```
In [1]:    1  import pandas as pd
           2  import numpy as np
           3  from sklearn.model_selection import train_test_split
           4  import re
           5  import contractions
           6  from unidecode import unidecode
           7  from nltk.tokenize import word_tokenize
           8  from nltk.corpus import stopwords
           9  from string import punctuation
          10  from rake_nltk import Rake
          11  import yake
          12  from nltk.stem import WordNetLemmatizer,LancasterStemmer
          13  from gensim.models import Word2Vec
          14  from sklearn.feature_extraction.text import CountVectorizer,TfidfVectorizer
          15  from nltk.util import ngrams
          16  from wordcloud import WordCloud
          17  import matplotlib.pyplot as plt
          18  from sklearn.linear_model import LogisticRegression
          19  from sklearn.neighbors import KNeighborsClassifier
          20  from sklearn.tree import DecisionTreeClassifier
          21  from sklearn.ensemble import RandomForestClassifier
          22  from sklearn.ensemble import AdaBoostClassifier
          23  from sklearn.naive_bayes import MultinomialNB
          24  import yake
          25  from rake_nltk import Rake
```

```
In [112]:   1  import warnings
            2  warnings.filterwarnings("ignore")
```

```
In [113]:   1  df=pd.read_csv(r"C:\Users\kannu\OneDrive\Desktop\Pandas_csvs\Datasets\sentim
```

```
In [114]:   1  df.shape
```

Out[114]: (515738, 17)

```
In [118]:   1  pos=df[["Positive_Review"]]
            2  pos.columns=["Review"]
            3  pos["Target"]=np.ones(len(pos),dtype=int)
```

```
In [119]:   1  for ind,data in enumerate (pos["Review"]):
            2      if ((len(data)<14) and ((data.lower().__contains__("nothing")) or (data.
            3          pos["Target"].iloc[ind]=2
```

```
In [120]:   1  a=" nothing"
            2  b=" Nothing"
            3  c="nothing"
            4  a.lower().__contains__("nothing")
            5  b.lower().__contains__("nothing")
            6  c.lower().__contains__("nothing")
```

Out[120]: True

```python
In [122]:    1  neg=df[["Negative_Review"]]
             2  neg.columns=["Review"]
             3  neg["Target"]=np.zeros(len(neg),dtype=int)
```

```python
In [123]:    1  # [word for word in neg["Review"] if len(word)<14]
```

```python
In [124]:    1  for ind,data in enumerate (neg["Review"]):
             2      if ((len(data)<14) and ((data.lower().__contains__("nothing")) or (data.
             3          neg["Target"].iloc[ind]=2
```

```python
In [125]:    1  neg["Target"].value_counts()
```

```
Out[125]:  0    364220
           2    151518
           Name: Target, dtype: int64
```

```python
In [127]:    1  pos_df=pos.sample(500)
```

```python
In [128]:    1  final_df=pd.concat([pos_df,neg_df],axis=0)
```

```python
In [130]:    1  final_df.reset_index(drop=True,inplace=True)
```

```python
In [131]:    1  final_df["Target"].value_counts()
             2  x=final_df.drop("Target",axis=1)
             3  y=final_df["Target"]
```

```python
In [132]:    1  x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.2,random_state
             2  x_train.shape,x_test.shape,y_train.shape
```

```
Out[132]:  ((800, 1), (200, 1), (800,))
```

```python
In [133]:    1  #Removing white spaces
             2
             3  def whitespace_removal(data):
             4      import re
             5      pattern=r"\s+"
             6      res=re.sub(pattern," ",data)
             7      return res
             8
             9  #contraction mapping
            10
            11  def expand_text(data):
            12      res=contractions.fix(data)
            13      return res
            14
            15  #accented character handling
            16
            17  def accent_handling(data):
            18      res=unidecode(data)
            19      return res
```

In [135]:
```python
x_train["Review"]=x_train["Review"].apply(whitespace_removal)
x_train["Review"]=x_train["Review"].apply(expand_text)
x_train["Review"]=x_train["Review"].apply(accent_handling)
```

In [136]:
```python
#datacleaning

stopword_list=stopwords.words("english")
stopword_list.remove("no")
stopword_list.remove("nor")
stopword_list.remove("not")


def cleaning(data):
    res=[word.lower() for word in word_tokenize(data) if (word.lower() not i
    return res
```

In [137]:
```python
x_train_cleaned=x_train["Review"].apply(cleaning)
x_train_cleaned.reset_index(drop=True,inplace=True)
```

In [138]:
```python
def ngram(data,n_range):
    ngram_list=[]
    op=ngrams(data,n_range)
    for ng in op:
        ngram_list.append(ng)

    return ngram_list
```

In [139]:
```python
y_train.reset_index(drop=True,inplace=True)
```

In [140]:
```python
#checking positive reviews

merged_data_for_ngrams=pd.concat([x_train_cleaned,y_train],axis=1)
pos=merged_data_for_ngrams.loc[merged_data_for_ngrams["Target"]==1]["Review"
```

In [141]:
```python
y_train.index
```

Out[141]: RangeIndex(start=0, stop=800, step=1)

In [144]:
```python
#checking negative reviews

neg=merged_data_for_ngrams.loc[merged_data_for_ngrams["Target"]==0]["Review"
neg.apply(lambda x: ngram(x,4)).iloc[7]
```

Out[144]: []

```
In [145]:  1  #checking neutral reviews
           2
           3  neutral=merged_data_for_ngrams.loc[merged_data_for_ngrams["Target"]==2]["Rev
           4  neutral.apply(lambda x: ngram(x,1)).iloc[17]
```
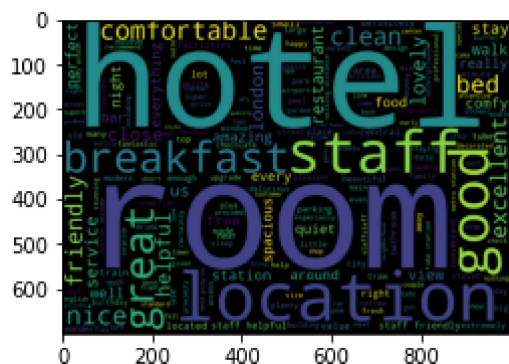
Out[145]:  [('no',), ('positive',)]

```
In [146]:  1  def text(data):
           2      return " ".join(data)
           3
           4  pos_string=pos.apply(text)
           5  neg_string=neg.apply(text)
           6  neutral_string=neutral.apply(text)
```
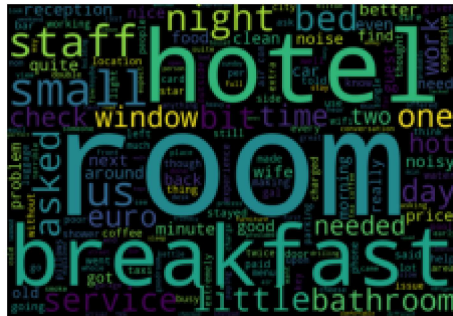
```
In [147]:  1  def df_2_string(df):
           2      final_text=[]
           3      list1=df.to_list()
           4      for sub_list in list1:
           5          val="".join(sub_list)
           6
           7          final_text.append(val)
           8      return "".join(final_text)
           9
```

```
In [148]:  1  pos_text=df_2_string(pos_string)
           2  neg_text=df_2_string(neg_string)
           3  neutral_text=df_2_string(neutral_string)
           4
```

```
In [217]:  1  #wordcloud of pos reviews
           2  word_cloud_pos=WordCloud(height=700,width=1000,background_color="black",min_
           3  plt.figure(figsize=(4,6))
           4  plt.imshow(word_cloud_pos)
           5  plt.show()
```

In [218]:
```python
#wordcloud of neg reviews
word_cloud_neg=WordCloud(height=700,width=1000,background_color="black",min_
plt.figure(figsize=(4,6))
plt.imshow(word_cloud_neg)
plt.axis("off")
plt.show()
```



In [152]:
```python
#lemmatization and stemming

def lemmatization(data):
    final_text=[]
    for word in data:
        lemmatizer=WordNetLemmatizer()
        lema=lemmatizer.lemmatize(word)
        final_text.append(lema)
    return " ".join(final_text)

def stemming(data):
    final_text=[]
    for word in data:
        stemmer=LancasterStemmer()
        stem=stemmer.stem(word)
        final_text.append(stem)
    return " ".join(final_text)
```

In [153]:
```python
lemmatized=x_train_cleaned.apply(lemmatization)
stemmed=x_train_cleaned.apply(stemming)
```

In [158]:
```python
#countvectorizer
count_vect=CountVectorizer()
bow=count_vect.fit_transform(lemmatized).A
```

In [160]:
```python
#tfidf
tfidf_vect=TfidfVectorizer()
tfidf=tfidf_vect.fit_transform(lemmatized).A
```

In [162]:
```python
#word to vect

#making in [[w1,w2,....],[w3,w4...]] format

def word_to_vec_ip(df):
    res=[word.split() for word in df]
    return res
```

In [163]:
```python
word_vec_data=word_to_vec_ip(lemmatized)
# word_vec_data
```

In [164]:
```python
word2vec=Word2Vec(word_vec_data,window=10,min_count=2)
word2vec.build_vocab(word_vec_data)
word2vec.train(word_vec_data,total_examples=word2vec.corpus_count,epochs=5)
```

Out[164]: (25212, 39070)

In [166]:
```python
def word2_vector(model,data):
    feature=[]
    zero_vector=np.zeros(model.vector_size) #a 1d array of 0.

    for vector in data:
        vectors=[]   #all vectors of a single doc will be here,[[vetor w1],[v
        for word in vector:

            if word in model.wv:
                try:
                    vectors.append(model.wv[word])
                except KeyError:
                    continue
        if vectors:
            vectors=np.array(vectors)    #it will have 2d array, means row an
            avg_vec=vectors.mean(axis=0)   # by axis=0 it will find mean roww
            feature.append(avg_vec)
        else:
            feature.append(zero_vector)


    return feature

feature=word2_vector(word2vec,word_vec_data)
```

In [167]:
```python
#modelling
```

## logistic_regression

In [168]:
```python
#logistic regression
#with_count_vectorizer
log_reg_count=LogisticRegression()
log_reg_count.fit(bow,y_train)
log_reg_count.score(bow,y_train)
```

Out[168]: 0.98125

In [169]:
```python
#with_tfidf
log_reg_tfidf=LogisticRegression()
log_reg_tfidf.fit(tfidf,y_train)
log_reg_tfidf.score(tfidf,y_train)
```

Out[169]: 0.975

In [170]:
```python
#with word2vec
log_reg_word2vec=LogisticRegression()
log_reg_word2vec.fit(feature,y_train)
log_reg_word2vec.score(feature,y_train)
```

Out[170]: 0.6475

## KNN

In [171]:
```python
#with_count_vectorizer
knn_count=KNeighborsClassifier()
knn_count.fit(bow,y_train)
knn_count.score(bow,y_train)
```

Out[171]: 0.78375

In [172]:
```python
#with_tfidf
knn_tfidf=KNeighborsClassifier()
knn_tfidf.fit(tfidf,y_train)
knn_tfidf.score(tfidf,y_train)
```

Out[172]: 0.35375

In [173]:
```python
#with_word2vec
knn_word2vec=KNeighborsClassifier()
knn_word2vec.fit(feature,y_train)
knn_word2vec.score(feature,y_train)
```

Out[173]: 0.80625

## DT

In [174]:
```python
#with count_vectorizer
dt_count=DecisionTreeClassifier()
dt_count.fit(bow,y_train)
dt_count.score(bow,y_train)
```

Out[174]: 0.99625

In [175]:
```python
#with tfidf
dt_tfidf=DecisionTreeClassifier()
dt_tfidf.fit(tfidf,y_train)
dt_tfidf.score(tfidf,y_train)
```

Out[175]: 0.99625

In [176]:
```python
# with word to vec
dt_word2vec=DecisionTreeClassifier()
dt_word2vec.fit(feature,y_train)
dt_word2vec.score(feature,y_train)
```

Out[176]: 0.98875

## Random Forest

In [177]:
```python
# with count vectorizer
rf_count=RandomForestClassifier()
rf_count.fit(bow,y_train)
rf_count.score(bow,y_train)
```

Out[177]: 0.99625

In [178]:
```python
#with tfidf
rf_tfidf=RandomForestClassifier()
rf_tfidf.fit(tfidf,y_train)
rf_tfidf.score(tfidf,y_train)
```

Out[178]: 0.99625

In [179]:
```python
#word2vec
rf_word2vec=RandomForestClassifier()
rf_word2vec.fit(feature,y_train)
rf_word2vec.score(feature,y_train)
```

Out[179]: 0.98875

## Adaboost

In [180]:
```python
#count_vectorizer
ada_count=AdaBoostClassifier()
ada_count.fit(bow,y_train)
ada_count.score(bow,y_train)
```

Out[180]: 0.68375

In [181]:
```python
#tfidf
ada_tfidf=AdaBoostClassifier()
ada_tfidf.fit(tfidf,y_train)
ada_tfidf.score(tfidf,y_train)
```

Out[181]: 0.8475

In [182]:
```python
#word2vec
ada_word2vec=AdaBoostClassifier()
ada_word2vec.fit(feature,y_train)
ada_word2vec.score(feature,y_train)
```

Out[182]: 0.5825

## Naive_bayes

In [183]:
```python
#count_vect
naive_count=MultinomialNB()
naive_count.fit(bow,y_train)
naive_count.score(bow,y_train)
```

Out[183]: 0.935

In [184]:
```python
naive_tfidf=MultinomialNB()
naive_tfidf.fit(tfidf,y_train)
naive_tfidf.score(tfidf,y_train)
```

Out[184]: 0.94875