	 Table Of Content:- 1- Import libraries 2- Read the dataset file 3- Show a brief of the dataframe 4- Data Exploration 5- Data cleaning 6- Exploratory Data Analysis EDA 7- Conclusion
In [2]:	1- import libraries:- # to make Data Manpulation import pandas as pd import numpy as np # to make Data Visuallization import matplotlib.pyplot as plt
In [3]:	<pre>import seaborn as sns # to ignore the Errors import warnings warnings.filterwarnings("ignore") 2- Read the dataset file:- # the Dataframe we will call as data data = pd.read_csv(r*C:\Users\del1\Desktop\powerbi projects\datasets\uber_eatscsv\uber_eat_rest.csv*)</pre>
In [4]:	3- Show abrief of dataframe:- # show top 5 row with all columns data.head() id position name score ratings category price_range full_address zip_code lat lng 0 1 19 PJ Fresh (224 Daniel Payne Drive) NaN NaN Burgers, American, Sandwiches \$ 224 Daniel Payne Drive, Birmingham, AL, 35207 35.06 36.830703
In [86]:	1 2 9 J'ti'z Smoothie-N-Coffee Bar NaN NaN Coffee and Tea, Breakfast and Brunch, Bubble Tea NaN 1521 Pinson Valley Parkway, Birmingham, AL, 35217 33.583640 -86.773330 2 3 6 Philly Fresh Cheesesteaks (541-B Graymont Ave) NaN NaN American, Cheesesteak, Sandwiches, Alcohol \$ 541-B Graymont Ave, Birmingham, AL, 35204 33.509800 -86.854640 3 4 17 Papa Murphy's (1580 Montgomery Highway) NaN NaN Pizza \$ 1580 Montgomery Highway, Hoover, AL, 35226 35.26 33.404439 -86.806614 4 5 162 Nelson Brothers Cafe (17th St N) 4.7 22.0 Breakfast and Brunch, Burgers, Sandwiches NaN 314 17th St N, Birmingham, AL, 35203 33.514730 -86.811700 # show last 5 row with all columns data.head() id position name score ratings category price_range full_address zip_code lat Ing
In [5]: Out[5]: In [88]:	1 19 PJ Fresh (224 Daniel Payne Drive) NaN NaN Burgers, American, Sandwiches \$ 224 Daniel Payne Drive, Birmingham, AL, 35207 35207 33.562365 -86.830703 1 2 9 J'ti'z Smoothie-N-Coffee Bar NaN NaN Coffee and Tea, Breakfast and Brunch, Bubble Tea NaN 1521 Pinson Valley Parkway, Birmingham, AL, 35217 33.583640 -86.773330 2 3 6 Philly Fresh Cheesesteaks (541-B Graymont Ave) NaN NaN American, Cheesesteak, Sandwiches, Alcohol \$ 541-B Graymont Ave, Birmingham, AL, 35204 33.509800 -86.854640 3 4 17 Papa Murphy's (1580 Montgomery Highway) NaN NaN Pizza \$ 1580 Montgomery Highway, Hoover, AL, 35226 33.404439 -86.806614 4 5 162 Nelson Brothers Cafe (17th St N) 4.7 22.0 Breakfast and Brunch, Burgers, Sandwiches NaN 314 17th St N, Birmingham, AL, 35203 35.03 33.514730 -86.811700 4- Data Exploration:- # Get the data shape data.shape (40227, 11)
Out[88]:	data.size 442497 # get more info about data data.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 40227 entries, 0 to 40226 Data columns (total 11 columns): # Column Non-Null Count Dtype</class>
	0 id 40227 non-null int64 1 position 40227 non-null int64 2 name 40227 non-null object 3 score 22254 non-null float64 4 ratings 22254 non-null float64 5 category 40204 non-null object 6 price_range 33581 non-null object 7 full_address 39949 non-null object 8 zip_code 39940 non-null object 9 lat 40227 non-null float64 10 lng 40227 non-null float64 dtypes: float64(4), int64(2), object(5) memory usage: 3.4+ MB That mean the dataset is consets of 40227 rows (instances) and 11 column (feature) There are more than datatypes (int64, Float64, object), so we have Numerical and Categorical Data
Out[90]: -	# Get more Statistical info about the Numerical data data.describe() id position score ratings lat lng
Out[91]:	# Find how many unique values data.nunique() id
In [92]: Out[92]:	<pre>lat</pre>
Out[6]:	data.isna().sum() id 0 position 0 name 0 score 17973 ratings 17973 category 23 price_range 6646 full_address 278 zip_code 287 lat 0 lng 0 dtype: int64 **There are Missing Values , We should make Data cleaning 5- Data cleaning:-
In [8]:	<pre># Check Duplicates print("\nDuplicates Rows:", data.duplicated().sum()) Duplicates Rows: 0 # check null values data.isna().sum() id</pre>
	score 17973 ratings 17997 category 23 price_range 6646 full_address 278 zip_code 287 lat 0 lng 0 dtype: int64 - We have missing Values in : - • 1- (score, ratings) which is numerical data • 2- (category, price_range, full_address, zip_code) which is Categorical data - Lets Handle This Missing values:-
	 1- numerical data> fill with median (beacuse of using mean affected by outliers) 2- Categorical data> fill (price_range) with Forward fill , drop (category, full_address, zip_code) # create Data Dictionary to fill the missing values values = {
n [11]:	'price_range': data['price_range'].fillna(method='ffill') # forward fill } # Use The Data Dictionary to fill missing Values And Save The Resualts data.fillna(value= values, inplace=True) We filled all Missing Values, Lets delete the others (category, full_address, zip_code) # drop null missing values and save the result data.dropna(inplace=True)
out[12]:	# check null values data.isna().sum() id
n [45]:	Now The Data is cleaned And Ready For Exploratory Data Analysis (EDA) 6- Exploratory Data Analysis:- data.columns Index(['id', 'position', 'name', 'score', 'ratings', 'category', 'price_range',
n [55]:	Index (['id', 'position', 'name', 'score', 'ratings', 'category', 'price_range',
n [70]: ut[70]: _	3 4 17 Papa Murphy's (1580 Montgomery Highway) 4.6 500 Pizza \$ 1580 Montgomery Highway, Hoover, AL, 35226 35226 33.40439 -86.80614 4 5 162 Nelson Brothers Cafe (17th St N) 4.7 500 Breakfast and Brunch, Burgers, Sandwiches \$ 314 17th St N, Birmingham, AL, 35203 35203 33.51470 -86.81170 data ['tprice_range'] == '\$\$\$') & (data ['score'] == 5.0)] .sort_values (by='id') tid position name score ratings category price_range full_address zip_code lat Ing 15 16 88 Jeni's Splendid Ice Cream (Pepper Place) 5.0 500 Ice Cream & amp; Frozen Yogurt, Comfort Food, D \$\$\$ 219 29th St S, Birmingham, AL, 35233 35233 33.51600 -86.789950 35230 35233 33.51600 -86.789950 35230 35233 33.51600 -86.789950 35230 35233 33.51600 -86.789950 35230 35233 33.51600 -86.789950 35230 35233 33.51600 -86.789950 35230
	13594 209 Vego Eatz 5.0 500 Vegetarian, Healthy \$\$\$ 203 W Pioneer Ave, Puyallup, WA, 98371 9871 47.190590 -122.295470 18947 18948 54 Jeni's Splendid Ice Creams (Old Town Alexandria) 5.0 500 Ice Cream & Cream & Frozen Yogurt, Comfort Food, D \$\$\$ 102 South Patrick Street, Alexandria, VA, 22314 22314 38.805220 -77.058314 20017 20018 219 Laporta's Restaurant 5.0 500 American, Burgers, Pasta \$\$\$ 1600 Duke St, Alexandria, VA, 22314 22314 38.803849 -77.058314 20619 20620 59 Bar Charley 5.0 500 Food American, Burgers, Pasta \$\$\$ 1825 18th St NW, Washington, DC, 20009 20009 38.915028 -77.058314 26444 26445 92 Jeni's Ice Cream Bethesda 5.0 500 Ice Cream & C
n [13]:	<pre># Calculate Correlation corr = data.drop(['ratings'], axis=1).corr() # make the figure size plt.figure(figsize=(8, 6)) sns.heatmap(corr, annot=True, cmap='coolwarm') # Add correlation values plt.title("Correlation Heatmap") # Add title plt.show()</pre>
	Correlation Heatmap 1 0.16 -0.016 -0.65 0.14 -0.8 -0.8 -0.66 -0.
	0.4 0.0 - 0.016 0.068 1 0.04 - 0.0 - 0.0 - 0.0
in [16]:	
	plt.show() score ratings 25000 20000 17500 15000
	10000
	position 2500 200 300 400 5
	8000
	2000 2000 250 300 0 10 20 30 40
	6000
	#Check Outlires For numerical columns by boxpolt num_cols = data.columns plt.figure(figsize=(10, 5)) pro_bemplet(datametra[sum_cols], enjoytemen)
	sns.boxplot(data=data[num_cols], orient='v') plt.title('Boxplot for Outlier Detection') plt.show() Boxplot for Outlier Detection 40000 35000
	25000 20000 15000
	10000 5000 id position score ratings lat lng # what is cateogries of resturants
ut[153]: n [169	<pre>data['category'].values array(['Burgers, American, Sandwiches',</pre>
ut[169]:	id position name score ratings category price_range full_address zip_code lat lng 551 552 23 Wild Burger (3400 Montgomery Hwy.) 4.1 500 Burgers \$ 3400 Montgomery Hwy, Dothan, AL, 36303 36303 31.25525 -85.43014 588 589 31 Songwriters Cafe (3320 Montgomery Highway) 4.6 500 American \$ 3320 Montgomery Hwy, Dothan, AL, 36303 36303 31.25448 -85.42925 664 665 15 Steak 'n Shake (5901 University Drive, Suite I) 4.4 500 American \$ 5901 University Drive, Suite I, Huntsville, AL 35806 34.73867 -86.66602 714 715 34 Wild Burger (7042 Highway 72 West) 4.6 500 Burgers \$ 7042 Highway 72 W, Huntsville, AL, 35806 35806 34.75441 -86.71043 1033 1034 58 Songwriters Cafe (3060 So. McKenzie Street) 4.6 500 American \$ 3060 S McKenzie St, Foley, AL, 36535 36535 30.36941 -87.68412
n [59]:	<pre>## change id to restaurant_id data.rename(columns={'id':'restaurant_id'}, inplace=True) - What are the top 10 highest-score restaurants? ## list the top 10 restaurants in score data[data['score']==np.max(data['score'])].head(10) restaurant_id position</pre>
	15 16 88 Jeni's Splendid Ice Cream (Pepper Place) 5.0 500.0 Ice Cream & Support, Comfort Food, D \$\$\$ 219 29th St. S, Birmingham, AL, 35233 3.51660 -86.7895 1.56
	356 357 49 Mr. Lin Chinese Restaurant 5.0 500.0 Chinese, Asian, Asian Fusion \$\$ 475, Helena, AL, 35080 3508 33.279000 -86.8511 410 411 94 Great American Cookies (Riverchase Galleria) 5.0 500.0 Bakery, Desserts, Comfort Food \$ 2000 Riverchase Galleria, Birmingham, AL, 35244 33.379202 -86.8087 615 616 1 Tropical Smoothie Cafe - 3230 Ross Clark Circl 5.0 500.0 Juice and Smoothies, Healthy, Fast Food \$ 3230 Ross Clark Circle, Suite 3, Dothan, AL, 3 36303 31.234995 -85.4312 632 633 10 Firehouse Subs (3255 South Oates Street. Suite 8) 5.0 500.0 Sandwich, Deli \$ 3255 South Oates Street. Suite 8, Dothan, AL, 36301 31.179849 -85.4010 62 Hunt Brothers Pizza 5.0 500.0 American, Italian, Wings \$\$ 6090 Old Madison Pike NW, Huntsville, AL, 35806 34.713470 -86.6580
n [78]:	- Which cuisine categories are most common in the dataset? # list the top 5 categories var = data['category'].value_counts() print(var.head(5)) Burgers, American, Sandwiches 1606 Mexican, Latin American, New Mexican 1161
	Fast Food, Sandwich, American 837
	Fast Food, Sandwich, American 837 Pizza, American, Italian 707 American, Burgers, Fast Food 685 Name: category, dtype: int64 category = ['Burgers', 'Mexican', 'Fast Food', 'Pizza', 'American'] nums = [1606,1161,837,707,685] # plot the pie chart plt.pie(nums, labels=category, autopct='%1.1f%%') plt.title('top 5 categories')
	Fast Food, Sandwich, American 837 Pizza, American, Italian 707 American, Burgers, Fast Food 685 Name: category, dtype: int64 category = ['Burgers', 'Mexican', 'Fast Food', 'Pizza', 'American'] nums = [1606,1161,837,707,685] # plot the pie chart plt.pie(nums, labels=category, autopct='%1.1f%%') plt.title('top 5 categories') plt.show() top 5 categories Burgers
n [83]:	Fast Food. Sandsich, kerkies
n [83]:	Fast Food **Rezican, Burgers, Past Pood servers and provided prov
n [83]:	Fig. 1 Fact Food, Standards, James Landards 1975 Landard Company and System (1985) L
n [93]:	Paut
n [93]:	Face Food Search