





```
from TrfidfVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer

#fit training and test data
tf_x_train = tv.fit_transform(x_train)
tf_x_test = tv.transform(x_test)

In (33):
#Using SVM for classification
from sklearn.svm import LinearSVC
svc = LinearSVC(random_state=0)

#fit training data into model
svc.fit(tf_x_train,y_train)

#Predicting test data
y_test_pred = svc.predict(tf_x_test)

In (34):
#analyze the results
from sklearn.metrics import classification_report

report=classification_report(y_test, y_test_pred,output_dict=True)
report

Out(34):
{'0': {'precision': 0.86154929357746479,
      'recall': 0.7806065377704505,
      'f1-score': 0.8115449915110357,
      'support': 305},
 '1': {'precision': 0.9643243243243244,
      'recall': 0.9753963914707491,
      'f1-score': 0.96986287578146236,
      'support': 1829},
 'accuracy': 0.947985004660356,
 'macro avg': {'precision': 0.90293681004948461,
               'recall': 0.879501744238991,
               'f1-score': 0.8908686746628296,
               'support': 2134},
 'weighted avg': {'precision': 0.9467768155578523,
                  'recall': 0.947985004660356,
                  'f1-score': 0.947206195157398,
                  'support': 2134}}

I got a 95% accuracy with SVM. My model did a pretty fair job at predicting both instances, we see that there are more positive reviews
overall and that's probably why my precision was higher than the first instances which shows me that the number of the negative reviews
is test. I am now going to try Logistic Regression modeling to see if it improves.

In (35):
#logReg
from sklearn.linear_model import LogisticRegression

log_reg = LogisticRegression(max_iter=1000,solver='saga')

In (36):
#fitting logistic regression model
log_reg.fit(tf_x_train,y_train)

Out(36):
LogisticRegression
LogisticRegression(max_iter=1000, solver='saga')

In (37):
#predicting the output of data
y_test_pred=log_reg.predict(tf_x_test)

In (38):
#analyzing with classification report
from sklearn.metrics import classification_report

report=classification_report(y_test, y_test_pred,output_dict=True)
report

Out(38):
{'0': {'precision': 0.9074074074074074,
      'recall': 0.6426229508196721,
      'f1-score': 0.752399322456914,
      'support': 305},
 '1': {'precision': 0.943169968717414,
      'recall': 0.98905628758885,
      'f1-score': 0.9655724579663731,
      'support': 1829},
 'accuracy': 0.939551405810584,
 'macro avg': {'precision': 0.9252868880624107,
               'recall': 0.815840058477603,
               'f1-score': 0.858985845106273,
               'support': 2134},
 'weighted avg': {'precision': 0.9380586373211854,
                  'recall': 0.939551405810584,
                  'f1-score': 0.9351048694730221,
                  'support': 2134}}

I got a 93.6% accuracy with LogReg which is better than the SVM. Also it did a better job with looking at the negative class instances and
identifying the precision of it. My model's recall was at 60% for the first instances which shows me that the number of actual positive
instances that were identified which is lower than I would have like.

Now I am going to repeat reviews with another dataframe that is similar to the original just at a later date to see how well my model is.
The dataset is 'hotel_reviews_june19th.csv', which is 1,400 hotels and about 10,000 reviews.

In (39):
df2 = pd.read_csv('hotel_reviews_June19th.csv')

In (40):
df2.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 26 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   id                    10000 non-null    object
 1   dateAdded             10000 non-null    object
 2   dateUpdated           10000 non-null    object
 3   address               10000 non-null    object
 4   categories            10000 non-null    object
 5   primaryCategories     10000 non-null    object
 6   city                 10000 non-null    object
 7   country              10000 non-null    object
 8   keys                 10000 non-null    object
 9   latitude              10000 non-null    float64
10  longitude             10000 non-null    float64
11  name                 10000 non-null    object
12  postalCode            10000 non-null    object
13  province              10000 non-null    object
14  reviews.date         10000 non-null    object
15  reviews.dateAdded    0 non-null       float64
16  reviews.dateSeen     10000 non-null    object
17  reviews.rating       10000 non-null    object
18  reviews.sourceURLs   10000 non-null    object
19  reviews.text         10000 non-null    object
20  reviews.title        9998 non-null     object
21  reviews.userCltly    10000 non-null    object
22  reviews.userPovince  9998 non-null     object
23  reviews.username     10000 non-null    object
24  sourceURLs           10000 non-null    object
25  websites             10000 non-null    object
dtypes: float64(3), int64(1), object(22)
memory usage: 2.0+ MB

In (41):
#dropping unnecessary columns because these are not essential in looking at the data.
df2 = df2.drop(['id', 'reviews.dateSeen','reviews.dateAdded', 'dateAdded', 'dateUpdated', 'reviews.sourceURLs',

In (42):
for i in range(0,len(df2)-1):
    if type(df2.iloc[i]['reviews.text']) !=str:
        df2.iloc[i]['reviews.text'] =str(df2.iloc[i])
    ['reviews.text']

Out(42):
['reviews.text']

In (43):
#Change ratings column to a int .
df2['reviews.rating'] = df2['reviews.rating'].astype(int)

#got rid of rating of 3
df2 = df2[df2['reviews.rating']!=3]

#seperate positive (4 or 5) and negative (1,2) and no middle
def sentiment (x):
    return 1 if x >= 4 else 0

df2['sentiment'] = df2['reviews.rating'].apply(sentiment)
df2.head()
```

	city	country	latitude	longitude	name	postalCode	province	reviews.date	reviews.rating	reviews.text	reviews.title	sentime
1	Carmel by the Sea	US	36.557220	-121.921940	Best Western Carmel's Town House Lodge	93921	CA	2016-04-02T00:00:00Z	4	We stayed in the king suite with the separate...	Clean rooms at solid rates in the heart of Carmel	
3	Carmel by the Sea	US	36.557220	-121.921940	Best Western Carmel's Town House Lodge	93921	CA	2016-08-22T00:00:00Z	5	Not cheap but excellent location. Price is som...	Very good	
4	Carmel by the Sea	US	36.557220	-121.921940	Best Western Carmel's Town House Lodge	93921	CA	2016-03-21T00:00:00Z	2	If you get the room that they advertised on th...	Low chance to come back here	
5	Lexington	US	38.047014	-84.497742	21c Museum Hotel Lexington	40507	KY	2016-03-15T00:00:00Z	4	This is such a fun, lovely hotel. The attentio...	Loved staying here	
6	Lexington	US	38.047014	-84.497742	21c Museum Hotel Lexington	40507	KY	2016-04-18T00:00:00Z	1	We recently stayed at this hotel on a trip to ...	Does not live up to its reputation	

```
df2['reviews'] = df2.apply(combined_features,axis=1 )
df2.head()
```

	city	country	latitude	longitude	name	postalCode	province	reviews.date	reviews.rating	reviews.text	reviews.title	sentime
1	Carmel by the Sea	US	36.557220	-121.921940	Best Western Carmel's Town House Lodge	93921	CA	2016-04-02T00:00:00Z	4	We stayed in the king suite with the separate...	Clean rooms at solid rates in the heart of Carmel	
3	Carmel by the Sea	US	36.557220	-121.921940	Best Western Carmel's Town House Lodge	93921	CA	2016-08-22T00:00:00Z	5	Not cheap but excellent location. Price is som...	Very good	
4	Carmel by the Sea	US	36.557220	-121.921940	Best Western Carmel's Town House Lodge	93921	CA	2016-03-21T00:00:00Z	2	If you get the room that they advertised on th...	Low chance to come back here	
5	Lexington	US	38.047014	-84.497742	21c Museum Hotel Lexington	40507	KY	2016-03-15T00:00:00Z	4	This is such a fun, lovely hotel. The attentio...	Loved staying here	
6	Lexington	US	38.047014	-84.497742	21c Museum Hotel Lexington	40507	KY	2016-04-18T00:00:00Z	1	We recently stayed at this hotel on a trip to ...	Does not live up to its reputation	

```
In (45):
#converting all letters to lower case
df2['reviews'] = df2['reviews'].apply(lambda x : " ".join(x.lower() for x in str(x).split()))

#remove non alpha characters
df2['reviews'] = df2['reviews'].apply(lambda x : " ".join([re.sub('[^A-Za-z]+',''),x] for x in nltk.word_tokenize(x)))

#remove extra spaces
df2['reviews'] = df2['reviews'].apply(lambda x : re.sub(' +', ' ',x))

#lemmatizer (takes away suffixes to bring out root word)
df2['reviews'] = df2['reviews'].apply(lambda x : " ".join([lemmatizer.lemmatize(w) for w in nltk.word_tokenize(x)]))

In (46):
#remove stop words
from nltk.corpus import stopwords
stop = stopwords.words('english')
df2['reviews'] = df2['reviews'].apply(lambda x : " ".join([x for x in x.split() if x not in stop]))

In (47):
#Using porter stemmer to stem words
df2['reviews'] = df2['reviews'].apply(lambda x : stem_words(x))

#tokenize
df2['reviews'] = df2['reviews'].apply(lambda x : tokenization(x))

df2['reviews'] = df2['reviews'].astype(str)

In (48):
df2['reviews']

Out(48):
1      ['clean room solid rate heart carmel stay king...
3      ['good cheap excel locat price somewhat standa...
4      ['low chanc come back get room advertis websit...
5      ['love stay fun love hotel attent detail impro...
6      ['doe live reput recent stay hotel trip lexing...

9995   ['accommod friendli staff friend took trip ham...
9996   ['comfort friendli clean profession check depar...
9997   ['great locat hampton locat quiet street across...
9998   ['great atmospher awesom wing favorit wa garli...
9999   ['health care appoint clean facil freeway staf...
Name: reviews, length: 8618, dtype: object

In (49):
x2 = df2['reviews']
y2 = df2['sentiment']

x2_train,x2_test,y2_train, y2_test = train_test_split(x2, y2, test_size = 0.25, random_state = 30)

print("Train: "x2_train.shape, y2_train.shape,Test: ", (x2_test.shape,y2_test.shape))

Train: (6463,) (6463,) Test: ((2155), (2155,))

In (50):
#fit training and test data
tf_x2_train = tv.fit_transform(x2_train)
tf_x2_test = tv.transform(x2_test)

In (51):
#fit training data into model
svc.fit(tf_x2_train,y2_train)

#Predicting test data
y2_test_pred = svc.predict(tf_x2_test)

In (52):
report2=classification_report(y2_test, y2_test_pred,output_dict=True)
report2

Out(52):
{'0': {'precision': 0.9095026248434389,
      'recall': 0.75563930977443609,
      'f1-score': 0.8254620232303285,
      'support': 266},
 '1': {'precision': 0.9663908996897621,
      'recall': 0.9894153875066173,
      'f1-score': 0.9777661522364635,
      'support': 1889},
 'accuracy': 0.9605568454575638,
 'macro avg': {'precision': 0.9379465810666006,
               'recall': 0.8792574262346893,
               'f1-score': 0.90161408278396,
               'support': 2155},
 'weighted avg': {'precision': 0.9593689147674781,
                  'recall': 0.9605568454575638,
                  'f1-score': 0.958966621122445,
                  'support': 2155}}

My accuracy was 96.15% which is better than what I did on my first data frame but also i see that my recall increased significantly here from
the previous 60%.

In (53):
df2['reviews']

Out(53):
1      ['clean room solid rate heart carmel stay king...
3      ['good cheap excel locat price somewhat standa...
4      ['low chanc come back get room advertis websit...
5      ['love stay fun love hotel attent detail impro...
6      ['doe live reput recent stay hotel trip lexing...

9995   ['accommod friendli staff friend took trip ham...
9996   ['comfort friendli clean profession check depar...
9997   ['great locat hampton locat quiet street across...
9998   ['great atmospher awesom wing favorit wa garli...
9999   ['health care appoint clean facil freeway staf...
Name: reviews, length: 8618, dtype: object

In (54):
#Example Review 1

print("Example Review: ", df2['reviews'][1])

sentiment = "Positive" if y2_test.iloc[1] == 1 else "Negative"
predicted_sentiment = "Positive" if y2_test_pred[1] == 1 else "Negative"

# Print the actual and predicted sentiments
print("Actual Sentiment:", sentiment)
print("Predicted Sentiment:", predicted_sentiment)

Example Review: ['clean room solid rate heart carmel stay king suit appar bedroom live space sofa bed and nt goo
d back discomfort day left three night stay room clean king bed comfort hotel locat within walk distanc w
ant']
Actual Sentiment: Positive
Predicted Sentiment: Positive

In (55):
#Example Review 2

print("Example Review: ", df2['reviews'][50])

sentiment = "Positive" if y2_test.iloc[50] == 1 else "Negative"
predicted_sentiment = "Positive" if y2_test_pred[50] == 1 else "Negative"

# Print the actual and predicted sentiments
print("Actual Sentiment:", sentiment)
print("Predicted Sentiment:", predicted_sentiment)

Example Review: ['disappoint area wa busi freeway describ stripper neighborhood need earli checkin arrang month
ago nt get abl check notic wet section carpet thought maybe clean someth next morn room satur went front desk sa
id probabl as problem said problem fix left hour later return problem fix water travel went front desk said som
eth like ca nt fix day expect odd expect fix least suck water ca nt hmmm get new room said fix time limit wa
soon said send somon look twenti minut bar nice bar staff went room mainten fix minut place towel floor stay l
eft thank feedback apoloy unabl fix air condit saw immedi sever interact earli day offer room move sever time
could get resourc fix leak air condit sorri staff spoke empathet offer altern accomod tri honor earli check re
quest depend deparur pattern quickli clean room new arriv hope doe prohibit join u']
Actual Sentiment: Positive
Predicted Sentiment: Positive

In (57):
# I made this code below because I was thinking my model and prediction model had only positive values.
# this is proof that there is negative and positive values.

import numpy as np

# Check if any predicted sentiment is negative (0)
if any(sentiment == 0 for sentiment in y2_test_pred):
    print("At least one predicted sentiment is negative (0.)")

# Filter negative values using list comprehension
negative_values = [sentiment for sentiment in y2_test_pred if sentiment == 0]

# Calculate the sum of the negative values using NumPy
sum_of_negatives = np.count_nonzero(y2_test_pred == 0)

# Print the sum of the negative values
print("Sum of negative values:", sum_of_negatives)
else:
    print("All predicted sentiments are positive (1).")

At least one predicted sentiment is negative (0).
Sum of negative values: 221
```