

Trip Advisor Hotel Reviews - Case Study

Explore Hotel aspects and Predict the rating of each review.

Hotels play a crucial role in traveling and with the increased access to information new pathways of selecting the best ones emerged. With this dataset, consisting of 20k reviews crawled from Tripadvisor, you can explore what makes a great hotel and maybe even use this model in your travels!

This Case Study is taken from : <https://www.kaggle.com/andrewmvd/trip-advisor-hotel-reviews>

Acknowledgements

Citation Alam, M. H., Ryu, W.-J., Lee, S., 2016. Joint multi-grain topic sentiment: modeling semantic aspects for online reviews. Information Sciences 339, 206–223. DOI https://zenodo.org/record/1219899#.X64_imgzblU

Dataset link: <https://www.kaggle.com/andrewmvd/trip-advisor-hotel-reviews>

In []:

In []:

Summary of my Idea and what I did.

The problem given is to predict whether the "review" given is either good or bad by using "rating" as class label. How good the model predicts is to be found based on the metrics given i.e MAE and RMSE.

We can generally understand that if the rating is 3 and above, it is good or else it is bad.

In this task, we have to predict the rating based on the review which should be automatically done by our model and that rating which in turn decides the review is either good or bad. As, we have 1-5 rating range for every review, we have to predict the rating in that range only by classification.

This means, the problem is a "Multi-Class Classification" problem.

We can use classification ML techniques like SVM, Decision Trees, Random Forest, Boosting techniques etc., and Deep Learning models too using certain type of loss and activation function.

Finally based on the metrics we have to decide which model to use. The metric we use is F1-Score and it's given to use RMSE and MAE.

I've covered 3 models which are Linear SVM, LGBMClassifier and a Deep Learning Model and analysed results.

In [1]:

```
# Importing the packages

%matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
```

```
from gensim.models import Word2Vec
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os
```

In [2]:

```
import matplotlib.pyplot as plt
from wordcloud import WordCloud
from plotly import graph_objs as go
import plotly.figure_factory as ff
import seaborn as sns
from bs4 import BeautifulSoup
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV
from sklearn.model_selection import cross_val_score, GridSearchCV
from sklearn.svm import SVC
from sklearn.preprocessing import StandardScaler
```

Reading the file

In [111]:

```
df = pd.read_csv("tripadvisor_hotel_reviews.csv")
df.head()
```

Out[111]:

	Review	Rating
0	nice hotel expensive parking got good deal sta...	4
1	ok nothing special charge diamond member hilt...	2
2	nice rooms not 4* experience hotel monaco seat...	3
3	unique, great stay, wonderful time hotel monac...	5
4	great stay great stay, went seahawk game aweso...	5

Exploratory Data Analysis

Shape of the file

In [4]:

```
df.shape
```

Out[4]:

(20491, 2)

There are 20491 rows and 2 columns.

In []:

Checking the type of columns

```
In [ ]:
```

```
df.dtypes
```

```
Out[ ]:
```

```
Review    object
Rating    int64
dtype: object
```

Checking for null values

```
In [ ]:
```

```
df[df.isnull().any(axis=1)]
```

```
Out[ ]:
```

Review Rating

```
In [ ]:
```

```
df.isnull().sum()
```

```
Out[ ]:
```

```
Review    0
Rating    0
dtype: int64
```

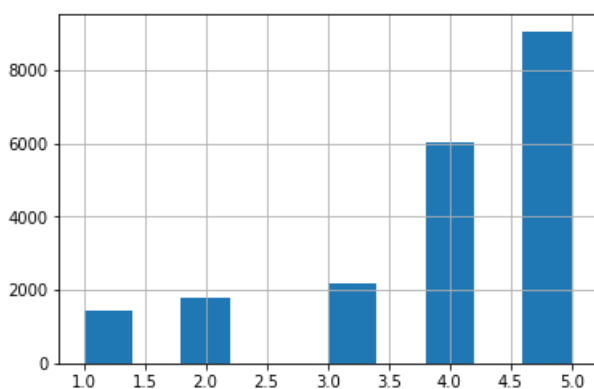
Histogram of frequency of data points according to the rating

```
In [ ]:
```

```
df['Rating'].hist()
```

```
Out[ ]:
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1fea304e668>
```



```
In [ ]:
```

Getting the count of data points of individual rating

```
In [ ]:
```

```
df['Rating'].value_counts()
```

Out[]:

```
5    9054
4    6039
3    2184
2    1793
1    1421
Name: Rating, dtype: int64
```

In []:

The distribution of the above values

In [5]:

```
class_dist = df['Rating'].value_counts()

def ditribution_plot(x,y,name):
    fig = go.Figure([
        go.Bar(x=x, y=y)
    ])

    fig.update_layout(title_text=name)
    fig.show()
```

In [6]:

```
ditribution_plot(x= class_dist.index, y= class_dist.values, name= 'Class Distribution')
```

In []:

Top words in the Reviews corpus

In [104]:

```
from wordcloud import WordCloud
```

In [105]:

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

In []:

In [112]:

```
vectorizer = TfidfVectorizer()  
vectorizer.fit_transform(df['Review'])
```

Out[112]:

```
<20491x52923 sparse matrix of type '<class 'numpy.float64'>'  
  with 1705964 stored elements in Compressed Sparse Row format>
```

In [113]:

```
word_list = vectorizer.get_feature_names()
```

In [116]:

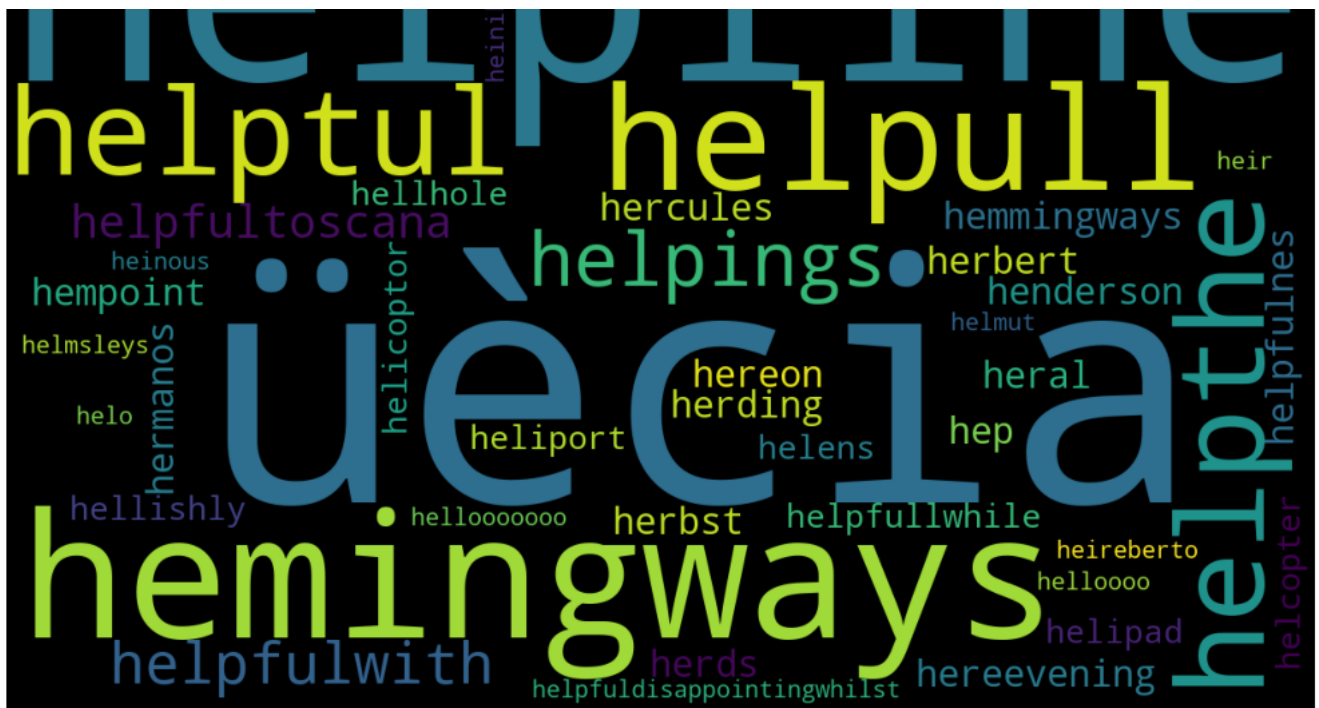
```
words2=[]  
idf2=vectorizer.idf_  
features=np.argsort(idf2)[::-1]  
for i in features[0:60]:  
    words2.append(word_list[i])  
print(words2)
```

```
['üecia', 'helpline', 'hemming', 'hemingways', 'helpull', 'helptul', 'helpthe', 'helpsthe',  
'helpings', 'hemp', 'helping_çî_', 'helpfulwith', 'helpfulwhen', 'helpfultoscana', 'helpfulthe',  
'helpfulrestaurants', 'hemmingways', 'hempel', 'hermanos', 'herculean', 'hereon', 'herefordshire',  
'hereevening', 'herds', 'herding', 'hercules', 'herbst', 'hempoint', 'herbivores', 'herbert',  
'heralded', 'heral', 'hep', 'henderson', 'helpfulone', 'helpfulnes', 'helpfullwhile', 'helens', 'h  
ellishly', 'hellhole', 'heliport', 'helipad', 'helicopter', 'helpfull', 'helicopter',  
'helpfuldisappointingwhilst', 'heist', 'heireberto', 'heir', 'heinz', 'heinous', 'heiniken',  
'hellllllllloooo', 'helloooo', 'helloooooooo', 'hellostayed', 'helmsleys', 'helmsly', 'helmut', 'hel  
o']
```

In [117]:

```
from wordcloud import WordCloud  
from wordcloud import WordCloud  
wordcloud = WordCloud(width = 1200, height = 1000).generate(" ".join(words2))  
plt.figure(figsize = (20, 15))  
plt.imshow(wordcloud)  
plt.axis("off")  
plt.show()
```





In []:

In []:

Data Preprocessing

In [4]:

```
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"\ 're", " are", phrase)
    phrase = re.sub(r"\ 's", " is", phrase)
    phrase = re.sub(r"\ 'd", " would", phrase)
    phrase = re.sub(r"\ 'll", " will", phrase)
    phrase = re.sub(r"\ 't", " not", phrase)
    phrase = re.sub(r"\ 've", " have", phrase)
    phrase = re.sub(r"\ 'm", " am", phrase)
    return phrase
```

In [5]:

```
import nltk
nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] /home/cloudportal456/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

Out [5] :

True

In [6]:

```
from nltk.corpus import stopwords
stopwords = set(stopwords.words('english'))
```

In [7]:

```
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(df['Review'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\S*\d\S*", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', ' ', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
    preprocessed_reviews.append(sentence.strip())
```

100%|██████████| 20491/20491 [00:08<00:00, 2393.44it/s]

In []:

In [8]:

```
preprocessed_reviews[232]
```

Out[8]:

'friendly staff comfortable beds stayed hotel tour included hotel coast california enjoyed hotel g
reat location tourists close market place attractions desk staff accommodating thanks christopher
rest hotel staff friendly loved beds comfortable breakfast included tour enjoyed want included buf
fet fruit fresh hi morgan wait staff hotel bright clean looking fresh flowers lobby area seattle tr
y make short trip boeing aircraft factory loved tour people seattle portland willing help tourists
directions suggestions'

In []:

Perform Modelling

In [9]:

```
le = LabelEncoder()
Y = le.fit_transform(df[r'Rating'])
```

In [10]:

```
X = np.array(preprocessed_reviews)
```

Performing Train_Test Split

In [11]:

```
X_train,X_cv,y_train,y_cv = train_test_split(X,Y,test_size=0.2)
X_train,X_test,y_train,y_test = train_test_split(X_train,y_train,test_size=0.2)
```

In [12]:

```
print(X_train.shape)
print(X_cv.shape)
```

```
print(X_test.shape)
```

```
(13113,)
(4099,)
(3279,)
```

Using Count Vectorizer with Bi-Grams

In [13]:

```
count_vect = CountVectorizer(lowercase=True, ngram_range=(1, 2))
X_train=count_vect.fit_transform(X_train)
X_cv=count_vect.transform(X_cv)
X_test=count_vect.transform(X_test)

scalar = StandardScaler(with_mean = False)
X_train_bow = scalar.fit_transform(X_train)
X_test_bow= scalar.transform(X_test)
X_cv_bow=scalar.transform(X_cv)
```

In []:

In [51]:

```
'''alphas = [10**-3,10**-2, 10**-1,1,10,10**2,10**3,10**4,10**6]
scores = []
for i in alphas:
    model = SGDClassifier(alpha = i,loss="hinge")
    model.fit(X_train_bow, y_train)
    scores.append(model.score(X_cv, y_cv))'''
```

In [52]:

```
#optimal_alpha= alphas[scores.index(max(scores))]
#optimal_alpha
```

Out[52]:

0.01

Fine tuning Linear SVM Model

In [58]:

```
#https://stackoverflow.com/questions/55893734/how-can-i-use-sgdclassifier-hinge-loss-with-gridsearchcv-using-log-loss-metric
```

```
%%time
grid_params = {'base_estimator__alpha': [ 10**-3,10**-2,10**-1,1,10,10**2,10**3]}
clf = SGDClassifier(loss='hinge')
calibrated_clf = CalibratedClassifierCV(base_estimator=clf, method='sigmoid', cv=3)
svm_model = GridSearchCV(calibrated_clf, param_grid=grid_params, cv=5)
svm_model.fit(X_train_bow, y_train)
```

Out[58]:

```
GridSearchCV(cv=5,
             estimator=CalibratedClassifierCV(base_estimator=SGDClassifier(),
                                              cv=3),
             param_grid={'base_estimator__alpha': [0.001, 0.01, 0.1, 1, 10, 100,
                                                    1000]}))
```

In []:

In [60]:

```
print(svm_model.best_params_)

{'base_estimator__alpha': 1}
```

Predicting the results after Fine Tuning

In [61]:

```
clf = SGDClassifier(alpha=1, fit_intercept=True,
                    learning_rate='optimal', loss='hinge',
                    verbose=0, warm_start=False)
clf.fit(X_train_bow, y_train)
tr_pred = clf.predict(X_train_bow)
cv_pred = clf.predict(X_cv_bow)
```

In [64]:

```
te_pred = clf.predict(X_test_bow)
```

In [62]:

```
print(confusion_matrix(y_cv, cv_pred))
```

```
[[ 127   42   11   29   70]
 [   50   49   38   84  130]
 [   23   19   34  159  220]
 [    4    9    6  264  928]
 [    2    2    2  130 1667]]
```

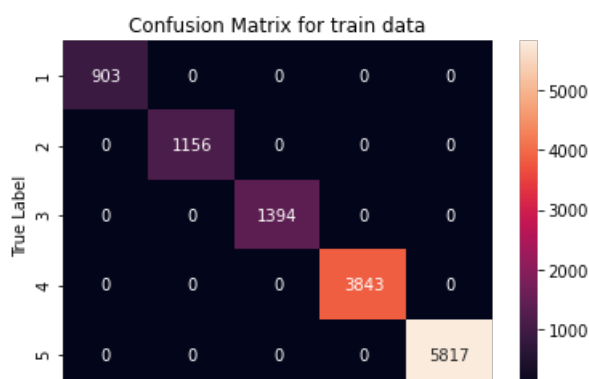
In []:

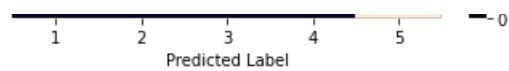
Confusion Matrix of Train data

In [92]:

```
from sklearn.metrics import confusion_matrix
import seaborn as sns

conf_mat = confusion_matrix(y_train, tr_pred)
class_label = ["1", "2", "3", "4", "5"]
df = pd.DataFrame(conf_mat, index = class_label, columns = class_label)
sns.heatmap(df, annot = True, fmt="d")
plt.title("Confusion Matrix for train data")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



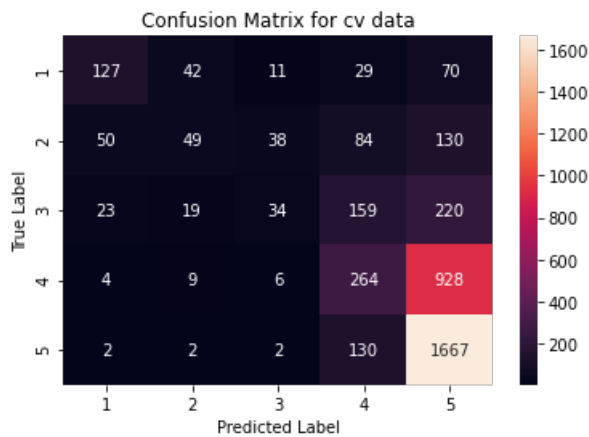


Confusion Matrix of CV data

In [91]:

```
from sklearn.metrics import confusion_matrix
import seaborn as sns

conf_mat = confusion_matrix(y_cv, cv_pred)
class_label = ["1", "2", "3", "4", "5"]
df = pd.DataFrame(conf_mat, index = class_label, columns = class_label)
sns.heatmap(df, annot = True, fmt="d")
plt.title("Confusion Matrix for cv data")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

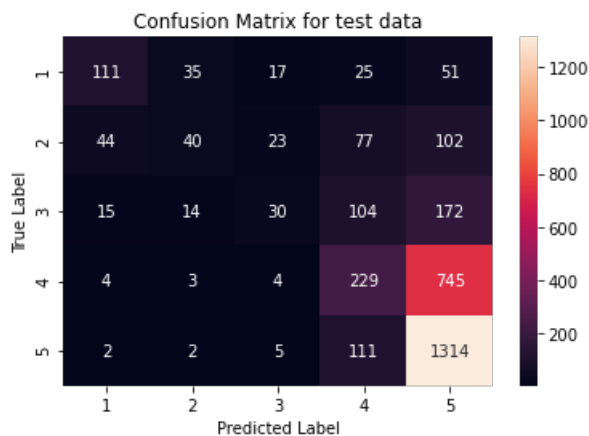


Confusion Matrix of Test data

In [66]:

```
from sklearn.metrics import confusion_matrix
import seaborn as sns

conf_mat = confusion_matrix(y_test, te_pred)
class_label = ["1", "2", "3", "4", "5"]
df = pd.DataFrame(conf_mat, index = class_label, columns = class_label)
sns.heatmap(df, annot = True, fmt="d")
plt.title("Confusion Matrix for test data")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



In []:

F1-Score

In [68]:

```
from sklearn.metrics import precision_score, recall_score, f1_score
```

In [99]:

```
print("Train ", f1_score(y_train, tr_pred, average='micro'))
print("cv ", f1_score(y_cv, cv_pred, average='micro'))
print("Test ", f1_score(y_test, te_pred, average='micro'))
```

```
Train  1.0
cv    0.5223225176872408
Test  0.5257700518450747
```

RMSE

In [76]:

```
rms = sqrt(mean_squared_error(y_test, te_pred))
print(rms)
```

```
1.142620221161812
```

Mean Absolute Deviation Error

In [79]:

```
from sklearn.metrics import mean_absolute_error
mae=mean_absolute_error(y_test, te_pred)
print(mae)
```

```
0.6919792619701128
```

In []:

In []:

In []:

Fine Tuning LGBMClassifier

In [82]:

```
from lightgbm import LGBMRegressor, LGBMClassifier
from sklearn.model_selection import GridSearchCV
```

In []:

```
grid={
```

```
'max_depth': [30,40,45],
'n_estimators': [4000,4500],
'learning_rate':[0.1,0.2]}
```

In []:

```
%time
clf= LGBMClassifier(colsample_bytree=0.8,subsample=0.9,min_child_samples=50,num_leaves=20)
rf_random = GridSearchCV(estimator = clf, param_grid = grid,
                        cv=3, verbose=1)
rf_random.fit(X_train_bow,y_train,verbose=True)
```

In []:

```
bestpar=rf_random.best_params_
bestpar
```

Out[]:

```
{'learning_rate': 0.1, 'max_depth': 30, 'n_estimators': 4000}
```

Predicting the results after Fine Tuning

In [83]:

```
clf=LGBMClassifier(num_leaves=10,colsample_bytree=0.8,subsample=0.9,min_child_samples=50,
                  learning_rate= 0.1, max_depth= 20, n_estimators= 4000)
clf.fit(X_train_bow,y_train)
predt=clf.predict(X_train_bow)
```

In [84]:

```
predcv=clf.predict(X_cv_bow)
```

In [85]:

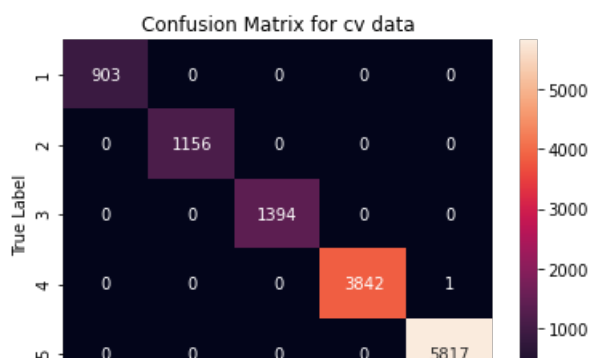
```
predte = clf.predict(X_test_bow)
```

Confusion Matrix of Train data

In [94]:

```
from sklearn.metrics import confusion_matrix
import seaborn as sns

conf_mat = confusion_matrix(y_train, predt)
class_label = ["1", "2", "3", "4", "5"]
df = pd.DataFrame(conf_mat, index = class_label, columns = class_label)
sns.heatmap(df, annot = True,fmt="d")
plt.title("Confusion Matrix for cv data")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



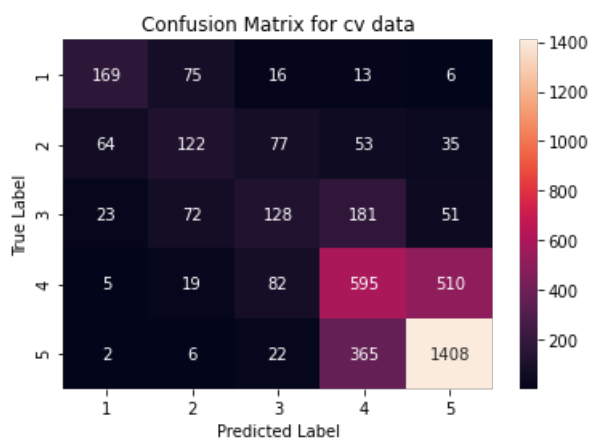


Confusion Matrix of CV data

In [90]:

```
from sklearn.metrics import confusion_matrix
import seaborn as sns

conf_mat = confusion_matrix(y_cv, predcv)
class_label = ["1", "2", "3", "4", "5"]
df = pd.DataFrame(conf_mat, index = class_label, columns = class_label)
sns.heatmap(df, annot = True, fmt="d")
plt.title("Confusion Matrix for cv data")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

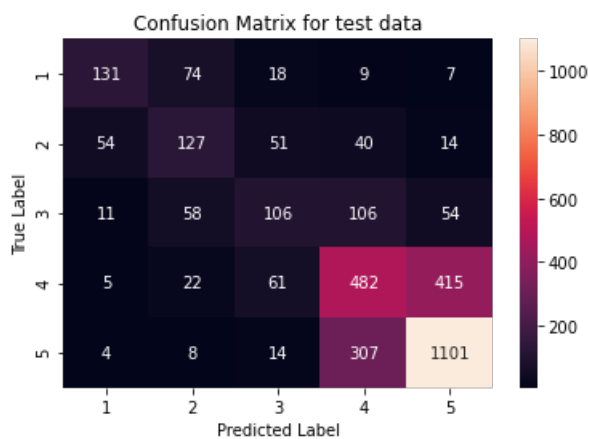


Confusion Matrix of Test data

In [89]:

```
from sklearn.metrics import confusion_matrix
import seaborn as sns

conf_mat = confusion_matrix(y_test, predte)
class_label = ["1", "2", "3", "4", "5"]
df = pd.DataFrame(conf_mat, index = class_label, columns = class_label)
sns.heatmap(df, annot = True, fmt="d")
plt.title("Confusion Matrix for test data")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```



In []:

In [86]:

```
predte
```

Out[86]:

```
array([1, 4, 4, ..., 0, 4, 4])
```

F1-Scores of the data

In [98]:

```
print("The train f1-score is", f1_score(y_train,predt,average='micro'))
print("The cv f1-score is", f1_score(y_cv,predcv,average='micro'))
print("The test f1-score is", f1_score(y_test,predte,average='micro'))
```

The train f1-score is 0.9999237398001983

The cv f1-score is 0.5908758233715541

The test f1-score is 0.5937785910338518

RMSE

In [95]:

```
rms = sqrt(mean_squared_error(y_test, predte))
print(rms)
```

0.830568759812122

Mean Absolute Error

In [96]:

```
from sklearn.metrics import mean_absolute_error
mae=mean_absolute_error(y_test, predte)
print(mae)
```

0.4867337602927722

In []:

In []:

In []:

Training Deep Learning Model

In []:

```
import plotly.figure_factory as ff
import os
```

```

import gc

from sklearn.model_selection import KFold
from sklearn.preprocessing import LabelEncoder
import json
from tensorflow.keras.preprocessing import text, sequence
from sklearn.feature_extraction.text import CountVectorizer
import numpy as np
%matplotlib inline
import matplotlib.pyplot as plt
from tensorflow.keras import layers
from keras.layers import Reshape, Concatenate
from tensorflow.keras.layers import Reshape
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.layers import Flatten
from tensorflow.keras.layers import Embedding
from tensorflow.keras import regularizers
from tensorflow.keras.layers import LeakyReLU

```

In [38]:

```

for i in X_train[:3]:
    print(i)

```

amazing room view hyatt hesitated bit giving hotel excellent rating difficult check experience balance really enjoyed hotel highly recommend good hotel great location near shopping action downtown walking distance major attractions stayed hotel week needed wished car importantly room spectacular tough hotel critic simply way complain room high end finishes beautiful foyer area large sitting area tvs automatic shades plush king bed views glorious sweeping view downtown bay south unobstructed view space needle west mention booked called emerald suite stayed floor splurge little highly recommend room size views higher floors worth price paid night lot kept thinking comparable room comparable views nyc probably cost nearly thousand suspect lower floors quite views enjoy impressive finishes furnishings bedding bad check process inefficient actually comical arrived cruise understandably way early desk staff said come room noon headed took time respectful returned room told cleaning service room ready shortly staff asked cell phone number said soon went lunch hour later room phone stopped told minutes waited view desk minutes later approached staff member dealing gone grand scheme things end world based postings tripadvisor clearly aberration staying excellent rating curious future reviewers similar experiences overall notwithstanding odd check experience location room quality views hotel simply can't beat best stayed wife weeks jan amazing resort best went animation team amazing year clean resort nice food lot animation read bad rating try price you regret going places caribbean europe africa best vacations enjoy hotels price quiet ifa villas animation ifa village ocean far beach problem hear disco night if stay ocean used went times disco nice ambiance brand new spa gym pizzeria ice cream parlour nice friendly place stay infos contact jacobberdugo canada com okay great location right faneuil hall unfortunately building bit old maintained room superior room view quincy market nice little balcony view hear

Tokenizing the data

In [43]:

```

tokenizer = Tokenizer(lower=False, num_words=80000)

for text in tqdm(preprocessed_reviews):
    tokenizer.fit_on_texts(text.split(" "))

```

100%|██████████| 20491/20491 [16:39<00:00, 20.50it/s]

In [44]:

```

pickle.dump(tokenizer, open("tokenizertripadv.pickle", "wb"))

```

In [27]:

```

tokenizer = pickle.load(open("tokenizertripadv.pickle", "rb"))

```

In [28]:

```
In [20]:
```

```
max_length = max([len(x) for x in X])
vocab_size = len(tokenizer.word_index)+1
exp_sen = 1
```

```
In [29]:
```

```
max_length
```

```
Out[29]:
```

```
12610
```

```
In [30]:
```

```
encoding = {1: 0,
            2: 1,
            3: 2,
            4: 3,
            5: 4
            }

#labels = ['1', '2', '3', '4', '5']

y = df['Rating'].copy()
y.replace(encoding, inplace=True)
```

```
In [ ]:
```

```
In [31]:
```

```
X_train,X_cv,y_train,y_cv = train_test_split(X,y,test_size=0.2, random_state=42, stratify=y)

X_train,X_test,y_train,y_test = train_test_split(X_train,y_train,test_size=0.2)
```

Converting text to sequences

```
In [32]:
```

```
def compute_text(X_train,X_cv,X_test,tokenizer):

    #train_text = tokenizer.texts_to_sequences(X_train.text.values)
    train = tokenizer.texts_to_sequences(X_train)
    cv = tokenizer.texts_to_sequences(X_cv)
    test = tokenizer.texts_to_sequences(X_test)

    #train_text = sequence.pad_sequences(train_text, maxlen=300)
    train = sequence.pad_sequences(train,maxlen=max_length)
    cv = sequence.pad_sequences(cv,maxlen=max_length)
    test = sequence.pad_sequences(test,maxlen=max_length)

    return train,cv,test
```

```
In [33]:
```

```
train,cv,test = compute_text(X_train,X_cv,X_test,tokenizer)
```

Using Glove Vectors as Pre-trianed word vectors embedding

```
In [34]:
```

```
with open('glove_vectors', 'rb') as f:
    glove=pickle.load(f)
    glove_words=set(glove.keys())
```



```

embedd_matrix= np.zeros((len(tokenizer.word_index)+1,300))
for i,j in tokenizer.word_index.items():
    if i in glove_words:
        embed_vec=glove[i]
        embedd_matrix[j]=embed_vec
print(embed_vec.shape,embedd_matrix.shape)

```

```
(300,) (48904, 300)
```

In [35]:

```
cv.shape
```

Out[35]:

```
(4099, 12610)
```

In []:

In [36]:

```

from tensorflow.keras.layers import Input, Dense, Embedding, SpatialDropout1D, concatenate, Masking
from tensorflow.keras.layers import LSTM, Bidirectional, GlobalMaxPooling1D, Dropout
from tensorflow.keras.preprocessing import text, sequence
from tqdm import tqdm_notebook as tqdm
import tensorflow as tf
import tensorflow.keras
import pickle
import tensorflow.keras
from tensorflow.keras import layers
from tensorflow.keras.layers import Dropout
from tensorflow.keras.layers import BatchNormalization
from tensorflow.keras.models import Model

```

In [37]:

```
k = tf.keras.initializers.he_normal(seed=None)
```

In [38]:

```
train.shape[1:]
```

Out[38]:

```
(12610,)
```

Model

In [45]:

```

text_in = Input(shape=(None,),name='input1')
t= Embedding(*embedd_matrix.shape, weights=[embedd_matrix])(text_in)
#t = Embedding(vocab_size, embedding_dim)(text_in)

t= layers.Bidirectional(tf.keras.layers.LSTM(32, activation="tanh",recurrent_activation="sigmoid",
                                             return_sequences=True))(t)

t=tf.keras.layers.LeakyReLU(alpha=0.3)(t)
t= tensorflow.keras.layers.GlobalMaxPooling1D()(t)

hidden = Dense(100, activation='relu',kernel_initializer=k)(t)
hidden=tf.keras.layers.LeakyReLU(alpha=0.3)(hidden)
hidden = Dropout(0.5)(hidden)
hidden = Dense(96, activation='relu',kernel_initializer=k)(hidden)
hidden=tf.keras.layers.LeakyReLU(alpha=0.3)(hidden)
hidden = Dropout(0.5)(hidden)
hidden = Dense(100, activation='relu',kernel_initializer=k)(hidden)

```

```

out1 = Dense(5, activation='softmax',name='out1')(hidden)
model = Model(inputs=[text_in], outputs=[out1])

model.compile(loss="sparse_categorical_crossentropy",
optimizer=tf.keras.optimizers.Adam(learning_rate=0.03),
metrics=['accuracy'])
model.summary()

```

Model: "functional_3"

Layer (type)	Output Shape	Param #
input1 (InputLayer)	[(None, None)]	0
embedding_1 (Embedding)	(None, None, 32)	1564928
bidirectional_1 (Bidirection	(None, None, 64)	16640
leaky_re_lu_3 (LeakyReLU)	(None, None, 64)	0
global_max_pooling1d_1 (Glob	(None, 64)	0
dense_3 (Dense)	(None, 100)	6500
leaky_re_lu_4 (LeakyReLU)	(None, 100)	0
dropout_2 (Dropout)	(None, 100)	0
dense_4 (Dense)	(None, 96)	9696
leaky_re_lu_5 (LeakyReLU)	(None, 96)	0
dropout_3 (Dropout)	(None, 96)	0
dense_5 (Dense)	(None, 100)	9700
out1 (Dense)	(None, 5)	505
Total params: 1,607,969		
Trainable params: 1,607,969		
Non-trainable params: 0		

In [46]:

```

EPOCHS = 3
BATCH_SIZE = 100

```

Training the model

In []:

```

history = model.fit(train, y_train, epochs=EPOCHS, validation_split=0.25, batch_size=BATCH_SIZE,
                    verbose=2)

```

The training is taking more time than expected. So, I'm skipping for now.

In []:

In []:

Conclusion

In [119]:

```
from prettytable import PrettyTable

x=PrettyTable()

x.field_names=(['Model','Test F1-Score','MAE','RMSE'])
x.add_row(['Linear SVM','0.525','1.142','0.691'])
x.add_row(['LGBMClassifier','0.593','0.830','0.48'])
x.add_row(['LSTM Model','-','-','-'])
print(x)
```

Model	Test F1-Score	MAE	RMSE
Linear SVM	0.525	1.142	0.691
LGBMClassifier	0.593	0.830	0.48
LSTM Model	-	-	-

By the above results we can LGBMClassifier performs better compared to Linear SVM.

In []:

In []: