

Text Classification on Hotel Reviews

Excelr NLP project presented by **Group 2**





Group Members

Group 2

Chithra K V

Shilpa Chavan

Sajir J M

Rahul Tiwari

Sreshta Sharon

Sefin Francis

Ankita Singh

Vedavathi M





Business Objective

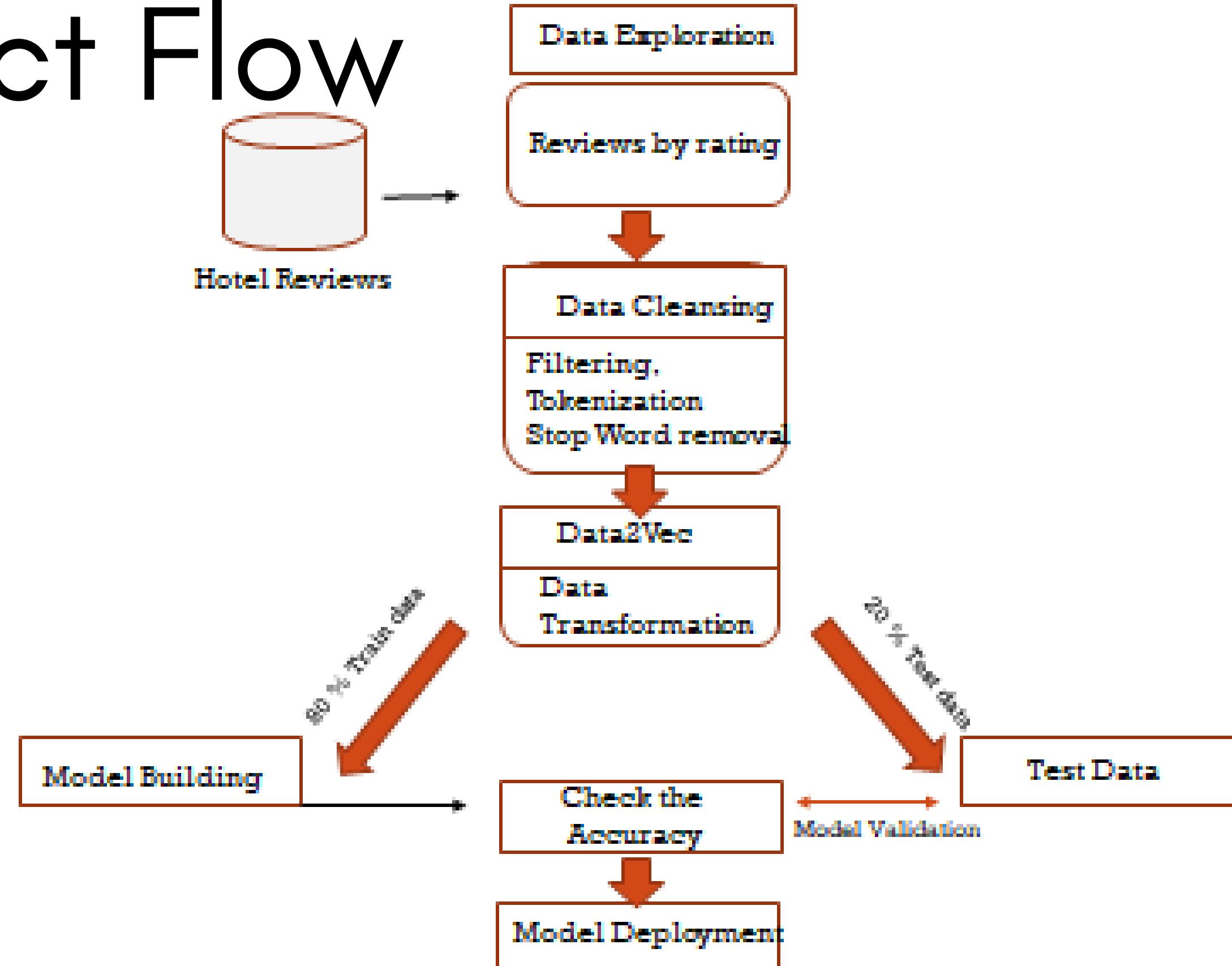


 This is a sample dataset that consists of 20,000 reviews and ratings for different hotels and our goal is to examine how travelers are communicating their positive and negative experiences in online platforms for staying in a specific hotel and major objective is what are the attributes that traveler are considering while selecting a hotel.

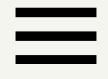
 This manager can understand which elements of their hotel influence more in forming a positive review or improves hotel brand image.



Project Flow

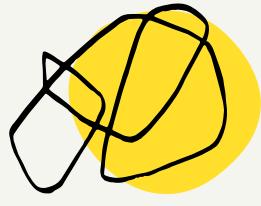


Sentiment
Prediction
→ →



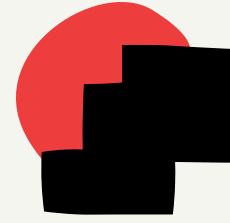


Data Sources & Fields



Source

Data has been
collected from Kaggle



Data Fields

ID - An Unique ID for each and
every review

Review - given by the
customers

Rating - given by customer (1
being the least and 5 being the best)





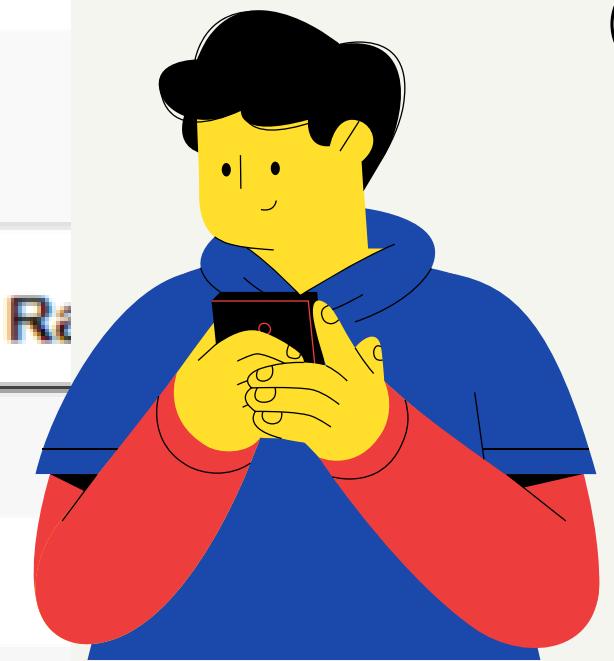
Exploratory Data Analysis



DATASET

```
] : train_df=pd.read_csv("train.csv")
train_df
```


	ID	Review	Rating
0	0	exceptional service nice all-around daughter s...	5.0
1	1	beautiful relaxing jw marriott desert ridge re...	5.0
2	2	great location great location 5 mins subway ta...	5.0
3	3	pleased nice safe hotel, flower market hotel v...	5.0
4	4	excellent hotel service great hotel excellent ...	5.0
...			
14338	14338	hotel madrid hotel perfect, location tiny quie...	5.0
14339	14339	excellent hotel stay florence hotel chosen tri...	5.0
14340	14340	great place relax know looking vacation book t...	5.0
14341	14341	better just got week seattle loved minute, pac...	5.0
14342	14342	stay clear, internet reservation friday rang h...	5.0
14343 rows × 3 columns			



```
8]: #check the information present in the train dataset  
train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 14343 entries, 0 to 14342  
Data columns (total 3 columns):  
 #   Column   Non-Null Count   Dtype     
 ---  --      --           --       --  
 0   ID       14343 non-null    int64  
 1   Review   14343 non-null    object  
 2   Rating   14343 non-null    int64  
dtypes: int64(2), object(1)  
memory usage: 336.3+ KB
```



The data set contains three field:

- **The ID and Rating are integer data types**
- **Review column is a object data type.**
- **The data set contain no null values.**



PREPROCESSING: DATA CLEANING

Remove

- Remove punctuation
- Remove numerical values
- Remove common non-sensical text (/n)
- Remove stop words

Text

- Make text all lower case
- Tokenize Text

After Tokenization

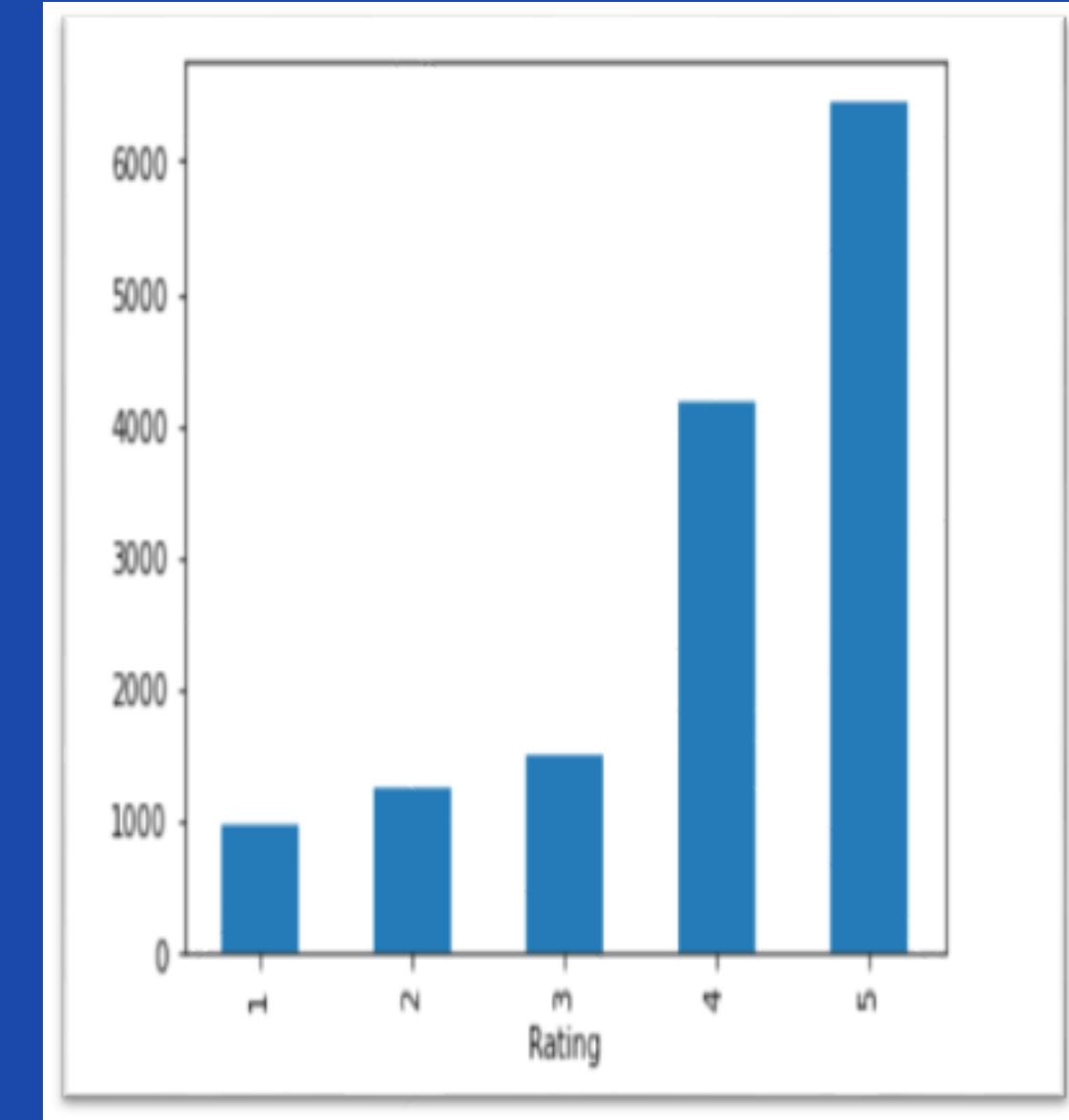
- Stemming / lemmatization
- Parts of speech tagging
- Create bi-grams or tri-grams etc.





Total Number of Customers and percentage who rated the Hotel as per rating category

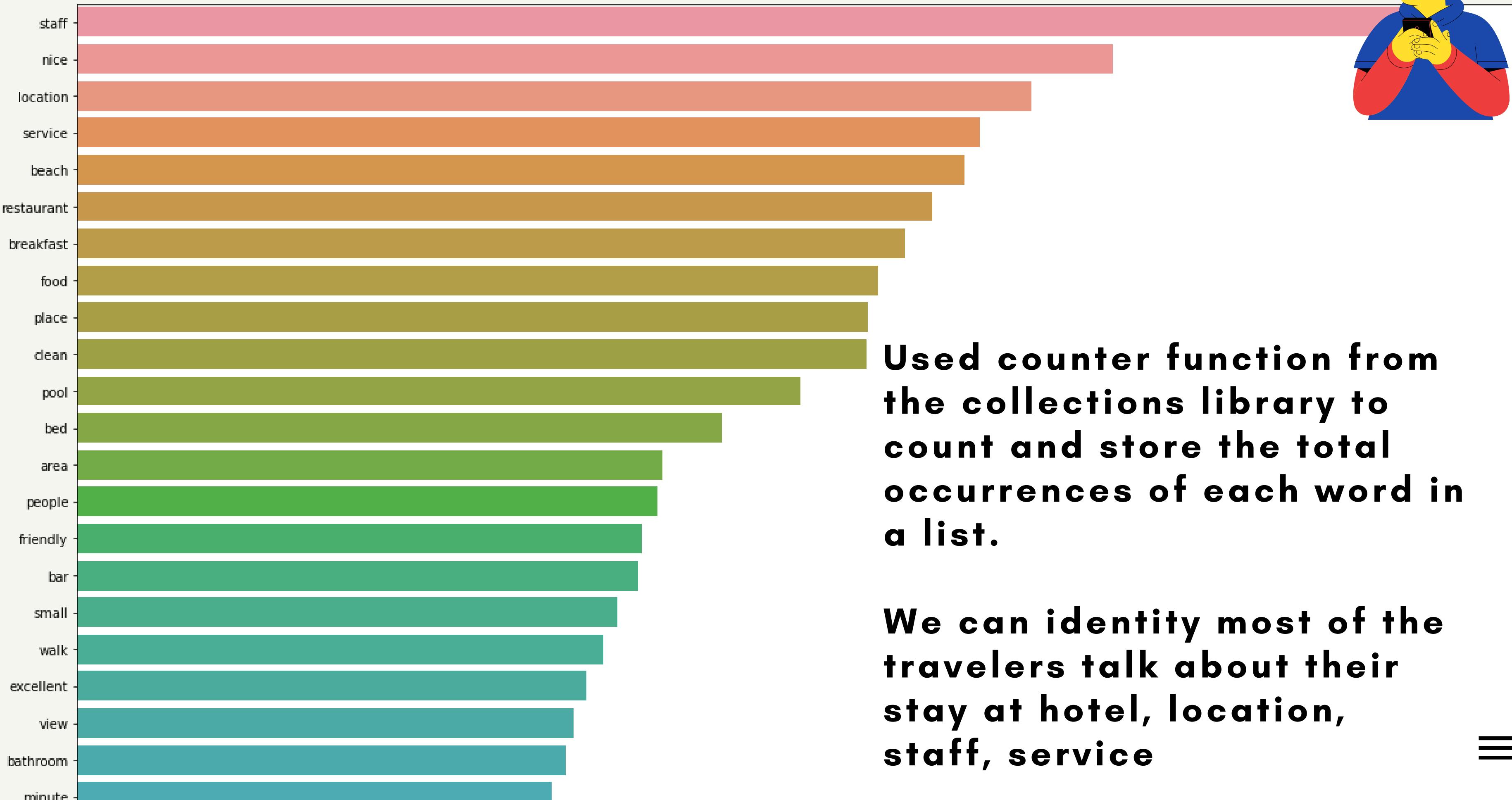
Rating	
1	977
2	1248
3	1510
4	4172
5	6436
Name:	ID, dtype: int64



Here we can see the distribution of reviews on basis of rating. Maximum customers have given positive review as feedback as compared to negative feedback.



Word frequency of Reviews



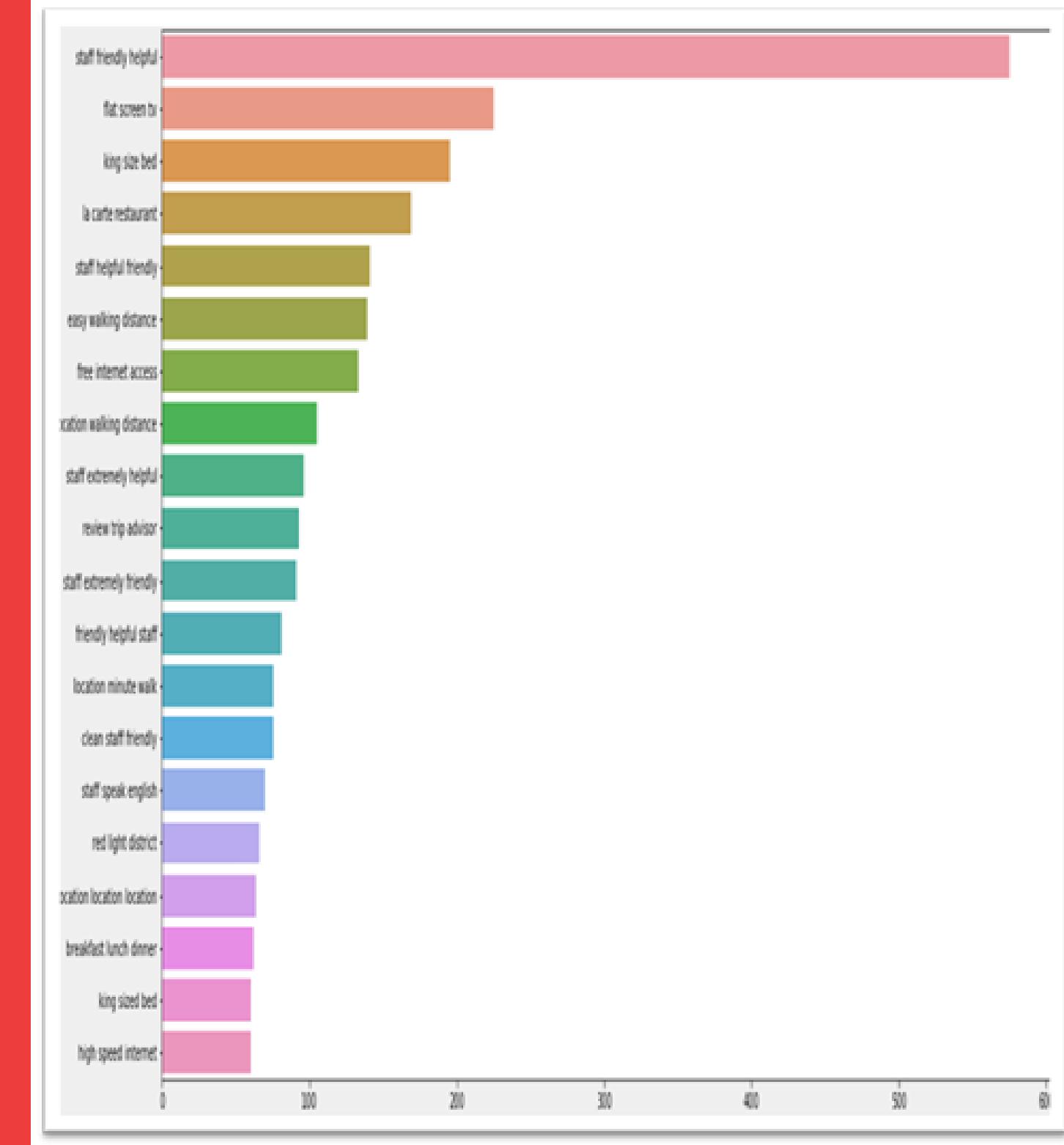
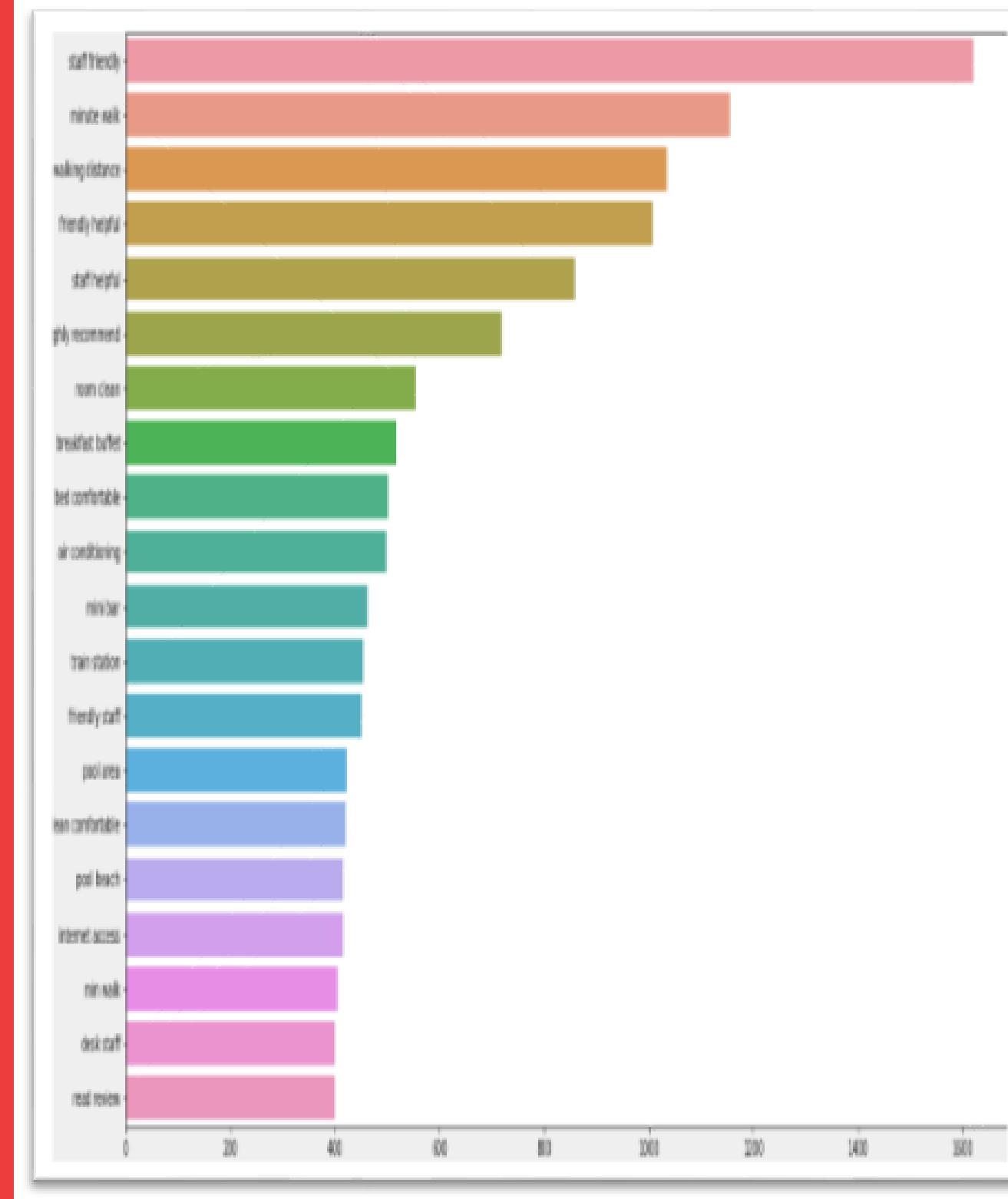


Ngram Exploration

- **Ngrams are simply sequences of n words.** If the number of words is two, it is called bigram. For 3 words it is called a trigram and so on.
- Looking at most frequent n-grams can give you a better understanding of the context in which the word was used.
- To implement n-grams we will use ngrams function from nltk.util.
- To build a representation of vocabulary we will use Countvectorizer.
- Countvectorizer is a simple method used to tokenize, vectorize and represent the corpus in an appropriate form.



visualization of top 20 bigram & Trigram



- Here , we observe that staff friendly, helpful, minute walk, room clean etc that are related to hotel dominates the hotel review.
- Hotel managers can easily identify what travelers are most interested with the help of ngrams.

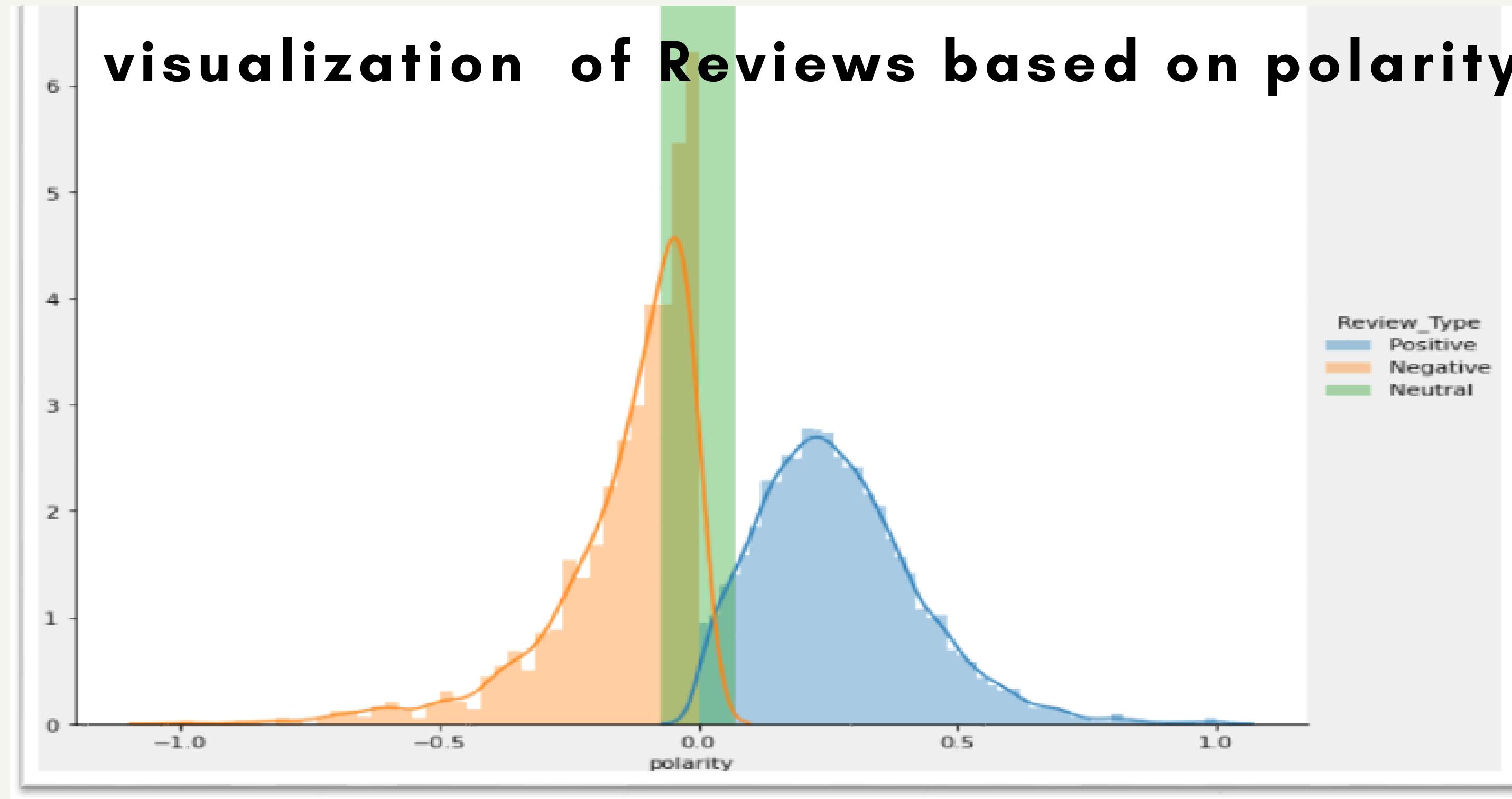
Sentimental Analysis of Reviews

- **Sentiment Analysis** Sentiment analysis is basically the process of determining the attitude or the emotion of the writer, i.e., whether it is positive or negative or neutral.
- The sentiment function of `textblob` returns two properties, **polarity**, and **subjectivity**.
- Polarity is float which lies in the range of [-1,1] where 1 means **positive statement** and -1 means **a negative statement**.
- Subjective sentences generally refer to personal opinion, emotion or judgment whereas objective refers to factual information.
- Subjectivity is also a float which lies in the range of [0,1].



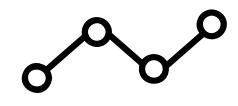


visualization of Reviews based on polarity

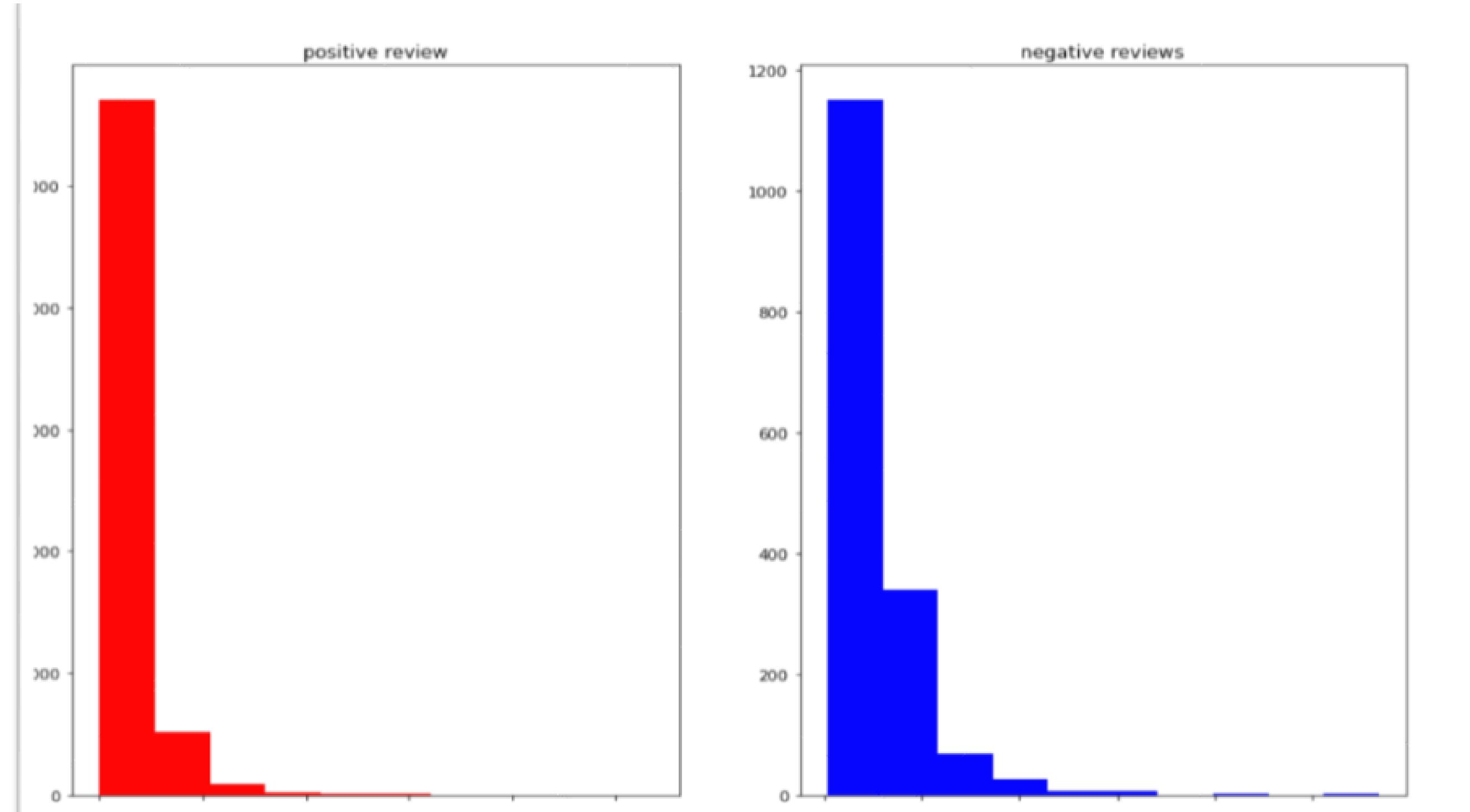


- Distribution of reviews based on polarity score :
- Most of positive reviews are distributed in the range between 0 to 0.5.
- Most negative reviews are distributed in the range of -0.5 to 0 .
- less distribution of neutral review





using histogram to visualize length of positive and negative reviews

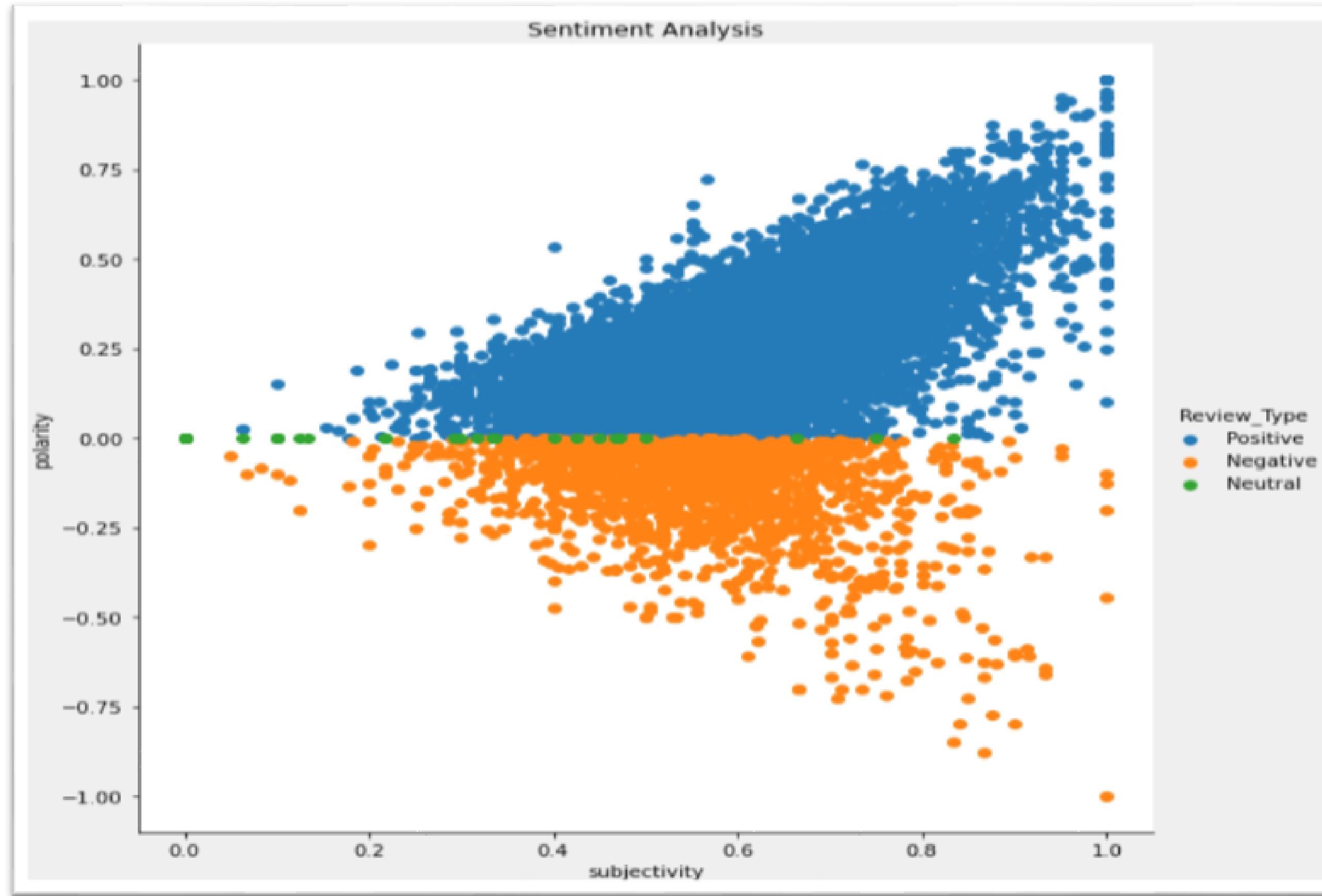


Both positive and negative reviews words length are positive skewed, all the reviews length in the range of 1000



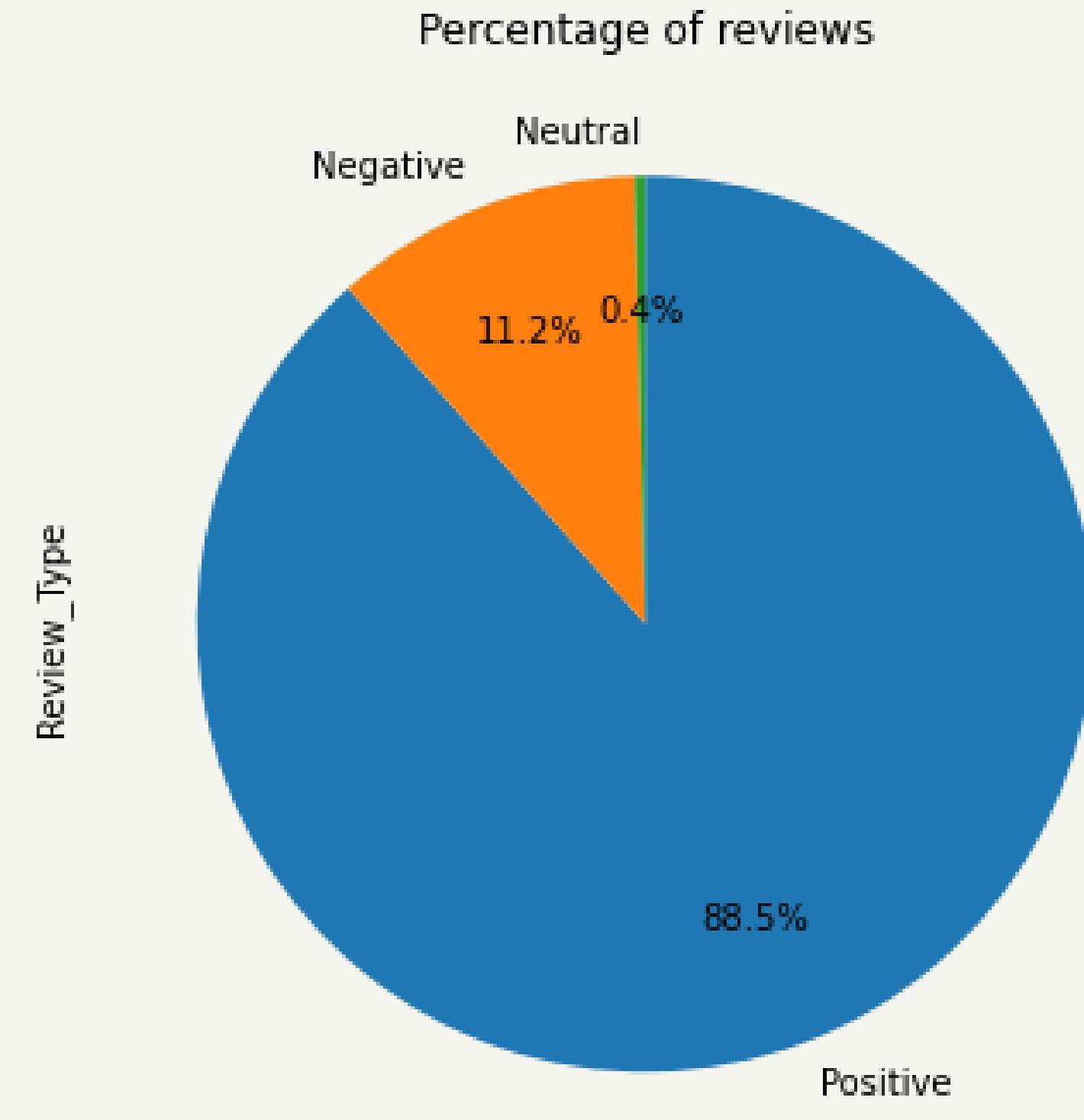


Scatter plot representation of reviews



Pie chart representation of reviews

- After Sentiment analysis, We can see that 88.5% reviews are positive.
- 11.2% reviews are negative.
- 0.4 % reviews are neutral.
- This imbalanced dataset issue does not give correct predictions.
- Inorder to solve this issue , we can do sampling technique such as over sampling and under sampling.



Highest Rating with 1 Rating

[28]:

```
#Reviews that have the highest polarity (most positive sentiment) but with a 1-star:  
df[(df.Rating == 1) & (df.polarity == 1)].head(5)
```

Out[28]:

ID		Review	Rating	polarity	subjectivity
53	53	bad ugly little good, wife booked trip majesti...	1	1.0	0.616495
104	104	think twice booking breezes, just got breezes ...	1	1.0	0.598214
156	156	dump, christmas vacation family brother family...	1	1.0	0.383333
195	195	wish stayed home, omg, wish brushed spanish le...	1	1.0	0.491643
223	223	not returning gone idapa year vacation time t...	1	1.0	0.548038

Lowest Rating with 5 Rating



```
#Reviews that have lowest polarity (most negative sentiment) but with a 5-star:  
df[(df.Rating == 5) & (df.polarity == -1)].head(5)
```

Out[27]:

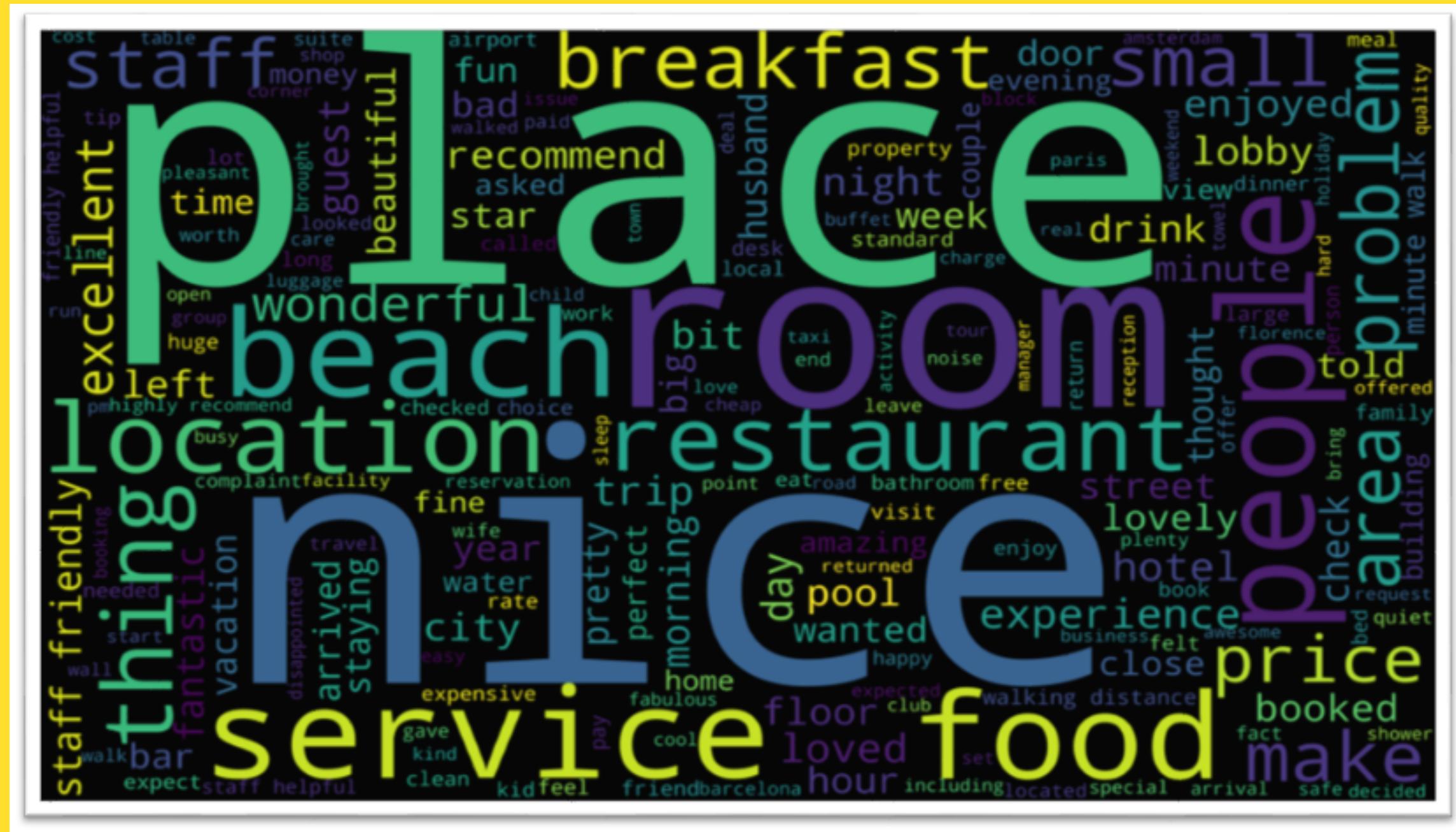
ID		Review	Rating	polarity	subjectivity
775	775	n't believe bad review hotel, 2nd time resort,...	5	-1.0	0.612982
1194	1194	lost madrid wife 10 month old recently stayed ...	5	-1.0	0.418750
5209	5209	disgusted comments just wanted make comments, ...	5	-1.0	0.531973
7644	7644	smooth sailing bad reviews worried soon went c...	5	-1.0	0.561111
8700	8700	oasis stayed near 3 nights march strength reco...	5	-1.0	0.747778

+ Code

+ Markdown

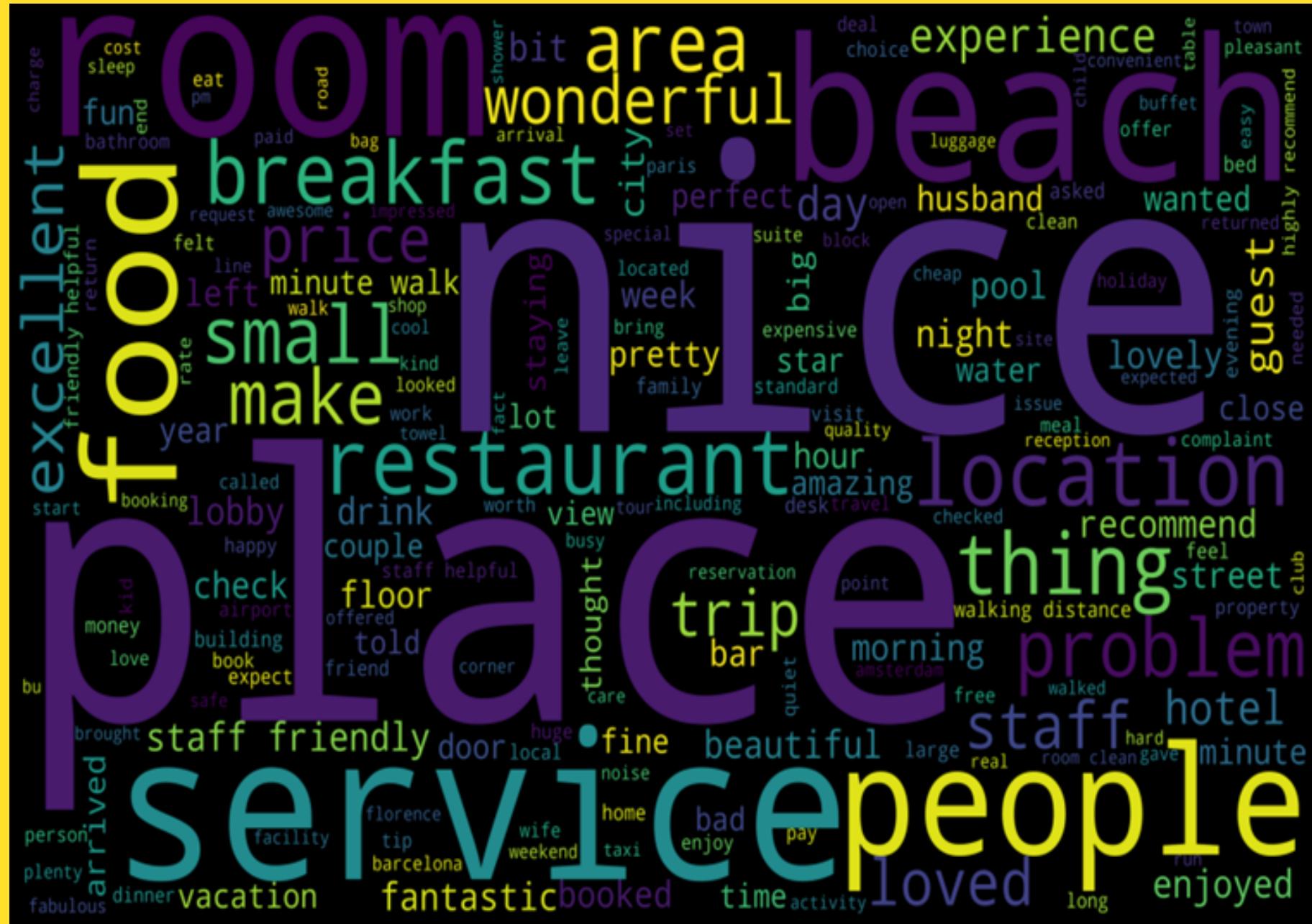
[28]:

Word cloud representation of all reviews



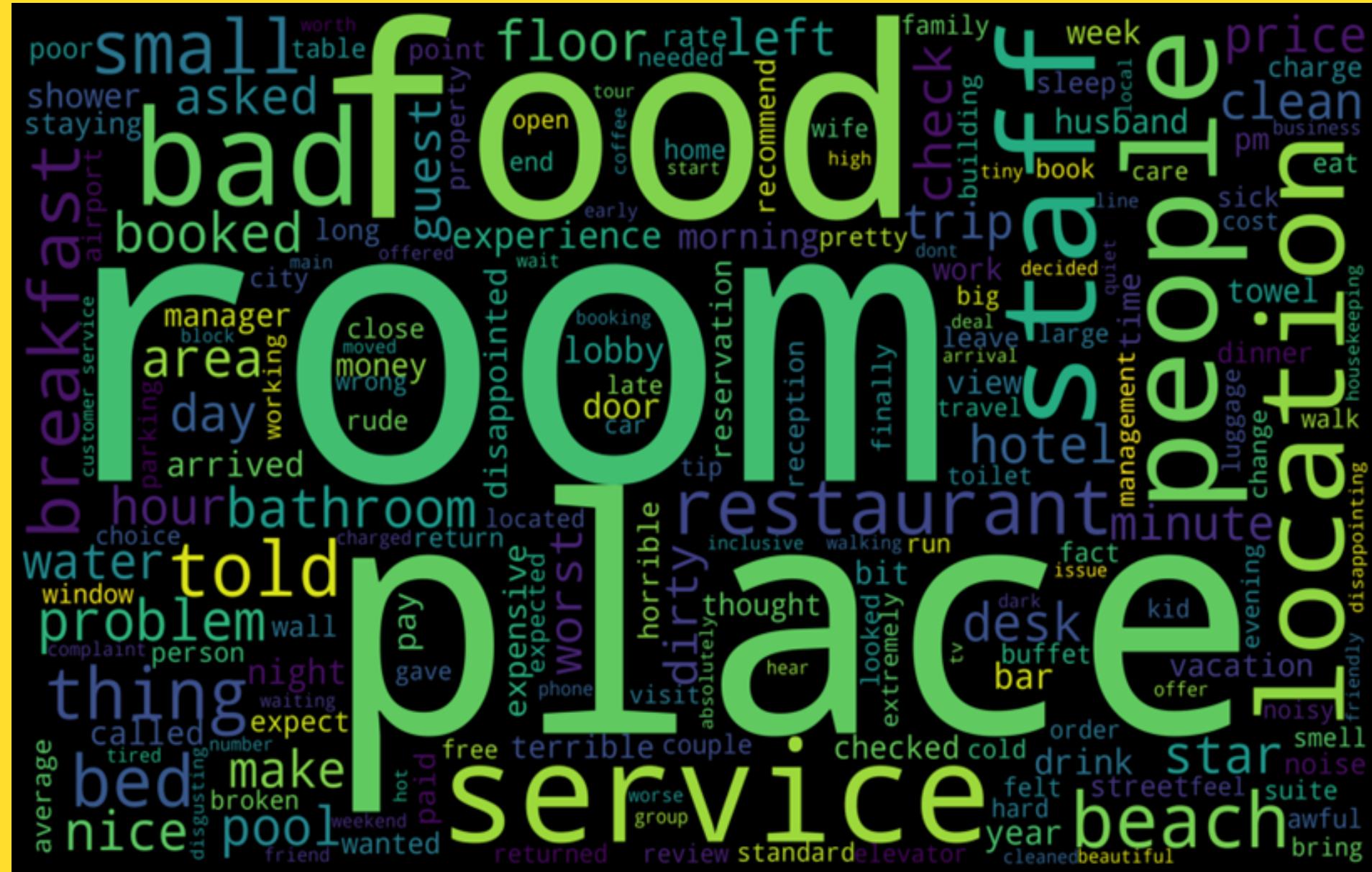
- Most of the words are related to the hotels: room, staff, breakfast, restaurant etc.
 - Some words are more related to the customer experience with the hotel stay: wonderful, beautiful, nice, small, excellent, quality etc.

Positive Word Cloud



- The most positive reviews correspond to some good feedbacks.
 - Here the customers in positive reviews related to hotel are place, service, people, food restaurant, staff, breakfast etc
 - Some words are more related to the customer experience with the hotel stay: wonderful, beautiful, nice, loved, excellent, quality etc.

Negative Word Cloud



- The most negative reviews correspond to some bad feedback . Here the customers in negative reviews related to hotel more complained about place, room,breakfast,staff,noise etc
 - Some words are more related to the customer experience with the hotel stay: wonderful, beautiful, nice, loved,excellent,quality etc.



Conclusions - Word Cloud

- Most of the customers talking about place, room, service, breakfast, beach, staff, reataurant, pool, vacation, walking distance
- Very few customers talking about price
- Except price the travelers would consider the place, room, service, breakfast, beach, staff, restaurant, pool, vacation, walking distance etc
- The hotel manager would consider these elements, which helpful more in forming a positive review or improves hotel brand image.





Sentiment Analysis Using Logistics Regression



FEATURE EXTRACTION



When dealing with classification of text documents, a first step is to convert text data into numerical format. For that we are using following methods

Bag Of Words

- A **bag of words** is a representation of text that describes the occurrence of words within a document.
- We just keep track of word counts and disregard the grammatical details and the word order.
- It is called a “bag” of words because any information about the order or structure of words in the document is discarded.
- With **count vectorizer**, we count the appearance of the words in each text.

TF-IDF

- Term Frequency-Inverse Document Frequency (TF-IDF) method. This evaluates how important a word is to a document within a large collection of documents (i.e. corpus).
- The importance increases proportionally based on the number of times a word appears in the document but is offset by the frequency of the word in the corpus.





Bag Of Words Approach (Using Bigram)

Extending Bag of words to Bigram:

Input (Reviews) :

Vector count (Sparse matrix of bigram words)

Output (Dependent Variable) :

Review Type (0: Negative , 1: Positive)

Performance metrics:

F1 score: 0.95





Bag Of Words Approach

Positive Bigram

	words	weights
190648	excellent location	1.211722
107921	clean comfortable	1.123184
228391	friendly staff	1.117777
117029	comfortable bed	1.055660
257734	highly recommend	1.032257
51040	bed comfortable	1.031388
227932	friendly helpful	1.031332
473192	service excellent	1.024022
509156	staff friendly	0.883065
359310	nice touch	0.876997

Negative Bigram

	words	weights
473931	service poor	-0.832528
45234	bathroom dirty	-0.865210
540628	terrible experience	-0.903880
52113	bed uncomfortable	-1.025504
263203	horrible experience	-1.027090
474424	service terrible	-1.058892
606671	worst experience	-1.145845
510030	staff rude	-1.196117
37106	bad service	-1.361182
36499	bad bad	-1.388531

From the above 10 Positive & Negative bigram we can see that Customers mainly talk about staff friendliness, Neatness, Service, Location etc



Word Cloud Representation using Bigram – BoW



Positive Wordcloud



Negative Worcloud

From the above top 100 Positive & Negative bigram wordcloud we can see that Customers mainly talk about staff, Neatness, Service, Location, Food, Amenities provided etc



Bag Of Words Approach (Using Trigram)

Extending Bag of words to Trigram:

Input (Reviews) :

Vector count (Sparse matrix of bigram words)

Output (Dependent Variable) :

Review Type (0: Negative , 1: Positive)

Performance metrics:

F1 score: 0.94





Bag Of Words Trigram

Positive Trigram

	words	weights
834822	staff friendly helpful	1.434760
349296	flat screen tv	0.772283
834214	staff extremely helpful	0.764992
367778	free internet access	0.694093
284755	easy walking distance	0.648718
375329	friendly helpful staff	0.644527
510500	location walking distance	0.597240
185678	clean staff friendly	0.567966
835552	staff helpful friendly	0.560707
834212	staff extremely friendly	0.540884

Negative Trigram

	words	weights
507613	location located minor	-0.569733
727936	resolved promptly manager	-0.569733
53887	bad bad renovation	-0.575901
664828	pool restaurant touch	-0.575901
721396	renovation beach pool	-0.575901
55007	bad renovation beach	-0.575901
666247	poor customer service	-0.627727
76265	beach food worst	-0.711389
362756	food worst cockroach	-0.711389
582561	nice beach food	-0.778789

From the above top 10 Positive & Negative trigram we can see that Customers mainly talk about staff, amenities, location distance, food, hotel interiors etc



Word Cloud Representation using Trigram – BoW



Positive Wordcloud



Negative Worcloud

From the above 100 words Positive & Negative wordcloud we can see that Customers mainly talk about staff, Neatness, Service, Location, Food, Amenities, Distance, hotel interior etc



TF - IDF Approach

Extending TF-IDF to Bigram:

Input (Reviews) :

Vector count (Sparse matrix of bigram words)

Output (Dependent Variable) :

Review Type (0: Negative , 1: Positive)

Performance metrics:

F1 score: 0.9432427435755223





TF - IDF Approach

Positive Bigram

	words	weights
227932	rer trains	2.587860
257734	sorts characters	2.145223
228391	reservations fellow	1.862854
190648	opportunities ton	1.853682
107921	filthy turned	1.768543
51040	checked expecting	1.662309
117029	fro visit	1.490482
311878	work internet	1.234567
312459	working loved	1.199356
224492	refreshments departed	1.118997

Negative Bigram

	words	weights
473931	service poor	-1.203204
474424	service terrible	-1.248694
260043	holiday inn	-1.402518
263203	horrible experience	-1.452493
36499	bad bad	-1.498981
37106	bad service	-1.667762
606671	worst experience	-1.706252
510030	staff rude	-1.729577
140247	customer service	-1.855088
137066	credit card	-2.042924

From the above top10 Negative bigram we can see that Customers mainly talk about staff and service



Word Cloud Representation using Trigram - TF IDF



Positive Wordcloud

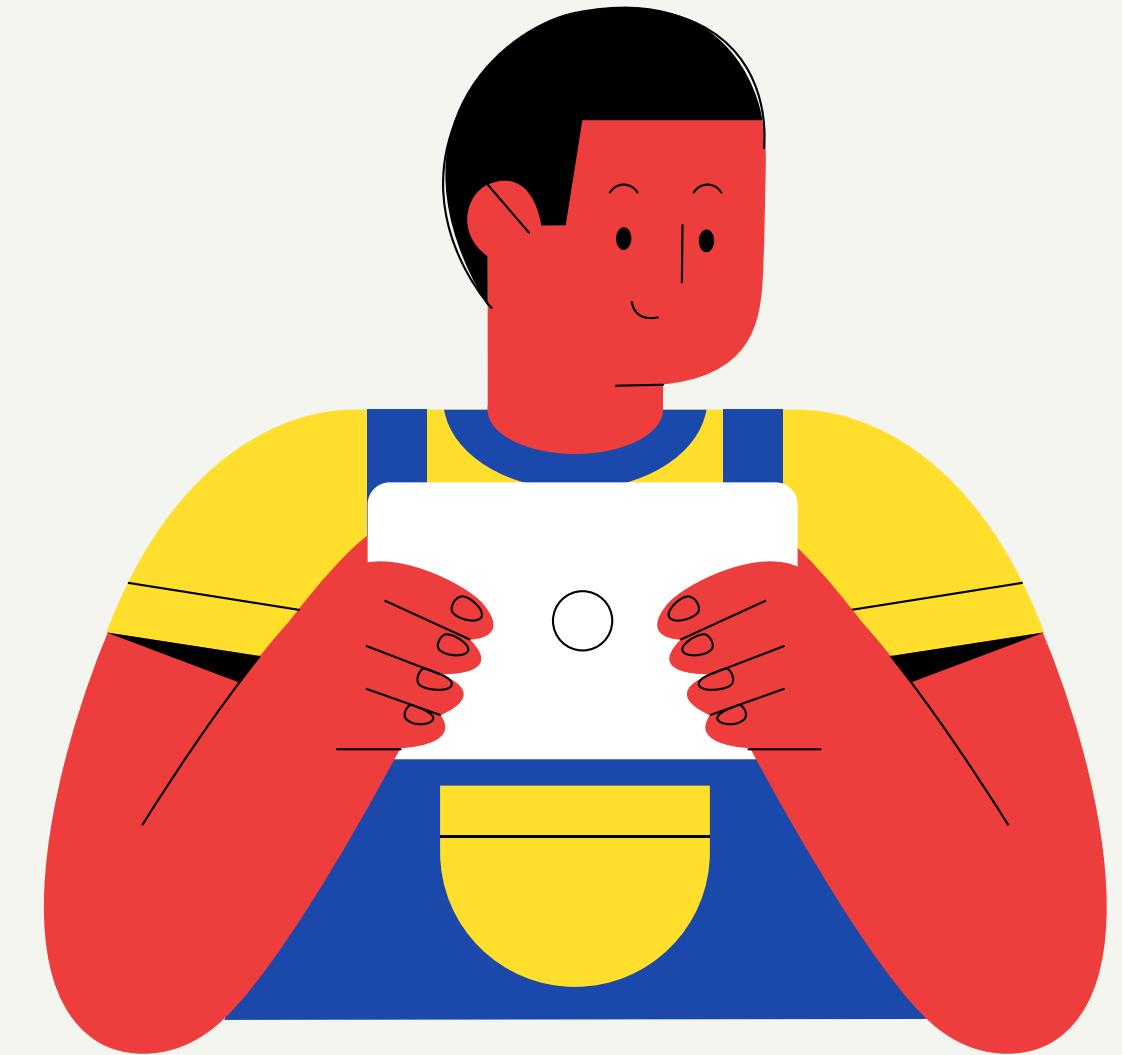


NegativeWordcloud

From the above top 100 Positive & Negative bigram wordcloud we can see that Customers mainly talk about staff, Neatness, Service, Location, Food, Amenities provided etc

Observations

- Bag of words approach give more useful information than tf-idf
- Bag of words bigram word cloud gives more meaningful insights as compared to trigram.
- In Positive reviews people are more talking about staff, hygiene, food, location, service, amenities, restaurant, parking etc
- In Negative reviews people are more talking about poor service, dirtybathroom, rude staff,food,noise,cockroach,water etc



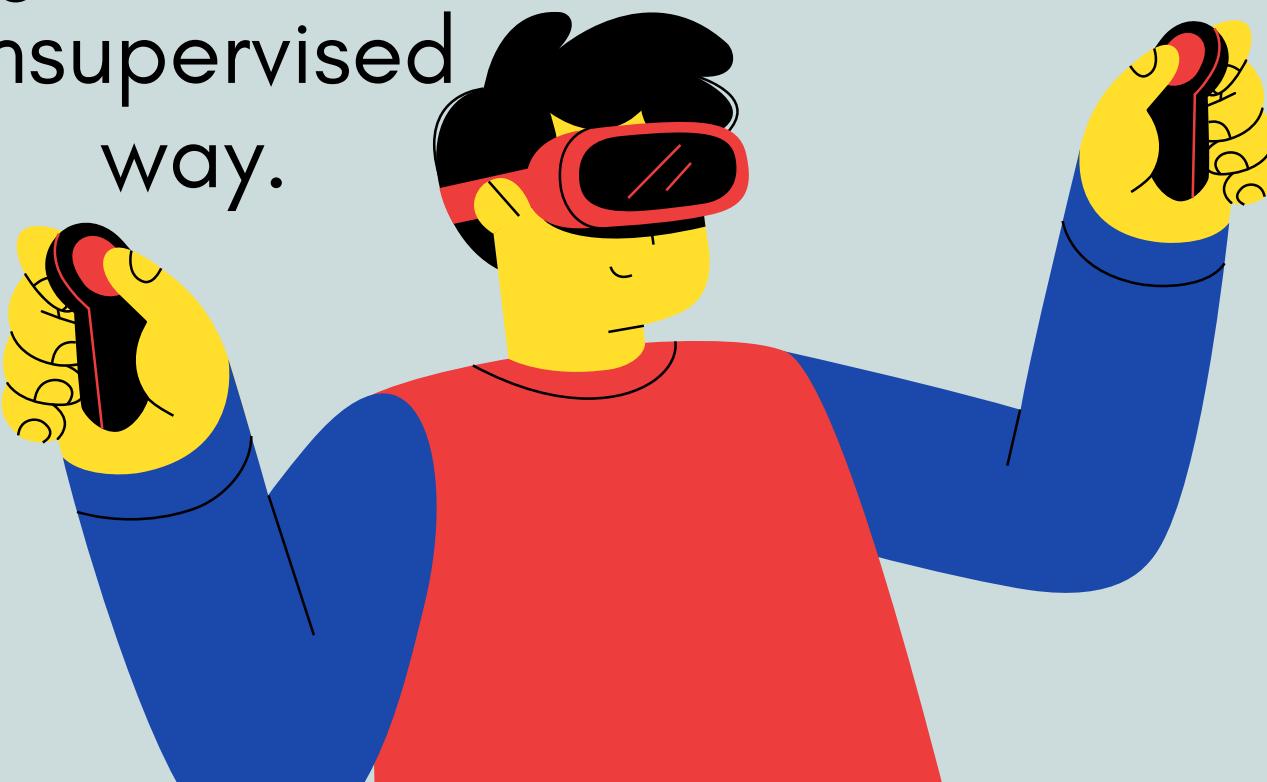
Topic Modelling



Topic Modeling in NLP seeks to find hidden semantic structure in documents.

They are probabilistic models that can help you comb through massive amounts of raw text and cluster similar groups of documents together in an unsupervised way.

Latent Dirichlet Allocation(LDA) is a popular algorithm for topic modeling.



For LDA we can use genism and scikit learn packages



Topic modelling

First take all reviews and generate five topics and print top 50 words



```
top 50 words in topics0
```

```
-----  
['thing', 'experience', 'car', 'suite', 'rate', 'window', 'recommend', 'problem', 'morning', 'lobby', 'internet', 'star', 'parking', 'large', 'bar', 'bit', 'review', 'noise', 'door', 'booked', 'minute', 'shower', 'check', 'walk', 'night', 'free', 'excellent', 'city', 'helpful', 'comfortable', 'desk', 'friendly', 'restaurant', 'hotel', 'street', 'price', 'area', 'view', 'floor', 'bathroom', 'place', 'clean', 'small', 'service', 'breakfast', 'bed', 'nice', 'staff', 'room', 'location']
```

```
top 50 words in topics1
```

```
-----  
['spanish', 'lot', 'tip', 'couple', 'chair', 'eat', 'island', 'area', 'friendly', 'activity', 'towel', 'entertainment', 'family', 'husband', 'dinner', 'loved', 'amazing', 'hour', 'year', 'clean', 'club', 'day', 'review', 'ground', 'bad', 'wonderful', 'problem', 'room', 'thing', 'ocean', 'fun', 'make', 'kid', 'week', 'trip', 'buffet', 'bar', 'staff', 'place', 'drink', 'nice', 'water', 'service', 'vacation', 'restaurant', 'beautiful', 'people', 'pool', 'food', 'beach']
```

```
top 50 words in topics2
```

```
-----  
['hostal', 'eileen', 'est', 'harvard', 'poland', 'yep', 'rink', 'ungherese', 'shaw', 'amandari', 'tel', 'pl', 'tout', 'costco', 'centralized', 'everywhere', 'leicester', 'santiago', 'bvi', 'chianti', 'recommended', 'abel', 'oberoi', 'thierry', 'mann', 'keeper', 'bencoolen', 'sophistication', 'sergio', 'vincents', 'tuilleries', 'equus', 'dai', 'ichi', 'adn', 'auto', 'bt', 'sapna', 'sabrina', 'hogar', 'mixup', 'safeco', 'dynasty', 'nouvel', 'ilikai', 'doumo', 'sedan', 'laura', 'uniquely', 'uma']
```

```
top 50 words in topics3
```

```
-----  
['aswell', 'du', 'cadran', 'kayser', 'kartika', 'hideaway', 'aubusson', 'kodak', 'adore', 'gina', 'campanile', 'ellie', 'onge', 'galeries', 'excelent', 'europa', 'disapoint', 'vaparetto', 'eric', 'atlantico', 'cellai', 'mouffetard', 'auberge', 'tx', 'rue', 'pendini', 'neri', 'sainte', 'za', 'perle', 'tourville', 'orto', 'walter', 'mocenigo', 'sandroni', 'sanjay', 'medici', 'brighton', 'colombia', 'monte', 'rosanna', 'artus', 'keppler', 'carlo', 'jardin', 'miriam', 'gassim', 'miramar', 'milena', 'pax']
```



Topic modelling

Generated top 50 Bigram



top 50 words in topics0

```
-----  
['bathroom clean', 'reception staff', 'bit small', 'centrally located', 'double bed', 'perfect location', 'size bed', 'met  
ro station', 'clean staff', 'easy walk', 'nice touch', 'helpful staff', 'central location', 'min walk', 'location excellen  
t', 'location close', 'screen tv', 'comfortable bed', 'trip advisor', 'air conditioning', 'bathroom small', 'french quarte  
r', 'flat screen', 'breakfast buffet', 'eiffel tower', 'breakfast included', 'excellent location', 'location perfect', 'st  
aff extremely', 'desk staff', 'room clean', 'continental breakfast', 'short walk', 'free internet', 'internet access', 'ca  
ble car', 'highly recommend', 'friendly staff', 'helpful friendly', 'bed comfortable', 'clean comfortable', 'time square',  
'train station', 'union square', 'room small', 'minute walk', 'walking distance', 'friendly helpful', 'staff helpful', 'st  
aff friendly']
```

top 50 words in topics1

```
-----  
['nice breakfast', 'darling harbour', 'la rambla', 'bus stop', 'flat screen', 'place eat', 'spent night', 'space needle',  
'breakfast included', 'reception staff', 'location walking', 'room large', 'free internet', 'nice touch', 'staff extremel  
y', 'room small', 'centrally located', 'metro station', 'desk staff', 'tea coffee', 'location perfect', 'extremely helpfu  
l', 'clean staff', 'location location', 'highly recommended', 'friendly staff', 'buffet breakfast', 'bed comfortable', 'ba  
r restaurant', 'short walk', 'internet access', 'air conditioning', 'clean comfortable', 'breakfast buffet', 'helpful staf  
f', 'mini bar', 'location excellent', 'comfortable bed', 'highly recommend', 'room clean', 'min walk', 'train station', 'l  
a ramblas', 'excellent location', 'friendly helpful', 'staff helpful', 'hong kong', 'staff friendly', 'walking distance',  
'minute walk']
```

top 50 words in topics2

```
-----  
['bed comfortable', 'location barcelona', 'nob hill', 'check told', 'star price', 'bed bug', 'view park', 'matter worse',  
'hot water', 'motor inn', 'extremely rude', 'nice bed', 'lobby nice', 'big disappointment', 'time year', 'mardi gras', 'ni  
ce view', 'make matter', 'service service', 'service terrible', 'parking parking', 'highly recommend', 'credit card', 'bad  
experience', 'room ready', 'millenium hilton', 'location center', 'star property', 'staff friendly', 'friendly staff', 'aq  
ua palm', 'worst experience', 'excellent service', 'breakfast buffet', 'battery park', 'star star', 'parking garage', 'pla  
ce armes', 'cow hollow', 'statue liberty', 'ritz carlton', 'worth money', 'terrible service', 'building site', 'air condit  
ioning', 'french quarter', 'nice place', 'ocean view', 'bad service', 'holiday inn']
```





Topic modelling
Interpret bigram & assign name to
each topics



- Topic 0: Hotel facility
- Topic 1: Food, bar, hotel restaurant
- Topic 2: Service
- Topic 3: Beach and food
- Topic 4: Staff



Topic modelling

Rating wise distribution of topics



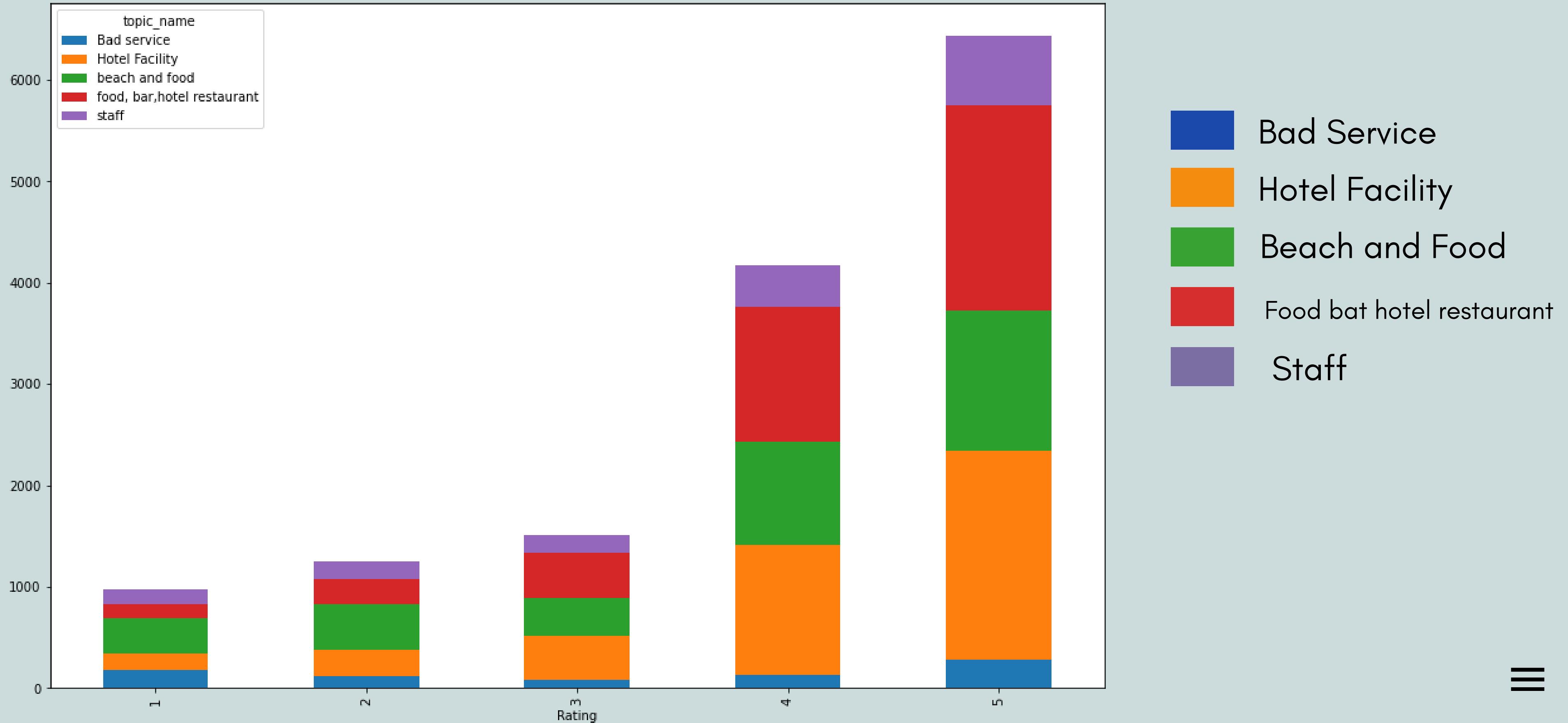
Rating	topic_name	
1	beach and food	346
	Bad service	174
	Hotel Facility	163
	staff	156
	food, bar,hotel restaurant	138
2	beach and food	440
	Hotel Facility	265
	food, bar,hotel restaurant	258
	staff	169
	Bad service	116
3	food, bar,hotel restaurant	444
	Hotel Facility	436
	beach and food	382
	staff	174
	Bad service	74
4	food, bar,hotel restaurant	1331
	Hotel Facility	1278
	beach and food	1016
	staff	417
	Bad service	130
5	Hotel Facility	2056
	food, bar,hotel restaurant	2030
	beach and food	1385
	staff	684
	Bad service	281

Name: topic_name, dtype: int64



Topic modelling

Rating wise distribution of topics



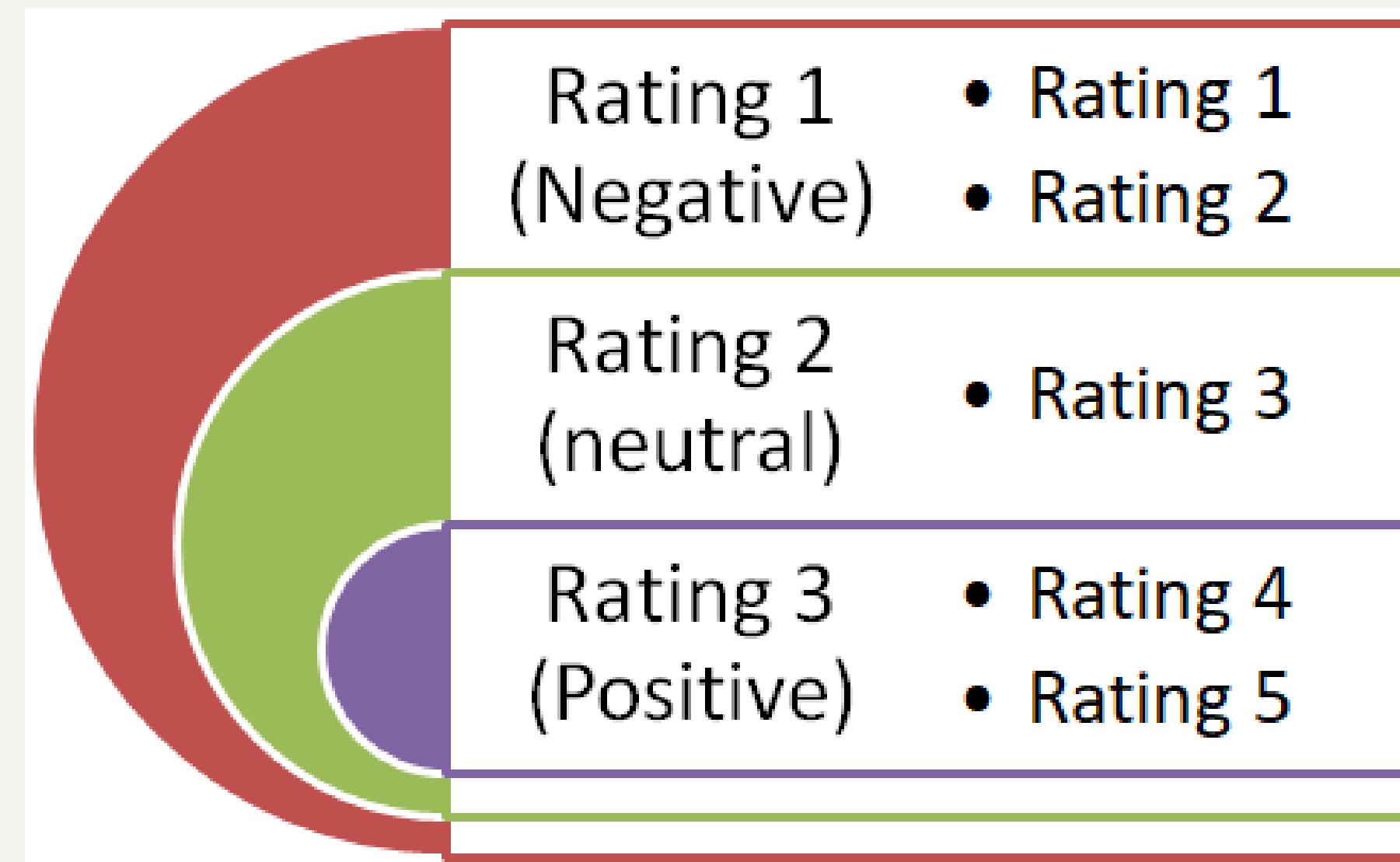
- Bad Service
- Hotel Facility
- Beach and Food
- Food bat hotel restaurant
- Staff



Rating Prediction



- 1 - Negative
- 2 - Neutral
- 3 - Positive





Rating prediction

For rating prediction take hotel rating as target variable

Build the Model

In this section, use several different models for our multiclass classification task.
There are a wide variety of techniques that can be used.

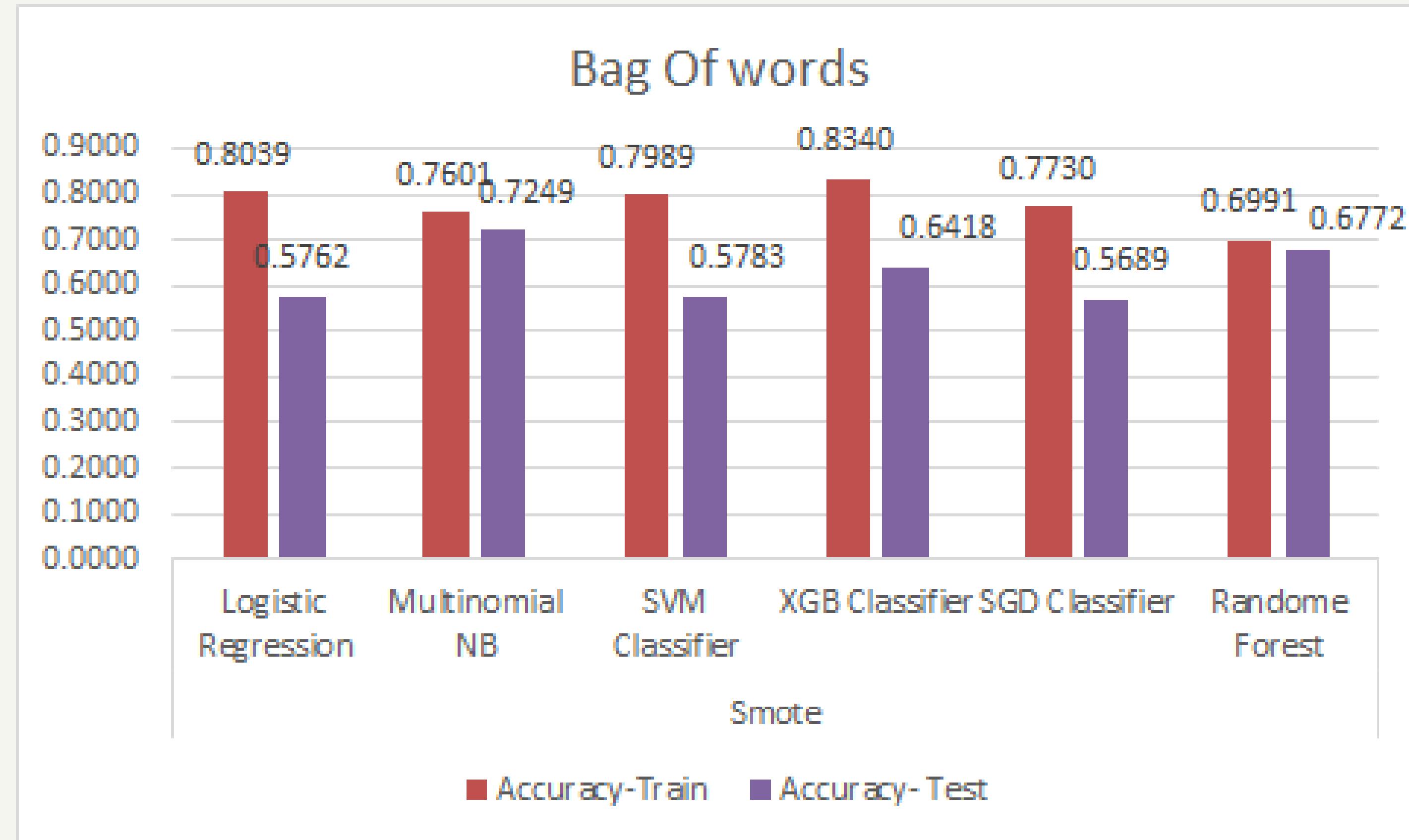
— 09

- Naive Bayes (NB)
- Logistic Regression
- Linear SVM
- Random Forest
- XGBoosting
- SGD classifier



Rating Prediction – imbalanced data

Using Bag of words approach



Logistic Regression is the best model



Classification Report of train and Test Data



	precision	recall	f1-score	support
1	0.61	0.78	0.68	435
2	0.36	0.28	0.32	488
3	0.91	0.90	0.91	2663
accuracy			0.80	3586
macro avg	0.63	0.66	0.64	3586
weighted avg	0.80	0.80	0.80	3586
0.8039598438371445				
LG classifier result:				
	precision	recall	f1-score	support
1	0.13	0.17	0.15	773
2	0.18	0.13	0.15	899
3	0.73	0.74	0.73	4476
accuracy			0.58	6148
macro avg	0.35	0.34	0.34	6148
weighted avg	0.57	0.58	0.57	6148
0.5762849707221861				
LG Classifier result:				
None				





Dealing with imbalanced data in classification

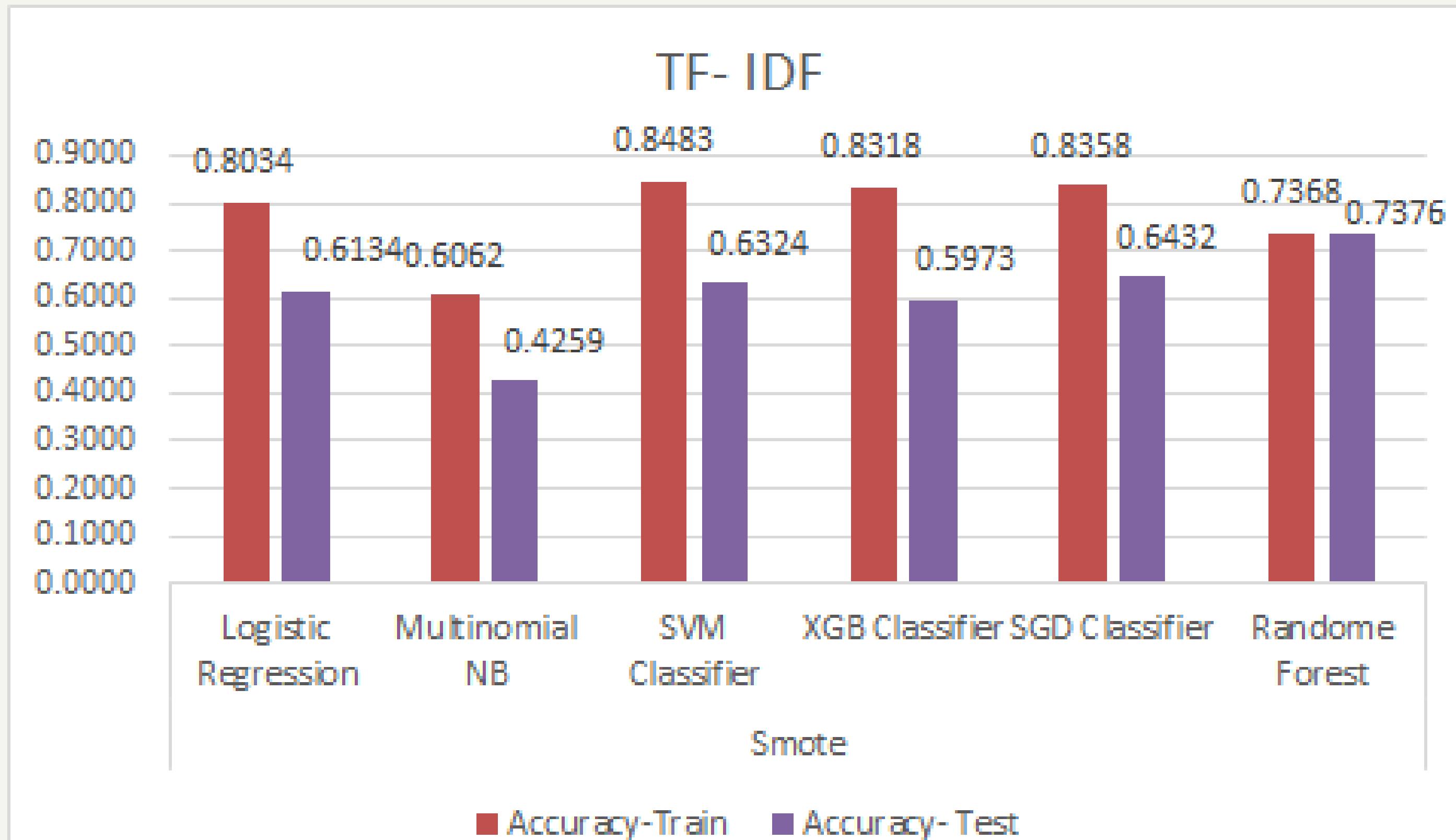


To deal with imbalanced data set using different resampling methods.

- Random oversampling
- SMOTE
- RandomUnderSampler
- TomekLinks(under sampling)
- ClusterCentroids(under sampling)
- SMOTETomek



Rating Prediction Balanced data - Smote Approach-Using TF_IDF



SVM Approach gives better accuracy in the combined rating prediction



Classification Report of SVM

	precision	recall	f1-score	support
1	0.68	0.82	0.74	464
2	0.20	0.49	0.29	160
3	0.98	0.87	0.92	2962
accuracy			0.85	3586
macro avg	0.62	0.73	0.65	3586
weighted avg	0.91	0.85	0.87	3586

0.8482989403234802

SVM Classifier result:
None

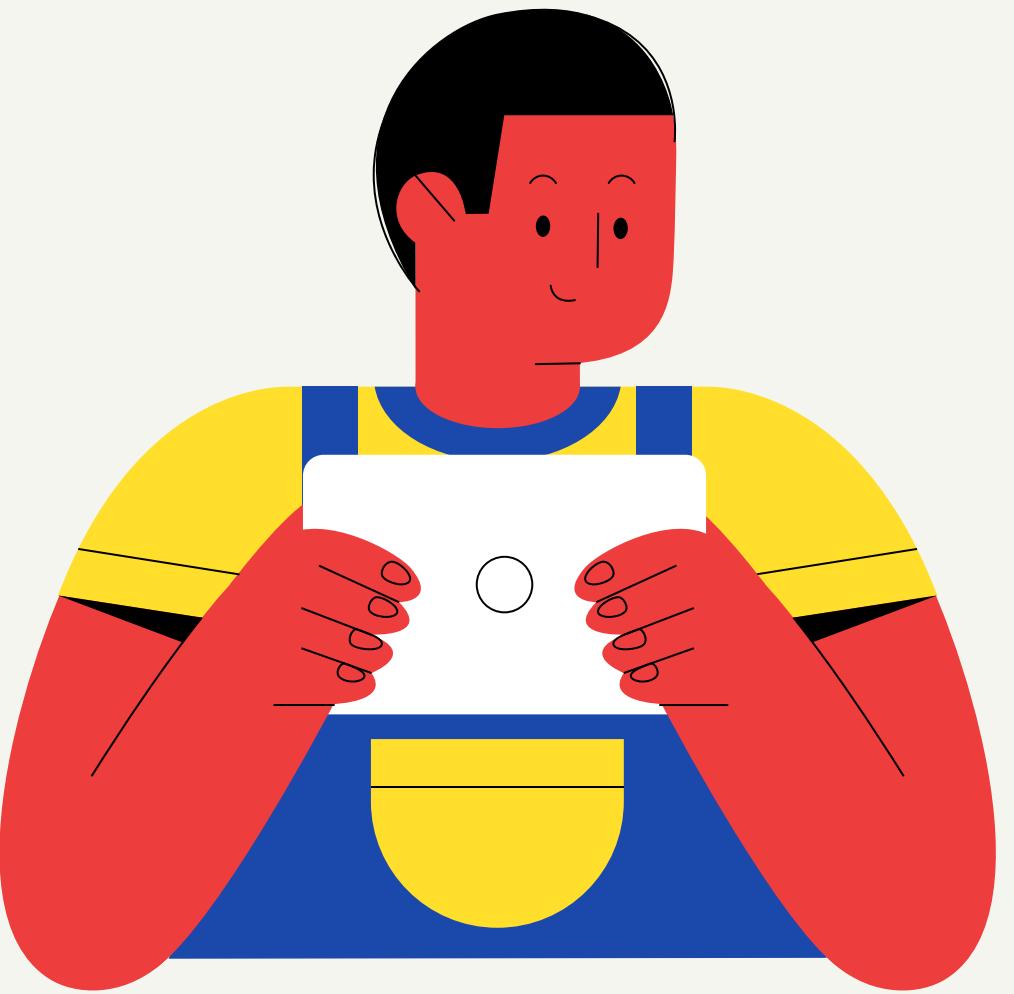
	precision	recall	f1-score	support
1	0.13	0.16	0.14	825
2	0.06	0.14	0.08	265
3	0.82	0.74	0.78	5058
accuracy			0.63	6148
macro avg	0.34	0.34	0.33	6148
weighted avg	0.70	0.63	0.66	6148

0.6322381262199089

SVM Classifier result:
None

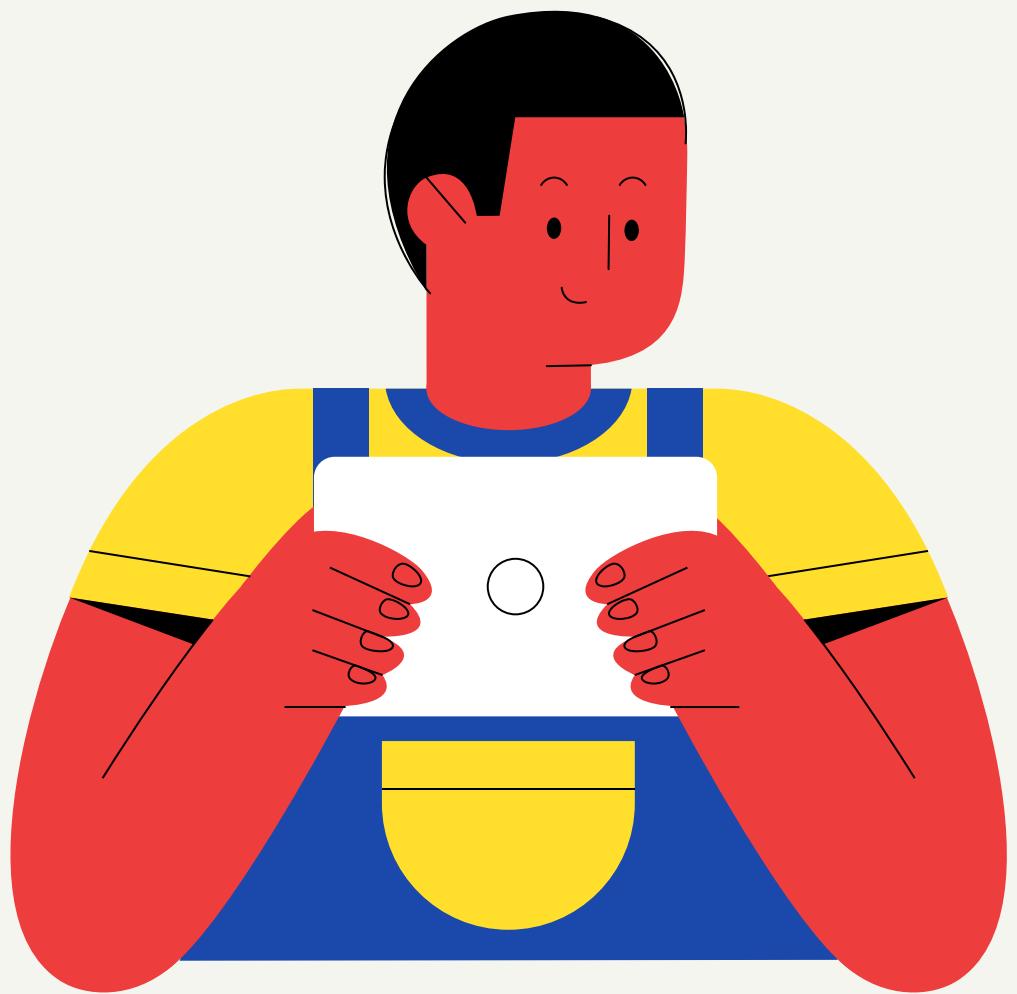
Best Model Information

- Use TF - IDF Approach
- Use Smote Sampling Technique
- Use SVM Classifier



DEPLOYMENT

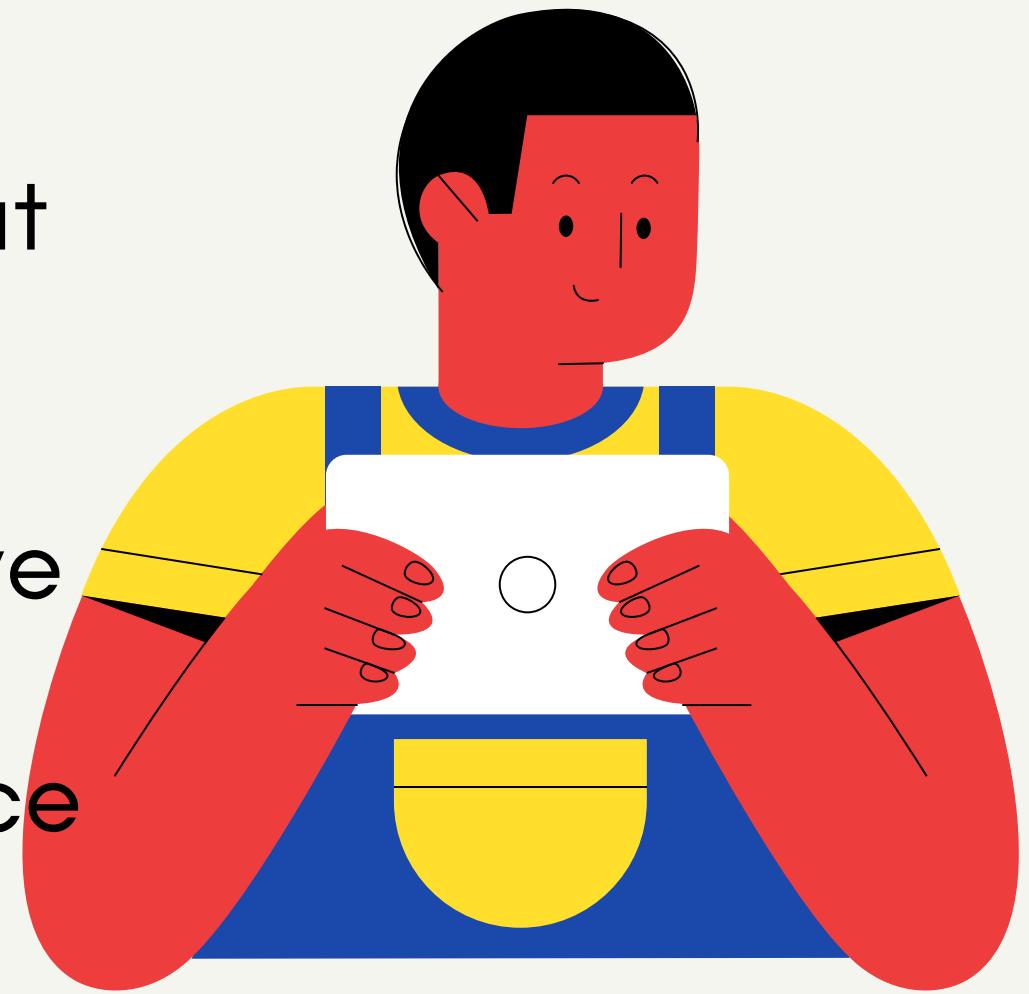
Here Model SVM is considered as best model and it is used for deployment



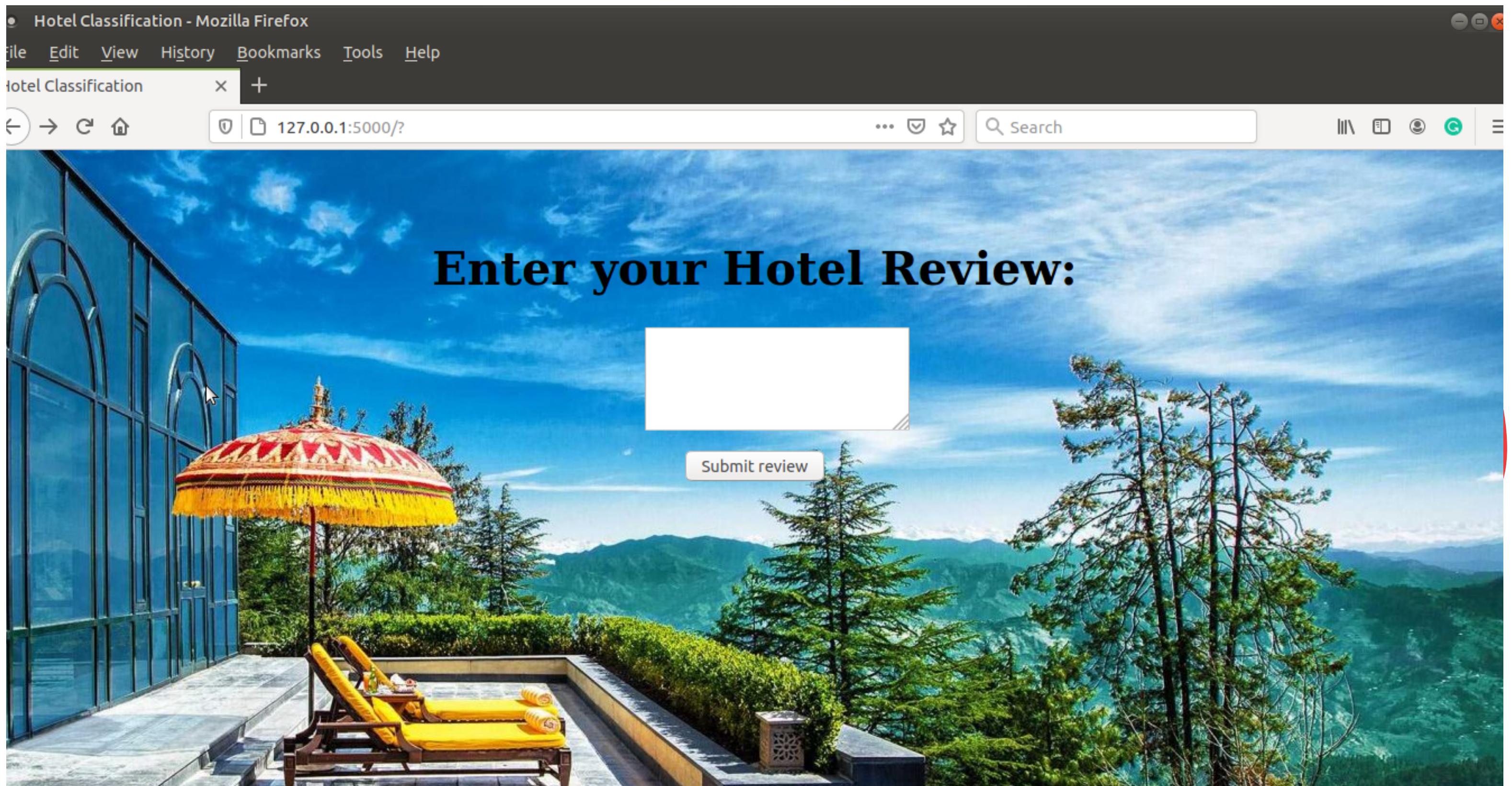
Deployment

We are Using Flask Method

- It makes the process of designing a web application simpler. Flask lets us focus on what the users are requesting and what sort of response to give back.
- To deploy our web application to the cloud, we will use Heroku
- heroku is an example of a Platform as a Service (PaaS)
- PaaS refers to the delivery of operating systems and associated services over the internet without downloads or installation

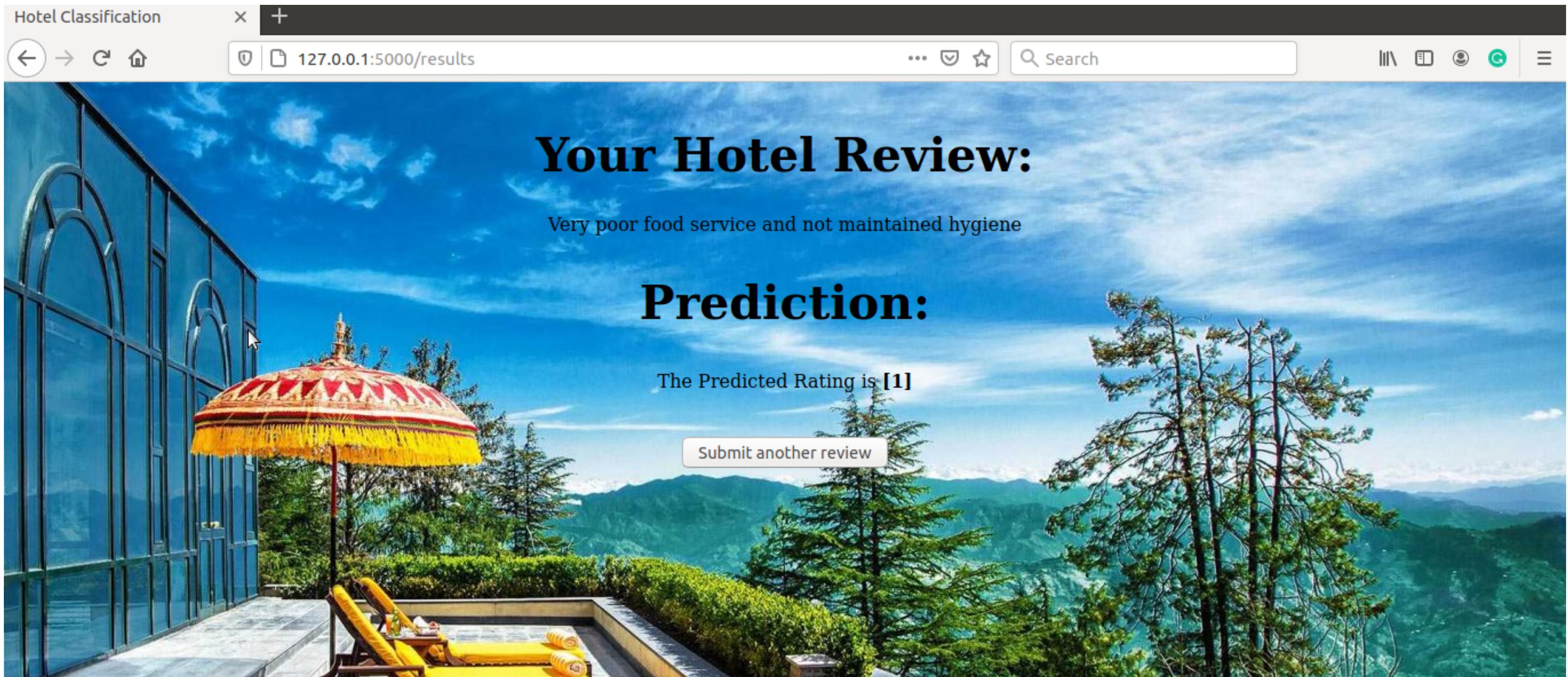


Deployment Screen shots-Index



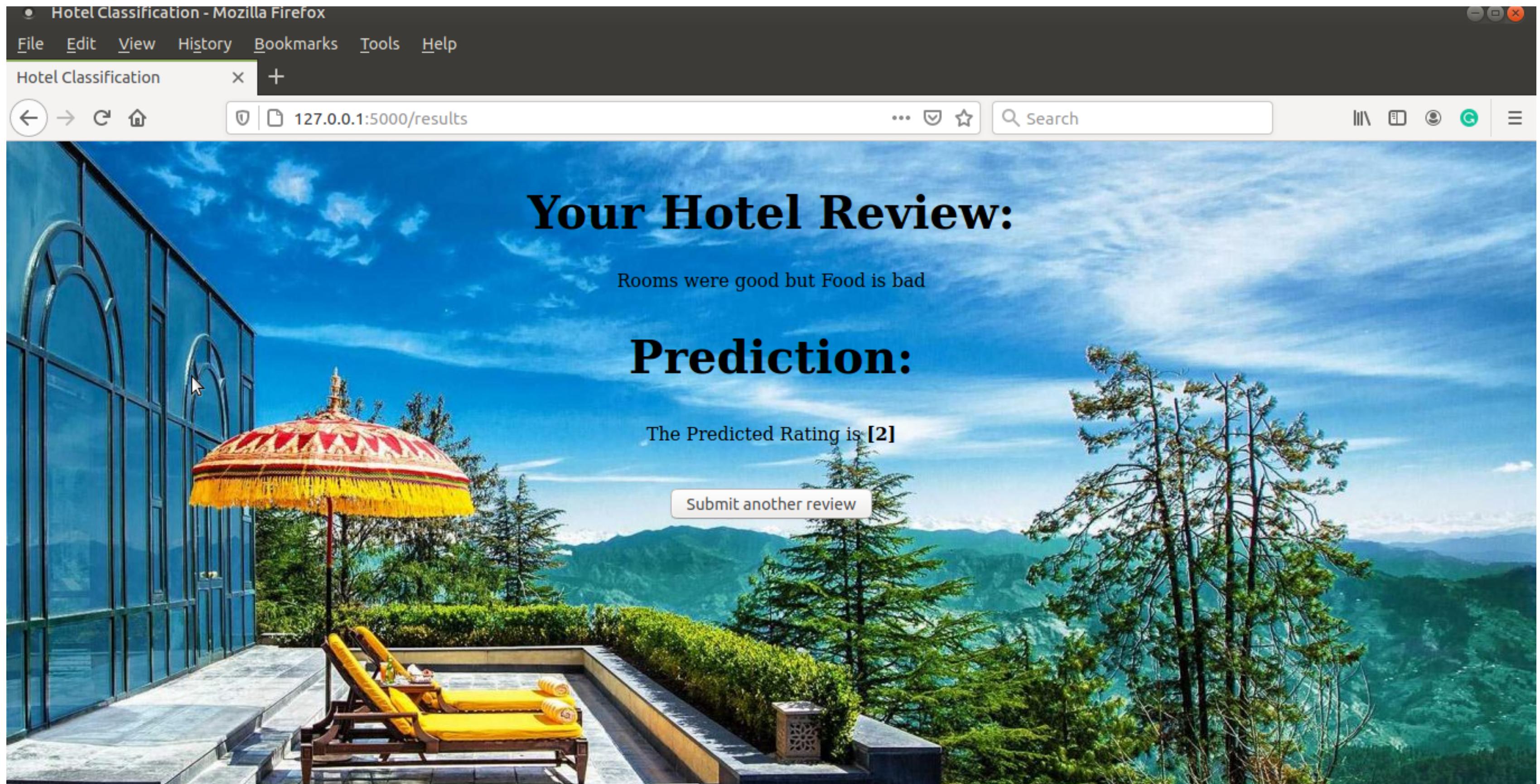
Deployment

Screen shots-Negative Rating

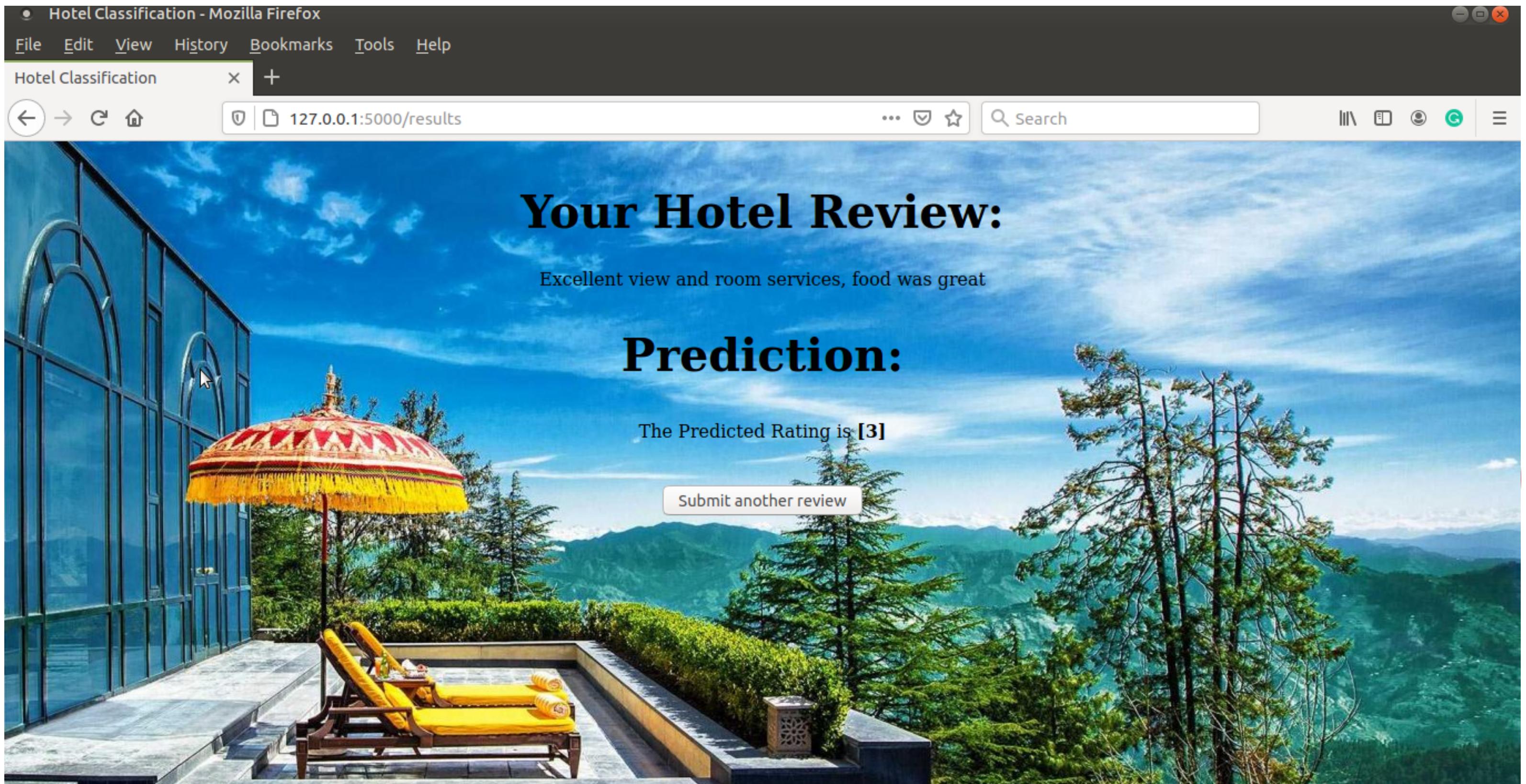


Deployment

Screen shots-Neutral Rating

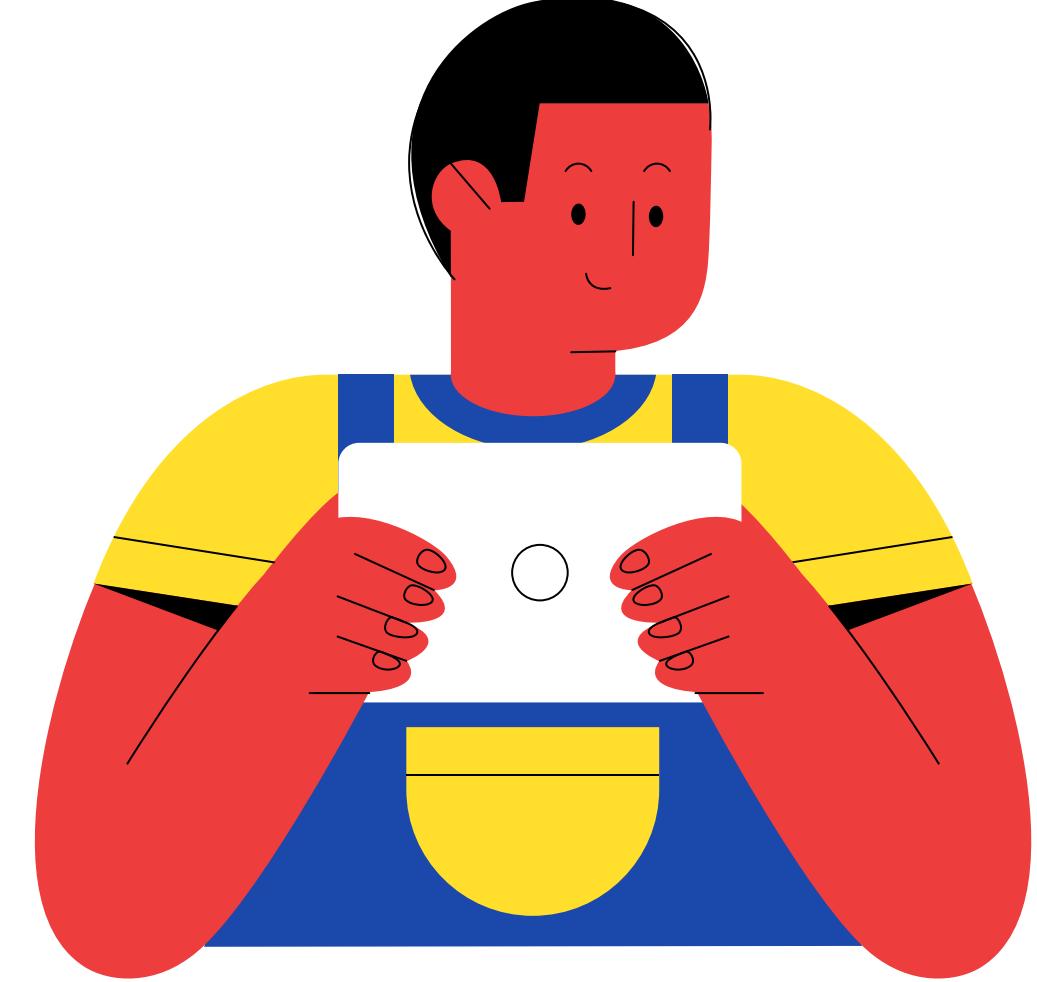


Deployment Screen shots-Positive Rating

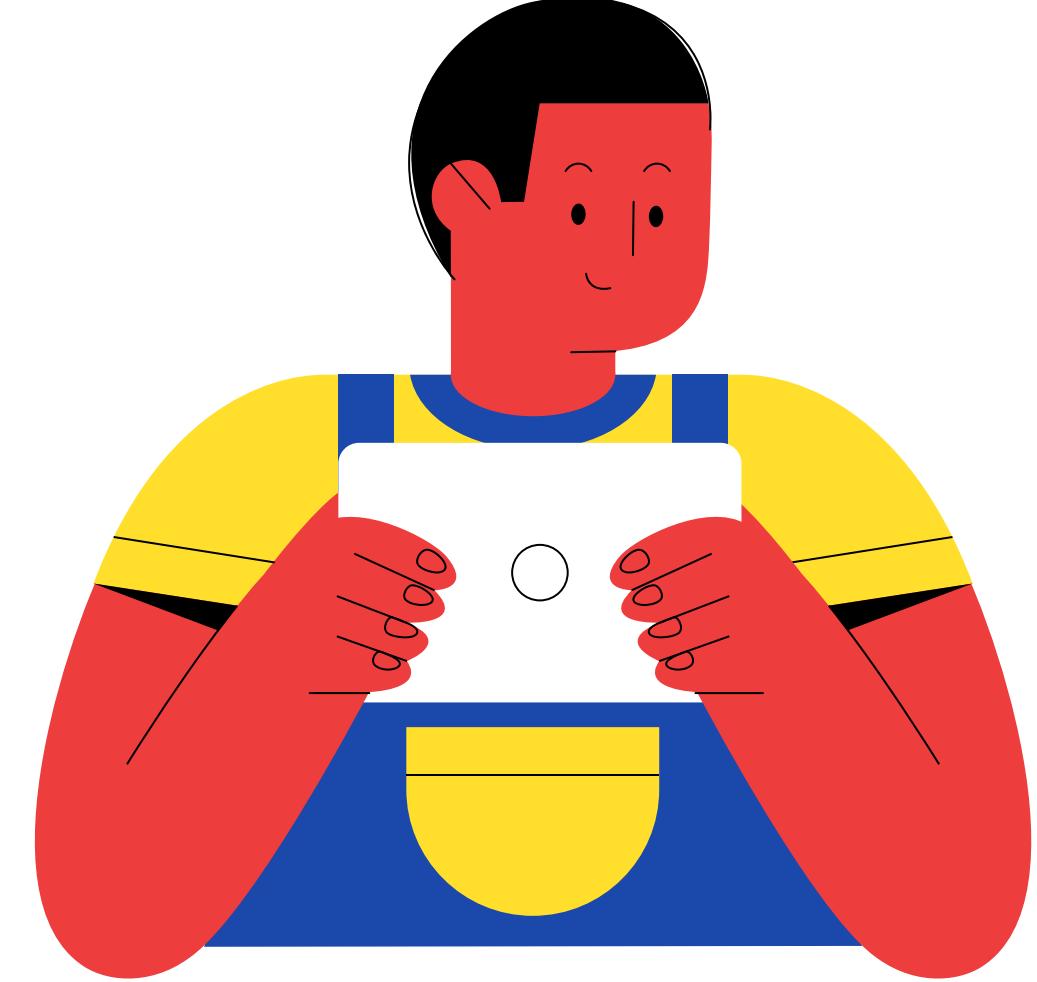


Challenges Faced

- Data Highly Imbalanced
- Preprocessed data.
- Provided few columns (No hotel names, Cities etc.)
- Has given both datasets separately ,
need to check accuracy for both sets separately



How did we overcome?



- Used ngrams and feature extraction to get relevant features from review text.
- Used various balancing techniques.
- Considered the review as positive, negative and neutral in the scale of 1-3 rating for increasing the predicting power(accuracy) of model.

