



INT 353 EDA PROJECT

BY RAJNISH BHARTI

REG. NO. - 12015883

ROLL NO. – B68

SECTION - K20RU





HOTELS ON MAKEMYTRIP

- DIL TO ROOMING HAI

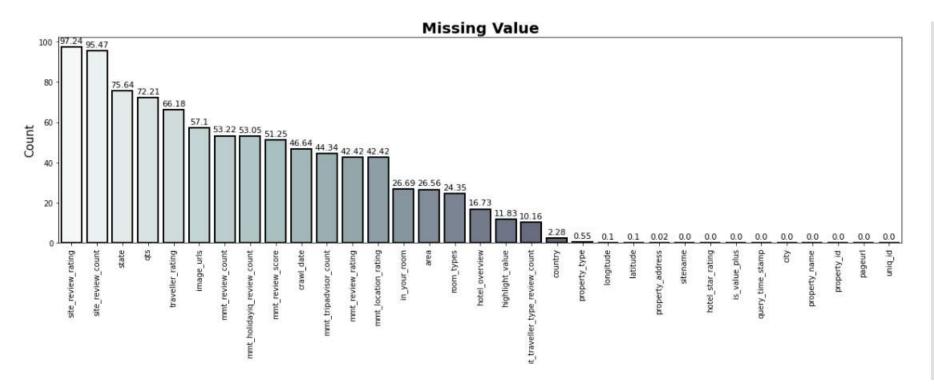
DATASET COLUMN INFO

• There are 20046 rows and 33 columns in the dataset

Here we can see that there are columns are of float type - 26 columns are of object type

#	Column	Non-Null Count	Dtype		
0	area	14722 non-null	object		
1	city	20046 non-null	object		
2	country	19588 non-null	object		
3	crawl date	10697 non-null	object		
4	highlight_value	17674 non-null	object		
5	hotel overview	16692 non-null	object		
6	hotel_star_rating	20046 non-null	object		
7	image_urls	8600 non-null	object		
8	in_your_room	14696 non-null	object		
9	is_value_plus	20046 non-null	object		
10	latitude	20025 non-null	float64		
11	longitude	20025 non-null	float64		
12	mmt_holidayiq_review_count	9412 non-null	float64		
13	mmt_location_rating	11543 non-null	object		
14	mmt review count	9378 non-null	float64		
15	mmt_review_rating	11543 non-null	object		
16	mmt_review_score	9772 non-null	float64		
17	mmt_traveller_type_review_count	18009 non-null	object		
18	mmt_tripadvisor_count	11158 non-null	float64		
19	pageurl	20046 non-null	object		
20	property_address	20042 non-null	object		
21	property_id	20046 non-null	object		
22	property_name	20046 non-null	object		
23	property_type	19936 non-null	object		
24	qts	5571 non-null	object		
25	query_time_stamp	20046 non-null	object		
26	room_types	15165 non-null	object		
27	site_review_count	908 non-null	object		
28	site_review_rating	554 non-null	float64		
29	sitename	20046 non-null	object		
30	state	4884 non-null	object		
31	traveller_rating	6780 non-null	object		
32	uniq_id	20046 non-null	object		
dtypes: float64(7), object(26)					

GRAPH FOR MISSING VALUES



- •Hai contains max 97.24% missing values.
- •We remove the column having null value more than 50 percent by taking threshold.
- •And for rest column we treat the missing values by each columns.



Here we observe that,

- area column having 5328 null values
- Country column having 458 null value
- And maximum null values in crawl_date column
- From here we can visualize the column one by one and clean the null values.
- Also we do further univariant and bivariant analysis

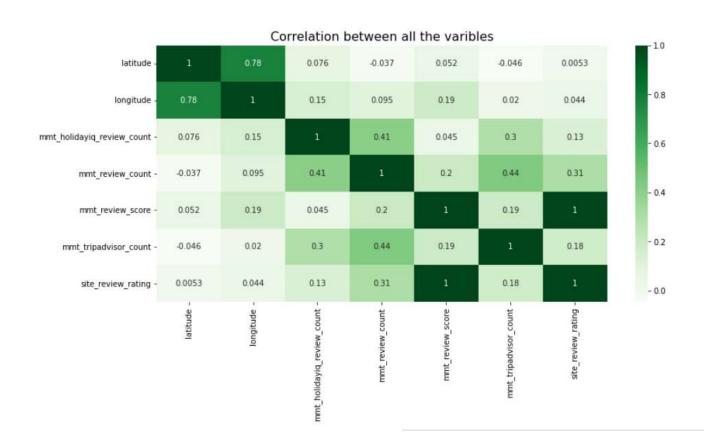
area	5324
city	0
country	458
crawl_date	9349
highlight_value	2372
hotel_overview	3354
hotel_star_rating	0
in_your_room	5350
is value plus	0
latitude	21
longitude	21
mmt_location_rating	8503
mmt_review_rating	8503
mmt_traveller_type_review_count	2037
mmt_tripadvisor_count	8888
pageurl	0
property_address	4
property_id	0
property_name	0
property_type	110
query_time_stamp	0
room_types	4881
sitename	0
uniq_id	0
dtype: int64	



CO- RELATION BETWEEN ALL COLUMNS

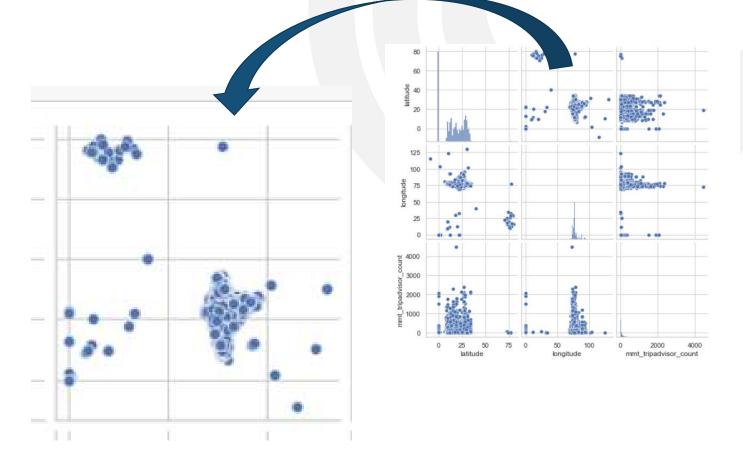


- Latitude and longitude as 0.78 co- related
- mmt_review_count and mmt_tripadvisor_count as 0.44
- Site_review_rating and mmt_review_score having most dependent.



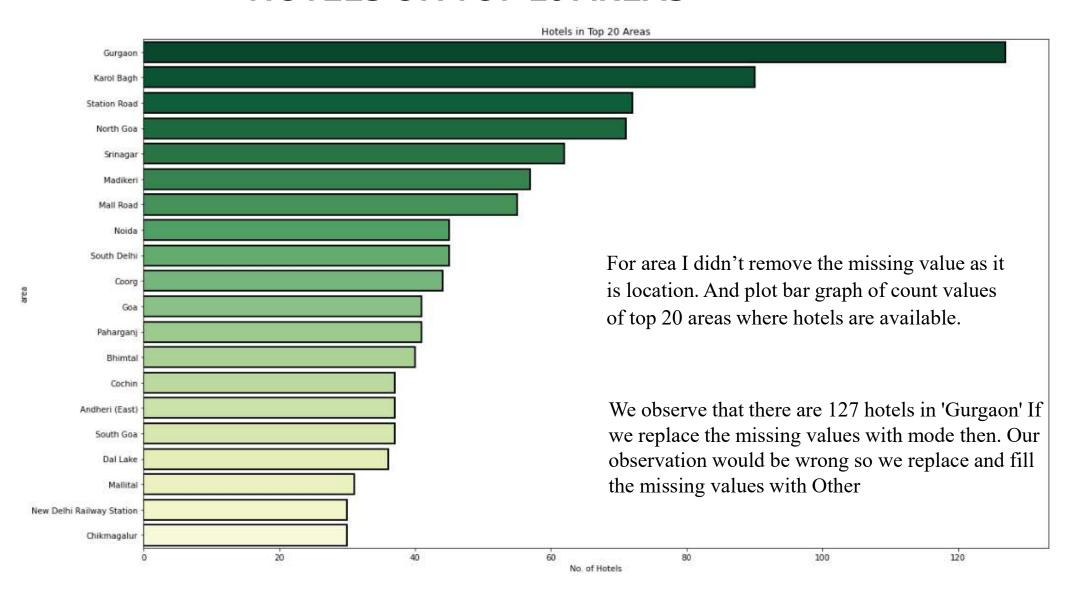


- Here I plot pair plot of MMT dataset. Where I get to know longitude and latitude are highly correlated
- Clusters showing no. of values present in the graph.
- So I decided to take both column and visualize further.

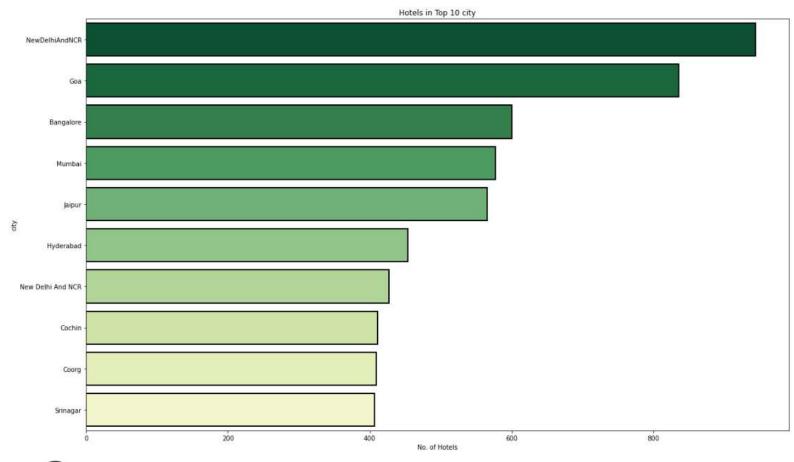




HOTELS ON TOP 20 AREAS



HOTELS ON TOP 10 CITIES



- NewDelhiAndNCR having max hotels 944
- Then Goa having 944 hotels



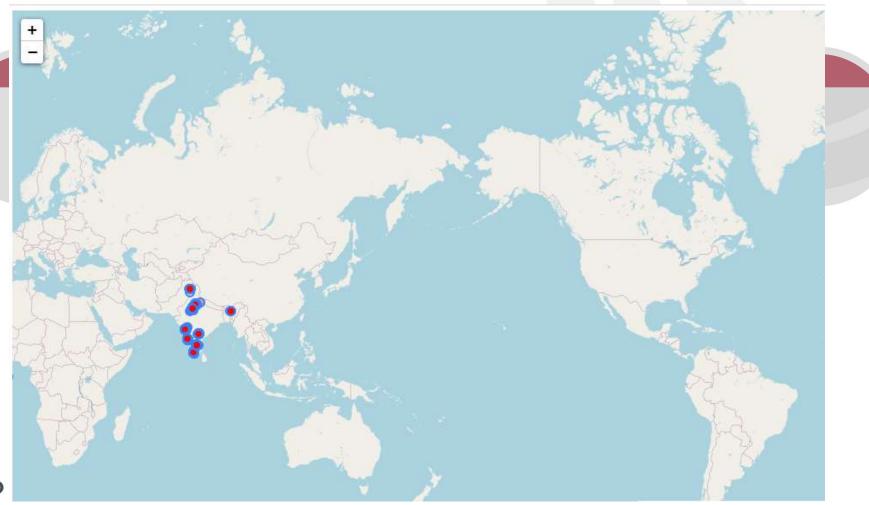
LETS VISUALIZE IN MAP WHERE HOTELS ARE PRESENT

- I had install folium library for plot map of world with the help of longitude and latitude columns
- I drop all the rows were null value are present
- Here we can find that some of the location are fetching out side the India.
- As Hotels on MMTdataset having only hotel related to India so we Drop that columns having outside location
- Some places are RMV Ext, Sanjaynagar, cochin, golden lake dal lake Srinagar Boulevard Road, Calangute, Boulevard Road, Dal Lake, MA Road, Besides Peddamma Temple, 201301, Dollars, Candolim, west Extension





CORRECT MAP OF HOTEL LOCATION ON MMT



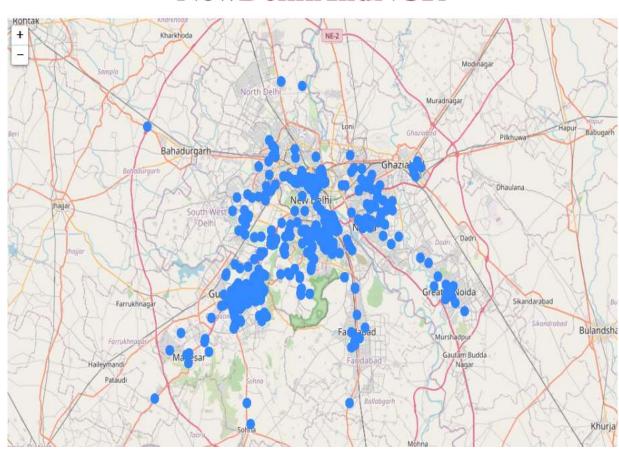


HOTEL LOCATION IN TOP CITY

- Mumbai

We observe that 2 hotels are out side Mumbai but its Area fetching to Mumbai so, I removed that.

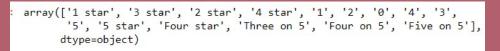
- NewDelhiAndNCR



HOTELS RATING BASED ON CITIES

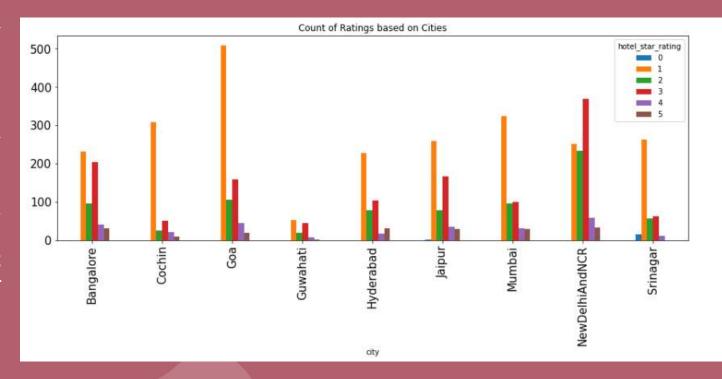


- We observe that hotel_star_rating are in the form of integer and string both. So we convert into integer.
- Also from above we get to know hotels in top 10 cities.
- Syntax: df['column_name']=df['column_name'].replace('value',1 to 5).astype(str)
- NewDelhiAndNCR having overall max hotels having 5 star.
- Goa having too much 1stars rating hotel. Here MMT can improve their facilities to grow their business.





1,2,3,4,5



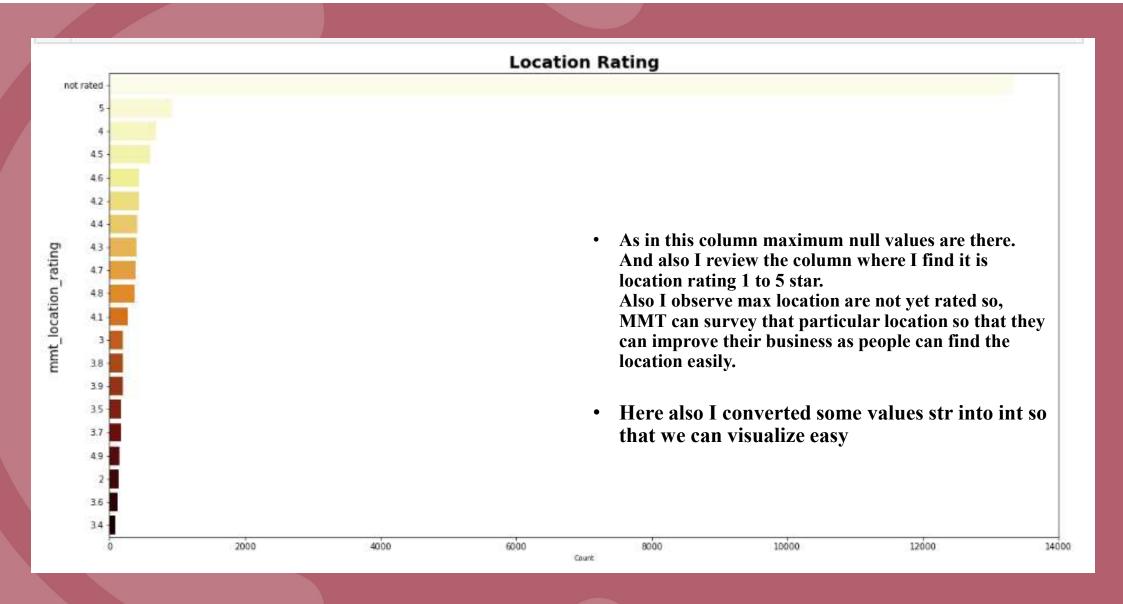


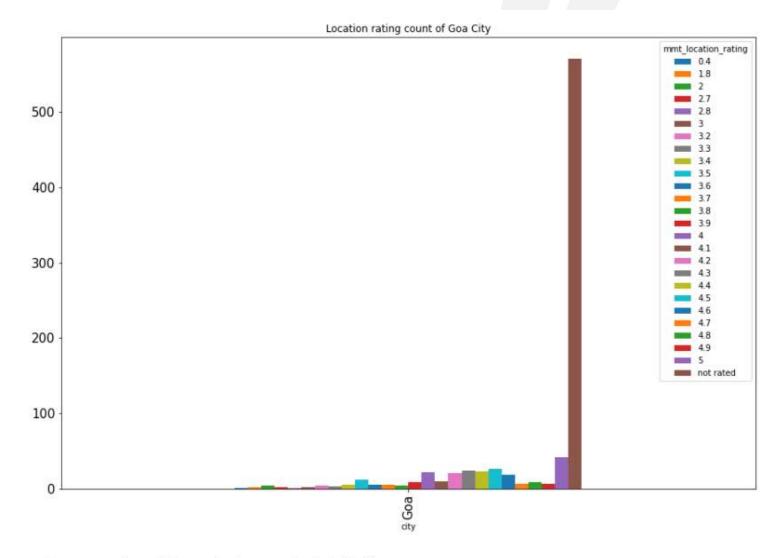
- **Crawl_date**: This columns having date format so replace the null values as mode. It just showing the no. of peoples are revisited the site. On 05-06-2016 max people are revisited
- In **highlight_value** column we observe that maximum row are filled with {{facilities}}. So I decided to fill the null values and replace {{facilities}} into the second facility max. i.e., **doctor on call.**
- In **Hotel_overview** column hotels overview is mentions as the max row of hotel_overview is ||less so in order to improve business MMT need to ask the hotel to mention their overview so that person can easily find the location and their needed hotels.
- In **column in_your_room**: The equipment's which are present in the rooms. As Here also "{{value}}" present in the max row so I replace it with needs of equipment present in a room.

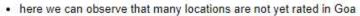
- In mmt_traveller-_type_review_count column there are description of type of visitors review of hotel.
- i.e., with family, friends, solo, Business, Couples.
- Here I decided to drop this column as this is of no use.

Cour	
483	amilies:{{ratingSummaryInfo.miscMap['family']}} Couples:{{ratingSummaryInfo.miscMap['couple']}} Business:{{ratingSummaryInfo.miscMap['business']}} Solo: {{ratingSummaryInfo.miscMap['solo']}} Friends:{{ratingSummaryInfo.miscMap['friends']}}
379	Families: Couples: Business: Solo: Friends:
219	Families:0 Couples:0 Business:0 Solo:0 Friends:0
97	Family:0 Couple:0 Solo:0 Friends:0 Business:0
30	Families:1 Couples:0 Business:0 Solo:0 Friends:0
	Family:7 Couple:7 Solo:2 Friends:0 Business:0
	Family:5 Couple:2 Solo:9 Friends:0 Business:0
	Family:4 Couple:1 Solo:5 Friends:1 Business:1
	Family:5 Couple:0 Solo:21 Friends:0 Business:0
	Families:160 Couples:41 Business:2 Solo:5 Friends:2









- after that 5 rated place are max
- . to improve business in Goa all the locations are to be rate so that people can visit in that location and book thier room



Room Types



- Maximum Hotels are having Standard Room, Non Ac room, Deluxe Room, Ac Room.
- Here nan showing the null values so we replace that as Other.

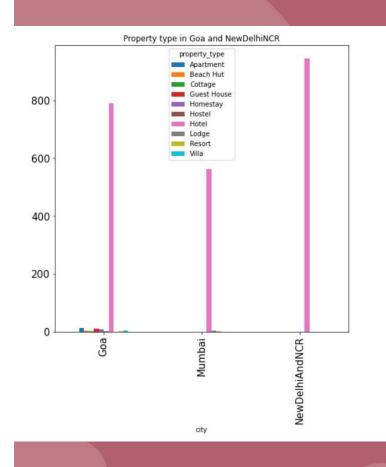


HOTELS PROPERTY TYPE



Count

	Count
Hotel	19587
Lodge	199
Homestay	28
Guest House	28
Houseboat	22
Apartment	16
Cottage	15
Resort	14
Camp	13
Villa	6
Beach Hut	4
Palace	2
Tree house	1
Hostel	1



WE CAN OBSERVE HERE THAT

- MAXIMUM PROPERTY ARE "HOTEL"
- ONLY ONE PROPERTY IS HOSTEL
- NEWDELHIANDNCR HAVING ALL PROPERTY TYPE ARE "HOTELS"

- We observe that after all analysis. There are some null values in latitude and longitude. Here we already drop the rows where null value are present. And observer the location of hotels in map.
- In property_address there are 4 null values. So we drop that row having address.
- By the help of "isna" we find all four column i.e., 686, 2896,15035, 19964.
- Or MMT ask that hotel address In order to customers will more book that hotel too.

