



Department of Electrical & Computer Engineering
ENCS3340 - Artificial Intelligence

Tweet Emotion Detection

Prepared by: Mohammad Abu-Shelbaia 1200198

Instructor: Dr. Aziz Qaroush

Date: February 17, 2023

Tweet Emotion Detection

Introduction

Our project consists of six parts:

1. Data Collection
2. Data Preprocessing and Cleaning
3. Feature Extraction
4. Model Training
5. Model Evaluation
6. Model Deployment

as for the data collection, we have a gathered and classified data provided by the instructor. The data is in the form of tweets, and each tweet is labeled with positive or negative.

The report is not enough to demonstrate everything, therefore we have a [Jupyter Notebook](#) or the [Results PDF](#) that contains all the code and the results, and we will back our report with some results from the notebook.

Data Processing and Cleaning

In this section, we will clean the data and make it ready for the next-step, as we work with arabic tweets, we went over these steps:

- remove mentions, hashtags, and links
 - remove stop words, that are words that are not important in the context of the sentence
 - remove punctuations
 - replace emojis with their meaning
- as for all the steps, we used the python library `nltk` to do the job.

Feature Extraction

In this part we will extract the features from the data, we will use the `TF-IDF` to extract the features, and what TF-IDF does is that it gives a score to each word in the sentence, and the score is based on how many times the word appears in the sentence, and how many times the word appears in the whole dataset, and the score is calculated by this formula:

$$Score = \frac{\text{number of times the word appears in the sentence}}{\text{number of words in the sentence}} \cdot \log \frac{\text{number of sentences}}{\text{number of sentences that contain the word}}$$

Additionally after extracting the features, we used `joblib` to save the features so we can use them later.

Model Training

In this part we will train 3 different models, on the data in two different ways, the first way is 75 - 25 split, and the second way is 5-fold cross validation, the models we will use are:

- Naive Bayes
- Random Forest
- SVM

after training the models, we saved them using `joblib` so we can use them later.

Regarding the two different ways of training the models, we found that the 5-fold cross validation gave us better results, but it took more time to train the models. Which is not a problem for us, as we evaluate the models once and then we deploy them.

e.g for the Naive Bayes model, the accuracy was 0.93 in the 75-25 split, and 0.93 in the 5-fold cross validation. But the 5-fold cross validation took `15m 20.1s` to train the model, while the 75-25 split took `2m 55.1s`.

```
# 75-25 split
Accuracy: 0.93
      precision    recall  f1-score   support

neg      0.93      0.94      0.93      5582
pos      0.94      0.93      0.94      5737

accuracy            0.93      11319
```

```

macro avg      0.93      0.93      0.93      11319
weighted avg    0.93      0.93      0.93      11319

Confusion Matrix:
[[5226  356]
 [ 382 5355]]

```

```

# 5-fold cross validation
Classification Report:
              precision    recall  f1-score   support

   neg         0.94         0.92         0.93     22514
   pos         0.93         0.94         0.93     22761

 accuracy              0.93         45275
macro avg         0.93         0.93         0.93     45275
weighted avg      0.93         0.93         0.93     45275

Confusion Matrix:
[[20800  1714]
 [ 1346 21415]]

```

As we see there is not much difference in the accuracy, but the 5-fold cross validation has trained and tested the model on more data, which means that it will be more accurate on unseen data.

For more details and examples, please check the [Jupyter Notebook](#).

Model Evaluation

As we know, accuracy is not the only metric we should use to evaluate a model, there are other metrics that we should use, such as:

- Accuracy: The proportion of correct predictions made by the model out of all predictions.
- Precision: The proportion of true positives (correctly predicted positive samples) out of all positive predictions (both true positives and false positives). Precision measures how often the model correctly predicts a positive sample.
- Recall: The proportion of true positives out of all actual positive samples (both true positives and false negatives). Recall measures how well the model is able to identify positive samples.
- Confusion Matrix: A table that shows the number of true positives, false positives, true negatives, and false negatives made by the model.
- F1-score: A weighted average of precision and recall. It takes into account both false positives and false negatives, making it a useful metric for imbalanced datasets where one class is more common than the other.

Precision is the fraction of true positive predictions out of all the positive predictions. It measures the accuracy of the positive predictions. High precision means that the model is making few false positive predictions. In the context of sentiment analysis, high precision means that the model is accurately predicting positive or negative sentiment.

Recall is the fraction of true positive predictions out of all the actual positive samples. It measures the completeness of the positive predictions. High recall means that the model is correctly identifying most of the positive samples. In the context of sentiment analysis, high recall means that the model is able to correctly identify most positive or negative tweets.

If the goal is to minimize false positives (i.e., avoid labeling a negative tweet as positive), then precision is more important. On the other hand, if the goal is to minimize false negatives (i.e., avoid labeling a positive tweet as negative), then recall is more important.

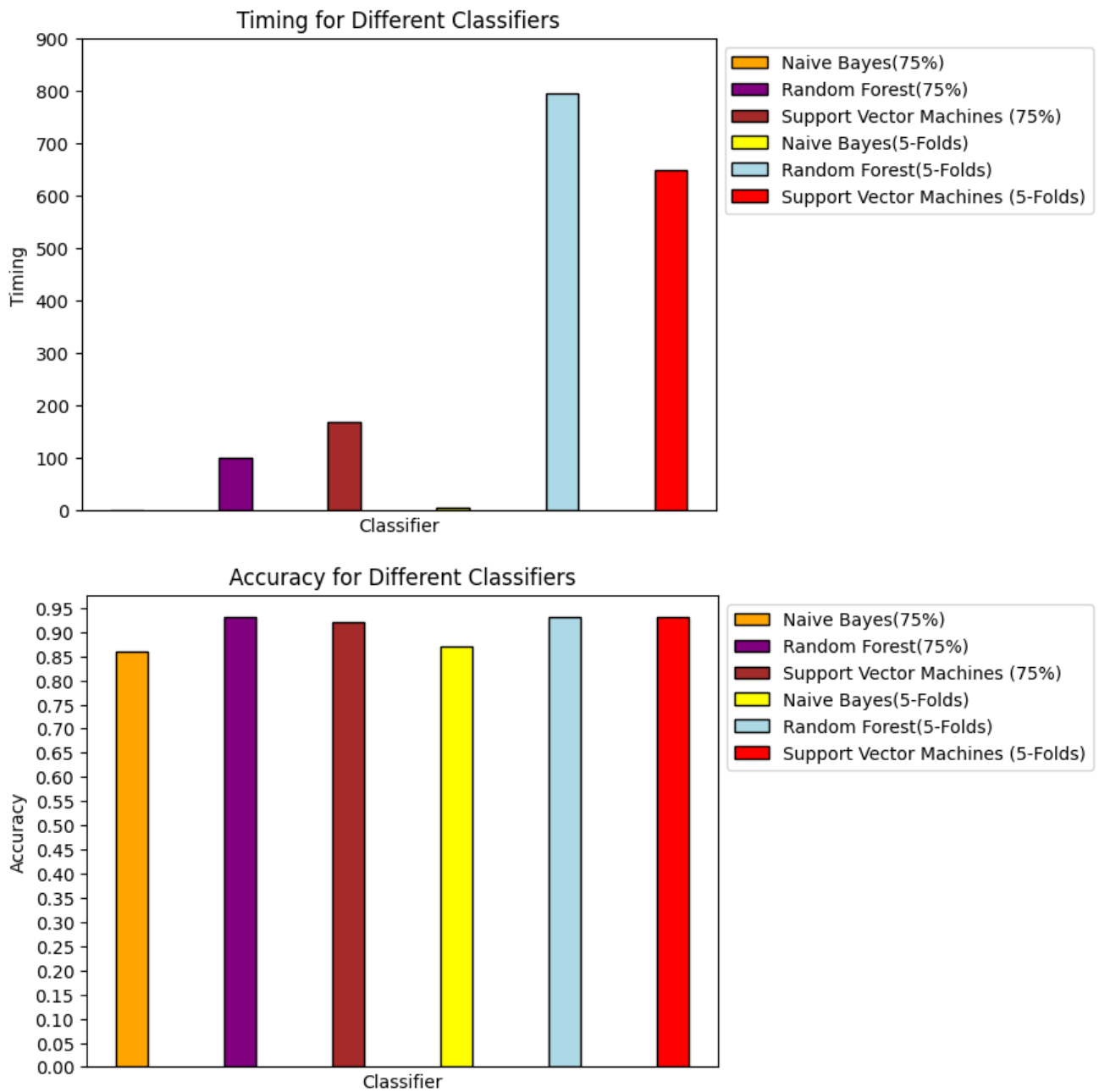
In general, it's important to consider both precision and recall together, as they are often in tension with each other. A common way to balance them is to use the F1 score, which is the harmonic mean of precision and recall.

Model Deployment

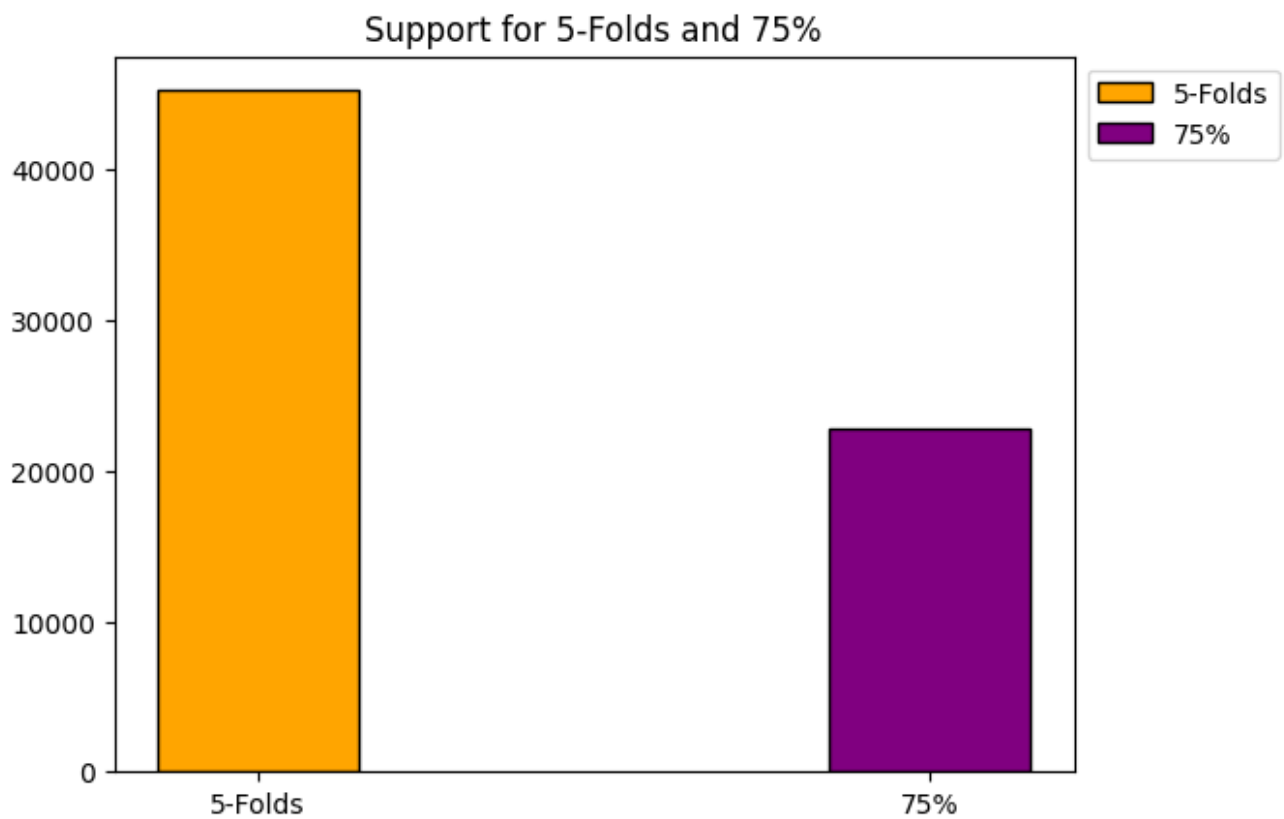
as for Model Deployment we used `joblib` to save the feature-extractor and the models, and then we ran the models using a script `main.py` that takes the tweet as an input and returns the classification of the tweet.

Results

In comparison between the models, we found that naive bayes classifier is the worst model; since it has the lowest accuracy, and the lowest f1-score, since it assumes that features are independent, which is not the case in our data, as we have a lot of repeated words, and the model is not able to capture the relationship between the words, but in contrast it is the fastest model to train since it only took a second to train the data while the other models took more than 10 minutes to train the data, as shown in the graphs below.



And as we could notice, we had a better accuracy using RandomForestClassifier, and SVM, but the difference is not that big, and the difference in the time it took to train the models is huge, as shown in the graphs above.



Furthermore, as we can see in the above graph the support: which is the number of occurrences of each class in the test set, in 5-folds cross validation it gets tested on more data, which means that it might be more accurate on unseed data, since it has been trained on more data.

Conclusion

In this project we built a model that can classify tweets as positive or negative, and we used 3 different models, and we trained them on 2 different ways, and we found that the 5-fold cross validation gave us better results, but it took more time to train the models. Which is not a problem for us, as we evaluate the models once and then we deploy them.

```
In [ ]: # Included Libraries
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import KFold
import re
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import SnowballStemmer
import emoji
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import *
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.naive_bayes import MultinomialNB
from sklearn.model_selection import cross_validate
from sklearn.model_selection import cross_val_score, train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.pipeline import make_pipeline
from sklearn.svm import SVC
import joblib
from sklearn.model_selection import cross_val_predict
```

```
In [ ]: positive_tweets = pd.read_csv('data/Positive_Tweets.tsv', sep='\t', header=None)
print("We have {} positive tweets in our dataset.".format(len(positive_tweets)))
negative_tweets = pd.read_csv('data/Negative_Tweets.tsv', sep='\t', header=None)
print("We have {} negative tweets in our dataset.".format(len(negative_tweets)))
all_tweets = pd.concat([positive_tweets, negative_tweets])
all_tweets.columns = ['Label', 'Tweet']
print("We have {} tweets in our dataset.".format(len(all_tweets)))
#A random sample of 5 tweets
print("A sample of the data\n", all_tweets.sample(5))
```

We have 22761 positive tweets in our dataset.

We have 22514 negative tweets in our dataset.

We have 45275 tweets in our dataset.

A sample of the data

	Label	Tweet
9954	pos	#... تمرين المقاومة بأنواعها : حدد العضله المراد
8791	neg	لا عاد الحلم يخوف مره 🤪
5059	pos	اروى حبيبتي نامي وش مسهرك
9483	pos	... اللهم اروي قبر ابي بفضلك ونعمتك بنورك و رحمتك
18815	pos	...البليهي أختصر مشاعر جمهور #الهلal تجاههك ي عب

Preprocessing Data

in this part we will preprocess the data and make it ready for the model

- remove mentions, hashtags, and links
- remove stop words, that are words that are not important in the context of the sentence
- remove punctuations
- replace emojis with their meaning

```
In [ ]: def proces_tweet(tweet: str) -> str:

    # Remove all mentions, hashtags, links, and special characters
    tweet = re.sub(r'http\S+|www\S+|https\S+|@\S+|#\S+', '', tweet, flags=re.MULTILINE)
    tweet = tweet.replace(' ', ' ').replace('!', '!').replace('!', '!')
    tweet = tweet.replace('ة', 'ة').replace('ي', 'ي').replace('و', 'و')

    # Remove all emojis and emoticons and replace them with unicode
    tweet = emoji.demojize(tweet)

    # Tokenize the tweet
    tokens = word_tokenize(tweet)

    # Remove all stop words, stop words are words that do not add any meaning to the sentence
    stop_words = set(stopwords.words('arabic'))
    filtered_tokens = [w for w in tokens if not w in stop_words]

    # Stemming, stemming is the process of reducing inflected (or sometimes derived) words to their word stem, base c
```

```

stemmer = SnowballStemmer('arabic')
stemmed_tokens = [stemmer.stem(w) for w in filtered_tokens]

return " ".join(stemmed_tokens)

all_tweets['Filtered_Tweet'] = all_tweets['Tweet'].apply(procces_tweet)

```

Feature Extraction

In this part we will extract the features from the data, we will use the TF-IDF to extract the features, and what TF-IDF does is that it gives a score to each word in the sentence, and the score is based on how many times the word appears in the sentence, and how many times the word appears in the whole dataset, and the score is calculated by this formula:

$$score = \frac{\text{number of times the word appears in the sentence}}{\text{number of words in the sentence}} \cdot \log\left(\frac{\text{number of sentences}}{\text{number of sentences that contain the word}}\right)$$

Additionally we are going to use 75% of the data for training, and 25% for testing, and we are going to use 5-fold cross validation to evaluate the model.

```

In [ ]: # Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(all_tweets['Filtered_Tweet'].values, all_tweets['Label'].values,

# Train on the training data and transform the training and testing data
feature_extraction = TfidfVectorizer(ngram_range=(1, 2), max_features=10000)
X_train_features = feature_extraction.fit_transform(X_train)
X_test_features = feature_extraction.transform(X_test)
joblib.dump(feature_extraction, 'models/features.joblib')

```

```

Out[ ]: ['models/features.joblib']

```

Model Training

We will use three models to train the data, a Naive Bayes model, a Search Tree model, and a not-assigned model, we will train the data in 5-fold cross validation, on a 75% training set and a 25% test set.

Naive Bayes Classifier

Naive Bayes is a probabilistic model that uses Bayes' theorem to predict the class of a given data point. It is a simple model that is easy to implement and is very fast. It is also very effective in text classification problems. The model is based on the assumption that the features are independent of each other, which is not true in most cases, but it still works well in practice.

75 - 25 Split

```
In [ ]: # Train a Naive Bayes classifier using 5-fold cross-validation
NaiveBayes_Classifier = MultinomialNB()
NaiveBayes_Classifier.fit(X_train_features, y_train)

# Evaluate the classifier using 5-fold cross-validation

# Evaluate the classifier on the test data
accuracy = NaiveBayes_Classifier.score(X_test_features, y_test)
y_pred = NaiveBayes_Classifier.predict(X_test_features)
print(f"Accuracy: {accuracy:.2f}")
print(classification_report(y_test, y_pred))
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)
print("True Positives (TP) = ", cm[0][0])
print("False Positives (FP) = ", cm[0][1])
print("False Negatives (FN) = ", cm[1][0])
print("True Negatives (TN) = ", cm[1][1])
joblib.dump(NaiveBayes_Classifier, 'models/NaiveBayes.joblib')
```

```

Accuracy: 0.86
      precision    recall  f1-score   support

   neg         0.83      0.90      0.86       5555
   pos         0.90      0.82      0.86       5764

 accuracy                   0.86       11319
 macro avg              0.86      0.86      0.86       11319
weighted avg              0.86      0.86      0.86       11319

```

Confusion Matrix:

```

[[5001  554]
 [1010 4754]]
True Positives (TP) = 5001
False Positives (FP) = 554
False Negatives (FN) = 1010
True Negatives (TN) = 4754

```

Out[]: ['models/NaiveBayes.joblib']

5 Fold Cross Validation

```

In [ ]: # Define the Random Forest model
naive_bayes_5f = MultinomialNB()

# Define the cross-validation splitter
kf = KFold(n_splits=5, shuffle=True, random_state=42)
# Define the evaluation metrics
metrics = {'accuracy': make_scorer(accuracy_score),
           'precision': make_scorer(precision_score, pos_label='pos'),
           'recall': make_scorer(recall_score, pos_label='pos'),
           'f1': make_scorer(f1_score, pos_label='pos')}

# Fit and evaluate the model with 5-fold cross-validation
X = all_tweets['Filtered_Tweet'].values
y = all_tweets['Label'].values
y_preds = cross_val_predict(naive_bayes_5f, feature_extraction.fit_transform(X), y, cv=kf)
print('Classification Report:')
print(classification_report(y, y_preds))
print('Confusion Matrix:')
print(confusion_matrix(y, y_preds))

```

```
# Fit the model on the full dataset
naive_bayes_5f.fit(feature_extraction.fit_transform(X), y)

# Save the model using joblib
joblib.dump(naive_bayes_5f, 'models/naive_bayes_5f.joblib')
```

```
Classification Report:
              precision    recall  f1-score   support

     neg         0.84         0.90         0.87         22514
     pos         0.90         0.83         0.86         22761

 accuracy                   0.87         45275
 macro avg         0.87         0.87         0.87         45275
 weighted avg         0.87         0.87         0.87         45275
```

```
Confusion Matrix:
[[20308  2206]
 [ 3835 18926]]
```

```
Out[ ]: ['models/naive_bayes_5f.joblib']
```

Decision Tree Classifier

We use Random Forest Classifier to train the data, Random Forest is an ensemble learning method for classification, regression, and other tasks, that operates by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests correct for decision trees' habit of overfitting to their training set.

75 - 25 Split

```
In [ ]: # Train a decision tree classifier on the features
RandomForest_Classifier = RandomForestClassifier(n_estimators=100)
RandomForest_Classifier.fit(X_train_features, y_train)

# Evaluate the classifier on the test data
accuracy = RandomForest_Classifier.score(X_test_features, y_test)
y_pred = RandomForest_Classifier.predict(X_test_features)
```

```

print(f"Accuracy: {accuracy:.2f}")

print(classification_report(y_test, y_pred))
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)
joblib.dump(RandomForest_Classifier, 'models/RandomFortress.joblib')

```

Accuracy: 0.93

	precision	recall	f1-score	support
neg	0.93	0.93	0.93	5555
pos	0.93	0.93	0.93	5764
accuracy			0.93	11319
macro avg	0.93	0.93	0.93	11319
weighted avg	0.93	0.93	0.93	11319

Confusion Matrix:

```

[[5162  393]
 [ 402 5362]]

```

Out[]: ['models/RandomFortress.joblib']

5 Fold Cross Validation

```

In [ ]: # Define the Random Forest model
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)

# Define the cross-validation splitter
kf = KFold(n_splits=5, shuffle=True, random_state=42)
# Define the evaluation metrics
metrics = {'accuracy': make_scorer(accuracy_score),
           'precision': make_scorer(precision_score, pos_label='pos'),
           'recall': make_scorer(recall_score, pos_label='pos'),
           'f1': make_scorer(f1_score, pos_label='pos')}

# Fit and evaluate the model with 5-fold cross-validation
X = all_tweets['Filtered_Tweet'].values
y = all_tweets['Label'].values
y_preds = cross_val_predict(rf_model, feature_extraction.fit_transform(X), y, cv=kf)
print('Classification Report:')

```

```

print(classification_report(y, y_preds))
print('Confusion Matrix:')
print(confusion_matrix(y, y_preds))

# Fit the model on the full dataset
rf_model.fit(feature_extraction.fit_transform(X), y)

# Save the model using joblib
joblib.dump(rf_model, 'models/RandomFortress_5f.joblib')

```

Classification Report:

	precision	recall	f1-score	support
neg	0.94	0.92	0.93	22514
pos	0.93	0.94	0.93	22761
accuracy			0.93	45275
macro avg	0.93	0.93	0.93	45275
weighted avg	0.93	0.93	0.93	45275

Confusion Matrix:

```

[[20800 1714]
 [ 1346 21415]]

```

Out[]: ['models/RandomFortress_5f.joblib']

SVMs (Support Vector Machines)

```

In [ ]: # Create a pipeline for the SVM classifier
svm_classifier = SVC(kernel='linear')
svm_classifier.fit(X_train_features, y_train)

# Evaluate the classifier on the test data
accuracy = svm_classifier.score(X_test_features, y_test)
print(f"Accuracy: {accuracy:.2f}")
print(classification_report(y_test, y_pred))
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)
joblib.dump(svm_classifier, 'models/SVM.joblib')

```

Accuracy: 0.92

	precision	recall	f1-score	support
neg	0.93	0.91	0.92	5555
pos	0.92	0.93	0.92	5764
accuracy			0.92	11319
macro avg	0.92	0.92	0.92	11319
weighted avg	0.92	0.92	0.92	11319

Confusion Matrix:

```
[[5061 494]
 [ 382 5382]]
```

Out[]: ['models/SVM.joblib']

5 Fold Cross Validation

```
In [ ]: # Define the Random Forest model
svc_model_5f = SVC(kernel='linear')

# Define the cross-validation splitter
kf = KFold(n_splits=5, shuffle=True, random_state=42)
# Define the evaluation metrics
metrics = {'accuracy': make_scorer(accuracy_score),
           'precision': make_scorer(precision_score, pos_label='pos'),
           'recall': make_scorer(recall_score, pos_label='pos'),
           'f1': make_scorer(f1_score, pos_label='pos')}

# Fit and evaluate the model with 5-fold cross-validation
X = all_tweets['Filtered_Tweet'].values
y = all_tweets['Label'].values
y_preds = cross_val_predict(svc_model_5f, feature_extraction.fit_transform(X), y, cv=kf)
print('Classification Report:')
print(classification_report(y, y_preds))
print('Confusion Matrix:')
print(confusion_matrix(y, y_preds))

# Fit the model on the full dataset
svc_model_5f.fit(feature_extraction.fit_transform(X), y)
```

```
# Save the model using joblib
joblib.dump(svc_model_5f, 'models/svc_model_5f.joblib')
```

Classification Report:

	precision	recall	f1-score	support
neg	0.93	0.91	0.92	22514
pos	0.91	0.94	0.92	22761
accuracy			0.92	45275
macro avg	0.92	0.92	0.92	45275
weighted avg	0.92	0.92	0.92	45275

Confusion Matrix:

```
[[20496 2018]
 [ 1440 21321]]
```

```
Out[ ]: ['models/svc_model_5f.joblib']
```