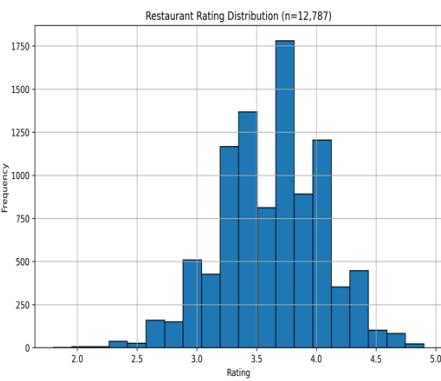
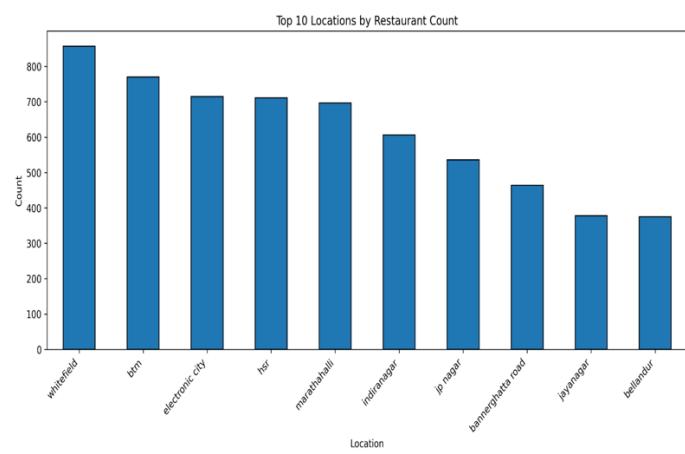
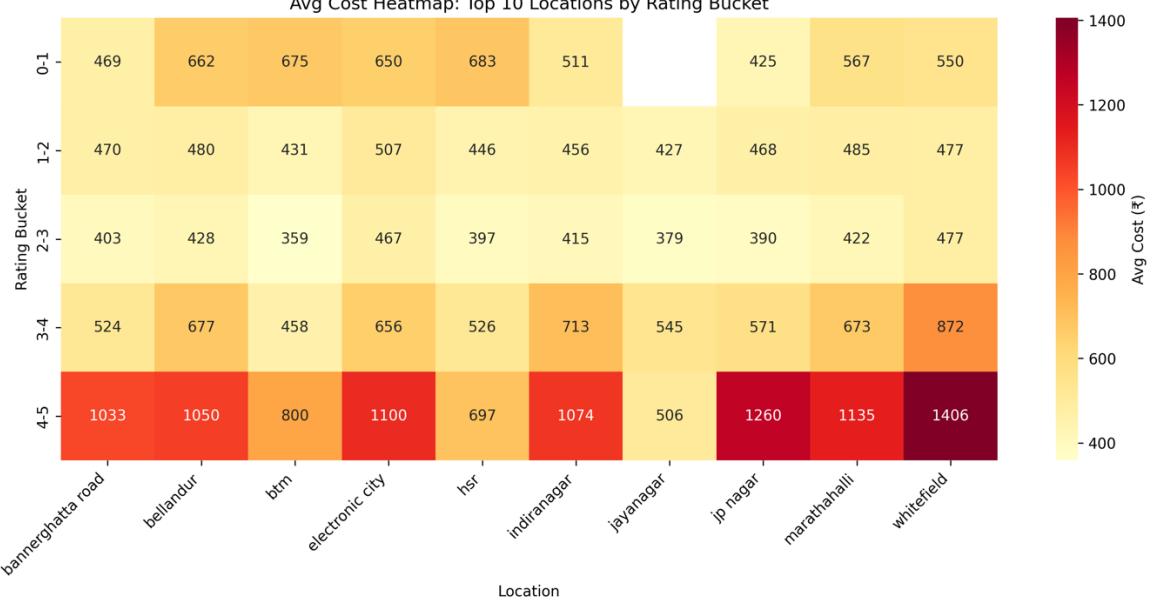


Data Collection and Preprocessing Phase

Date	3 Dec 2025
Team ID	
Project Title	Restaurant Recommendation System
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

Section	Description																																																																																																																														
Data Overview	<p>📊 DESCRIPTIVE STATISTICS TABLE - Dataset Shape: 51,717 rows × 17 columns</p> <hr/> <table> <thead> <tr> <th>Feature</th> <th>Data Type</th> <th>Count</th> <th>Missing</th> <th>Mean</th> <th>Median</th> <th>Std Dev</th> </tr> </thead> <tbody> <tr> <td>votes</td> <td>int64</td> <td>51717</td> <td>0</td> <td>283.7</td> <td>41.0</td> <td>803.84</td> </tr> <tr> <td>url</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>address</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>name</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>online_order</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>book_table</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>rate</td> <td>Categorical</td> <td>51717</td> <td>7775</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>phone</td> <td>Categorical</td> <td>51717</td> <td>1208</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>location</td> <td>Categorical</td> <td>51717</td> <td>21</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>rest_type</td> <td>Categorical</td> <td>51717</td> <td>227</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>dish_liked</td> <td>Categorical</td> <td>51717</td> <td>28078</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>cuisines</td> <td>Categorical</td> <td>51717</td> <td>45</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>approx_cost(for two people)</td> <td>Categorical</td> <td>51717</td> <td>346</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>reviews_list</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>menu_item</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>listed_in(type)</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>listed_in(city)</td> <td>Categorical</td> <td>51717</td> <td>0</td> <td>-</td> <td>-</td> <td>-</td> </tr> </tbody> </table> <p>TOTAL MISSING VALUES: 37700</p>	Feature	Data Type	Count	Missing	Mean	Median	Std Dev	votes	int64	51717	0	283.7	41.0	803.84	url	Categorical	51717	0	-	-	-	address	Categorical	51717	0	-	-	-	name	Categorical	51717	0	-	-	-	online_order	Categorical	51717	0	-	-	-	book_table	Categorical	51717	0	-	-	-	rate	Categorical	51717	7775	-	-	-	phone	Categorical	51717	1208	-	-	-	location	Categorical	51717	21	-	-	-	rest_type	Categorical	51717	227	-	-	-	dish_liked	Categorical	51717	28078	-	-	-	cuisines	Categorical	51717	45	-	-	-	approx_cost(for two people)	Categorical	51717	346	-	-	-	reviews_list	Categorical	51717	0	-	-	-	menu_item	Categorical	51717	0	-	-	-	listed_in(type)	Categorical	51717	0	-	-	-	listed_in(city)	Categorical	51717	0	-	-	-
Feature	Data Type	Count	Missing	Mean	Median	Std Dev																																																																																																																									
votes	int64	51717	0	283.7	41.0	803.84																																																																																																																									
url	Categorical	51717	0	-	-	-																																																																																																																									
address	Categorical	51717	0	-	-	-																																																																																																																									
name	Categorical	51717	0	-	-	-																																																																																																																									
online_order	Categorical	51717	0	-	-	-																																																																																																																									
book_table	Categorical	51717	0	-	-	-																																																																																																																									
rate	Categorical	51717	7775	-	-	-																																																																																																																									
phone	Categorical	51717	1208	-	-	-																																																																																																																									
location	Categorical	51717	21	-	-	-																																																																																																																									
rest_type	Categorical	51717	227	-	-	-																																																																																																																									
dish_liked	Categorical	51717	28078	-	-	-																																																																																																																									
cuisines	Categorical	51717	45	-	-	-																																																																																																																									
approx_cost(for two people)	Categorical	51717	346	-	-	-																																																																																																																									
reviews_list	Categorical	51717	0	-	-	-																																																																																																																									
menu_item	Categorical	51717	0	-	-	-																																																																																																																									
listed_in(type)	Categorical	51717	0	-	-	-																																																																																																																									
listed_in(city)	Categorical	51717	0	-	-	-																																																																																																																									

Univariate Analysis	 																																																																																																
Bivariate Analysis	 <table border="1"> <thead> <tr> <th colspan="2"></th> <th colspan="10">Rating Bucket</th> </tr> <tr> <th colspan="2"></th> <th>0-1</th> <th>1-2</th> <th>2-3</th> <th>3-4</th> <th>4-5</th> <th>0-1</th> <th>1-2</th> <th>2-3</th> <th>3-4</th> <th>4-5</th> </tr> <tr> <th colspan="2"></th> <th>469</th> <th>662</th> <th>675</th> <th>650</th> <th>683</th> <th>511</th> <th>425</th> <th>567</th> <th>550</th> <th>1033</th> </tr> </thead> <tbody> <tr> <th colspan="2"></th> <td>470</td> <td>480</td> <td>431</td> <td>507</td> <td>446</td> <td>456</td> <td>427</td> <td>468</td> <td>485</td> <td>1050</td> </tr> <tr> <th colspan="2"></th> <td>403</td> <td>428</td> <td>359</td> <td>467</td> <td>397</td> <td>415</td> <td>379</td> <td>390</td> <td>422</td> <td>677</td> </tr> <tr> <th colspan="2"></th> <td>524</td> <td>677</td> <td>458</td> <td>656</td> <td>526</td> <td>713</td> <td>545</td> <td>571</td> <td>673</td> <td>872</td> </tr> <tr> <th colspan="2"></th> <td>1033</td> <td>1050</td> <td>800</td> <td>1100</td> <td>697</td> <td>1074</td> <td>506</td> <td>1260</td> <td>1135</td> <td>1406</td> </tr> <tr> <th colspan="2"></th> <th>bannerghatta road</th> <th>bellandur</th> <th>btm</th> <th>electronic city</th> <th>hsr</th> <th>indiranagar</th> <th>jayanagar</th> <th>jp nagar</th> <th>marathahalli</th> <th>whitefield</th> </tr> </tbody> </table>			Rating Bucket												0-1	1-2	2-3	3-4	4-5	0-1	1-2	2-3	3-4	4-5			469	662	675	650	683	511	425	567	550	1033			470	480	431	507	446	456	427	468	485	1050			403	428	359	467	397	415	379	390	422	677			524	677	458	656	526	713	545	571	673	872			1033	1050	800	1100	697	1074	506	1260	1135	1406			bannerghatta road	bellandur	btm	electronic city	hsr	indiranagar	jayanagar	jp nagar	marathahalli	whitefield
		Rating Bucket																																																																																															
		0-1	1-2	2-3	3-4	4-5	0-1	1-2	2-3	3-4	4-5																																																																																						
		469	662	675	650	683	511	425	567	550	1033																																																																																						
		470	480	431	507	446	456	427	468	485	1050																																																																																						
		403	428	359	467	397	415	379	390	422	677																																																																																						
		524	677	458	656	526	713	545	571	673	872																																																																																						
		1033	1050	800	1100	697	1074	506	1260	1135	1406																																																																																						
		bannerghatta road	bellandur	btm	electronic city	hsr	indiranagar	jayanagar	jp nagar	marathahalli	whitefield																																																																																						
Outliers and Anomalies	Rate >5 or <1 (clipped); cost >10k (z-score >3, capped); anomalies: duplicate addresses (removed).																																																																																																
Data Preprocessing Code Screenshots																																																																																																	
Loading Data	<pre> raw_path = Path("/Users/rohit/Desktop/swayam_project/DATASET/zomato.csv") clean_path = Path("/Users/rohit/Desktop/swayam_project/DATASET/zomato_clean.csv") # Phase 2: EDA print(" PHASE 2: DATA ANALYSIS") df_raw = pd.read_csv(raw_path) analyze_data(df_raw) </pre>																																																																																																

Handling Missing Data	<pre>df.drop_duplicates(inplace=True); df['rate'].fillna(df['rate'].median(), inplace=True); df.dropna(subset=['cuisines'], inplace=True)</pre>
Data Transformation	<pre># Fill minor missing values df["rest_type"] = df["rest_type"].fillna("Casual Dining") df = df.dropna(subset=["cuisines", "location"]) # Critical for recommender # Text cleaning for TF-IDF df["cuisines"] = df["cuisines"].str.lower().str.strip() df["location"] = df["location"].str.lower().str.strip() df["name"] = df["name"].str.strip()</pre>
Feature Engineering	<pre>def create_features(self): """TF-IDF vectorization on cuisines + location + name.""" print("\nCreating TF-IDF features...") # Combine features into "soup" text self.df["soup"] = (self.df["cuisines"].fillna("") + " " + self.df["location"].fillna("") + " " + self.df["rest_type"].fillna("") + " " + self.df["name"].fillna("")) # TF-IDF (core algorithm) self.tfidf = TfidfVectorizer(max_features=5000, stop_words="english", ngram_range=(1, 2), lowercase=True)</pre>
Save Processed Data	<pre>def save_model(self, model_path: str = "../model/tfidf_model.pkl"): """Save trained model for Flask.""" Path(model_path).parent.mkdir(exist_ok=True) with open(model_path, "wb") as f: pickle.dump(self, f) print(f"Model saved: {model_path}")</pre>