

Faitus Jeline Joseph

9/12/2022

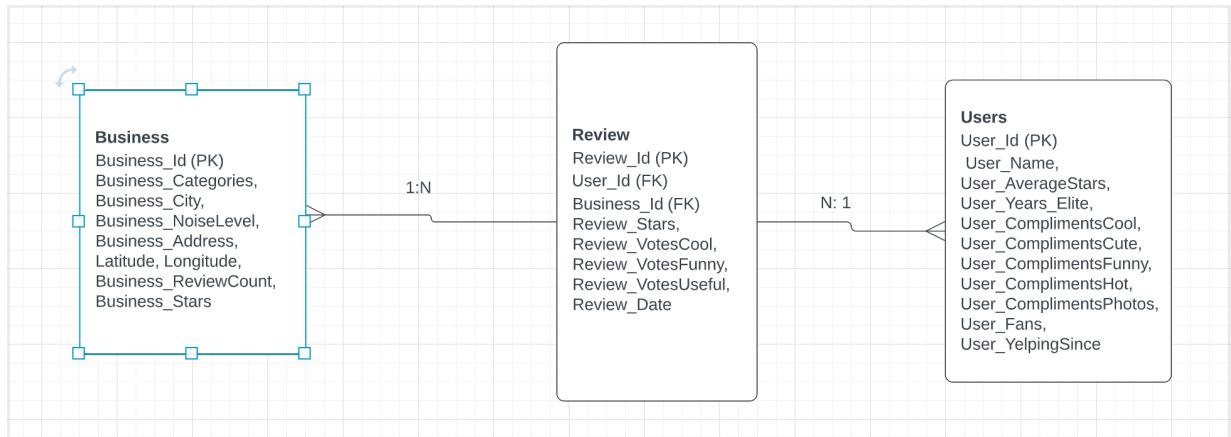
Approach:

1. Download the CSV
2. Export the CSV to a relational DB Mysql to analyze the data.
 1. yelp_data table
3. Separate the data into 3 different subject areas (Business, Users and Reviews)
 1. Business table (Dimension)
 2. Users table (Dimension)
 3. Reviews table (Fact)
4. Data Visualization using Python
 1. Top 10 Cities with the highest Yelp 5 Star reviewers
 2. Plot the number of Cool votes received based on the latitude and longitude in the US map
 3. Most popular star rating
 4. Yelp reviews over the years
 5. Relation between the restaurant's noise level and the Ratings

Tools and softwares used:

MySQL - Database
Python - For Data Ingestion and Visualization
VS Code - IDE
Jupyter Notebook - Data Visualization

Data Model:



Code used to load the CSV to MySQL:

The screenshot shows a Jupyter Notebook interface with the following details:

- File Explorer:** Shows a tree view of files and folders. The current file is `csv_mysql.py`. Other files include `test.py`, `app`, `compressors.py`, `processors.py`, and `setup.py`.
- Code Editor:** The main pane displays Python code for connecting to a MySQL database and reading CSV data into it. The code uses `pandas` and `sqlalchemy` libraries.
- Terminal:** The bottom terminal window shows command-line output from a Mac OS X terminal. It runs `python3 csv_mysql.py` and creates a table named `yelp_data` in the `testdb` database.
- Status Bar:** Shows the current file is `csv_mysql.py`, the line number is 9, the column number is 15, and the file size is 3.95 MB.

```
csv_mysql.py -- devskiller-code-FFYQ-QQF4-EY-S0F
Users > fatusjelinejoseph > Documents > Job > Cornerstone > csv_mysql.py > ...
2 import pandas as pd
3 import pymysql
4 from sqlalchemy import create_engine
5
6 yelp_df = pd.read_csv('/Users/fatusjelinejoseph/Documents/Job/Cornerstone/yelp_data.csv')
7
8 user = 'root'
9 passw = '*****'
10 host = 'localhost'
11 port = 3306
12 database = 'testdb'
13
14 mydb = create_engine('mysql+pymysql://'+ user + ':' + passw + '@' + host + ':' + str(port))
15
16 dbConnection = mydb.connect()
17 tablename = "yelp_data"
18 try:
19     frame = yelp_df.to_sql(tablename, dbConnection, if_exists='fail');
20 except ValueError as vx:
```

(base) fatusjelinejoseph@MacBook-Pro-2 Cornerstone % python3 csv_mysql.py
Table 'yelp' already exists.
(base) fatusjelinejoseph@MacBook-Pro-2 Cornerstone % python3 csv_mysql.py
Table yelp_data created successfully.
(base) fatusjelinejoseph@MacBook-Pro-2 Cornerstone %

DDL queries used:

```
select count(*) from yelp_data;
```

285764

```
Create table business as  
(select distinct Business_Id, Business_Categories, Business_City,  
Business_NoiseLevel,Business_Address,  
Latitude,Longitude,Business_ReviewCount, Business_Stars  
from yelp_data);
```

Create table users as

```
(select distinct User_Id,  
User_Name,User_AverageStars,User_Years_Elite,User_ComplimentsCool,  
User_ComplimentsCute, User_ComplimentsFunny, User_ComplimentsHot,  
User_ComplimentsPhotos, User_Fans, User_YelpingSince  
from yelp_data)  
;
```

Create table review as

```
(select distinct Review_Id, User_Id,Business_Id,Review_Stars,  
Review_VotesCool,Review_VotesFunny, Review_VotesUseful, Review_Date  
from yelp_data)  
;
```

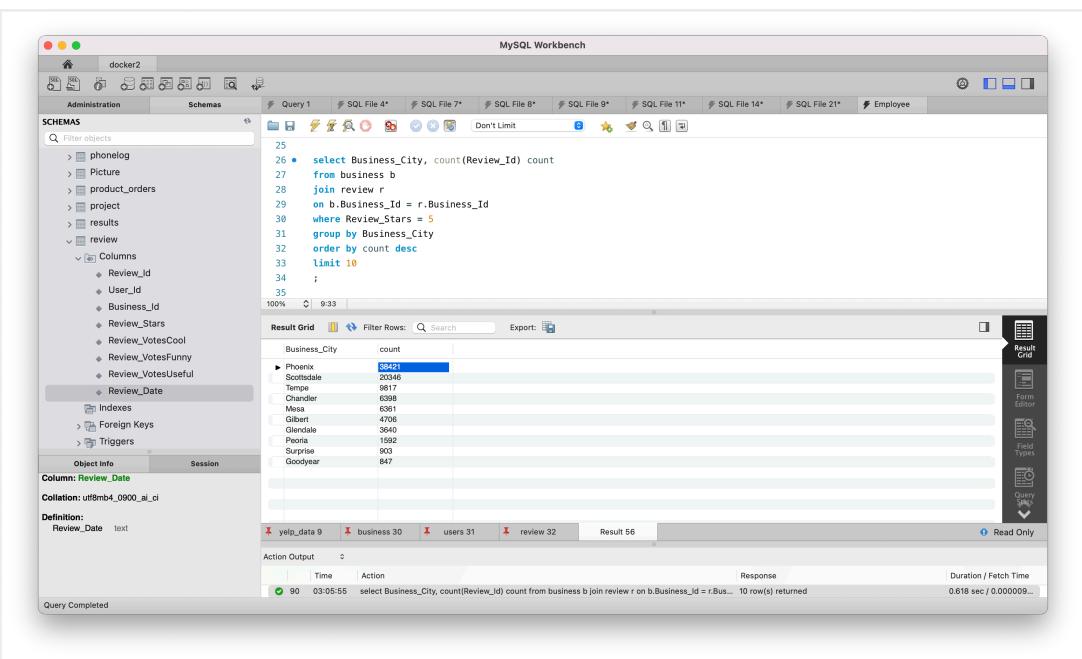
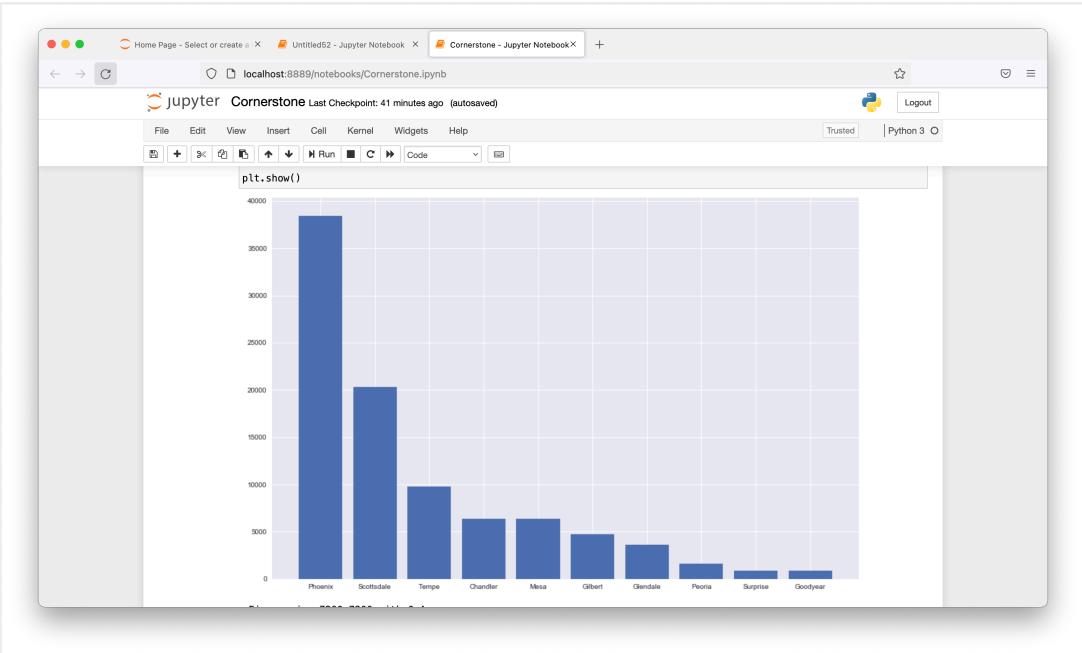
```
select * from business;  
select * from users;  
select * from review;
```

Data Visualization:

Note: For all the data visualizations, I have attached the screenshots of data visualization, query used and the python code

1. Top 10 Cities with the highest Yelp 5 Star reviewers

Key observations : Phoenix has the highest Yelp reviewers with a total of 38421 reviews



```

In [82]: import mysql.connector
import matplotlib.pyplot as plt

mydb = mysql.connector.connect(host='localhost',
                                user='root',
                                password = '*****',
                                db='testdb')

mycursor = mydb.cursor()

# Fetching Data From mysql to my python programe
mycursor.execute("select Business_City, count(Review_Id) count
from business b
join review r
on b.Business_Id = r.Business_Id
where Review_Stars = 5
group by Business_City
order by count desc
limit 10")
result = mycursor.fetchall()

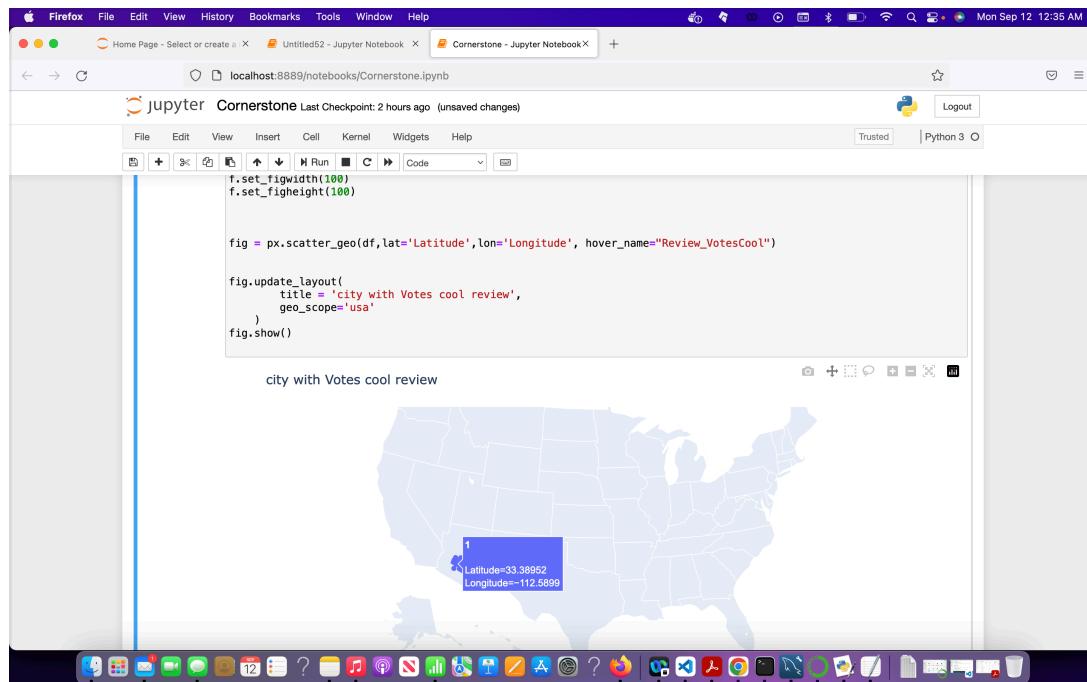
Business_City = []
count = []

for i in mycursor:
    Business_City.append(i[0])
    count.append(i[1])

f = plt.figure()
f.set_figwidth(15)
f.set_figheight(10)
plt.bar(Business_City, count)
plt.figure(figsize=(100,100))
plt.show()

```

2. Plot the number of Cool votes received based on the latitude and longitude in the US map





MySQL Workbench

```

36
37 • select distinct Latitude,Longitude, Review_VotesCool
38   from business b
39   join review r
40   on b.Business_Id = r.Business_Id
41   order by Review_VotesCool desc;
42
43
44
45
46 #3

```

Latitude	Longitude	Review_VotesCool
33.4014048761	-112.07134580	65
33.46585	-112.0691813	65
33.44925	-112.065450	58
33.50242	-111.995515	51
33.553044	-111.995008	51
33.6194983	-111.8983914	51
33.5972864554	-111.977340353	49
33.5592224	-112.0081614	48
33.5577545211	-112.0081614	47
33.4071625343	-111.909640747	45
33.5822419729	-111.883346098	45
33.5822419731	-111.883346098	45
33.5822419731	-111.883346098	44
33.6020429	-111.8767215	44
33.4054537	-112.019776	42

Object Info Session

Column: Review_Date

Collation: utf8mb4_0900_ai_ci

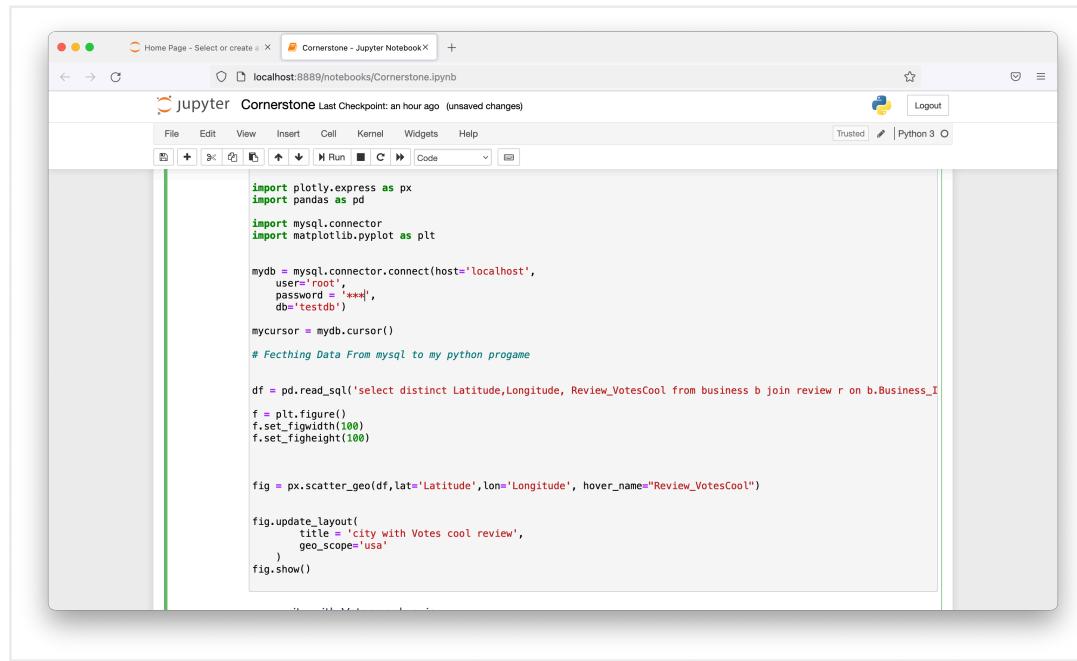
Definition:

Review_Date text

Action Output Response Duration / Fetch Time

91 03:14:47 select distinct Latitude,Longitude, Review_VotesCool from business b join review r on b.Business_Id... 21742 row(s) returned 0.872 sec / 0.055 sec

Query Completed



A screenshot of a Jupyter Notebook interface. The title bar says "Cornerstone - Jupyter Notebook". The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, Help. The toolbar has icons for New, Open, Save, Run, Cell, Code, and Help. The status bar shows "Trusted" and "Python 3". The code cell contains the following Python code:

```
import plotly.express as px
import pandas as pd

import mysql.connector
import matplotlib.pyplot as plt

mydb = mysql.connector.connect(host='localhost',
                                user='root',
                                password = '***',
                                db='testdb')
mycursor = mydb.cursor()

# Fetching Data From mysql to my python programme

df = pd.read_sql('select distinct Latitude,Longitude, Review_VotesCool from business b join review r on b.Business_ID = r.Business_ID')

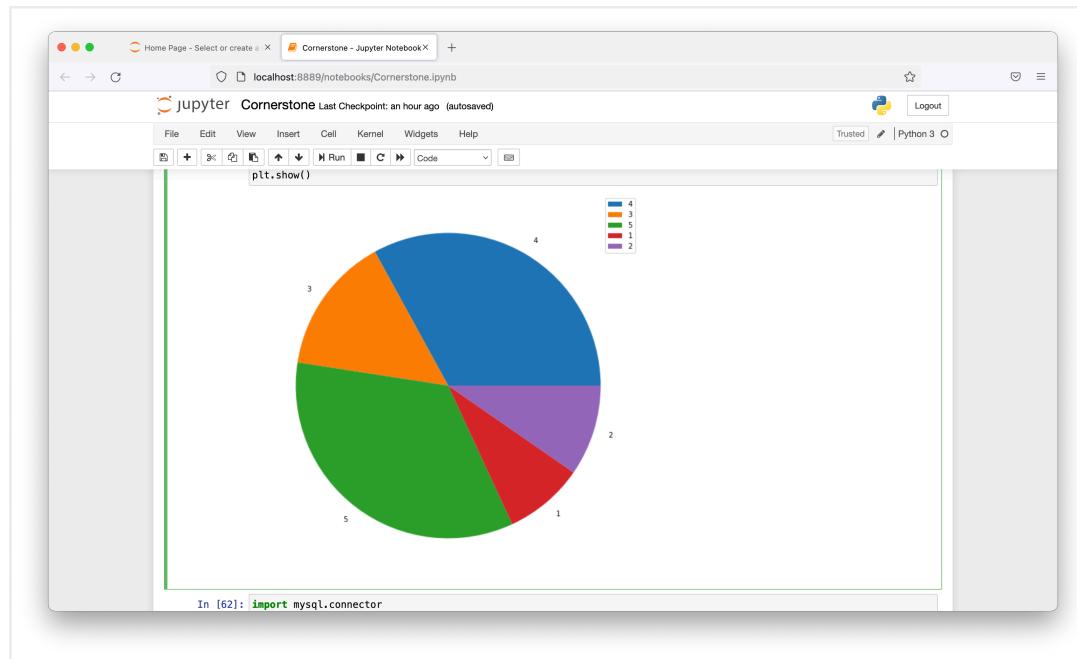
f = plt.figure()
f.set_figwidth(100)
f.set_figheight(100)

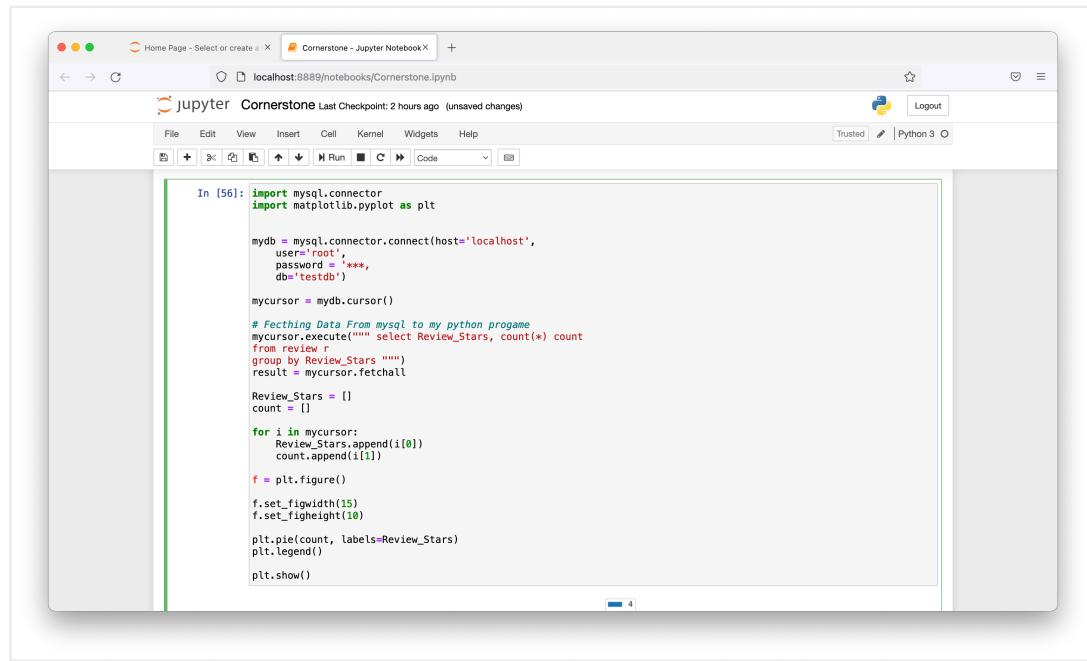
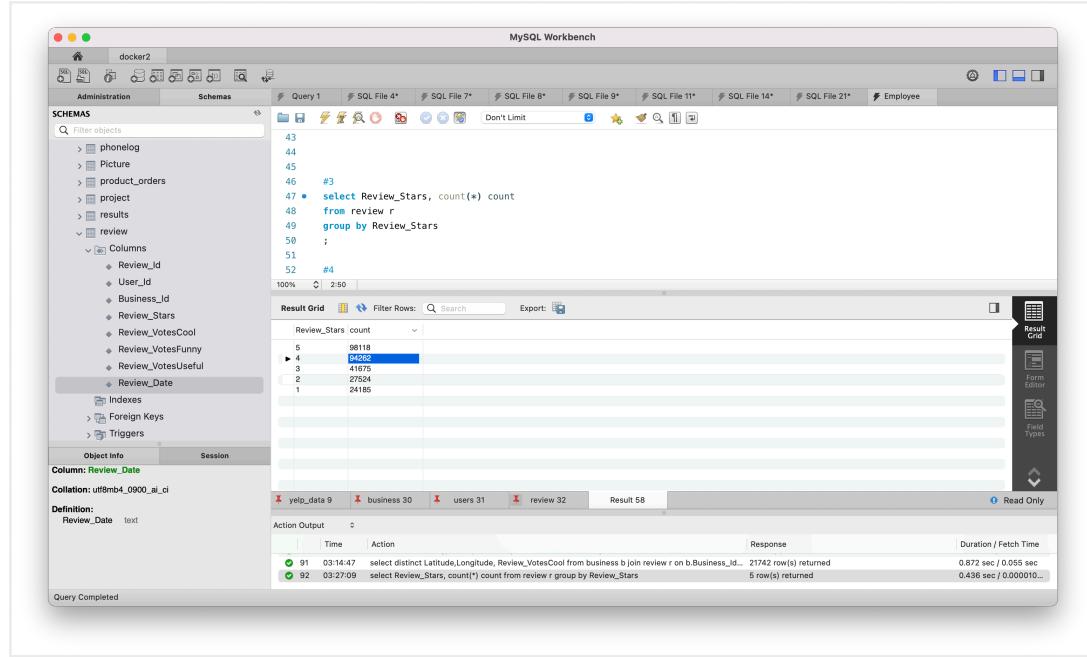
fig = px.scatter_geo(df,lat="Latitude",lon="Longitude", hover_name="Review_VotesCool")

fig.update_layout(
    title = 'city with Votes cool review',
    geo_scope='usa'
)
fig.show()
```

3. Most popular star rating

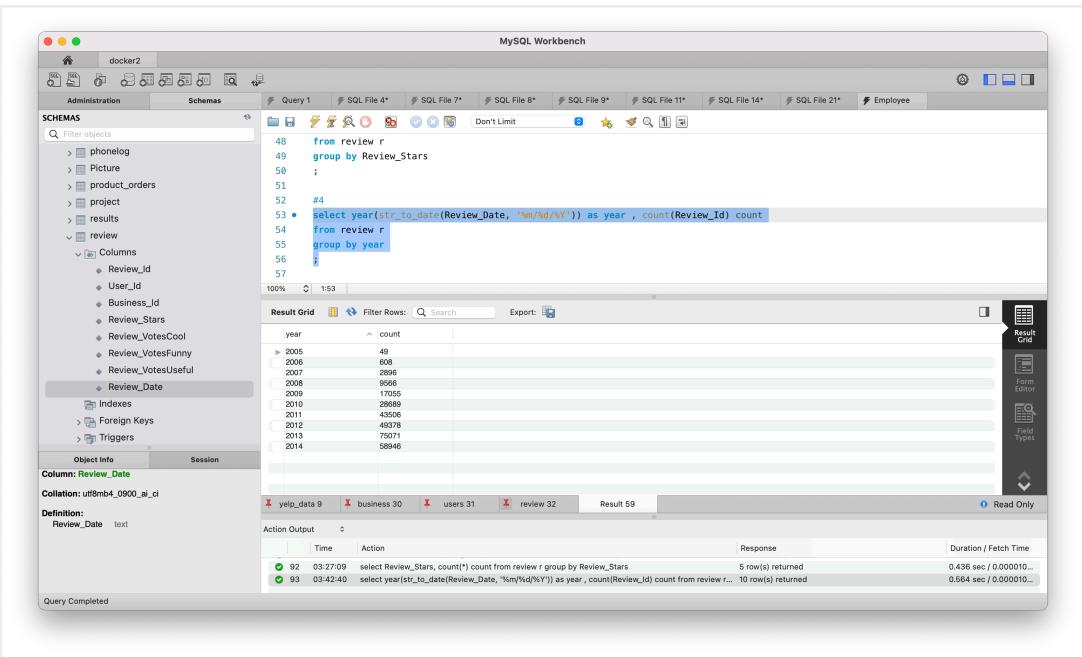
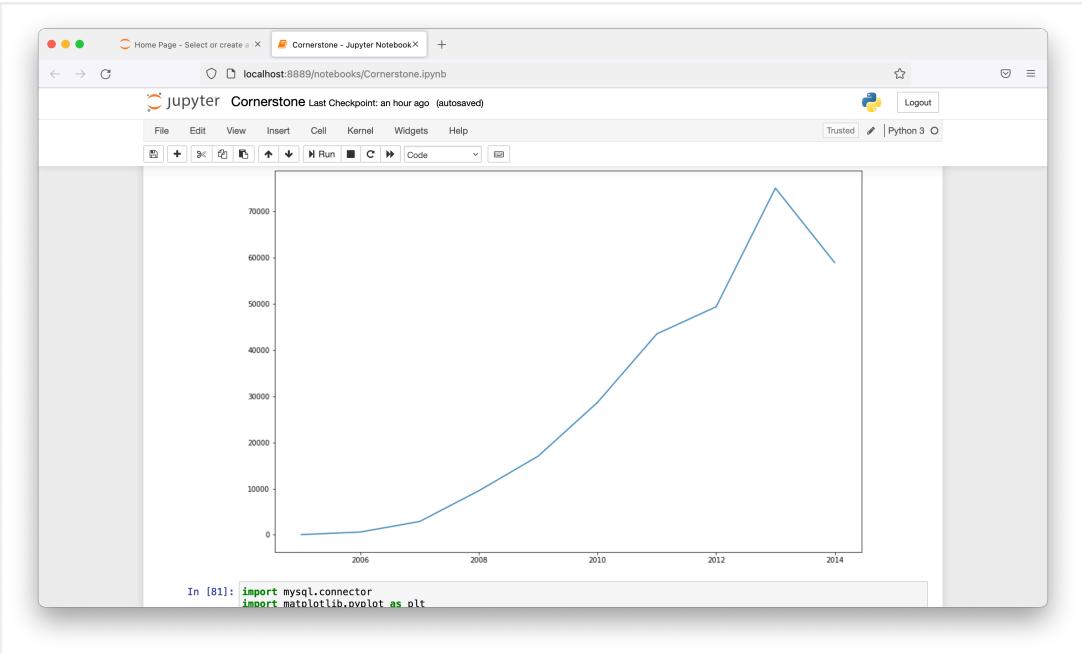
Key Observation: Star 5 rating is the most popular with the count of '98118', followed by 4 Star rating with the count of '94262'





4. Yelp reviews over the years

Key Observations: There is a steady increase in the number of Yelp reviews from 2005 to 2013. But there is a slight drop in 2014.



A screenshot of a Jupyter Notebook interface. The title bar says "Cornerstone - Jupyter Notebook". The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, Help. The toolbar has icons for New, Open, Save, Run, Cell, Kernel, Help, and Code. The status bar shows "Trusted" and "Python 3".

```
In [62]: import mysql.connector
import matplotlib.pyplot as plt

mydb = mysql.connector.connect(host='localhost',
                                user='root',
                                password='*****',
                                db='testdb')

mycursor = mydb.cursor()

# Fetching Data From mysql to my python progame
mycursor.execute("select year(str_to_date(Review_Date, '%m/%d/%Y')) as year , count(Review_Id) count
from review
group by year ")
result = mycursor.fetchall()

year = []
count = []

for i in mycursor:
    year.append(i[0])
    count.append(i[1])

f = plt.figure()
f.set_figwidth(15)
f.set_figheight(10)
plt.plot(year, count)
plt.show()
```

5. Relation between the restaurant's noise level and the Ratings

Key Observation: Majority of the users rated 5 for the Average noisy restaurants when compared to the quiet and loud noisy restaurants

