# Data Collection and Preprocessing Phase
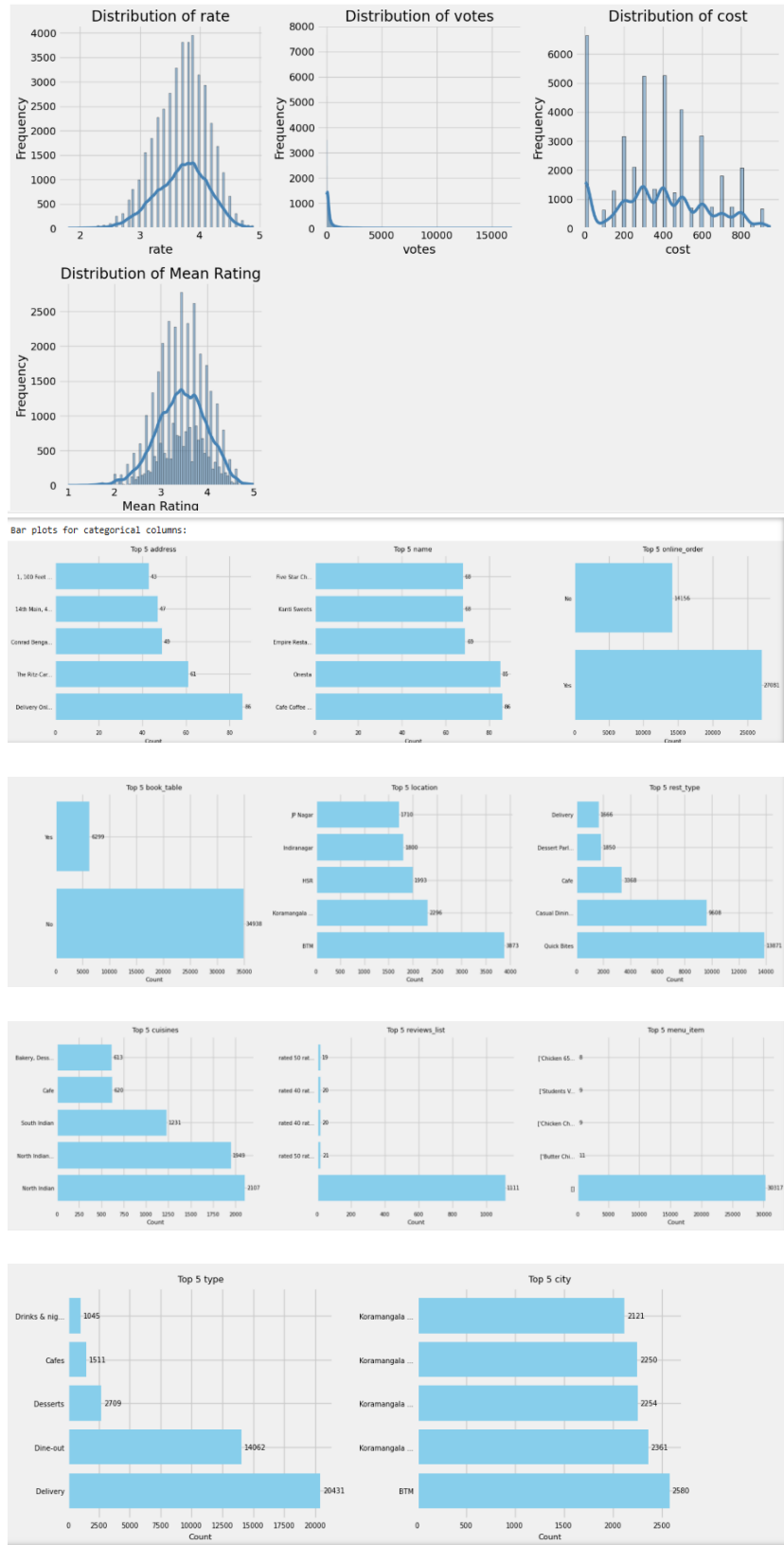
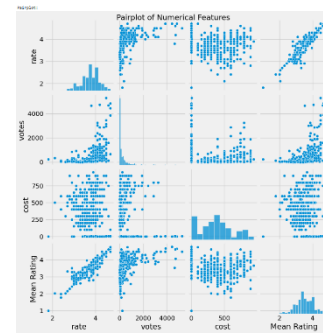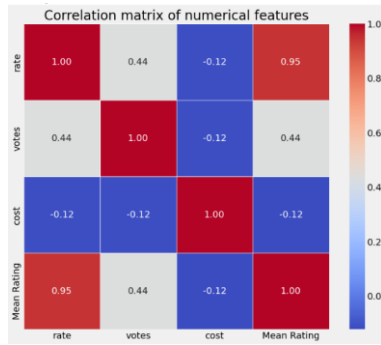| | |
|---|---|
| Date | 18 June 2025 |
| Team ID | xxxxxx |
| Project Title | Restaurant Recommendation System |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Report**

Data set variables will be statistically analyzed to identify patterns and outliers, with Python Employed for pre-processing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

| Section | Description |
|---|---|
| Data Overview | **Dimensions:**<br>51717 rows x 17 columns<br>**Descriptive statistics:**<br><br>**votes**<br>count 51717.000000<br>mean 283.697527<br>std 803.838853<br>min 0.000000<br>25% 7.000000<br>50% 41.000000<br>75% 198.000000<br>max 16832.000000 |

Univariate Analysis



**Distribution of rate**

**Distribution of votes**

**Distribution of cost**

**Distribution of Mean Rating**

Bar plots for categorical columns:

Top 5 address

Top 5 name

Top 5 online_order

Top 5 book_table

Top 5 location

Top 5 rest_type

Top 5 cuisines

Top 5 reviews_list

Top 5 menu_item

Top 5 type

Top 5 city

Bivariate Analysis


Correlation matrix of numerical features


Pairplot of Numerical Features

**Numerical vs Categorical Analysis:**


Rating vs Online Order


Rating vs City

**Categorical vs Categorical Analysis:**


Online order availability vs type of meal

Heatmap of Table booking option vs location



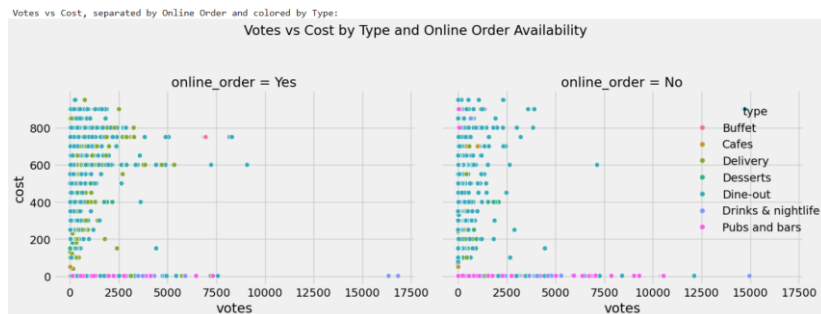restaurant type vs type of meal

**Multivariate Analysis**



Cost by Rest Type, broken down by Online Order:

Cost by Rest Type and Online Order Availability



Votes vs Cost, separated by Online Order and colored by Type:

Votes vs Cost by Type and Online Order Availability

| | |
|---|---|
| Outliers and Anomalies | <br>Outlier Detection (using Box Plots)<br>Box Plot of rate, Box Plot of votes, Box Plot of cost |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data |  |
| Handling Missing Data | ```python<br># Drop unnecessary columns<br>zomato_df = zomato_df.drop(['phone', 'dish_liked', 'url'], axis=1, errors='ignore')<br><br># Remove rows with missing values<br>zomato_df.dropna(how='any', inplace=True)<br>``` |
| Data Transformation | ```python<br># Rename columns for clarity<br>zomato_df = zomato_df.rename(columns={<br>    'approx_cost(for two people)': 'cost',<br>    'listed_in(type)': 'type',<br>    'listed_in(city)': 'city'<br>})<br><br># Filter out rows where 'rate' is 'NEW' or '-' (corrected typo: 'NEM' → 'NEW')<br>zomato_df = zomato_df[~zomato_df['rate'].isin(['NEW', '-'])].reset_index(drop=True)<br>```<br><br>```python<br>#We will be using the 'Review' and 'Cuisines' feature in order to create a recommender system<br>## Lower Casing<br>zomato_df['reviews_list'] = zomato_df['reviews_list'].str.lower()<br><br>## Removal of Punctuations<br>import string<br>PUNCT_TO_REMOVE = string.punctuation<br>def remove_punctuation(text):<br>    """custom function to remove the punctuation"""<br>    return text.translate(str.maketrans('', '', PUNCT_TO_REMOVE))<br>zomato_df['reviews_list'] = zomato_df['reviews_list'].apply(lambda text: remove_punctuation(text))<br>``` |

| Feature Engineering | ```python
# Clean and convert 'rate' to float (remove '/5' and whitespace)
zomato_df['rate'] = (
    zomato_df['rate']
    .str.replace('/5', '')   # Remove '/5' if present
    .str.strip()             # Remove whitespace
    .astype(float)           # Convert to float
)

# Clean and convert 'cost' to float (handle commas/decimals)
zomato_df['cost'] = (
    zomato_df['cost']
    .astype(str)
    .str.replace(',', '.')   # Replace commas with periods (decimal separator)
    .astype(float)
)
```

prepares **reviews_list** to be used as a textual feature for TF-IDF vectorization

```python
#We will be using the 'Review' and 'Cuisines' feature in order to create a recommender system
## Lower Casing
zomato_df['reviews_list'] = zomato_df['reviews_list'].str.lower()

## Removal of Punctuations
import string
PUNCT_TO_REMOVE = string.punctuation
def remove_punctuation(text):
    """custom function to remove the punctuation"""
    return text.translate(str.maketrans('', '', PUNCT_TO_REMOVE))
zomato_df['reviews_list'] = zomato_df['reviews_list'].apply(lambda text: remove_punctuation(text))
``` |
| Save Processed Data | - |