

DATA ANALYSIS TASK LIST

Level 1

Task 1

Task: Top Cuisines

Determine the top three most common cuisines in the dataset.

Calculate the percentage of restaurants that serve each of the top cuisines.

```
In [3]: import pandas as pd
import matplotlib.pyplot as plt
import json, express as px

In [4]: df = pd.read_csv('C:/Users/Aravind/Desktop/project.py/ps3405.py/Dataset_1.csv')
df.columns = df.columns.str.strip()

In [5]: df.head()
```

Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude	Latitude	Cuisines ...	Currency	Has Table booking	Has Online delivery	Is delivering now	Switch to order menu	Price range	Aggregate rating	Rating color	Rating text	Votes	
0	6317637	Le Funt Suifu	162	Makati City	Third Floor, Century City Mall, Kataysan Avenue...	Century City Mall, Poblacion, Makati City, Mak...	Century City Mall, Poblacion, Makati City, Mak...	121.027235	14.565443	French, Japanese, ... Desserts	Botswana Pula(P)	Yes	No	No	No	3	4.6	Dark Green	Excellent	314
1	6304287	Izaya Kajufu	162	Makati City	Little Tokyo, 2277 Chino Rome Avenue, Legaspi...	Little Tokyo, Legaspi Village, Makati City, Ma...	Little Tokyo, Legaspi Village, Makati City, Ma...	121.014101	14.553708	Japanese ...	Botswana Pula(P)	Yes	No	No	No	3	4.5	Dark Green	Excellent	591
2	6300002	Heel - Edes Shang-Li	162	Mandakyong City	Edsa Shangri-La 1 Garden Way, Origas, Mandal...	Edsa Shangri-La, Origas, Mandakyong City	Edsa Shangri-La, Origas, Mandakyong City	121.056831	14.561404	Szechuan, Asian, Filipino, Indian ...	Botswana Pula(P)	Yes	No	No	No	4	4.4	Green	Very Good	270
3	6316506	Ooma	162	Mandakyong City	Third Floor, Mega Fashion Hall, SM Megamall, O...	SM Megamall, Origas, Mandakyong City	SM Megamall, Origas, Mandakyong City	121.066475	14.585318	Japanese, Sushi	Botswana Pula(P)	No	No	No	No	4	4.8	Dark Green	Excellent	365
4	6314302	Sambo Kojin	162	Mandakyong City	Third Floor, Mega Atrium, SM Megamall, Origas...	SM Megamall, Origas, Mandakyong City	SM Megamall, Origas, Mandakyong City	121.037508	14.564450	Japanese, Korean ...	Botswana Pula(P)	Yes	No	No	No	4	4.8	Dark Green	Excellent	229

5 rows * 21 columns

Task 2

Task: City Analysis

Identify the city with the highest number of restaurants in the dataset.

Calculate the average rating for restaurants in each city.

Determine the city with the highest average rating.

```
In [6]: city_counts = df['City'].value_counts()
print("City with highest number of restaurants:\n", city_counts.head(1))

City with highest number of restaurants:
City
New Delhi    5473
Name: count, dtype: int64

In [7]: avg_rating_per_city = df.groupby('City')['Aggregate rating'].mean()
print("Average rating per city:\n", avg_rating_per_city)

Average rating per city:
City
Asia Dhabl    4.300000
Agra          3.950000
Ahmedabad     4.613903
Alibay        3.355000
Alibahad      3.350000
...
Wilmington Bay 4.250000
Winchester Bay 3.200000
Yokotom       3.300000
Zanzibar      4.292857
Name: Aggregate rating, Length: 141, dtype: float64

In [8]: top_avg_city = avg_rating_per_city.sort_values(ascending=False)
print("\nCity with highest average rating:\n", top_avg_city.head(1))

City with highest average rating:
City
Inner City    4.9
Name: Aggregate rating, dtype: float64

Level Task 3 Task: Price Range Distribution
Create a histogram or bar chart to visualize the distribution of price ranges among the restaurants.
Calculate the percentage of restaurants in each price range category.
```

price_range	price_counts
1	4444
2	3113
3	1408
4	586

Restaurant count by price range:

Percentage by price range:

price_range	percentage
1	46.529159
2	32.539466
3	14.741912
4	6.114882

Name: count, dtype: float64

Bar chart visualization of Price Range Distribution of Restaurants

Level 1 Task 4

Task: Online Delivery

Determine the percentage of restaurants that offer online delivery.

Compare the average ratings of restaurants with and without online delivery.

```
In [14]: online_counts = df['Has Online delivery'].value_counts()
print("Online delivery counts:\n", online_counts)

Online delivery counts:
Has Online delivery
Yes    7100
No     2451
Name: count, dtype: int64

In [15]: total = len(df)
online_percentage = (online_counts / total) * 100
print("Percentage by price range:\n", online_percentage)

Percentage:
Has Online delivery
No     74.337166
Yes    25.662834
Name: count, dtype: float64

In [16]: avg_with_delivery = df[(df['Has Online delivery'] == 'Yes')]['Aggregate rating'].mean()
avg_without_delivery = df[(df['Has Online delivery'] == 'No')]['Aggregate rating'].mean()

In [17]: print("Average rating with online delivery:", round(avg_with_delivery, 2))
print("Average rating without online delivery:", round(avg_without_delivery, 2))

Average rating with online delivery: 5.25
Average rating without online delivery: 5.47

Level 2 Task 1
Task: Restaurant Ratings
Analyze the distribution of aggregate ratings and determine the most common rating range.
Calculate the average number of votes received by restaurants.
```

Aggregate rating	rating_counts
5.0	2148
4.8	1
4.9	2
4.0	7
4.1	15
4.2	27
4.3	47
4.4	67
4.5	110
4.6	151
4.7	250
4.8	315
4.9	381
3.0	468
3.1	519
3.2	592
3.3	483
3.4	498
3.5	460
3.6	458
3.7	427
3.8	400
3.9	335
4.0	246
4.1	274
4.2	274
4.3	174
4.4	144
4.5	95
4.6	78
4.7	42
4.8	25
4.9	41

Name: count, dtype: int64

Most common rating is: 5.0

Most common rating is: 5.0

Average votes = 156.91

Average number of votes per restaurant: 156.91

Bar chart visualization of Aggregate Rating Distribution

Level 2 Task 2 Task: Cuisine Combination

Identify the most common combinations of cuisines in the dataset.

Determine if certain cuisine combinations tend to have higher ratings.

```
In [22]: top_combo = df['Cuisines'].value_counts().head(10)
print("Top Cuisine Combinations:\n", top_combo)

Top Cuisine Combinations:
Cuisines
North Indian 936
North Indian, Chinese 512
Chinese 354
Fast Food 354
North Indian, Mughlai 334
Cafe 299
Bakery 228
North Indian, Mughlai, Chinese 197
Bakery, Desserts 170
Street Food 149
Name: count, dtype: int64

In [23]: combo_rating = df.groupby('Cuisines')['Aggregate rating'].mean().sort_values(ascending=False)
print("\nCuisine Combinations with Highest Average Ratings:\n", combo_rating.head(10))

Cuisine Combinations with Highest Average Ratings:
Cuisines
Italian, Deli 4.9
Bavarian, Beef Cook 4.9
American, Sandwich, Tea 4.9
Continental, Indian 4.9
European, Indian, Indian 4.9
European, Contemporary 4.9
European, German 4.9
BBQ, Breakfast, Southern 4.9
American, Coffee and Tea 4.9
Sunda, Indonesian 4.9
Name: Aggregate rating, dtype: float64

Level Task 3 Task: Geographic Analysis
Plot the locations of restaurants on a map using longitude and latitude coordinates.
Identify any patterns or clusters of restaurants in specific areas.
```

Map visualization of restaurant locations

Level 2 Task 4 Task: Restaurant Chains

Identify if there are any restaurant chains present in the dataset.

Analyze the ratings and popularity of different restaurant chains.

```
In [24]: chains_counts = df['Restaurant Name'].value_counts()
restaurant_chains = chain_counts[chain_counts > 1]
print("Top Restaurant Chains:\n", restaurant_chains.head(10))

Top restaurant chains:
Restaurant Name
Cafe Coffee Day    83
Domino's Pizza    79
Subway            65
Green Chick Chop  51
McDonald's        34
Keweenaw          30
Pizza Hut         29
Gianni            28
Baskin Robbins    26
Bachchan Nation  25
Name: count, dtype: int64

In [25]: chains_df = df[df['Restaurant Name'].isin(restaurant_chains.index)]

In [26]: chains_analysis = chains_df.groupby('Restaurant Name').agg(
    {'Aggregate rating': 'mean',
     'Votes': 'sum',
     'City': 'count'})

[[Restaurant Name]
Aggregate rating: 'Average Rating',
Votes: 'Average Votes',
City: 'Branch Count']

[[sort_values(by='Branch Count', ascending=False)

print("\nTop Restaurant Chains (with ratings & votes):\n")
print(chains_analysis.head(10))

Top Restaurant Chains (with ratings & votes):
Average Rating Average Votes Branch Count
Restaurant Name
Cafe Coffee Day    2.419277    29.253012    83
Domino's Pizza    2.740506    84.086608    79
Subway            2.507937    87.006349    65
Green Chick Chop  2.672549    18.301461    51
McDonald's        3.339583    110.239167    48
Keweenaw          2.870388    37.147039    34
Pizza Hut         2.690555    165.364667    29
Gianni            2.689555    29.488276    29
Baskin Robbins    1.680714    15.289714    26
Bachchan Nation  4.333846    1982.384615    26

Level 3 Task 1 Task: Restaurant Reviews
Analyze the text reviews to identify the most common positive and negative keywords.
Calculate the average length of reviews and explore if there is a relationship between review length and rating.
```

```
In [29]: import pandas as pd
import re
import matplotlib.pyplot as plt
import random as rd
from nltk.tokenize import word_tokenize

In [30]: rating_counts = df['Rating text'].value_counts()
print("Number of Restaurants in each Rating Category:")

Number of Restaurants in each Rating Category:
Rating text
Average    2157
Not rated  2148
Good       2100
Very Good  1079
Excellent   701
Poor        186
Name: count, dtype: int64

In [31]: avg_votes = df.groupby('Rating text')['Votes'].mean().sort_values(ascending=False)
print("\nAverage Votes per Rating Category:")

Average Votes per Rating Category:
Rating text
Excellent    851.770764
Very Good    520.158758
Good         229.354219
Poor          80.715054
Average      46.149120
Not rated     0.870312
Name: Votes, dtype: float64

In [32]: plt.figure(figsize=(15,5))
sns.barplot(x=avg_votes.index, y=avg_votes.values, palette='viridis')
plt.title("Average Votes by Rating Category")
plt.xlabel("Rating Category")
plt.ylabel("Average Number of Votes")
plt.xticks(rotation=50)
plt.tight_layout()
plt.show()

C:/Users/Aravind/AppData/Local/Temp/ipykernel_13504/3009967702.py:2: FutureWarning:
    Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'x' variable to 'hue' and set 'legend=False' for the same effect.
```

Bar chart visualization of Average Votes by Rating Category

Level 3 Task 2 Task: Votes Analysis

Identify the restaurants with the highest and lowest number of votes.

Analyze if there is a correlation between the number of votes and the rating of a restaurant.

```
In [33]: top_votes = df.sort_values(by='Votes', ascending=False).head(5)
low_votes = df[(df['Votes'] > 0).sort_values(by='Votes', ascending=True).head(5)]
print("\nTop 5 Most Voted Restaurants:\n", top_votes[['Restaurant Name', 'Votes', 'Aggregate rating']])
print("\nBottom 5 (non-zero votes):\n", low_votes[['Restaurant Name', 'Votes', 'Aggregate rating']])

Top 5 Most Voted Restaurants:
Restaurant Name    Votes    Aggregate rating
728              Toit    10594          4.8
728              Triffid    868          4.7
3994             Rauz Khao Social    7931          4.3
2412             Peter Cat    7574          4.3
739             Absolute Barchinese    4907          4.4

Bottom 5 (non-zero votes):
Restaurant Name    Votes    Aggregate rating
1014             Abata    1          0.0
5653             The Cake Basket    1          0.0
5447             Ramatik Restaurant    1          0.0
5446             RT's Shik-Shack    1          0.0
5836             Approval Dweets    1          0.0

In [34]: correlation = df[['Votes', 'Aggregate rating']].corr().iloc[0,1]
print("Pearson Correlation between Votes and Rating: (round(correlation, 2))*")

Pearson Correlation between Votes and Rating: 0.31

In [35]: plt.figure(figsize=(6,4))
sns.scatterplot(data=df, x='Votes', y='Aggregate rating', alpha=0.6)
plt.title("Votes vs Aggregate Rating")
plt.xlabel("Number of Votes")
plt.ylabel("Aggregate Rating")
plt.tight_layout()
plt.show()
```

Scatter plot visualization of Votes vs Aggregate Rating

Level 3 Level 3 Task: Price Range vs. Online Delivery and Table Booking

Analyze if there is a relationship between the price range and the availability of online delivery and table booking.

Determine if higher-priced restaurants are more likely to offer these services.

```
In [36]: df[['Has Online delivery']] = df[['Has Online delivery']].map({'Yes': 1, 'No': 0})
df[['Has Table booking']] = df[['Has Table booking']].map({'Yes': 1, 'No': 0})

In [37]: grouped = df.groupby('Price range')[['Has Online delivery', 'Has Table booking']].mean() * 100
grouped = grouped.reset_index()

In [38]: print("\nService Availability by Price Range:\n")
print(grouped)

Service Availability by Price Range:
Price range    Has Online delivery (%)    Has Table booking (%)
1             15.774077             0.022502
2             41.110833             7.277482
3             29.180741             45.738636
4             9.044369             46.757679

In [39]: grouped.plot(kind='bar', figsize=(7,5), color=['orange', 'green'])
plt.title("Online Delivery & Table Booking % DV Price Range")
plt.xlabel("Percentage of Restaurants")
plt.ylabel("Price Range (1 = Lowest, 4 = Highest)")
plt.xticks(rotation=50)
plt.tight_layout()
plt.show()
```

Bar chart visualization of Service Availability by Price Range

While completing these data analysis tasks, I explored real-world restaurant data, performed data cleaning, visualization, and gained hands-on experience with tools like Python (pandas, matplotlib) and Excel. This project helped me build confidence in solving real datasets and improved my overall understanding of data analytics.