



Middle East Technical University



Department of Computer Engineering

**CENG 352**  
Database Management Systems  
Spring 2020–2021  
Project 1

---

## 1 Introduction

In this project you are supposed to write many SQL queries on a relational database which will be created by using a dataset. The dataset is taken from Yelp Open Dataset which is available for access to users for personal and non-commercial use. You can reach it from [here](#). Also, you can find the documentation on the same link if you are interested in the dataset. Yelp Dataset contains too much information for this task, so we have done some cleaning to reduce the size of the dataset and converted json files to csv files. You can reach the data which will be used on this project from [here](#). This project has three parts:

- Create the database using the given 'csv' files.
- Write proper SQL queries for certain problems.
- Create triggers and views.

Note that all tasks should be completed using **PostgreSQL**.

## 2 Database Schema

Business(business\_id, business\_name, address, state, is\_open, stars)

Users(user\_id, user\_name, review\_count, yelping\_since, useful, funny, cool, fans, average\_stars)

Friend(user\_id1, user\_id2)

Review(review\_id, user\_id, business\_id, stars, date, useful, funny, cool)

Tip(tip\_id, business\_id, user\_id, date, tip\_text, compliment\_count)

Yelp is a system in which you can find any kind of service such as restaurants, dentists or vet. You can see businesses in your area, their stars (measure for appreciation given by other users) and also you can vote on such businesses that you've visited. The main purpose of Yelp is create an environment for users to share information among themselves.

## 2.1 Foreign Key Constraints

Friend's user\_id1 and user\_id2 references Users' user\_id.

Review's user\_id references Users' user\_id.

Review's business\_id references Business's business\_id.

Tip's business\_id references Business's business\_id.

Tip's user\_id references Users' user\_id.

## 2.2 Explanation of Attributes

Business	<ul style="list-style-type: none"><li>• <b>business_id</b>: Primary key of the business table that stores id's of each business.</li><li>• <b>business_name</b>: Name of that business.</li><li>• <b>address</b>: Address of the business.</li><li>• <b>state</b>: State code that the business is located at.</li><li>• <b>is_open</b>: Boolean value, that indicates whether the business still working or closed completely.</li><li>• <b>stars</b>: Float value, average stars given to that business by users.</li></ul>
Users	<ul style="list-style-type: none"><li>• <b>user_id</b>: Primary key of the Users table that stores id's of each user.</li><li>• <b>user_name</b>: Name of the user.</li><li>• <b>review_count</b>: Indicates how many reviews have been made by that user.</li><li>• <b>yelping_since</b>: The date user started using Yelp system.</li><li>• <b>useful</b>: Count of this useful votes that sent by this user.</li><li>• <b>funny</b>: Count of this funny votes that sent by this user.</li><li>• <b>cool</b>: Count of this cool votes that sent by this user.</li><li>• <b>fans</b>: Number of fans of that user.</li><li>• <b>average_stars</b>: Float value, average stars given by that user.</li></ul>
Friend	<ul style="list-style-type: none"><li>• <b>user_id1</b>: First user's id.</li><li>• <b>user_id2</b>: Friend of the first user.</li></ul>
Review	<ul style="list-style-type: none"><li>• <b>review_id</b>: Primary key of the Review table that stores id's of each review.</li><li>• <b>user_id</b>: User that gave the review.</li><li>• <b>business_id</b>: Business that is reviewed by the user.</li><li>• <b>stars</b>: Given stars to the business by the user.</li><li>• <b>date</b>: Review date.</li><li>• <b>useful</b>: Count of useful tags given by other users to this review.</li><li>• <b>funny</b>: Count of funny tags given by other users to this review.</li><li>• <b>cool</b>: Count of cool tags given by other users to this review.</li></ul>
Tip	<ul style="list-style-type: none"><li>• <b>tip_id</b>: The id of the tip. <b>It is NOT</b> provided in the .csv file, it must be an auto incremented value.</li><li>• <b>business_id</b>: Business that is given a tip.</li></ul>

- **user\_id**: User that gave the tip about the business.
- **date**: Date of the given Tip.
- **compliment\_count**: Compliments that the tip is received.
- **tip\_text**: String, text of the tip.

## 2.3 General Information About Relations

- Reviews are supposed to be detailed reports created by a user for a business. However, since Review texts are too large in terms of bytes we've ignored them.
- Tips are shorter reviews like 'this restaurant's pizza is really good.'. They are similar to the reviews but less detailed.

## 3 Tasks

### 3.1 Task 1 - Creating the Database - 15 pts

Using the given 'csv' files and considering the database schema above create a database using PostgreSQL. For this task you should create a file named 'task1.sql' that contains SQL statements that you've used to create the database.

### 3.2 Task 2 - Advanced SQL Queries - 70 pts

For this task you should be able to write SQL queries for given problems. Please order your queries (from query1 to query14) and create a file named 'task2.sql'. Each question is 5 pts.

1. Find the users whose review\_count is higher than its fans and reviewed at least a business which has more than 3.5 stars. List user\_id, user\_name, difference between review\_count and fans (Ordered by difference between review\_count and fans, user\_id DESC) (1504630 rows)
2. Find the users who have tipped a currently open business located in 'TX' and get compliments more than 2. List all the user-business pairs (this means if there is two businesses that a user get more than 2 compliments list both tips for that user) with user\_name, business\_name, tip\_date, compliment\_count (Ordered by compliment\_count, tip\_date DESC) (38 rows)
3. Find the top 20 users by their friend\_count (more friends are better). List user\_names and count of friends (Ordered by friend\_count, user\_name DESC) (20 rows)
4. Find the users who have given lower stars to a business than its current stars. List distinct users' user\_name (not distinct user\_name distinct user), average\_stars and yelping\_since (Ordered by average\_stars, yelping\_since DESC) (1139971 rows)
5. List open "good" businesses who received the highest number of "good" tips in 2020. "Good" business is the one which has word 'good' in its tip text. List the business\_name, state and the stars (Ordered by stars, business\_name DESC). (4 rows)
6. Find the users who have lower average\_stars than **all** of his/her friends' average\_stars(Consider the average\_stars of Users table.). List user\_name, yelping\_since and average\_stars of such users.(Ordered by average\_stars, yelping\_since DESC) (105952 rows)

7. The average stars of businesses in a state gives us Average of the State. Find top 10 state by highest average stars. List state code and average stars (Ordered by average stars DESC) (10 rows).
8. A tip is a GOOD tip if its compliment\_count is higher than 0 (at least 1 compliment needed). For each year calculate the percentages of GOOD tips if the percentage is higher than 1 percent, this year is called a TIP YEAR. Find all the TIP YEARS. List date of the year and average compliment\_count (Ordered by year ASC) (5 rows).
9. List the names of the user's who **only** reviewed businesses who have stars more than 3.5. (Ordered by user\_name ASC) (832155 rows)
10. Popular businesses are those with more than 1000 reviews. For each such business find the business\_name and average stars for each year. List only those with average stars greater than 3 (in ascending order of years, business\_name) (5601 rows).
11. Find the users who have got more useful votes than cool votes (Note that Users table contains votes sent by the users, now we need to find votes s/he got by reviews). List user\_name, useful, cool and difference of useful and cool. (Order by difference, user\_name DESC) (1003635 rows)
12. List pairs of friends and business names if both friends reviewed the same business with the same stars. List the business id, friends as user\_id\_1 user\_id\_2 and stars (Multiple pair reviews should also returned). You need to return all businesses for each (u1,u2) pair (There might be more than one review that two friends gave to a business. The stars were same before and even both changed their stars they meet in the same point again). However, if (u1,u2) is listed (u2,u1) should not be listed for the same business (Ordered by business\_id, stars DESC). (77157/79393 or 126629 rows)

Example to Clarification:

user_id1	user_id2	business_id	stars
user1	user2	business1	5
user1	user2	business1	4

This is a valid return.

13. **Cross tabulations** Write a single SQL query that computes the statistics present in a cross tabulation over stars and state on Business table. The aggregate reported in the crosstab should simply count businesses, but restrict your attention to the open businesses. You only need to list stars, states and count. (147 rows)
14. **Window functions** List the top 3 users (user\_id, review\_count and fans) and their rank (i.e. 1, 2, or 3) when ranked by review\_count for each grouping formed by number\_of\_fans. Restrict your attention to those with a number\_of\_fans between 50 and 60, inclusive. (33 rows)

### 3.2.1 Task 2 Specifications

Your submission file format (for task2) should be as follows:

```
/* Question 1 */  
SELECT...  
/* Question 2 */  
SELECT...  
...  
/* Question 14 */  
SELECT...
```

You don't have to write your queries on single lines. You are allowed (and encouraged) to write them in multiple lines for better readability. You should not write anything for the unsolved questions. You should have lines for only the solved questions.

## 3.3 Task 3 - Triggers and Views - 15 pts

**Triggers.** It is often useful to have the DBMS perform some actions automatically in response to operations on the database. Write necessary triggers to achieve the following functionalities:

1. We want to keep review count for each user consistent with the number of reviews they have in Users table. Every time a new user review is inserted into Review table, the review\_count in Users table must be incremented. Write a trigger to keep review\_count of the users up-to-date in this manner.
2. Our system does not want Users who write reviews with 0 stars. If there is a review with 0 stars that review should not be inserted into the Review table. Further, delete all reviews and tips of that user.

**Views.** It is inefficient to calculate the number of reviews (review\_count) of a business for each related query. Create a view called 'BusinessCount' that contains the columns business\_id, business\_name, and review\_count. The review\_count should be the total number of reviews made for that business.

## 4 Submission

Send a 'tar.gz' file with your 7 digit student id and starting with 'e' like 'e1234567.tar.gz' that contains all 3 '.sql' file completed in previous tasks. You should submit 'tar.gz' file to Odtuclass before the deadline.