```
In [1]:    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 unidecode
          16  import html
          17
          18  from sklearn.model_selection import train_test_split, GridSearchCV, Ran
          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, conf
          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

This project it available on github here: link

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

Blog: blog link

**Classification Metric Understanding**

**ACTUAL**

|  | **True Positive Tweet** | **False Negative** |
|---|---|---|
| **Positive Tweet** | A tweet that is correctly predicted to contain positive sentiment by model. | 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. |
| **Predicted** | Positive Tweet | Negative Tweet |

## Confusion Matrix Description

A true positive in the current context would be when the model correctly identifies a tweet with positive sentiment as positive. A true negative would be when the model correctly identifies a tweet with negative sentiment as containing negative sentiment. Both are important and both can be described by the overall accuracy of the model.

True negatives are really at the heart of the model, as this is the situation in which Apple 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 better understanding at what is upsetting customers involved with our brand and competitor's brands. Bigrams, quadgrams, and other word frequency analysis can help Apple to address brand concerns.

True positives are 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.

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.

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 buisness problem: making informed decisions from information directly from consumers in the form of social media text. Minimizing false positives is very 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 analyisis model.

## Function Definition

In [2]:

```python
#list of all functions for modeling
#and processing

#force lowercase of text data
def lower_case_text(text_series):
    text_series = text_series.apply(lambda x: str.lower(x))
    return text_series

#remove URL links from text
def strip_links(text):
    link_regex = re.compile('((https?):((\/\/)|(\\\\))+([\w\d:#@%\/;$(
    links = re.findall(link_regex, text)
    for link in links:
        text = text.replace(link[0], ', ')
    return text

#remove '@' and '#' symbols from text
def strip_all_entities(text):
    entity_prefixes = ['@','#']
    for separator in  string.punctuation:
        if separator not in entity_prefixes:
            text = text.replace(separator,' ')
    words = []
    for word in text.split():
        word = word.strip()
        if word:
            if word[0] not in entity_prefixes:
                words.append(word)
    return ' '.join(words)

#tokenize text and remove stopwords
def process_text(text):
    tokenizer = TweetTokenizer()

    stopwords_list = stopwords.words('english') + list(string.punctuat
    stopwords_list += ["'''", '"""', '...', '``']
    my_stop = ["#sxsw",
                "sxsw",
                "sxswi",
                "#sxswi's",
                "#sxswi",
                "southbysouthwest",
                "rt",
                "tweet",
                "tweet's",
                "twitter",
                "austin",
                "#austin",
                "link",
                "1/2",
                "southby",
                "south",
                "texas",
                "@mention",
                "ï",
                "ï",
```

```python
57                   "½ï",
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 sto
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_entit
80       text_series = text_series.apply(unidecode.unidecode).apply(html.un
81       text_series =text_series.apply(process_text)
82       lemmatizer = WordNetLemmatizer()
83       text_series = text_series.apply(lambda x: [lemmatizer.lemmatize(wo
84       return text_series.str.join(' ').copy()
85
86
87  #function for intepreting results of models
88  #takes in a pipeline and training data
89  #and prints cross_validation scores
90  #and average of scores
91
92
93  def cross_validation(pipeline, X_train, y_train):
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
97
98
99  #function to fit pipeline
100 #and return subplots
101 #that show normalized and
102 #regular confusion matrices
103 #to easily intepret results
104 def plot_confusion_matrices(pipe, pathway):
105
106      pipe.fit(X_train, y_train)
107      y_true = y_test
108      y_pred = pipe.predict(X_test)
109
110      matrix_norm = confusion_matrix(y_true, y_pred, normalize='true')
111      matrix = confusion_matrix(y_true, y_pred)
112
113      fig, (ax1, ax2) = plt.subplots(ncols = 2,figsize=(10, 5))
```

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

In [3]:
```python
#import cleaned dataset
df = pd.read_csv('data/clean_df.csv')
df.drop(columns=['Unnamed: 0'], inplace=True)

# X = df['tweet'].str.join(' ').copy()
X = df['tweet'].copy()
y = df['target'].copy()

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2

#clean and prepare data
#for TF_IDF vector transformation
X_train = Master_Pre_Vectorization(X_train)
X_test = Master_Pre_Vectorization(X_test)

#vecorize text data
vectorizer = TfidfVectorizer()
tf_idf_X_train = vectorizer.fit_transform(X_train)
tf_idf_X_test = vectorizer.transform(X_test)
```

In [4]:
```python
X_train
```

Out[4]:
```
1429    google building real time engine monitor illeg...
67                                          haz ifrom gr8
1703    yup apple come cool technology one ever heard ...
3819                               u used b cool happened
221     think fell bit love today thanks throwing nerd...
                              ...
1238                      pollak much trouble sell apple duh
466     new whrrl app live iphone app store android ma...
3092    wait anymore google launch major new social ne...
3772    dear ever lovin heck put navigation arrow itun...
860                        awesome traded last year ipad 1
Name: tweet, Length: 3280, dtype: object
```

In [5]:
```python
vectorizer = TfidfVectorizer()

#multinomial bayes classifier
nb_classifier = MultinomialNB()
NB_pipe = make_pipeline(vectorizer, nb_classifier)
cross_validation(NB_pipe, X_train, y_train)
#random forest classifier
rf_classifier = RandomForestClassifier(n_estimators=100)
rf_pipe = make_pipeline(vectorizer, rf_classifier)
cross_validation(rf_pipe, X_train, y_train)
#balanced random forest classifier
balanced_rf_classifier = BalancedRandomForestClassifier(n_estimators=10
balanced_rf_pipe = make_pipeline(vectorizer, balanced_rf_classifier)
cross_validation(balanced_rf_pipe, X_train, y_train)
#XGBoosted classifier
xgb_classifier = XGBClassifier()
xgb_pipe = make_pipeline(vectorizer, xgb_classifier)
cross_validation(xgb_pipe, X_train, y_train)
```

```
MultinomialNB(): Average cross validation score is 0.8134146341463415.
RandomForestClassifier(): Average cross validation score is 0.86402439024
39024.
BalancedRandomForestClassifier(): Average cross validation score is 0.839
0243902439025.
XGBClassifier(): Average cross validation score is 0.8033536585365854.
```
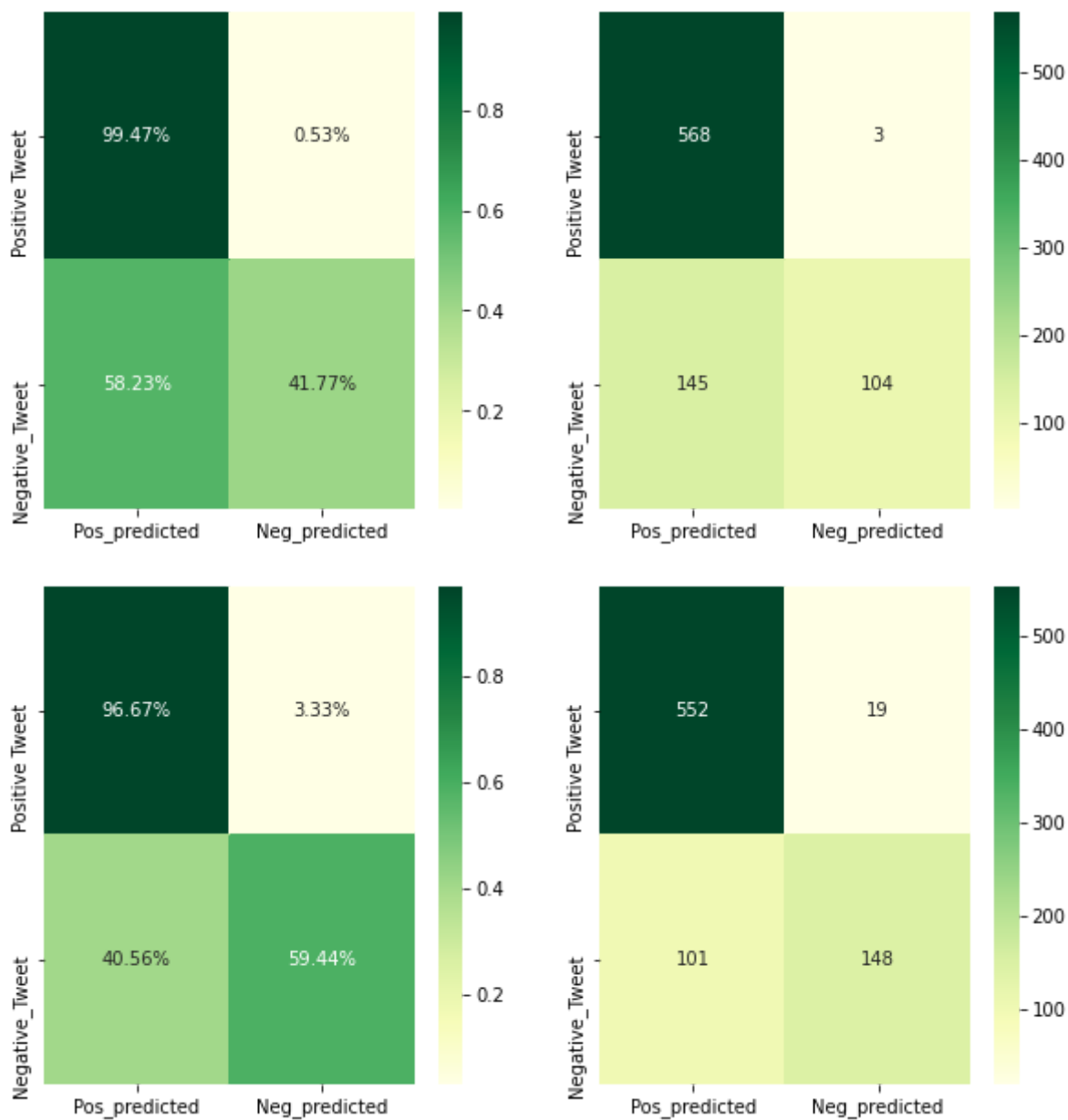
In [6]:
```python
pipes = [NB_pipe,
         rf_pipe,
         balanced_rf_pipe,
         xgb_pipe
        ]
```

```
In [7]:   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)
```
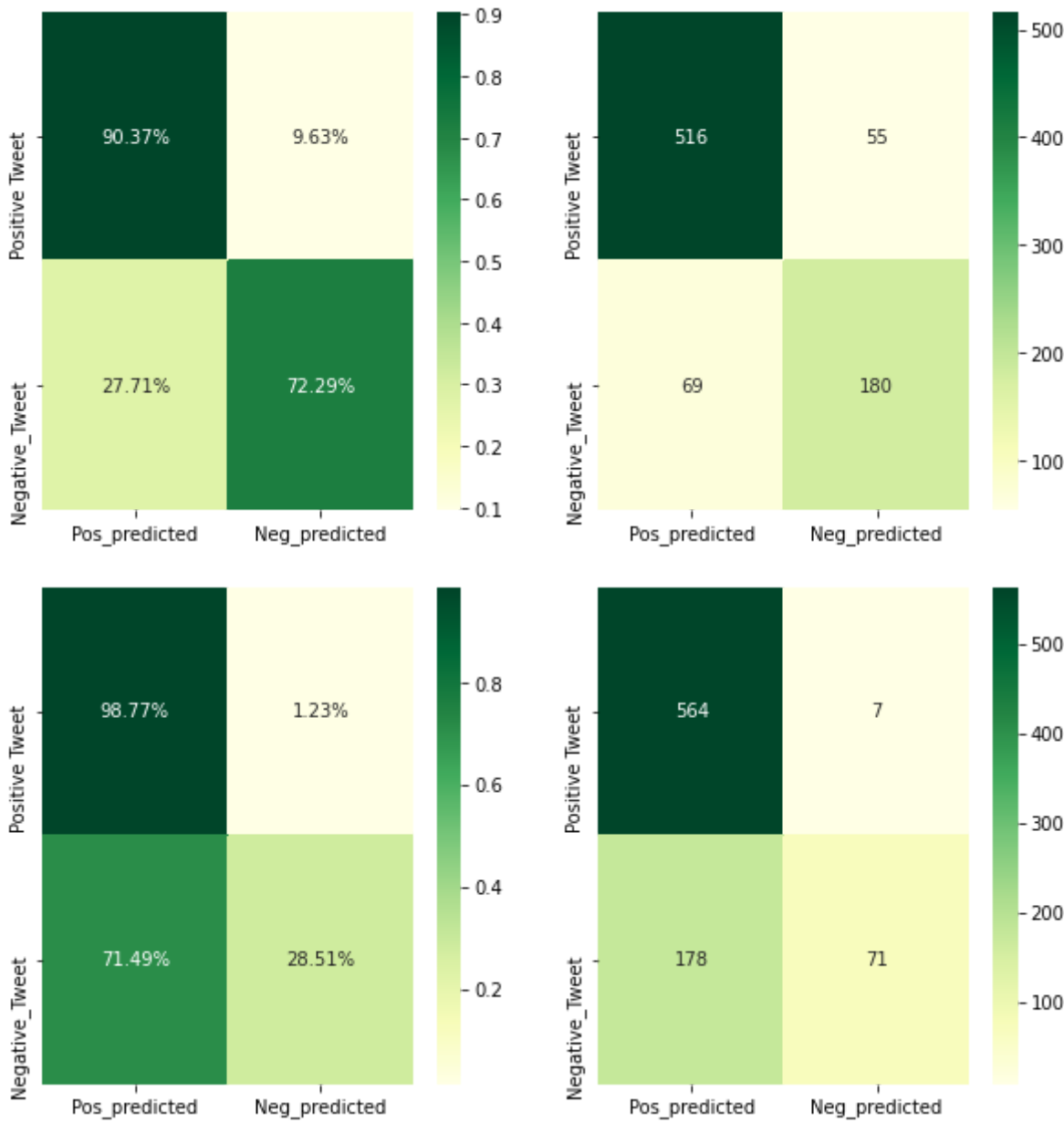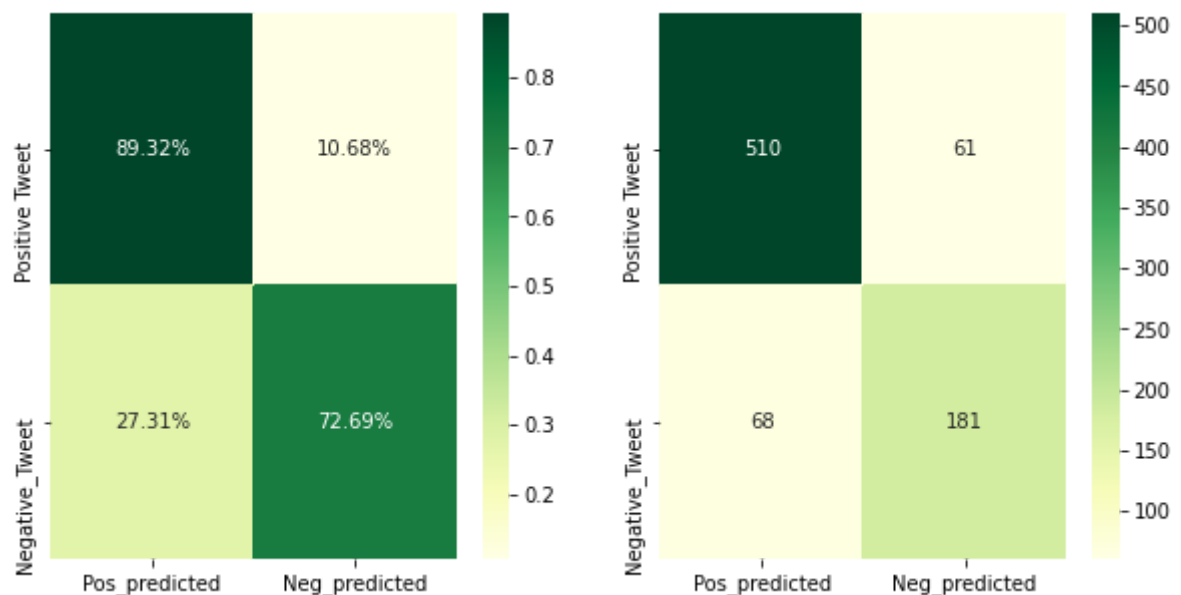
```
In [8]:   1  #initialize grid search variables
          2  n_estimators = [int(x) for x in np.linspace(start = 10, stop = 200, num
          3  criterion = ["gini", "entropy"]
          4  min_samples_split = [8, 10, 12]
          5  max_depth = [int(x) for x in np.linspace(10, 1000, num = 10)]
          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                                        n_iter = 1667,
         20                                        cv = 3,
         21                                        verbose=2,
         22                                        random_state=11,
         23                                        n_jobs = -1
         24                                       )
         25
         26
         27  #fit random grid search and determine best_estimator_
         28  balanced_rfc_rs.fit(tf_idf_X_train, y_train)
         29
         30  #create pipeline for best result from random grid search
         31  balanced_rfc_rs_pipe = make_pipeline(vectorizer,
         32                                        balanced_rfc_rs.best_estimator_)
         33
         34  plot_confusion_matrices(balanced_rfc_rs_pipe, 'images/best_balanced_rf_
```

Fitting 3 folds for each of 1667 candidates, totalling 5001 fits

Now that supervised learning models have been built, trained, and tuned without any pre-training, 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.

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

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

Sckit-learn wrapper provided by Charles Nainan. GitHub of Scikit Learn BERT wrapper (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.

In [9]:

```python
"""
The first model was fitted as seen commeted out below
after some trial and error to determine an appropriate
max_seq_length given my computer's capibilities.

"""


# bert_1 = BertClassifier(do_lower_case=True,
#                         train_batch_size=32,
#                         max_seq_length=50
#                         )



"""
My second model contains 2 hidden layers with 600 neurons.
It only passes over the corpus one time when learning.
It trains fast and gives impressive results.

"""


# bert_2 = BertClassifier(do_lower_case=True,
#                         train_batch_size=32,
#                         max_seq_length=50,
#                         num_mlp_hiddens=500,
#                         num_mlp_layers=2,
#                         epochs=1
#                         )

"""
My third bert model has 600 neurons still but
only one hidden layer. However, the model
passes over the corpus 4 times in total
while learning.

"""

# bert_3 = BertClassifier(do_lower_case=True,
#                         train_batch_size=32,
#                         max_seq_length=50,
#                         num_mlp_hiddens=600,
#                         num_mlp_layers=1,
#                         epochs=4
#                         )

"""
My fourth bert model has 750 neurons and
two hidden layers. The corpus also gets
transversed four times in total while
learning.

"""

# 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 corp us 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)
```

```
  "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,
  "vocab_size": 30522
}


Defaulting to linear classifier/regressor
Building sklearn text classifier...

Predicting: 100%|██████████| 103/103 [01:27<00:00,  1.18it/s]
```
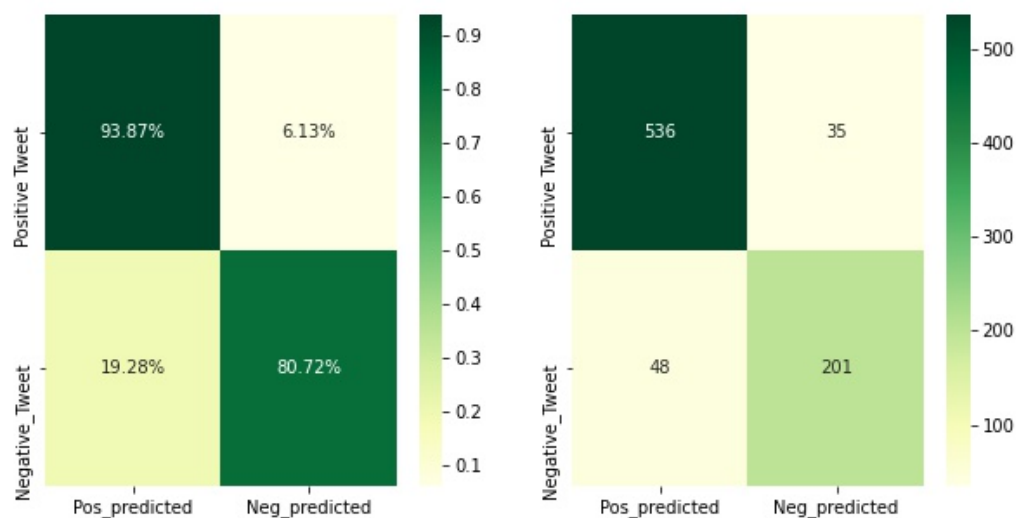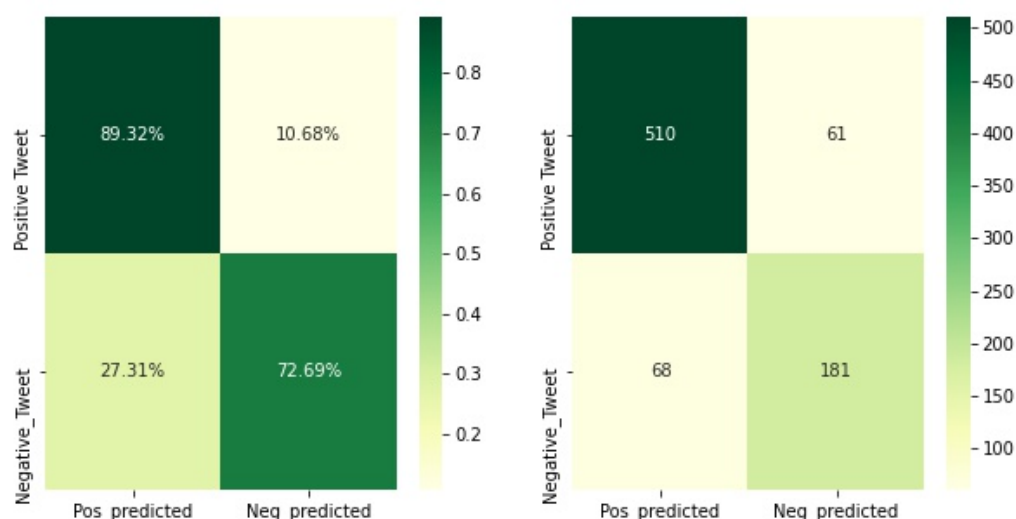
# Evaluation

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 exisiting dataset and retrain the other models.