

OLA RIDES ANALYSIS PROJECT

- Project Overview

This project analyses OLA ride data with specific focus on ****ride cancellation patterns****, identifying root causes, and providing data-driven solutions to reduce cancellation rates and improve overall service efficiency.

- Dataset Summary

The dataset contains over 20,000+ rows and 19 columns. It includes data related to ride bookings, vehicle types, payment methods, customer and driver ratings, ride distances, and cancellation reasons.

- Data Cleaning And Standardization Using Python

- **Data Loading** : Imported Dataset Using Pandas
- **Initial Exploration** : Used `df.info()` to check structure and `.describe()` for summary statistics.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20407 entries, 0 to 20406
Data columns (total 20 columns):
 #   Column                                Non-Null Count  Dtype  
---  -
 0   Date                                20407 non-null  object 
 1   Time                                20407 non-null  object 
 2   Booking_ID                           20407 non-null  object 
 3   Booking_Status                       20407 non-null  object 
 4   Customer_ID                         20407 non-null  object 
 5   Vehicle_Type                         20407 non-null  object 
 6   Pickup_Location                     20407 non-null  object 
 7   Drop_Location                       20407 non-null  object 
 8   V_TAT                               12652 non-null  float64 
 9   C_TAT                               12652 non-null  float64 
10   Canceled_Rides_by_Customer           2081 non-null   object 
11   Canceled_Rides_by_Driver             3654 non-null   object 
12   Incomplete_Rides                     12652 non-null  object 
...
25%      3.500000      3.500000      NaN
50%      4.000000      4.000000      NaN
75%      4.500000      4.500000      NaN
max       5.000000      5.000000      NaN
```

- **Standardization Of Column Names & Removing ALL Whitespaces** : using pandas string methods (.str.lower() , .str.replace() , .str.strip()) .

```
Index(['date', 'time', 'booking_id', 'booking_status', 'customer_id',
      'vehicle_type', 'pickup_location', 'drop_location', 'v_tat', 'c_tat',
      'canceled_rides_by_customer', 'canceled_rides_by_driver',
      'incomplete_rides', 'incomplete_rides_reason', 'booking_value',
      'payment_method', 'ride_distance', 'driver_ratings', 'customer_rating',
      'vehicle_images'],
      dtype='object')
```

- **Data Consistency Check** : Checked the datatypes of all numerical columns and changed for the required data type .
- **Database Integration**: Connected Python Script to MYSQL Workbench and located the cleaned dataset into the database for further inspection.

```
# Connecting to mysql
!pip install pymysql sqlalchemy
!-m pip install --upgrade pip
```

```
1 from sqlalchemy import create_engine
2 import mysql.connector
3 # Setting up connection with MYSQL
4 username = "root"
5 password = "Secure$4u"
6 host = "localhost"
7 port = "3306"
8 database = "OLA_RIDES_ANALYSIS"
9
10 engine = create_engine(f"mysql+pymysql://{username}:{password}@{host}:{port}/{database}")
11
```

```
1 # Write DataFrame to MySQL ( Loaded this dataframe under table name "TABLE1")
2 table_name = "TABLE1" # choose any table name
3 df.to_sql(table_name, engine , if_exists="replace", index=False)
```

• Data Analysis & Finding Insights Using Python

We performed structured analysis in Python (pandas) to answer key business questions:

- Retrieve all successful bookings :

```
1 ## TOTAL SUCCESSFUL or UNSUCCESSFUL BOOKINGS
2 var=df[df['booking_status']=="Success"]
3 print(len(var))
```

12652

- Find the average ride distance for each vehicle type:

```
1 ## FINDING AVERAGE DISTANCE RIDE FOR EACH VEHICLE ( let's do it for prime_suv)
2 var1=df[df['vehicle_type']=="prime_suv"]
3 var2=var1['ride_distance'].mean()
4 print(var2)
```

15.187096774193549

- Get the total number of cancelled rides by customers:

```
1 ## CANCELLED RIDES BY THE CUSTOMERS or drivers
2 var3=df[df['canceled_rides_by_customer'].notna()]
3 print(len(var3))
```

2081

- List the top 5 customers who booked the highest number of rides:

```
1 ## TOP 5 CUSTOMERS WHO BOOKED HIGHEST NUMBER OF RIDES or rating wise or booking value wise ( just adjust the filter condition & see the
2 var4 =df.groupby('customer_id').agg(number_of_bookings=('booking_id','count')).sort_values('number_of_bookings',ascending=False).head(
3 print(var4)
```

customer_id	number_of_bookings
CID940408	3
CID393964	2
CID143850	2
CID190536	2
CID126952	2

- Get the number of rides cancelled by drivers due to personal and car-related issues:

```
1 ## Number of rides cancelled by drivers due to personal and car-related issues or any other reasons as
2 var5 =df[df['canceled_rides_by_driver']=="Personal & Car related issue"]
3 print(len(var5))
```

1263

- Find the maximum and minimum driver ratings for Prime Sedan bookings:

```
1 ## MAX & MIN DRIVER RATING ON PRIME SUV or in any vehicle
2 #print(df['driver_ratings'].dtype) # firstly checking the datatype of this column so that numerical op
3 var6= df[df['vehicle_type']=="prime_suv"]
4 #print(var6[['vehicle_type','driver_ratings']])
5 print(var6['driver_ratings'].min())
6 print(var6['driver_ratings'].max())
```

✓ 0.0s

3.0
5.0

- Retrieve all rides where payment was made using UPI:

```

1 ## RIDES WHERE PAYMENT METHOD IS UPI or any other
2 #First cleaing the strings of payment methoid column
3 df['payment_method']=df['payment_method'].str.lower()
4 df['payment_method']=df['payment_method'].str.replace(' ','_')
5 df['payment_method']=df['payment_method'].str.strip()
6 print(df['payment_method'])
7 var7= df[df['payment_method']=="upi"]
8 print(len(var7))

```

Python

```

0      NaN
1     cash
2      upi
3      NaN
4  credit_card
...
20402    NaN
20403    cash
20404    cash
20405    upi
20406    NaN
Name: payment_method, Length: 20407, dtype: object
5113

```

- Find the average customer rating per vehicle type:

```

1 ## Average customer rating or driver ratings per vehicle type ( let's take mini vehicle here)
2 #print(df['customer_rating'].dtype) # Checking the Data type of cyustomer RATING COLUMN
3 var8= df[df['vehicle_type']=="mini"]
4 var9=var8['customer_rating'].mean()
5 print(var9)

```

✓ 0.0s Python

4.019822320932816

- Calculate the total booking value of rides completed successfully:

```

1 ## TOTAL value OF rides COMPLETED successfully or unsuccessfully
2 print(df['incomplete_rides'].unique()) # checking the different values inside the incomplete_rides column
3 var10= df[df['incomplete_rides']=="No"]
4 var11= var10['booking_value'].sum()
5 print(var11)

```

Python

```

[nan 'No' 'Yes']
6429897

```

- List all incomplete rides along with the reason:

```

1 ## ALL incomplete rides along with reason :
2 #print(df[['incomplete_rides_reason','incomplete_rides']].tail(30)) # just doing basic inspection
3 var12=df[df['incomplete_rides']=="Yes"]
4 #print(var12[['booking_id','customer_id','incomplete_rides','incomplete_rides_reason']])
5 print(var12[['booking_id','customer_id','incomplete_rides','incomplete_rides_reason']])

```

Python

	booking_id	customer_id	incomplete_rides	incomplete_rides_reason
38	CNR5176704322	CID296026	Yes	Customer Demand
49	CNR9312632867	CID649563	Yes	Vehicle Breakdown
70	CNR7924302885	CID517661	Yes	Customer Demand
93	CNR1640228587	CID190281	Yes	Other Issue
101	CNR7623690602	CID526261	Yes	Other Issue
...
20197	CNR3500429121	CID108258	Yes	Other Issue
20275	CNR9654855291	CID923733	Yes	Other Issue
20289	CNR6348178557	CID948352	Yes	Vehicle Breakdown
20378	CNR7902243107	CID514958	Yes	Customer Demand
20390	CNR7459607546	CID845808	Yes	Customer Demand

[795 rows x 4 columns]

- Dashboard In Power BI

Finally, we built an interactive dashboard in Power BI to present insights visually.

OLA

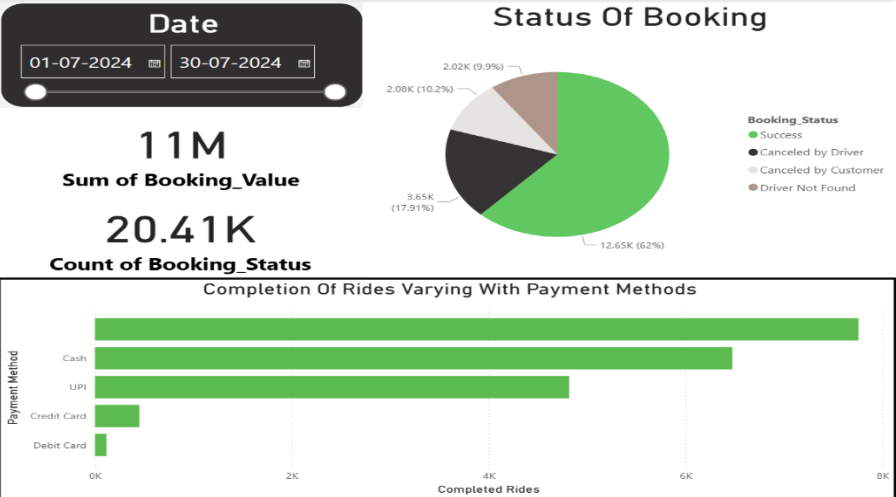
Overall

Vehicle Type

Revenue

Cancellation

Ratings



OLA

Overall

Vehicle Type

Revenue

Cancellation

Ratings

01-07-202430-07-2024

Vehicle Type	Total Booking Value	Success Booking Value	Avg. Distance Travelled	Total Distance Travelled
<div>Prime Sedan</div>	1.67M	1.06M	15.27	44.97K
<div>Prime SUV</div>	1.58M	962.25K	15.19	44.73K
<div>Prime Plus</div>	1.54M	934.88K	14.87	41.46K
<div>Mini</div>	1.57M	973.65K	15.72	45.06K
<div>Auto</div>	1.59M	992.78K	6.21	18.12K
<div>Bike</div>	1.58M	985.07K	16.16	47.64K
<div>E-Bike</div>	1.62M	994.76K	15.68	46.99K

OLA

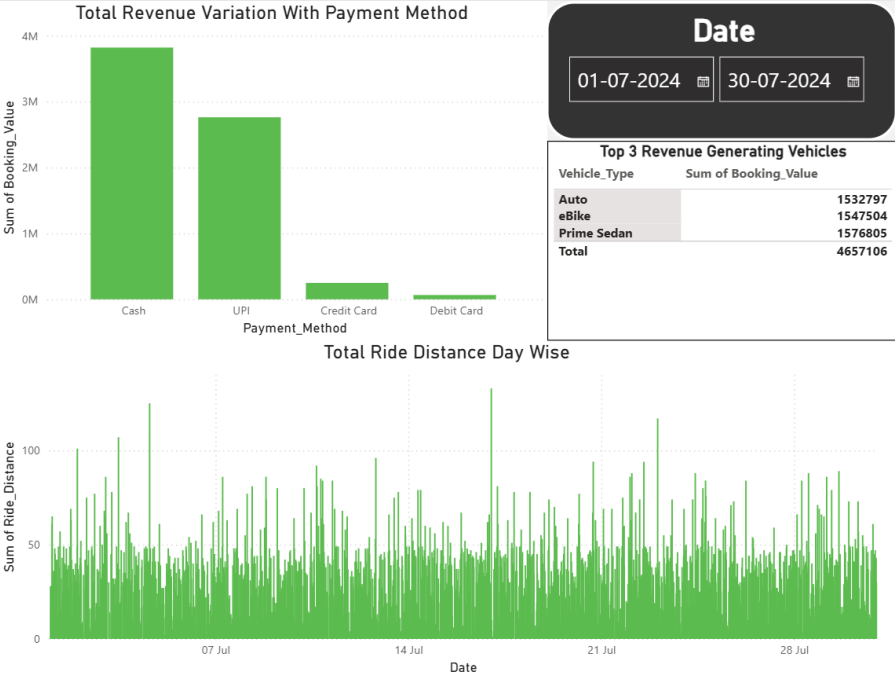
Overall

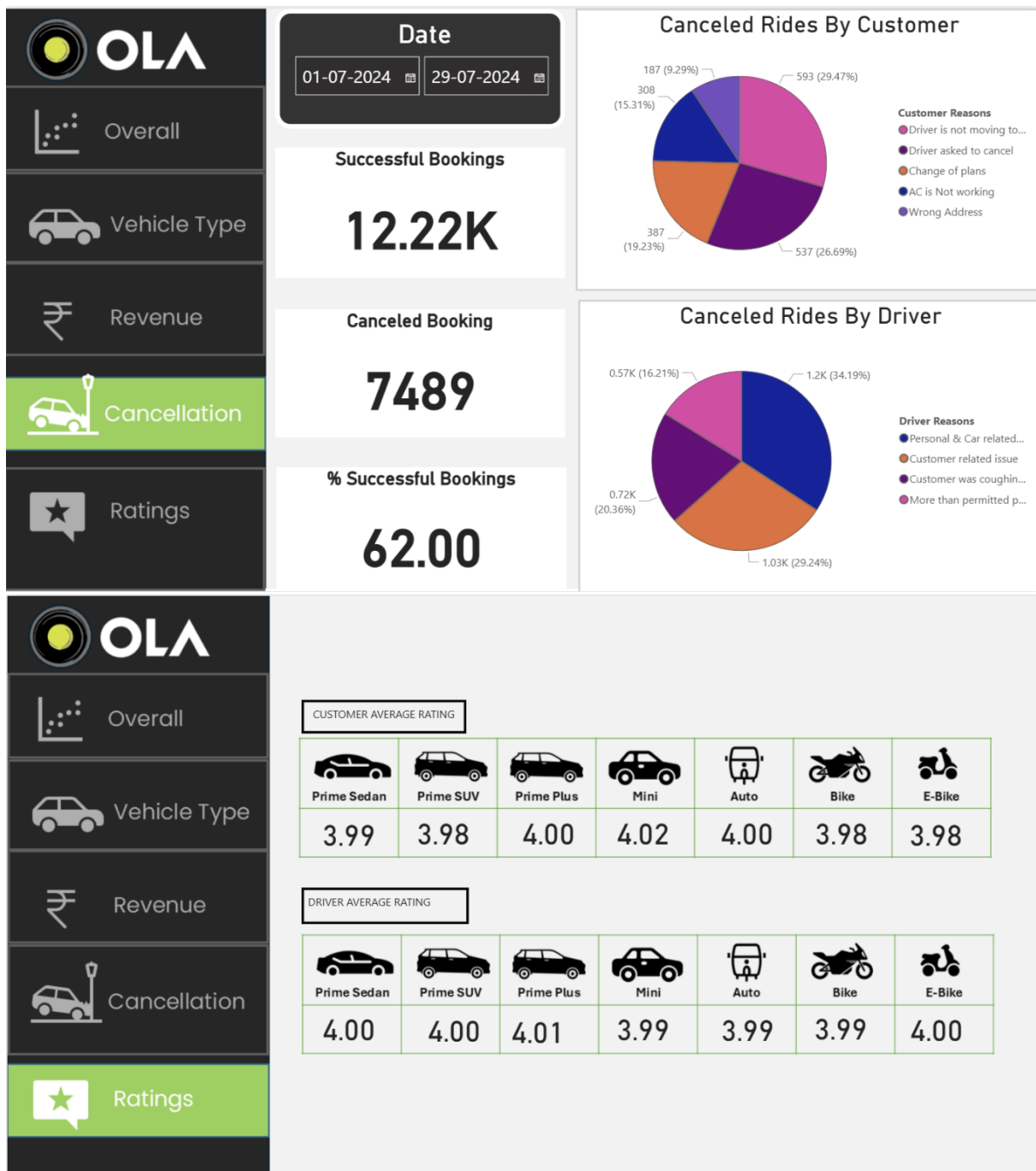
Vehicle Type

Revenue

Cancellation

Ratings





• Business Insights

We found several high impactful insights which upon implementation can significantly reduce the cancellation rate and will increase overall efficiency .

- Payment Method Impact on Cancellations** : UPI payments show significantly lower cancellation rates compared to cash payments
- Vehicle Wise Cancellation Analysis** : Premium vehicles have higher cancellation rates due to availability issues

3. Cancellation Reason Analysis :

Customer Side : Driver is not moving (29.47%) & Change of plans (26.69%)

Driver Side : Personal & Car related (34.19%) & Customer related issue (29.24%)

4. Rating Analysis :

Customer Side : Customers are happy with mini vehicles

Driver Side : Drivers are happy with premium vehicles (like prime SUV)

● Business Recommendations

- 1. Charge a fee for last-minute “change of plans “ cancellations : this reason makes up over a quarter (26.69 %) of customer cancellations.**
- 2. Fix the “Driver not moving issue “ : it’s the biggest reason (29.47 %) of customer canceled rides.**
- 3. Help drivers to keep their cars in good condition : “personal & car related issue “ is the major reason (34.19 &) for driver canceled rides.**
- 4. Make a strict rule against drivers asking customers to cancel : this bad practice is responsible for (15.31%) of customer canceled rides.**
- 5. Give better cars and training to prime drivers : these have generally the lower ratings in comparison to the mini type vehicles.**