

```
In [3]: 1 #import necessary libraries
2
3 import pandas as pd
4 import numpy as np
5 import string
6 import re
7 from matplotlib import pyplot as plt
8 import seaborn as sns
9 import nltk
10 from nltk.corpus import stopwords
11 from nltk import FreqDist, word_tokenize
12 from nltk.tokenize import TweetTokenizer
13 from nltk.stem import WordNetLemmatizer
14 from sklearn.feature_extraction.text import TfidfVectorizer
15 import unicode
16 import html
17
18 from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV
19
20 from imblearn.ensemble import BalancedRandomForestClassifier
21 from sklearn.ensemble import RandomForestClassifier
22 from xgboost import XGBClassifier
23 from sklearn.naive_bayes import MultinomialNB
24
25 from imblearn.pipeline import make_pipeline
26 from sklearn.model_selection import cross_val_score
27 from sklearn.metrics import accuracy_score, plot_confusion_matrix, confusion_matrix
28
29
30 from bert_sklearn import BertClassifier
31 from bert_sklearn import BertRegressor
32 from bert_sklearn import load_model
```

Apple Tweet Sentiment Analysis

Modeling Notebook

Author: Dylan Dey

The Author can be reached at the following email: ddey2985@gmail.com
(<mailto:ddey2985@gmail.com>)

Blog: [Quick BERT Pre-Trained Model for Sentiment Analysis with Scikit Wrapper](https://dev.to/ddey117/quick-bert-pre-trained-model-for-sentiment-analysis-with-scikit-wrapper-3jcp)
(<https://dev.to/ddey117/quick-bert-pre-trained-model-for-sentiment-analysis-with-scikit-wrapper-3jcp>)

Classification Metric Understanding

ACTUAL	Predicted	
	Positive Tweet	Negative Tweet
Positive Tweet	True Positive Tweet A tweet that is correctly predicted to contain positive sentiment by model.	False Negative The model incorrectly identifies a tweet that contains positive sentiment as one that contains negative sentiment. Given the context of the business problem, this would mean extra noise added when trying to isolate for negative sentiment of brand/product.
Negative Tweet	False Positive A false positive would occur when the model incorrectly identifies a tweet containing negative sentiment as a tweet that contains positive sentiment. Given the context of the business model, this would mean more truly negative sentiment will be left out of analyzing key word pairs for negative tweets.	True Negative Tweet A tweet that is correctly predicted to contain negative sentiment by model.

Confusion Matrix Description

There will always be some error involved in creating a predictive model. The model will incorrectly identify positive tweets as negative and vice versa. That means the error in any classification model in this context can be described by ratios of true positives or negatives vs false positives or negatives.

Correctly predicting a tweet to have negative sentiment is at the heart of the model, as this is the situation in which a company would have a call to action. An appropriately identified tweet with negative sentiment can be properly examined using some simple NLP techniques to get a quick buy effective way to view what is upsetting customers about the company it's products.

Correctly predicting a tweet to have positive sentiment is also important. Word frequency analysis can be used to summarize what consumers think Apple is doing right and also what consumers like about Apple's competitors.

A false positive would occur when the model incorrectly identifies a tweet containing negative sentiment as a tweet that contains positive sentiment. Given the context of the business model, this would mean more truly negative sentiment will be left out of analyzing key word pairs for negative tweets. This could be interpreted as loss in analytical ability for what we care about most given the business problem: making informed decisions from information directly from consumers in the form of social media text. Minimizing false positives is important.

False negatives are also important to consider. A false negative would occur when the model incorrectly identifies a tweet that contains positive sentiment as one that contains negative sentiment. Given the context of the business problem, this would mean extra noise added to the data when trying to isolate for negative sentiment of brand/product.

In summary, overall accuracy of the model and a reduction of both false negatives and false positives are the most important metrics to consider when developing the sentiment analysis model.

MVP Metric

[balanced_accuracy_score \(https://scikit-learn.org/stable/modules/model_evaluation.html#scoring-parameter\)](https://scikit-learn.org/stable/modules/model_evaluation.html#scoring-parameter)

From the documentation:

"The balanced accuracy in binary and multiclass classification problems to deal with imbalanced datasets. It is defined as the average of recall obtained on each class."

This is a great metric for this problem as optimizing for the average of recall for each class will give the best performance given the context of the business problem.

Function Definition

```

In [35]: 1 #list of all functions for modeling
2 #and processing
3
4 #force lowercase of text data
5 def lower_case_text(text_series):
6     text_series = text_series.apply(lambda x: str.lower(x))
7     return text_series
8
9 #remove URL links from text
10 def strip_links(text):
11     link_regex = re.compile('((https?):(\\//\\/)|(\\\\\\\\))+(\\w\\d:##%\\/;$)')
12     links = re.findall(link_regex, text)
13     for link in links:
14         text = text.replace(link[0], ' ')
15     return text
16
17 #remove '@' and '#' symbols from text
18 def strip_all_entities(text):
19     entity_prefixes = ['@', '#']
20     for separator in string.punctuation:
21         if separator not in entity_prefixes:
22             text = text.replace(separator, ' ')
23     words = []
24     for word in text.split():
25         word = word.strip()
26         if word:
27             if word[0] not in entity_prefixes:
28                 words.append(word)
29     return ' '.join(words)
30
31 #tokenize text and remove stopwords
32 def process_text(text):
33     tokenizer = TweetTokenizer()
34
35     stopwords_list = stopwords.words('english') + list(string.punctuation)
36     stopwords_list += ['"', "'", '...', '`']
37     my_stop = ["#sxsw",
38               "sxsw",
39               "sxswi",
40               "#sxswi's",
41               "#sxswi",
42               "southbysouthwest",
43               "rt",
44               "tweet",
45               "tweet's",
46               "twitter",
47               "austin",
48               "#austin",
49               "link",
50               "1/2",
51               "southby",
52               "south",
53               "texas",
54               "@mention",
55               "i",
56               "i",

```

```

57         "½",
58         "¿",
59         "½",
60         "link",
61         "via",
62         "mention",
63         "quot",
64         "amp",
65         "austin"
66     ]
67
68     stopwords_list += my_stop
69
70     tokens = tokenizer.tokenize(text)
71     stopwords_removed = [token for token in tokens if token not in stopwords_list]
72     return stopwords_removed
73
74
75
76 #master preprocessing function
77 def Master_Pre_Vectorization(text_series):
78     text_series = lower_case_text(text_series)
79     text_series = text_series.apply(strip_links).apply(strip_all_entities)
80     text_series = text_series.apply(unidecode.unidecode).apply(html.unescape)
81     text_series = text_series.apply(process_text)
82     lemmatizer = WordNetLemmatizer()
83     text_series = text_series.apply(lambda x: [lemmatizer.lemmatize(word) for word in x])
84     return text_series.str.join(' ').copy()
85
86
87 #function for interpreting results of models
88 #takes in a pipeline and training data
89 #and prints cross_validation scores
90 #and average of scores
91 #It also returns balanced_accuracy_score
92
93 def cross_validation_plus(pipeline, X_train, y_train, X_test, y_test):
94     scores = cross_val_score(pipeline, X_train, y_train)
95     agg_score = np.mean(scores)
96     print(f'{pipeline.steps[1][1]}: Average cross validation score is {agg_score}')
97     pipeline.fit(X_train, y_train)
98     y_pred = pipeline.predict(X_test)
99     balanced_accuracy = balanced_accuracy_score(y_test, y_pred)
100    print(f'{pipeline.steps[1][1]}: Balanced accuracy score is {balanced_accuracy}')
101    print("-----")
102
103
104 #function to fit pipeline
105 #and return subplots
106 #that show normalized and
107 #regular confusion matrices
108 #to easily interpret results
109 def plot_confusion_matrices(pipe, pathway):
110
111     pipe.fit(X_train, y_train)
112
113     y_pred = pipe.predict(X_test)

```

```

114
115 matrix_norm = confusion_matrix(y_test, y_pred, normalize='true')
116 matrix = confusion_matrix(y_test, y_pred)
117
118 fig, (ax1, ax2) = plt.subplots(ncols = 2,figsize=(10, 5))
119 sns.heatmap(matrix_norm,
120             annot=True,
121             fmt='.2%',
122             cmap='YlGn',
123             xticklabels=['Pos_predicted', 'Neg_predicted'],
124             yticklabels=['Positive Tweet', 'Negative_Tweet'],
125             ax=ax1)
126 sns.heatmap(matrix,
127             annot=True,
128             cmap='YlGn',
129             fmt='d',
130             xticklabels=['Pos_predicted', 'Neg_predicted'],
131             yticklabels=['Positive Tweet', 'Negative_Tweet'],
132             ax=ax2)
133
134 plt.savefig(pathway)
135
136 plt.show();
137
138
139
140
141 #loads a fitted model from memory
142 #returns confusion matrix and
143 #returns normalized confusion matrix
144 #calculated using given test data
145 def confusion_matrix_bert_plots(model_path, X_test, y_test, fig_pathwa
146
147     model = load_model(model_path)
148
149     y_pred = model.predict(X_test)
150
151     matrix_norm = confusion_matrix(y_test, y_pred, normalize='true')
152
153     matrix = confusion_matrix(y_test, y_pred)
154
155     fig, (ax1, ax2) = plt.subplots(ncols = 2,figsize=(10, 5))
156     sns.heatmap(matrix_norm,
157                 annot=True,
158                 fmt='.2%',
159                 cmap='YlGn',
160                 xticklabels=['Pos_predicted', 'Neg_predicted'],
161                 yticklabels=['Positive Tweet', 'Negative_Tweet'],
162                 ax=ax1)
163     sns.heatmap(matrix,
164                 annot=True,
165                 cmap='YlGn',
166                 fmt='d',
167                 xticklabels=['Pos_predicted', 'Neg_predicted'],
168                 yticklabels=['Positive Tweet', 'Negative_Tweet'],
169                 ax=ax2)
170

```

```
171 plt.savefig(fig_pathway);  
172 plt.show();
```

```
In [36]: 1 #import cleaned dataset  
2 df = pd.read_csv('data/clean_df.csv')  
3 df.drop(columns=['Unnamed: 0'], inplace=True)  
4  
5 # X = df['tweet'].str.join(' ').copy()  
6 X = df['tweet'].copy()  
7 y = df['target'].copy()  
8  
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  
10  
11 #clean and prepare data  
12 #for TF_IDF vector transformation  
13 X_train = Master_Pre_Vectorization(X_train)  
14 X_test = Master_Pre_Vectorization(X_test)  
15  
16 #vecorize text data  
17 vectorizer = TfidfVectorizer()  
18 tf_idf_X_train = vectorizer.fit_transform(X_train)  
19 tf_idf_X_test = vectorizer.transform(X_test)
```

```
In [37]: 1 vectorizer = TfidfVectorizer()
2
3 #multinomial bayes classifier
4 nb_classifier = MultinomialNB()
5 NB_pipe = make_pipeline(vectorizer, nb_classifier)
6 cross_validation_plus(NB_pipe, X_train, y_train, X_test, y_test)
7
8 #random forest classifier
9 rf_classifier = RandomForestClassifier(n_estimators=100)
10 rf_pipe = make_pipeline(vectorizer, rf_classifier)
11 cross_validation_plus(rf_pipe, X_train, y_train, X_test, y_test)
12
13 #balanced random forest classifier
14 balanced_rf_classifier = BalancedRandomForestClassifier(n_estimators=100)
15 balanced_rf_pipe = make_pipeline(vectorizer, balanced_rf_classifier)
16 cross_validation_plus(balanced_rf_pipe, X_train, y_train, X_test, y_test)
17
18 #XGBoosted classifier
19 xgb_classifier = XGBClassifier()
20 xgb_pipe = make_pipeline(vectorizer, xgb_classifier)
21 cross_validation_plus(xgb_pipe, X_train, y_train, X_test, y_test)
```

MultinomialNB(): Average cross validation score is 0.8134146341463415.

MultinomialNB(): Balanced accuracy score is 0.7062083711377911.

RandomForestClassifier(): Average cross validation score is 0.8689024390243902.

RandomForestClassifier(): Balanced accuracy score is 0.7883266867821549.

BalancedRandomForestClassifier(): Average cross validation score is 0.8408536585365853.

BalancedRandomForestClassifier(): Balanced accuracy score is 0.8146737563212569.

XGBClassifier(): Average cross validation score is 0.8033536585365854.

XGBClassifier(): Balanced accuracy score is 0.6364406839265996.

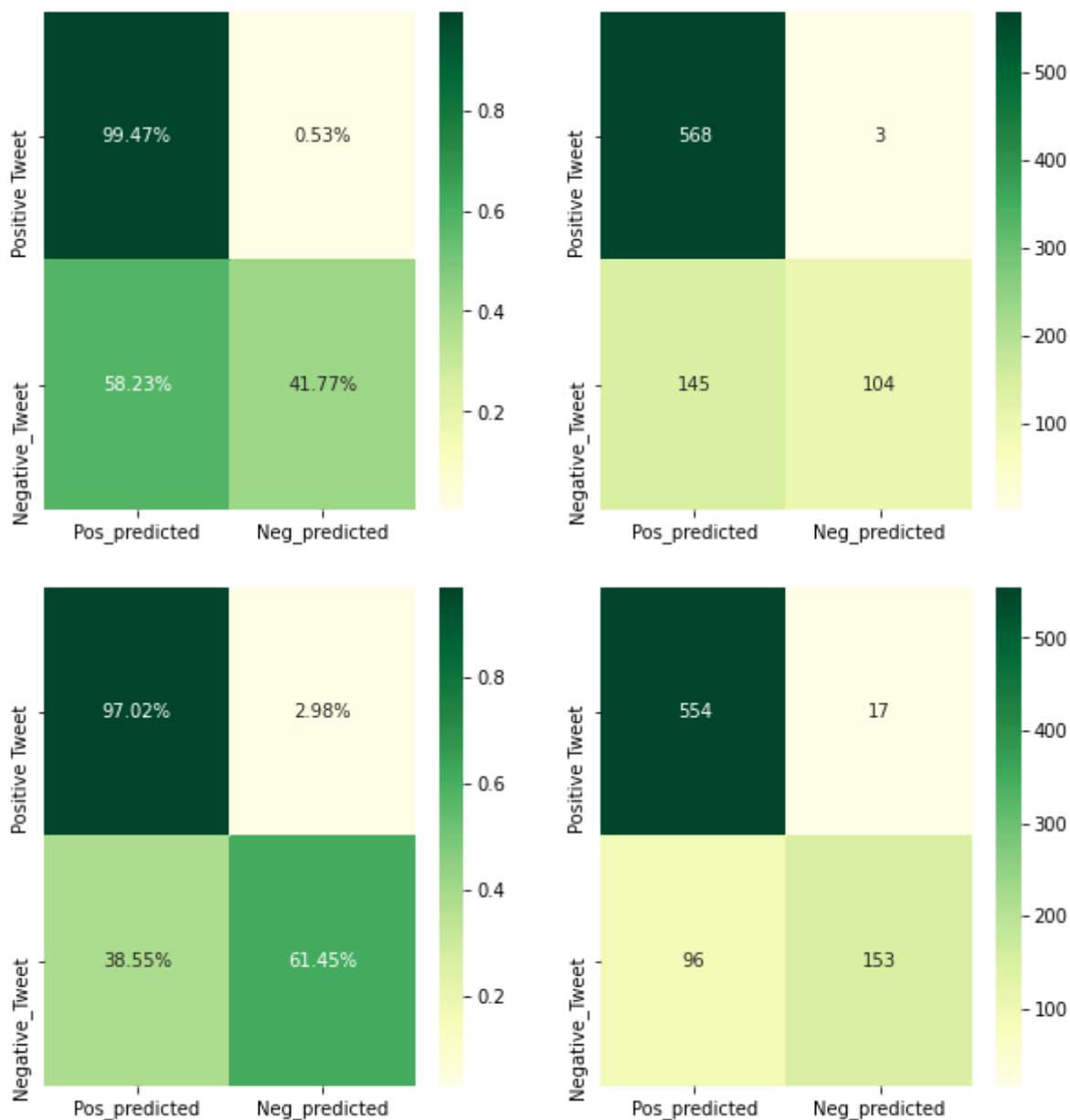
```
In [38]: 1 pipes = [NB_pipe,
2           rf_pipe,
3           balanced_rf_pipe,
4           xgb_pipe
5           ]
```

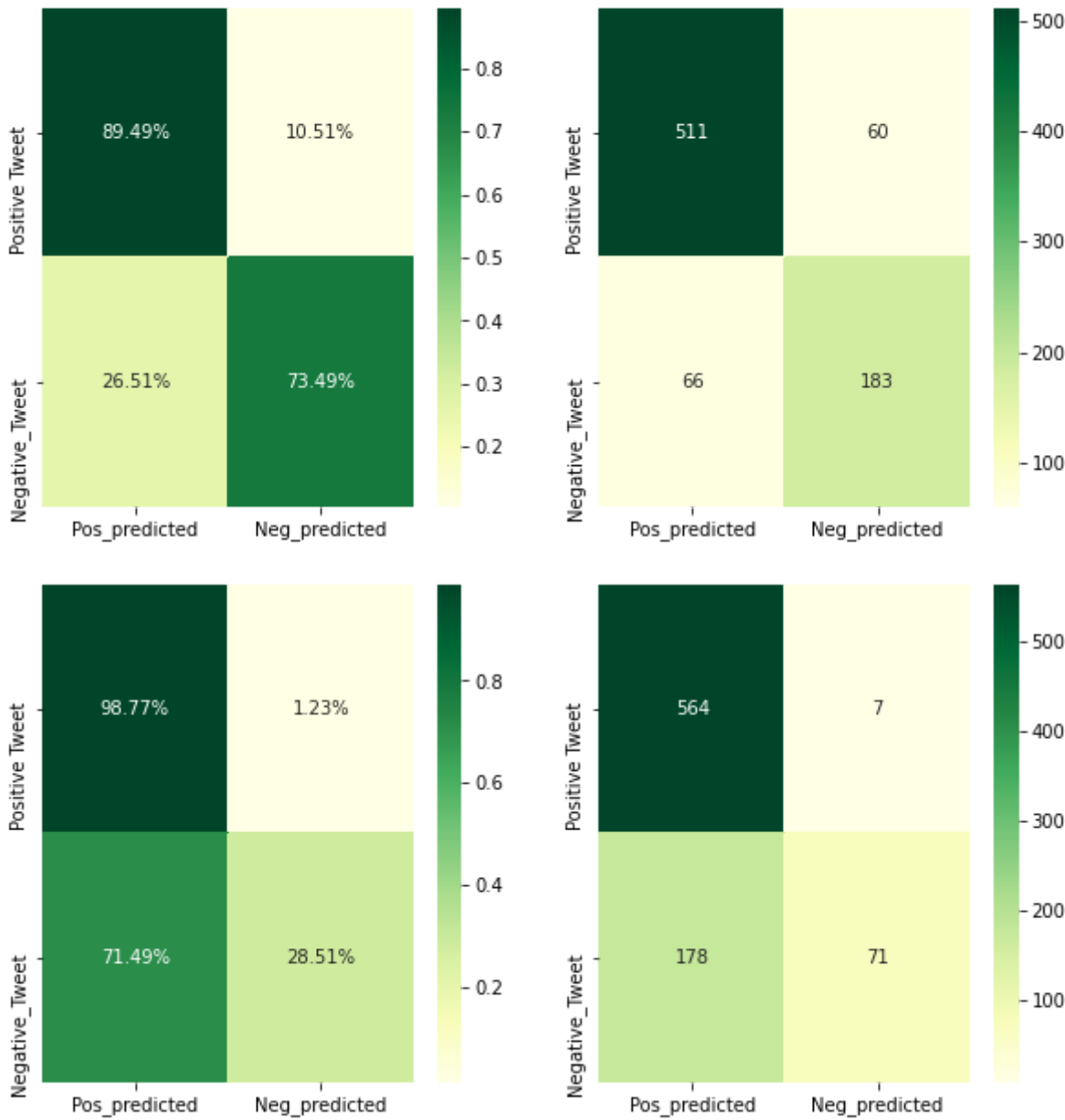


```

In [39]: 1 pathways = ['images/XGBoosted_matrix',
2               'images/balanced_RF_matrix',
3               'images/RF_matrix',
4               'images/NB_matrix'
5               ]
6
7
8 for pipe in pipes:
9     pathway = pathways.pop()
10    plot_confusion_matrices(pipe, pathway)

```



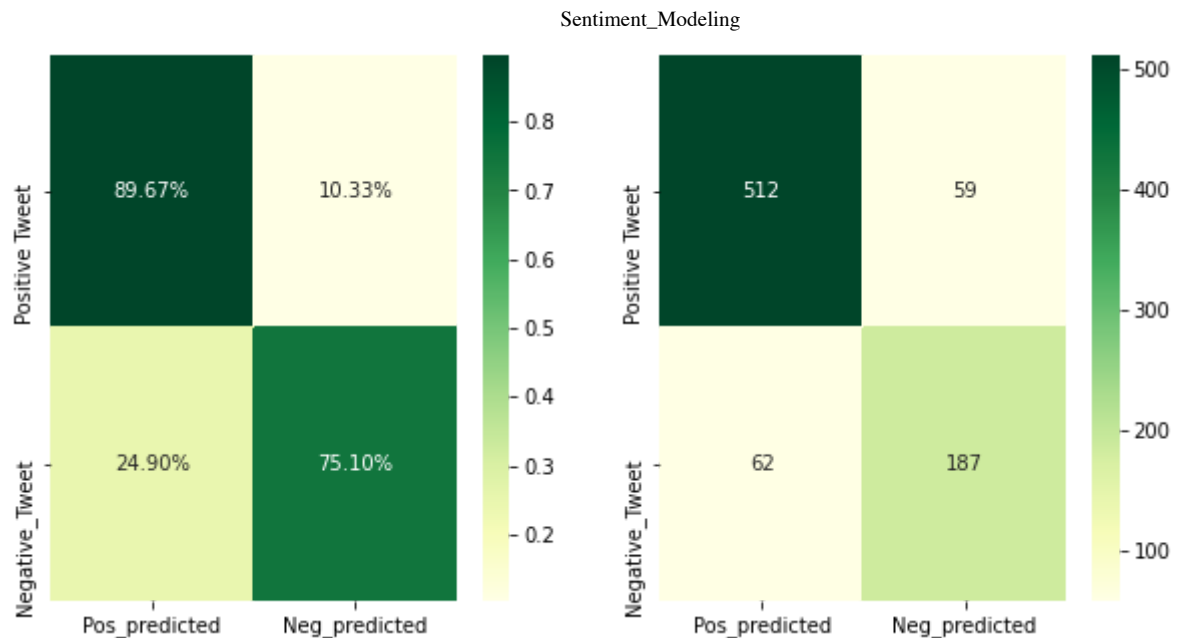


```

In [41]:
1  #initialize grid search variables
2  n_estimators = [int(x) for x in np.linspace(start = 10, stop = 200, num
3  criterion = ["gini", "entropy"]
4  min_samples_split = [8, 10, 12]
5  max_depth = [int(x) for x in np.linspace(10, 1000, num = 10)]
6  min_samples_leaf = [0.01, 0.1, 1, 2, 4]
7
8  # Create the random grid
9  random_grid = {'n_estimators': n_estimators,
10                 'criterion': criterion,
11                 'max_depth': max_depth,
12                 'min_samples_split': min_samples_split,
13                 'min_samples_leaf': min_samples_leaf
14                 }
15
16  #rrandomly iterate 1667*3 times through the grid
17  balanced_rfc_rs = RandomizedSearchCV(estimator = BalancedRandomForestCl
18                                     param_distributions = random_grid,
19                                     scoring = 'balanced_accuracy',
20                                     n_iter = 1667,
21                                     cv = 3,
22                                     verbose=2,
23                                     random_state=11,
24                                     n_jobs = -1
25                                     )
26
27
28  #fit random grid search and determine best_estimator_
29  balanced_rfc_rs.fit(tf_idf_X_train, y_train)
30
31  #create pipeline for best result from random grid search
32  balanced_rfc_rs_pipe = make_pipeline(vectorizer,
33                                     balanced_rfc_rs.best_estimator_)
34
35  cross_validation_plus(balanced_rfc_rs_pipe, X_train, y_train, X_test, y
36  plot confusion matrices(balanced rfc rs pipe, 'images/best balanced rfc

```

```
Fitting 3 folds for each of 1667 candidates, totalling 5001 fits
BalancedRandomForestClassifier(criterion='entropy', max_depth=340,
                                min_samples_split=8, n_estimators=94): Ave
rage cross validation score is 0.8408536585365853.
BalancedRandomForestClassifier(criterion='entropy', max_depth=340,
                                min_samples_split=8, n_estimators=94): Bal
anced accuracy score is 0.81194128528123.
```



Now that supervised learning models have been built, trained, and tuned without being pre-trained on any other data, our focus will now turn to transfer learning using Bidirectional Encoder Representations from Transformers (BERT), developed by Google. BERT is a transformer-based machine learning technique for natural language processing pre-training. BERTBASE models are pre-trained from unlabeled data extracted from the BooksCorpus with 800M words and English Wikipedia with 2,500M words. For this project, a BERT_base model will be trained (110 million parameters in the core of the network). Results of tuning a BERT_large model will be added as future work.

A really interesting and succinct blog by Author Rani Horev about [BERT](https://towardsdatascience.com/bert-explained-state-of-the-art-language-model-for-nlp-f8b21a9b6270) (<https://towardsdatascience.com/bert-explained-state-of-the-art-language-model-for-nlp-f8b21a9b6270>) briefly covers the architecture and capabilities of this massive pre-trained model.

[Click Here for more from Wikipedia](https://en.wikipedia.org/wiki/BERT_(language_model)) ([https://en.wikipedia.org/wiki/BERT_\(language_model\)](https://en.wikipedia.org/wiki/BERT_(language_model))).

[GitHub for BERT release code](https://github.com/google-research/bert) (<https://github.com/google-research/bert>).

Scikit-learn wrapper provided by Charles Nainan. [GitHub of Scikit Learn BERT wrapper](https://github.com/charles9n/bert-sklearn) (<https://github.com/charles9n/bert-sklearn>).

This scikit-learn wrapper is used to finetune Google's BERT model and is built on the huggingface pytorch port.

Below is the documentation for Hugging Face Transformers. Taken from the website: "Transformers (formerly known as pytorch-transformers and pytorch-pretrained-bert) provides thousands of pretrained models to perform tasks on different modalities such as text, vision, and audio." [Hugging Face Transformers](https://huggingface.co/docs/transformers/index) (<https://huggingface.co/docs/transformers/index>).

```
In [9]: 1 """
2 The first model was fitted as seen commented out below
3 after some trial and error to determine an appropriate
4 max_seq_length given my computer's capabilities.
5
6 """
7
8
9 # bert_1 = BertClassifier(do_lower_case=True,
10 #                         train_batch_size=32,
11 #                         max_seq_length=50
12 #                         )
13
14
15
16 """
17 My second model contains 2 hidden layers with 600 neurons.
18 It only passes over the corpus one time when learning.
19 It trains fast and gives impressive results.
20
21 """
22
23
24 # bert_2 = BertClassifier(do_lower_case=True,
25 #                         train_batch_size=32,
26 #                         max_seq_length=50,
27 #                         num_mlp_hiddens=500,
28 #                         num_mlp_layers=2,
29 #                         epochs=1
30 #                         )
31
32 """
33 My third bert model has 600 neurons still but
34 only one hidden layer. However, the model
35 passes over the corpus 4 times in total
36 while learning.
37
38 """
39
40 # bert_3 = BertClassifier(do_lower_case=True,
41 #                         train_batch_size=32,
42 #                         max_seq_length=50,
43 #                         num_mlp_hiddens=600,
44 #                         num_mlp_layers=1,
45 #                         epochs=4
46 #                         )
47
48 """
49 My fourth bert model has 750 neurons and
50 two hidden layers. The corpus also gets
51 transversed four times in total while
52 learning.
53
54 """
55
56 # bert_4 = BertClassifier(do_lower_case=True,
```

```

57 # train_batch_size=32,
58 # max_seq_length=50,
59 # num_mlp_hiddens=750,
60 # num_mlp_layers=2,
61 # epochs=4
62 # )

```

Out[9]: '\nMy fourth bert model has 750 neurons and \ntwo hidden layers. The corpus also gets\ntransversed four times in total while \nlearning.\n\n'

```

In [11]: 1 #Review confusion matrix plots
2 #For all bert models saved in memory
3
4 bert_paths= ['data/bert_model_1.bin',
5              'data/bert_model_2.bin',
6              'data/bert_model_3.bin',
7              'data/bert_model_4.bin'
8              ]
9
10 figure_paths = ['images/bert4_matrix.jpg',
11                 'images/bert3_matrix.jpg',
12                 'images/bert2_matrix.jpg',
13                 'images/bert1_matrix.jpg',
14                 ]
15
16 for bert_path in bert_paths:
17     figure_path = figure_paths.pop()
18     confusion_matrix_bert_plots(bert_path, X_test, y_test, figure_path)

```

Loading model from data/bert_model_1.bin...

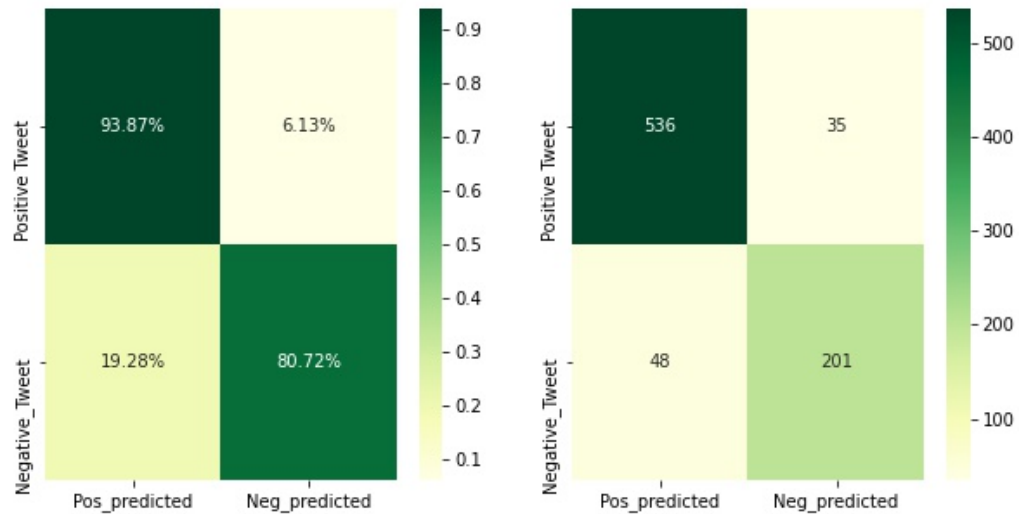
```

02/09/2022 21:58:35 - INFO - bert_sklearn.model.pytorch_pretrained.modeling - Model config {
  "architectures": [
    "BertForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "bert",
  "num_attention_heads": 12,
  "num_hidden_layers": 12,
  "pad_token_id": 0,
  "type_vocab_size": 2
}

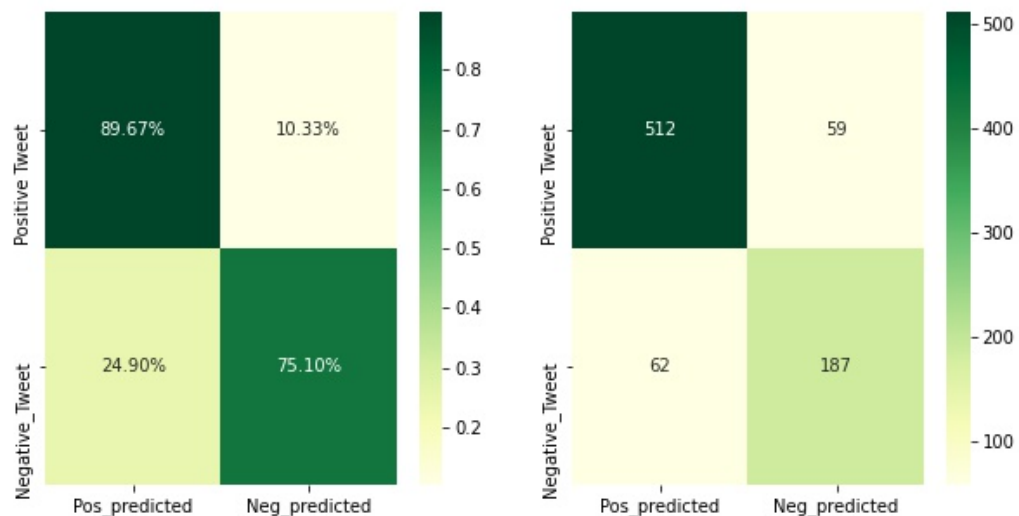
```

Evaluation

The best performing model was the BERT Classifier with 4 epochs, one hidden layer, and 600 neurons. This classifier was able to correctly predict over 80% of negative tweets correctly, which is really impressive given the imbalance in the original data. It also correctly identifies positive tweets nearly 94% of the time.

Balanced Random Forest Confusion Matrix

While the BERT classifier performed the best, the balanced random forest classifier has moderate predictive abilities using sparse vectors.

Balanced Random Forest Confusion Matrix**Conclusions**

- Either classifier could be used to predict sentiment on new brand-centric social media data for the company's own products or that of a competitor.

Future Work

- Use the BERT classifier to predict the sentiment on new unlabeled twitter data filtered for product or brand of interest (Apple/Google) from another source to find more actionable

insights to further proof of concept.

- Use the BERT classifier to predict the sentiment on new twitter data to help balance existing dataset and retrain the other models.
- leverage a state-of-the-art early stopping algorithm (ASHA) using Ray Tune and PyTorch.(1)(2)

(1)Author Amog Kamsetty explores the importance of hyperparameter tuning in his blog [Hyperparameter Optimization for Transformers: A guide \(https://medium.com/distributed-computing-with-ray/hyperparameter-optimization-for-transformers-a-guide-c4e32c6c989b\)](https://medium.com/distributed-computing-with-ray/hyperparameter-optimization-for-transformers-a-guide-c4e32c6c989b). This [Colabrative Notebook \(https://colab.research.google.com/drive/1tQgAKgcKQzheoh503OzhS4N9NtfFgmjF?usp=sharing\)](https://colab.research.google.com/drive/1tQgAKgcKQzheoh503OzhS4N9NtfFgmjF?usp=sharing) shared in the blog is a good starting point to try optimize with Ray Tune.

(2)Author Richard Liaw shares a blog that shows how simple it is to leverage all of the cores and GPUs on your machine to perform parallel asynchronous hyperparameter tuning and how to launch a massive distributed hyperparameter search on the cloud (and automatically shut down hardware after completion). [Ray Tune: a Python library for fast hyperparameter tuning at any scale \(https://towardsdatascience.com/fast-hyperparameter-tuning-at-scale-d428223b081c\)](https://towardsdatascience.com/fast-hyperparameter-tuning-at-scale-d428223b081c) also showcases a lot of exciting algorithms to explore when tuning models.