

## CptS -451 Introduction to Database Systems

### Project Milestone-2

Online

Due Date: Thursday March 24, 11:59pm

#### Summary:

In this milestone you will:

- ✓ design the database schema for your application and provide the ER diagram for your database design,
- ✓ translate your entity relationship model into relations and produce DDL SQL statements for creating the corresponding tables in a relational DBMS,
- ✓ populate your database with the Yelp data and get to practice generating INSERT statements and running those to insert data into your DB,
- ✓ write triggers to enforce additional constraints,
- ✓ start developing your application GUI . In Milestone3 you will develop the full final application with all required features.

#### Milestone Description:

You need to complete the following in milestone-2:

- 1) **Revise the database schema that you created in milestone-1. Make sure that your database schema is complete** (i.e., stores all data necessary for the application) and appropriate (i.e., all queries/data retrievals on/from the database can be run efficiently and effectively).

Additional notes about the database schema:

- ✓ **Your business table should include the following attributes** (in addition to the business attributes provided in the business JSON objects):
  - “numCheckins” which stores the number of total check-ins for each business.
  - “reviewrating” which stores the avg star rating for each business.
- 2) (4%) Translate your revised ER model into relations and produce DDL SQL (CREATE TABLE) statements for creating the corresponding tables in a relational DBMS. Note the constraints, including primary key constraints, foreign key constraints, not NULL constraints, etc. needed for the relational schema to capture and enforce the semantics of your ER design. Write your CREATE TABLE statements to a file named “<your-name>\_RELATIONS\_v2.sql”.
- 3) (46%) **Populate your database with the Yelp data.**
  - Generate INSERT statements for your tables and run those to insert data onto your DB. You will use your JSON parsing code from Milestone-1 and use the data you extracted from JSON objects to generate the INSERT statements for your tables.

**Note:** In your business table, initialize “numCheckins” attribute to 0, and “reviewrating” attribute to 0.0 for all businesses (i.e., assign “numCheckins” to 0 in the INSERT statements). In task-5, you will

write an UPDATE statement where you calculate the sum of all check-ins and avg star rating for each business and update the value of “*numCheckins*” and “*reviewrating*” attributes in the business table.

- You may populate your DB with data in 2 different ways:
  - i. You may embed your INSERT statement inside your JSON parsing code and execute them one by one as you generate them. (Sample Python code for connecting to PostgreSQL database and executing SQL statements will be available on Blackboard).
  - ii. Alternatively, you may write the INSERT statements to a SQL script file and then run this (large) script file. (You will find some information about how to generate and run SQL scripts in Appendix-A of this document).

Please note that due to foreign key constraints, the order of tables you insert data to matters. The INSERTs to referenced tables should be run before INSERTs to referencing tables.

Please do not create any INDEXES for your tables until you insert all the data. Indexes can slow down the data insertion.

- 4) After you insert all your data, retrieve the number of tuples in each of your tables and copy/paste that information into a text file. (Simply run the query “SELECT COUNT(\*) FROM tablename” for each table in your database.) Copy the results to “<your-name>\_TableSizes.txt file.
- 5) (6%) Calculate and update the “*numCheckins*”, “*reviewcount*”, and “*reviewrating*” attributes for each business.
  - “*numCheckins*” value for a business should be updated to the sum of all check-in counts for that business. Similarly, “*reviewcount*” should be updated to the number of reviews provided for that business (Note that you will overwrite the values extracted from the JSON data). “*reviewrating*” is the average of the review star ratings provided for each business. You should query the review table to calculate the number of reviews and avg review rating for each business. Similarly, you should query the check-in table to calculate the total number of check-ins. In grading, points will be deducted if you don’t update these values.
  - Write your UPDATE statements to a file named “<your-name>\_UPDATE.sql”.

**(Note:** Running the update statement for the *reviewcount* and *reviewrating* may take a long time. To speed up the process, you may calculate and store the calculated *reviewcount*, *reviewrating*, and *numCheckins* into a temporary table and then update the business table using the values from this temporary table.)

- 6) (24%) Write a 2 page paper where you describe your proposed metrics for classifying businesses as “popular” and “successful”, i.e.,
  - Popular businesses that seem to attract more customers compared to other businesses in the same category.
  - Successful businesses which have been serving the community for a long time and which have loyal customers.

In your paper you need to propose and formulate your own metrics for classifying the businesses into these three groups. You should include the following in your paper:

- Provide a brief description about each of your metrics. Your description should identify all data items you will use in your metrics and how you will combine that information to argue whether a business is expensive, popular, or successful.

- If you made some pre-processing on the data, extracted some information, and used those information in your metrics, make sure to explain them in your paper. For example:
  - You may query your check-in data and calculate the average number of check-ins for each “zipcode , business category” pair. Then you can use this information in your metric for finding successful businesses. If a business has more check-ins than the average number of check-ins within its own zipcode for the same category, then it may be considered successful. The larger the difference is (with respect to the average), the more successful the business is.
- Include all queries you used to process and analyze the data. Also include a brief description that summarizes the goal of each query and what information it extracts.
- If you stored the results of your data-analysis in any tables, include the schemas of those tables.

7) (20%) Start implementing your user interface. You should complete the following features in milestone2.

- Retrieve all distinct states that appear in the Yelp business data and list them (e.g. in a QComboBox).
- When user selects a state, retrieve and display the cities that appear in the Yelp business data in the selected state (e.g. in a list a QListWidget.)
- When a city is selected, retrieve the zipcodes that appear in the Yelp business data in the selected city (e.g. in a QListWidget.)
- When a zipcode is selected, retrieve and display all the business categories for the businesses that appear in that zipcode. (e.g. in a QListWidget.)
- When the user searches for businesses, all the businesses in the selected zipcode will be displayed (e.g. in a QTableWidget).

(Note: The mentioned interface components are only suggested ways to display the data. As long as it is functional and easy to use, any design is acceptable.)

#### Milestone-2 Deliverables - Checklist:

1. The revised E-R diagram for your database design. To create your ER diagram, you can use an ER modeling tool that you get from the web (ERDraw), or your favorite drawing tool (e.g., Visio, Word, PowerPoint). **Should be submitted in .pdf format.** Name this file “<your-name>\_ER.pdf”
2. SQL script file containing all CREATE TABLE statements. Name this file “<your-name>\_RELATIONS\_v2.sql”
3. Text file containing all table sizes. Name this file “<your-name>\_TableSizes.txt”
4. SQL script file containing all UPDATE TABLE statements. Name this file “<your-name>\_UPDATE.sql”
5. The paper explaining your proposed metrics for classifying businesses.
6. The source code of your application. Create a zip file including all your application files (excluding binaries and any data files). Name this file “<your-name>\_application.zip”

Create a zip archive “<your-name>\_milestone2.zip” that includes all the 6 items above. Upload your milestone-2 submission on Blackboard until the deadline.

**You will demonstrate your Milestone-2 to the instructor and the TA after spring break.**

#### References:

1. Yelp Dataset Challenge, [http://www.yelp.com/dataset\\_challenge/](http://www.yelp.com/dataset_challenge/)
2. Samples for users of the Yelp Academic Database, <https://github.com/Yelp/dataset-examples>

3. Yelp Challenge, University of Washington Student Paper 1  
<http://courses.cs.washington.edu/courses/cse544/13sp/final-projects/p08-fants.pdf>
4. Yelp Challenge, University of Washington Student Paper 2,  
<http://courses.cs.washington.edu/courses/cse544/13sp/final-projects/p10-michelmj.pdf>

## Appendix A – How to create and run a SQL script file in PostgreSQL

Simply open a text editor and write all of your queries, separating them with empty lines. Make sure that each query is terminated by a ‘;’.

As an example, suppose that you have the following two queries, and your database name is yelpDB:

**Q1:**

```
select * from reviewTable;
```

**Q2:**

```
select name from businessTable  
where state>'AZ';
```

Your script file should then look like as follows:

```
select * from reviewTable;  
select name from businessTable  
where star>3;
```

You should save this file with a “.sql” extension.

### Running The Script

In the command line, run the following :

```
psql -d yelpdb -U postgres
```

(on Windows: run cmd to open command line window)

(if `psql` is not recognized, you need to add the PostgreSQL installation path to the PATH environment variable. Alternatively, you may browse to the installation directory of PostgreSQL and then run the above command).

- You have to supply a database name to connect to. The above statement assumes your database name is “yelpdb”.

- If you would be running postgresQL with another username (other than `postgres`), replace `postgres` with that username. You will be asked to enter your password for the username you specify.

Assuming that you have saved the script file in the folder `c:\myfolder`, run the following in command line:

```
yelpdb=#> \i ./myscript.sql
```

(update the path of the file if your script file is not in the current directory. )

(The “yelpdb=#” here is the command prompt. Yours will look different depending on your database name.)

The above command will execute all the queries in the `myscript.sql` file .

Check <http://www.postgresqlforbeginners.com/2010/11/interacting-with-postgresql-psql.html> for a brief tutorial about interacting with PostgreSQL in the command line.