C=1, gamma=0.001, kernel=rbf, total= 10.1s  
[CV] C=1, gamma=0.001, kernel=rbf .....
```

```
[CV] ..... C=1, gamma=0.001, kernel=rbf, total= 10.0s
[CV] C=1, gamma=0.001, kernel=rbf .....
[CV] ..... C=1, gamma=0.001, kernel=rbf, total= 9.7s
[CV] C=1, gamma=0.001, kernel=rbf .....
[CV] ..... C=1, gamma=0.001, kernel=rbf, total= 9.6s
[CV] C=1, gamma=0.001, kernel=poly .....
[CV] ..... C=1, gamma=0.001, kernel=poly, total= 6.7s
[CV] C=1, gamma=0.001, kernel=poly .....
[CV] ..... C=1, gamma=0.001, kernel=poly, total= 6.9s
[CV] C=1, gamma=0.001, kernel=poly .....
[CV] ..... C=1, gamma=0.001, kernel=poly, total= 7.1s
[CV] C=1, gamma=0.001, kernel=poly .....
[CV] ..... C=1, gamma=0.001, kernel=poly, total= 8.0s
[CV] C=1, gamma=0.001, kernel=poly .....
[CV] ..... C=1, gamma=0.001, kernel=poly, total= 6.8s
[CV] C=1, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=1, gamma=0.001, kernel=sigmoid, total= 10.0s
[CV] C=1, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=1, gamma=0.001, kernel=sigmoid, total= 10.2s
[CV] C=1, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=1, gamma=0.001, kernel=sigmoid, total= 10.1s
[CV] C=1, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=1, gamma=0.001, kernel=sigmoid, total= 10.0s
[CV] C=1, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=1, gamma=0.001, kernel=sigmoid, total= 10.1s
[CV] C=10, gamma=1, kernel=rbf .....
[CV] ..... C=10, gamma=1, kernel=rbf, total= 1.5min
[CV] C=10, gamma=1, kernel=rbf .....
[CV] ..... C=10, gamma=1, kernel=rbf, total= 1.5min
[CV] C=10, gamma=1, kernel=rbf .....
[CV] ..... C=10, gamma=1, kernel=rbf, total= 1.5min
[CV] C=10, gamma=1, kernel=rbf .....
[CV] ..... C=10, gamma=1, kernel=rbf, total= 1.6min
[CV] C=10, gamma=1, kernel=rbf .....
[CV] ..... C=10, gamma=1, kernel=rbf, total= 1.6min
[CV] C=10, gamma=1, kernel=poly .....
[CV] ..... C=10, gamma=1, kernel=poly, total= 21.7s
[CV] C=10, gamma=1, kernel=poly .....
[CV] ..... C=10, gamma=1, kernel=poly, total= 27.2s
[CV] C=10, gamma=1, kernel=poly .....
[CV] ..... C=10, gamma=1, kernel=poly, total= 28.8s
[CV] C=10, gamma=1, kernel=poly .....
[CV] ..... C=10, gamma=1, kernel=poly, total= 21.2s
[CV] C=10, gamma=1, kernel=poly .....
[CV] ..... C=10, gamma=1, kernel=poly, total= 22.7s
[CV] C=10, gamma=1, kernel=sigmoid .....
[CV] ..... C=10, gamma=1, kernel=sigmoid, total= 8.9s
```

```
[CV] C=10, gamma=1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=1, kernel=sigmoid, total= 8.9s  
[CV] C=10, gamma=1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=1, kernel=sigmoid, total= 8.9s  
[CV] C=10, gamma=1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=1, kernel=sigmoid, total= 9.2s  
[CV] C=10, gamma=1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=1, kernel=sigmoid, total= 8.9s  
[CV] C=10, gamma=0.1, kernel=rbf .....  
[CV] ..... C=10, gamma=0.1, kernel=rbf, total= 38.1s  
[CV] C=10, gamma=0.1, kernel=rbf .....  
[CV] ..... C=10, gamma=0.1, kernel=rbf, total= 44.0s  
[CV] C=10, gamma=0.1, kernel=rbf .....  
[CV] ..... C=10, gamma=0.1, kernel=rbf, total= 35.5s  
[CV] C=10, gamma=0.1, kernel=rbf .....  
[CV] ..... C=10, gamma=0.1, kernel=rbf, total= 35.3s  
[CV] C=10, gamma=0.1, kernel=rbf .....  
[CV] ..... C=10, gamma=0.1, kernel=rbf, total= 39.6s  
[CV] C=10, gamma=0.1, kernel=poly .....  
[CV] ..... C=10, gamma=0.1, kernel=poly, total= 35.1s  
[CV] C=10, gamma=0.1, kernel=poly .....  
[CV] ..... C=10, gamma=0.1, kernel=poly, total= 32.9s  
[CV] C=10, gamma=0.1, kernel=poly .....  
[CV] ..... C=10, gamma=0.1, kernel=poly, total= 31.7s  
[CV] C=10, gamma=0.1, kernel=poly .....  
[CV] ..... C=10, gamma=0.1, kernel=poly, total= 33.4s  
[CV] C=10, gamma=0.1, kernel=poly .....  
[CV] ..... C=10, gamma=0.1, kernel=poly, total= 32.1s  
[CV] C=10, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=0.1, kernel=sigmoid, total= 3.7s  
[CV] C=10, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=0.1, kernel=sigmoid, total= 3.8s  
[CV] C=10, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=0.1, kernel=sigmoid, total= 3.8s  
[CV] C=10, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=0.1, kernel=sigmoid, total= 3.7s  
[CV] C=10, gamma=0.1, kernel=sigmoid .....  
[CV] ..... C=10, gamma=0.1, kernel=sigmoid, total= 3.6s  
[CV] C=10, gamma=0.01, kernel=rbf .....  
[CV] ..... C=10, gamma=0.01, kernel=rbf, total= 7.0s  
[CV] C=10, gamma=0.01, kernel=rbf .....  
[CV] ..... C=10, gamma=0.01, kernel=rbf, total= 7.4s  
[CV] C=10, gamma=0.01, kernel=rbf .....  
[CV] ..... C=10, gamma=0.01, kernel=rbf, total= 7.2s  
[CV] C=10, gamma=0.01, kernel=rbf .....  
[CV] ..... C=10, gamma=0.01, kernel=rbf, total= 7.2s  
[CV] C=10, gamma=0.01, kernel=rbf .....
```

```
[CV] ..... C=10, gamma=0.01, kernel=rbf, total= 7.5s
[CV] C=10, gamma=0.01, kernel=poly .....
[CV] ..... C=10, gamma=0.01, kernel=poly, total= 20.7s
[CV] C=10, gamma=0.01, kernel=poly .....
[CV] ..... C=10, gamma=0.01, kernel=poly, total= 22.7s
[CV] C=10, gamma=0.01, kernel=poly .....
[CV] ..... C=10, gamma=0.01, kernel=poly, total= 22.1s
[CV] C=10, gamma=0.01, kernel=poly .....
[CV] ..... C=10, gamma=0.01, kernel=poly, total= 20.4s
[CV] C=10, gamma=0.01, kernel=poly .....
[CV] ..... C=10, gamma=0.01, kernel=poly, total= 23.8s
[CV] C=10, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.01, kernel=sigmoid, total= 6.6s
[CV] C=10, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.01, kernel=sigmoid, total= 6.7s
[CV] C=10, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.01, kernel=sigmoid, total= 6.5s
[CV] C=10, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.01, kernel=sigmoid, total= 6.5s
[CV] C=10, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.01, kernel=sigmoid, total= 6.6s
[CV] C=10, gamma=0.001, kernel=rbf .....
[CV] ..... C=10, gamma=0.001, kernel=rbf, total= 7.8s
[CV] C=10, gamma=0.001, kernel=rbf .....
[CV] ..... C=10, gamma=0.001, kernel=rbf, total= 7.9s
[CV] C=10, gamma=0.001, kernel=rbf .....
[CV] ..... C=10, gamma=0.001, kernel=rbf, total= 8.1s
[CV] C=10, gamma=0.001, kernel=rbf .....
[CV] ..... C=10, gamma=0.001, kernel=rbf, total= 8.2s
[CV] C=10, gamma=0.001, kernel=rbf .....
[CV] ..... C=10, gamma=0.001, kernel=rbf, total= 8.2s
[CV] C=10, gamma=0.001, kernel=poly .....
[CV] ..... C=10, gamma=0.001, kernel=poly, total= 7.3s
[CV] C=10, gamma=0.001, kernel=poly .....
[CV] ..... C=10, gamma=0.001, kernel=poly, total= 7.1s
[CV] C=10, gamma=0.001, kernel=poly .....
[CV] ..... C=10, gamma=0.001, kernel=poly, total= 7.1s
[CV] C=10, gamma=0.001, kernel=poly .....
[CV] ..... C=10, gamma=0.001, kernel=poly, total= 7.2s
[CV] C=10, gamma=0.001, kernel=poly .....
[CV] ..... C=10, gamma=0.001, kernel=poly, total= 7.3s
[CV] C=10, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.001, kernel=sigmoid, total= 9.2s
[CV] C=10, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.001, kernel=sigmoid, total= 9.0s
[CV] C=10, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=10, gamma=0.001, kernel=sigmoid, total= 8.9s
```

```
[CV] C=10, gamma=0.001, kernel=sigmoid .....  
[CV] ..... C=10, gamma=0.001, kernel=sigmoid, total= 8.9s  
[CV] C=10, gamma=0.001, kernel=sigmoid .....  
[CV] ..... C=10, gamma=0.001, kernel=sigmoid, total= 8.9s  
[CV] C=100, gamma=1, kernel=rbf .....  
[CV] ..... C=100, gamma=1, kernel=rbf, total= 1.5min  
[CV] C=100, gamma=1, kernel=rbf .....  
[CV] ..... C=100, gamma=1, kernel=rbf, total= 1.5min  
[CV] C=100, gamma=1, kernel=rbf .....  
[CV] ..... C=100, gamma=1, kernel=rbf, total= 1.4min  
[CV] C=100, gamma=1, kernel=rbf .....  
[CV] ..... C=100, gamma=1, kernel=rbf, total= 1.5min  
[CV] C=100, gamma=1, kernel=rbf .....  
[CV] ..... C=100, gamma=1, kernel=rbf, total= 1.5min  
[CV] C=100, gamma=1, kernel=poly .....  
[CV] ..... C=100, gamma=1, kernel=poly, total= 19.7s  
[CV] C=100, gamma=1, kernel=poly .....  
[CV] ..... C=100, gamma=1, kernel=poly, total= 25.1s  
[CV] C=100, gamma=1, kernel=poly .....  
[CV] ..... C=100, gamma=1, kernel=poly, total= 27.0s  
[CV] C=100, gamma=1, kernel=poly .....  
[CV] ..... C=100, gamma=1, kernel=poly, total= 19.9s  
[CV] C=100, gamma=1, kernel=poly .....  
[CV] ..... C=100, gamma=1, kernel=poly, total= 21.1s  
[CV] C=100, gamma=1, kernel=sigmoid .....  
[CV] ..... C=100, gamma=1, kernel=sigmoid, total= 6.8s  
[CV] C=100, gamma=1, kernel=sigmoid .....  
[CV] ..... C=100, gamma=1, kernel=sigmoid, total= 7.2s  
[CV] C=100, gamma=1, kernel=sigmoid .....  
[CV] ..... C=100, gamma=1, kernel=sigmoid, total= 6.7s  
[CV] C=100, gamma=1, kernel=sigmoid .....  
[CV] ..... C=100, gamma=1, kernel=sigmoid, total= 7.1s  
[CV] C=100, gamma=1, kernel=sigmoid .....  
[CV] ..... C=100, gamma=1, kernel=sigmoid, total= 6.6s  
[CV] C=100, gamma=0.1, kernel=rbf .....  
[CV] ..... C=100, gamma=0.1, kernel=rbf, total= 31.5s  
[CV] C=100, gamma=0.1, kernel=rbf .....  
[CV] ..... C=100, gamma=0.1, kernel=rbf, total= 32.0s  
[CV] C=100, gamma=0.1, kernel=rbf .....  
[CV] ..... C=100, gamma=0.1, kernel=rbf, total= 32.0s  
[CV] C=100, gamma=0.1, kernel=rbf .....  
[CV] ..... C=100, gamma=0.1, kernel=rbf, total= 31.9s  
[CV] C=100, gamma=0.1, kernel=rbf .....  
[CV] ..... C=100, gamma=0.1, kernel=rbf, total= 30.9s  
[CV] C=100, gamma=0.1, kernel=poly .....  
[CV] ..... C=100, gamma=0.1, kernel=poly, total= 27.2s  
[CV] C=100, gamma=0.1, kernel=poly .....
```

```
[CV] ..... C=100, gamma=0.1, kernel=poly, total= 27.9s
[CV] C=100, gamma=0.1, kernel=poly .....
[CV] ..... C=100, gamma=0.1, kernel=poly, total= 28.3s
[CV] C=100, gamma=0.1, kernel=poly .....
[CV] ..... C=100, gamma=0.1, kernel=poly, total= 28.9s
[CV] C=100, gamma=0.1, kernel=poly .....
[CV] ..... C=100, gamma=0.1, kernel=poly, total= 27.9s
[CV] C=100, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.1, kernel=sigmoid, total= 3.4s
[CV] C=100, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.1, kernel=sigmoid, total= 3.3s
[CV] C=100, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.1, kernel=sigmoid, total= 3.5s
[CV] C=100, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.1, kernel=sigmoid, total= 3.4s
[CV] C=100, gamma=0.1, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.1, kernel=sigmoid, total= 3.5s
[CV] C=100, gamma=0.01, kernel=rbf .....
[CV] ..... C=100, gamma=0.01, kernel=rbf, total= 6.2s
[CV] C=100, gamma=0.01, kernel=rbf .....
[CV] ..... C=100, gamma=0.01, kernel=rbf, total= 6.2s
[CV] C=100, gamma=0.01, kernel=rbf .....
[CV] ..... C=100, gamma=0.01, kernel=rbf, total= 6.2s
[CV] C=100, gamma=0.01, kernel=rbf .....
[CV] ..... C=100, gamma=0.01, kernel=rbf, total= 6.3s
[CV] C=100, gamma=0.01, kernel=rbf .....
[CV] ..... C=100, gamma=0.01, kernel=rbf, total= 6.3s
[CV] C=100, gamma=0.01, kernel=poly .....
[CV] ..... C=100, gamma=0.01, kernel=poly, total= 26.2s
[CV] C=100, gamma=0.01, kernel=poly .....
[CV] ..... C=100, gamma=0.01, kernel=poly, total= 26.0s
[CV] C=100, gamma=0.01, kernel=poly .....
[CV] ..... C=100, gamma=0.01, kernel=poly, total= 25.9s
[CV] C=100, gamma=0.01, kernel=poly .....
[CV] ..... C=100, gamma=0.01, kernel=poly, total= 25.7s
[CV] C=100, gamma=0.01, kernel=poly .....
[CV] ..... C=100, gamma=0.01, kernel=poly, total= 26.0s
[CV] C=100, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.01, kernel=sigmoid, total= 5.8s
[CV] C=100, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.01, kernel=sigmoid, total= 5.9s
[CV] C=100, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.01, kernel=sigmoid, total= 5.8s
[CV] C=100, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.01, kernel=sigmoid, total= 5.8s
[CV] C=100, gamma=0.01, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.01, kernel=sigmoid, total= 5.7s
```

```

[CV] C=100, gamma=0.001, kernel=rbf .....
[CV] ..... C=100, gamma=0.001, kernel=rbf, total= 5.9s
[CV] C=100, gamma=0.001, kernel=rbf .....
[CV] ..... C=100, gamma=0.001, kernel=rbf, total= 5.9s
[CV] C=100, gamma=0.001, kernel=rbf .....
[CV] ..... C=100, gamma=0.001, kernel=rbf, total= 5.8s
[CV] C=100, gamma=0.001, kernel=rbf .....
[CV] ..... C=100, gamma=0.001, kernel=rbf, total= 5.8s
[CV] C=100, gamma=0.001, kernel=rbf .....
[CV] ..... C=100, gamma=0.001, kernel=rbf, total= 5.8s
[CV] C=100, gamma=0.001, kernel=poly .....
[CV] ..... C=100, gamma=0.001, kernel=poly, total= 8.6s
[CV] C=100, gamma=0.001, kernel=poly .....
[CV] ..... C=100, gamma=0.001, kernel=poly, total= 9.0s
[CV] C=100, gamma=0.001, kernel=poly .....
[CV] ..... C=100, gamma=0.001, kernel=poly, total= 9.4s
[CV] C=100, gamma=0.001, kernel=poly .....
[CV] ..... C=100, gamma=0.001, kernel=poly, total= 8.4s
[CV] C=100, gamma=0.001, kernel=poly .....
[CV] ..... C=100, gamma=0.001, kernel=poly, total= 8.7s
[CV] C=100, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.001, kernel=sigmoid, total= 6.1s
[CV] C=100, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.001, kernel=sigmoid, total= 6.1s
[CV] C=100, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.001, kernel=sigmoid, total= 6.1s
[CV] C=100, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.001, kernel=sigmoid, total= 6.1s
[CV] C=100, gamma=0.001, kernel=sigmoid .....
[CV] ..... C=100, gamma=0.001, kernel=sigmoid, total= 6.1s
[Parallel(n_jobs=1)]: Done 240 out of 240 | elapsed: 79.0min finished
SVC(C=10, gamma=0.01, kernel='sigmoid')

```

From the GridSearch, we can see that we should use SVC(C=10, gamma=0.01, kernel='sigmoid'). Now we will make that so that we have the best possible SVC.

```
In [19]: grid_predictions = grid.predict(X_test)
print(confusion_matrix(y_test,grid_predictions))
print(classification_report(y_test,grid_predictions))
```

	precision	recall	f1-score	support
0	0.84	0.88	0.86	1239
1	0.98	0.97	0.97	6196

accuracy			0.95	7435
macro avg	0.91	0.92	0.91	7435
weighted avg	0.95	0.95	0.95	7435

3.2.4.1 Creating the hypertuned version

```
In [20]: svc_1 = SVC(C=10, gamma=0.01, kernel='sigmoid')
X = countVectorizer.transform(df['tweet'].values)
y = df['is_hate_speech'].astype(int)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state=46)
svc_1.fit(X_train, y_train)
y_pred = svc_1.predict(X_test)
print(classification_report(y_test,y_pred))
print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
print("Precision:",metrics.precision_score(y_test, y_pred))
print("Recall:",metrics.recall_score(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.84	0.88	0.86	1239
1	0.98	0.97	0.97	6196
accuracy			0.95	7435
macro avg	0.91	0.92	0.91	7435
weighted avg	0.95	0.95	0.95	7435

Accuracy: 0.9515803631472763
 Precision: 0.9763303950375449
 Recall: 0.9653001936733376

3.2.5 Save as a pickle

For the deployment part of the project, we will be having a Django web application where the user can enter the Tweet they want to get checked and then return if the Tweet is hate or not. To do so, we need the already trained model. For this, I have decided to use pickle because it is very easy to integrate with Django. Now we will extract this SVC model

```
In [21]: import pickle
pickl = {
    'vectorizer': countVectorizer,
    'regressor': svc_1
}
pickle.dump( pickl, open( "HateBotModels.p", "wb" ) )
```

3.3 Training using a SVC with a HashingVectorizer

This strategy has several advantages:

- it is very low memory scalable to large datasets as there is no need to store a vocabulary dictionary in memory
- it is fast to pickle and un-pickle as it holds no state besides the constructor parameters
- it can be used in a streaming (partial fit) or parallel pipeline as there is no state computed during fit.

There are also a couple of cons (vs using a CountVectorizer with an in-memory vocabulary):

- there is no way to compute the inverse transform (from feature indices to string feature names) which can be a problem when trying to introspect which features are most important to a model.
- there can be collisions: distinct tokens can be mapped to the same feature index. However in practice this is rarely an issue if n_features is large enough (e.g. 2^{18} for text classification problems).
- no IDF weighting as this would render the transformer stateful.

Source: https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.HashingVectorizer.html

3.3.1 Create the HashingVectorizer

```
In [22]: hashingVectorizer = HashingVectorizer(n_features=20)
vector = hashingVectorizer.transform(df.tweet)
df['vectorised_words'] = hashingVectorizer.transform(df.tweet.values)
```

3.3.2 Create the needed variables and split the dataset into test and train

```
In [23]: svc_2 = SVC(kernel='linear')
```

```
In [24]: X = hashingVectorizer.transform(df.tweet.values)
y = df['is_hate_speech'].astype(int)
```

```
In [25]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state=46)
```

3.3.3 Fit the SVC

```
In [26]: svc_2.fit(X_train, y_train)
```

```
Out[26]: SVC(kernel='linear')
```

3.3.4 Evaluate the results

```
In [27]: y_pred = svc_2.predict(X_test)
print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	1239
1	0.83	1.00	0.91	6196
accuracy			0.83	7435
macro avg	0.42	0.50	0.45	7435
weighted avg	0.69	0.83	0.76	7435

```
C:\Users\karin\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1221: 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 [28]: print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
print("Precision:",metrics.precision_score(y_test, y_pred))
print("Recall:",metrics.recall_score(y_test, y_pred))
```

```
Accuracy: 0.8333557498318762
Precision: 0.8333557498318762
Recall: 1.0
```

We can see that when using the SVC with a HashingVectorizer, the accuracy of the model goes down. this can be explained by the fact that distinct tokens can be mapped to the same feature index using the HashingVectorizer in contrast to using a CountVectorizer. However, this is still a very good result (because the accuracy is 0.83) even though it is worse than our previous one.

4. Machine learning using Naive Bayes

In statistics, Naive Bayes classifiers are a family of simple "probabilistic classifiers" based on applying Bayes' theorem with strong (naive) independence assumptions between the features. They are among the simplest Bayesian network model but coupled with kernel density estimation, they can achieve higher accuracy levels.

Naive Bayes classifiers are highly scalable, requiring a number of parameters linear in the number of variables (features/predictors) in a learning problem.

4.1 Imports

```
In [29]: from sklearn.naive_bayes import GaussianNB  
from sklearn import metrics
```

4.2 Split the dataset into test and train

```
In [30]: X = countVectorizer.transform(df.tweet.values).toarray()  
y = df['is_hate_speech'].astype(int)
```

```
In [31]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state=46)
```

4.3 Create the Gaussian Naive Bayes

```
In [32]: gnb = GaussianNB()
```

4.4 Fit the model

```
In [33]: gnb.fit(X_train, y_train)
```

```
Out[33]: GaussianNB()
```

4.5 Predict

```
In [34]: y_pred = gnb.predict(X_test)
```

4.6 Calculate the metrics (accuracy, precision, recall)

```
In [35]: print("Accuracy:", metrics.accuracy_score(y_test, y_pred))  
print("Precision:", metrics.precision_score(y_test, y_pred))  
print("Recall:", metrics.recall_score(y_test, y_pred))
```

```
Accuracy: 0.655413584398117  
Precision: 0.9037777777777778  
Recall: 0.6563912201420271
```

5. Conclusions

We can see that from the three types of classifications we tried, the one with the highest accuracy was the SVC with a CountVectorizer so this is going to be the one we use. That is not very surprising because SVM Classifiers offer good accuracy and perform faster prediction compared to the Naive Bayes algorithm. They also use less memory because they use a subset of training points in the decision phase.

References

How to Encode Text Data for Machine Learning with scikit-learn (article)- <https://machinelearningmastery.com/prepare-text-data-machine-learning-scikit-learn/>

Text Classification (article) - <https://monkeylearn.com/text-classification/>

Credits

The following Jupyter notebook was made by Karina Kozarova as part of the Artificial Intelligence specialization at Fontys University of Applied Sciences