

**Group Name:** CA

**Name:** Angel Huang, Cecilia Huynh

**Github Link:** [https://github.com/angelxhuang/final\\_project206](https://github.com/angelxhuang/final_project206)

**Course:** SI 206: Final Project Report

## **1. Project Goals**

Our goal for the project was to determine if there was an effect on Yelp ratings of restaurants and cafes based on the population in California cities. By calculating the average Yelp ratings of restaurants and the average Yelp ratings of cafes in select cities, we compared this to the population sizes. We grabbed 25 data points each from San Francisco, Silicon Valley, Southern California, Bay Area, Los Angeles, and overall California for restaurants and did the same for cafes for a total of 300 ratings. We matched these data points with their respective cities to examine our findings. We also compared the ratings of cafes to restaurants in a city. To visualize this, we planned to have four charts: one bar graph for average ratings of cafes and restaurants by city, one scatterplot for city population versus their average cafe and restaurant ratings, one bar graph for average price levels of cafes and restaurants by city, and one scatterplot for city population versus their average cafe and restaurant price levels.

## **2. Achieved Project Goals**

In the beginning, we planned to look at Yelp and TripAdvisor APIs to compare ratings of restaurants and the city location. However, we revised our plan due to complications with getting the latter API and looked at a website ([www.california-demographics.com/cities\\_by\\_population](http://www.california-demographics.com/cities_by_population)) that had California city population instead. The Yelp API provided valuable data of restaurant and cafe ratings and price levels, while the website gave us information about California city populations. By extracting information from these two sources, we were able to calculate the average ratings and price levels, collect the population sizes, and make graphs to represent our findings. We found that cafes generally have better ratings than restaurants do in the same city, and there is not much correlation between a city's population and their average cafe and restaurant ratings. Our conclusion comes from our two graphs about average ratings: a bar graph and a scatterplot. We also discovered that average restaurant price levels are higher than the average cafe price level in a city, and there was no correlation between the city population versus average cafe and restaurant price levels. Overall, this project allowed us to learn more about web scraping and APIs, understand that there is a very small relation between restaurant and cafe ratings in California cities, and lastly, visualize how there is little to no correlation between the price levels and ratings of restaurants and cafes in addition to their city population.

### **3. Problems We Faced**

*Problem 1:* Unable to access TripAdvisor API for free.

*Solution:* Revised our project plan to find a different API/website instead which also meant we had to change our goals for the project. Found a California city population website to compare restaurants and cafes with the population rather than Yelp ratings with overall TripAdvisor ratings.

*Problem 2:* Extracting randomized data from Yelp API gave cities that did not appear on California cities population website.

*Solution:* Specifically choose parts/regions of California for the majority of data points (25 from San Francisco, 25 from Silicon Valley, 25 from Southern California, 25 from Bay Area, 25 from Los Angeles, 25 California) to find ratings that correspond with population website. The latter of the 25 points are to account for some randomization.

*Problem 3:* Running code without duplicates in the data.

*Solution:* Used an offset query that allowed us to grab the data once without having repeated data over and over again.

#### 4. Calculations from the Database File

##### Average Ratings of Cafes and Restaurants vs. Population of California Cities

cafes\_population

Average Ratings of California Cities Places vs. Population	
Population	Average Ratings
117145	4.0
122989	4.5
544510	4.5
3849297	4.4
433823	4.24
148338	4.75
815201	4.3
127151	4.33
152258	4.14

res\_population

Average Ratings of California Cities Places vs. Population	
Population	Average Ratings
117145	4.0
122989	4.33
544510	4.5
3849297	4.4
433823	4.17
148338	4.5
815201	4.26
983489	4.5
127151	4.33
152258	4.11

##### Average Price Levels of Cafes and Restaurants vs. Population of California Cities

cafe\_price\_pop

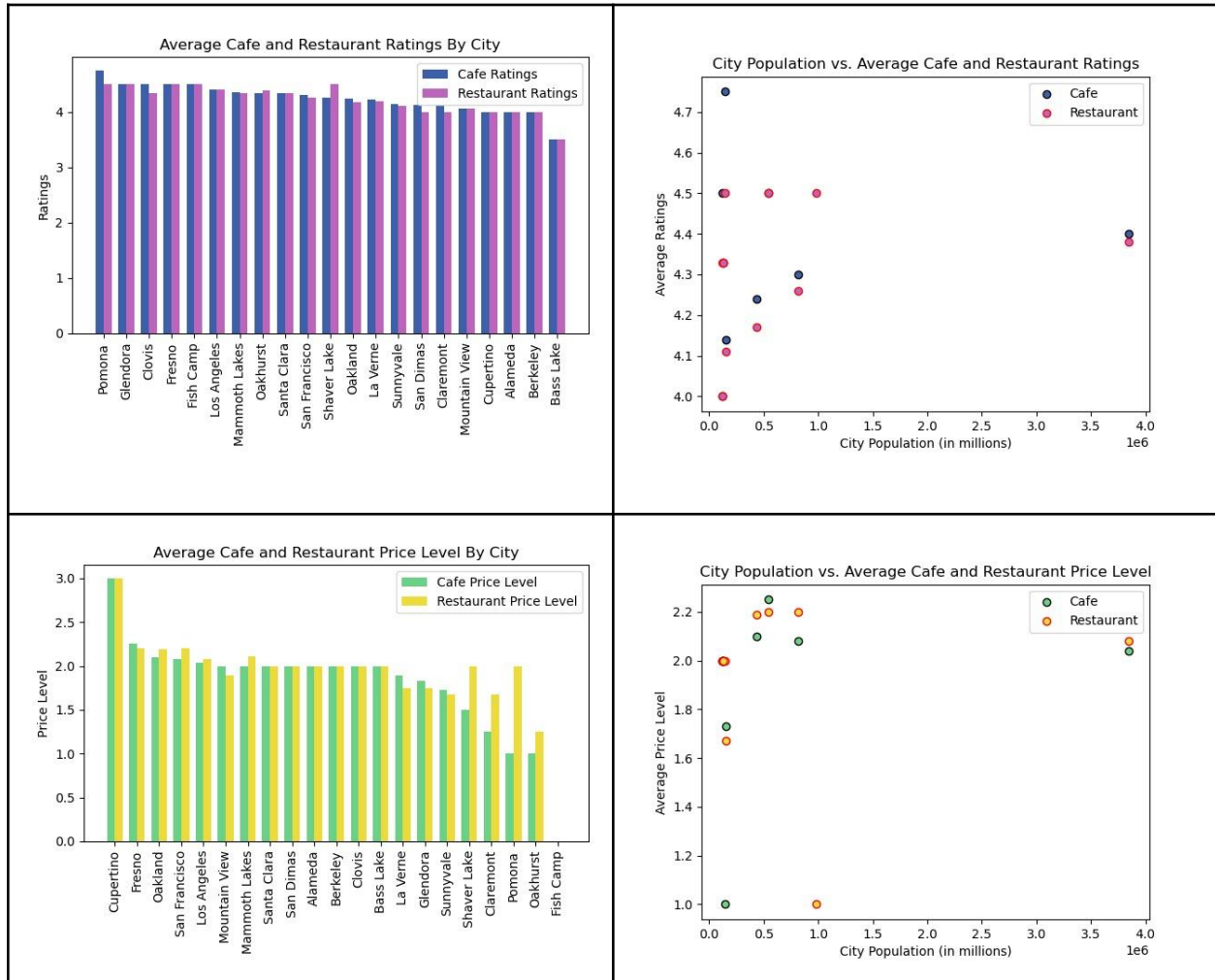
Average Price Levels of California Cities Places vs. Population	
Population	Average Price Levels
117145	2.0
122989	2.0
544510	2.25
3849297	2.04
433823	2.1
148338	1.0
815201	2.08
127151	2.0
152258	1.73

res\_price\_pop

Average Price Levels of California Cities Places vs. Population	
Population	Average Price Levels
117145	2.0
122989	2.0
544510	2.2
3849297	2.08
433823	2.19
148338	2.0
815201	2.2
983489	1.0
127151	2.0
152258	1.67

## 5. Created Data Visualizations

### [Enlarged Versions of the Images](#)



## 6. Instructions For Running Code

Our code does not require any special instructions to run. Compile and run the code to have the CSV files and data visualizations appear. Make sure to use your DB Browser for SQLite to open your database.

## 7. Documentation For Each Function (input and output)

```
# Input: Database name
# Output: Creates a database
def setUpDatabase(db_name):

# Input: Dictionaries that you want to combine
# Output: List
def combine(dict1, dict2, dict3, dict4, dict5, dict6):
```

```

# Input: Location, Offset
# Output: Dictionary
# About: Retrieves cafe data from the Yelp API
def yelp_cafes(location, offset):

# Input: Dictionary, Cur, Conn
# Output: A table based off Yelp cafes data, sorted by city, cafe name, rating, price
level
def create_yelp_cafes(cafes, cur, conn):

# Input: Location, Offset
# Output: Dictionary
# About: Retrieves restaurant data from the Yelp API
def yelp_restaurants(location, offset):

# Input: Dictionary, Cur, Conn
# Output: A table based off Yelp restaurants data, sorted by city, restaurant name,
rating, price level
def create_yelp_restaurants(restaurants, cur, conn):

# Input: None
# Output: Dictionary
# About: We extracted the first 100 ranked california cities & their population
def cali_ratings():

# Input: Dictionary, Cur, Conn
# Output: A table based off our california city population data, sorted by city and
population
def population_table(pop_dict, cur, conn):

# Input: Cur, Conn
# Output: Dictionary
# About: We selected city + cafe rating from our cafe data and calculated the average
ratings of cafes in each city
def avg_yelp_cafes(cur, conn):

# Input: Cur, Conn
# Output: Dictionary
# About: We selected city + restaurant rating from our restaurant data and calculated
the average ratings of restaurants in each city
def avg_yelp_restaurants(cur, conn):

```

```

# Input: Cur, Conn
# Output: Dictionary
# About: We selected city + cafe price level from our cafe data and calculated the
average price level of cafes in each city
def avg_yelp_cafes_price(cur, conn):

# Input: Cur, Conn
# Output: Dictionary
# About: We selected city + restaurant price level from our restaurant data and
calculated the average price level of restaurants in each city
def avg_yelp_restaurants_price(cur, conn):

# Input: Cur, Conn
# Output: List
# About: Joins restaurants data and population data in a tuple format: (City name,
Population, Average restaurant ratings)
def res_pop_join(cur, conn):

# Input: Cur, Conn
# Output: Dictionary
# About: Joins restaurants data and population data in a dictionary format: {City
name: (Population, Average restaurant price level)
def res_price_pop_join(cur, conn):

# Input: Cur, Conn
# Output: List
# About: Joins cafe data and population data in a tuple format: (City name,
Population, Average cafe ratings)
def cafe_pop_join(cur, conn):

# Input: Cur, Conn
# Output: Dictionary
# About: Joins cafe data and population data in a dictionary format: {City name:
(Population, Average cafe price level)
def cafe_price_pop_join(cur, conn):

# Input: Data table, File name
# Output: CSV File
# About: Writes the csv file for dot plot of "Average Ratings of California Cities
Places vs. Population"
def write_csv_dot(data1, filename):

```

```

# Input: Data table, File name
# Output: CSV File
# About: Writes the csv file for dot plot of "Average Price Levels of California
Cities Places vs. Population"
def write_csv_dot_price(data1, filename):

# Input: Data table, File name
# Output: CSV File
# About: Writes the csv file for bar plot of "Average Ratings of California Cities
Cafes vs. Restaurants"
def write_csv_bar(data1, data2, filename):

# Input: Data table, File name
# Output: CSV File
# About: Writes the csv file for bar plot of "Average Price Levels of California
Cities Cafes vs. Restaurants"
def write_csv_bar_price(data1, data2, filename):

# Input: File
# Output: Bar graph
# About: Creates a bar graph of 'Average Cafe and Restaurant Ratings By City'
def cali_bar_graph(file):

# Input: Cafe file and Restaurant file
# Output: Scatter plot
# About: Creates a scatterplot of 'City Population vs. Average Cafe and Restaurant
Ratings'
def cali_dot_plot(file1, file2):

# Input: File
# Output: Bar graph
# About: Creates a bar graph of 'Average Cafe and Restaurant Price Level By City'
def cali_price_bar_graph(file):

# Input: Cafe file and Restaurant file
# Output: Scatter plot
# About: Creates a scatterplot of 'City Population vs. Average Cafe and Restaurant
Price Level'
def cali_price_dot_plot(file1, file2):

# Calls all of our main functions

```

```
def main():
```

## 8. Resources Used

Date	Issue Description	Location of Resource	Result (did it solve the issue?)
12/2	Finding an API/website to replace our initial plan of using TripAdvisor	<a href="https://www.california-demo-graphics.com/cities_by_population">https://www.california-demo-graphics.com/cities_by_population</a>	Yes, we found and used this website to extract population size data using BeautifulSoup
12/3	We realized that we were getting restaurants/cafes from the same cities in California (that did not match the top 100 California cities we had in our other data table). We wanted a more spread-out set of data points that include restaurants/cafes from major cities in California	<a href="https://docs.developer.yelp.com/reference/v3_business_search">https://docs.developer.yelp.com/reference/v3_business_search</a>	Yes, we realized that we can set the location to specific locations, such as “Silicon Valley” or “Bay Area”, which included more major cities such as Berkeley or Santa Clara. This is why we decided to extract 25 data points each from different areas within California.
12/5	Because we could only retrieve 25 items each time we run our code, we wanted to combine all of the data (150 items) for restaurants and cafes.	<a href="https://www.geeksforgeeks.org/python-merging-two-dictionaries/">https://www.geeksforgeeks.org/python-merging-two-dictionaries/</a>	Yes, we ended up combining our cafes/restaurants data into one big dictionary that we can utilize for the rest of our functions.
12/7	We had trouble running the code without duplicating existing data (ex. Trying to not get the same 25 restaurants/cafes over and over again)	<a href="https://docs.developer.yelp.com/reference/v3_business_search">https://docs.developer.yelp.com/reference/v3_business_search</a>	Yes, we realize that there was an offset query that allows us to get unique restaurants/cafes every time we run our code. This is also why we had offset as an input for our yelp_restaurants and yelp_cafes functions so we can control the uniqueness of our data.
12/8	We did not know how to set up database	Discussion 11	Yes, it worked. We grabbed the code from Discussion 11. <pre>def setUpDatabase(db_name):     path =     os.path.dirname(os.path.abspath(__file__))     conn =     sqlite3.connect(path+'/'+db_name)     cur = conn.cursor()     return cur, conn</pre>
12/8	We had trouble joining our city population data and our calculated average ratings/price level data.	<a href="https://www.w3schools.com/python/python_mysql_join.asp">https://www.w3schools.com/python/python_mysql_join.asp</a>	Yes, we figured out how to join the two data tables together. Because the price level data was in string forms (ex. “\$” or “\$\$”), we had to convert the data to integer first (ex. “\$” = 1, “\$\$” = 2, etc.) before calculating the



			average price levels and joining it with the population data into one table.
12/8	We had trouble making a grouped bar graph with labels.	<a href="https://matplotlib.org/stable/gallery/lines_bars_and_markers/barchart.html">https://matplotlib.org/stable/gallery/lines_bars_and_markers/barchart.html</a>	Yes, we followed the code format for creating a grouped bar chart from this website and adjusted titles/data points based on our own project (ex. x-axis/y-axis labels + x-ticks/y-ticks).
12/9	We had trouble plotting multiple data from two different files onto the same scatterplot.	<a href="https://www.scaler.com/topics/matplotlib/scatter-plot-matplotlib/">https://www.scaler.com/topics/matplotlib/scatter-plot-matplotlib/</a> <a href="https://www.geeksforgeeks.org/visualize-data-from-csv-file-in-python/">https://www.geeksforgeeks.org/visualize-data-from-csv-file-in-python/</a>	Yes, we followed the code format for implementing multiple scatterplots on the same graph from this website and adjusted titles/data points based on our own project (ex. x-axis/y-axis labels + x-scatter/y-scatter).
12/11	We were not able to change the color or edgcolor of the bar graphs.	<a href="https://www.python-graph-gallery.com/3-control-color-of-barplots">https://www.python-graph-gallery.com/3-control-color-of-barplots</a> <a href="https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.scatter.html">https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.scatter.html</a>	Yes, we realized that we were using “c=” instead of “color=”. We also had to put the color code in [] for it to work.