

CARNEGIE MELLON UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE
15-415/615 - DATABASE APPLICATIONS
C. FALOUTSOS & A. PAVLO, FALL 2015

Homework 2 (by Jinliang Wei)

Due: hard copy, in class at 3:00pm, on Monday, Sep. 28

Due: tarball, BlackBoard at 3:00pm, on Monday, Sep. 28

Reminders:

- *Plagiarism*: Homework is to be completed *individually*.
- *Typeset* all of your answers whenever possible. Illegible handwriting may get zero points, at the discretion of the graders.
- *Late homeworks*: in that case, please email it
 - to all TAs
 - with subject line: 15-415 Homework Submission (HW 2)
 - and the count of slip-days you are using.

For your information:

- Graded out of **100** points; **2** questions total
- Rough time estimate: *1-2 hours for setting up postgres; approx. 4-8 hours for completing the questions*

Revision : 2015/09/24 00:44

Question	Points	Score
Yelp Reviews: Part 1 - Statistics, etc.	55	
Yelp Reviews - Part 2 - Fraud Detection	45	
Total:	100	

Preliminaries

Cluster machine assignments

Each student is assigned with an andrew cluster machine, and a specific port, for this homework (you might be sharing the machine with other 415/615 students). Your machine and port assignment is on Blackboard, under “grade center”. You use some other machines in the range ghc25..86 (say, if your machine is down) but you **MUST** use your assigned port on any machine that you try, to avoid collisions with other students.

Postgres set-up

Before starting the homework, follow the instructions for setting up PostgreSQL and importing the data we’ll be working with, available at <http://www.cs.cmu.edu/~christos/courses/dbms.F15/hws/HW2/postgresql-setup.html>.

What to deliver: Check-list

Both hard copy, and soft copy:

1. **Hard copy:**

- What: hard copy of your **SQL queries**, plus their **output**.
- When: Sep. 28, 3:00pm
- Where: in class

Contrary to HW1, keep all your answers in one document, but still provide (course#, Homework#, Andrew ID, name).

2. **Soft copy: tar-file:**

- What: A tar file (<your-andrew-id>.tar) with all your SQL code. See next paragraph on how to generate the tarball.
- When: Sep. 28, 3:00pm
- Where: on *Blackboard*, under ‘Assignments’/‘Homework #2’

Details of the tar file. Obtain <http://www.cs.cmu.edu/~christos/courses/dbms.F15/415F15HW2/hw2.tar.gz> - after `tar xvfz`, check the directory `./hw2`: replace the content of each place-holder `hw2/queries/*.sql` file, with your SQL code.

- Your SQL queries will be auto-graded by comparing their outputs to the correct outputs. When comparing each query’s output, the grading script prints “**checking qnn; correct if nothing below ----**”, where *qnn* is the number of the question. Only when your answer is wrong, the grading script prints the difference (computed by `diff`) between your query output and the correct output. The goal is that, with your answers in the directory `queries`, when the grader loads the correct outputs in the directory `outputs` and types `make`, he/she should see no difference.
- For your queries, the **order** of the output columns is important; their names, are not. (our grading script turns off the column headers.)

Naming Convention. The place-holder `queries/.sql` files have the obvious naming convention: For example, the place-holder file for Question 1(a) is named as `q1a.sql`. Each

place-holder file contains a trivial query `SELECT 'hi'`; . Except for `outputs/q1a.txt`, all other outputs just contain `hi`. The only exception is `outputs/q1a.txt`, which contains the output for the (“Sample”) question.

Sanity Check. Before doing any other question, check your answers to the (“Sample”) question (Question 1(a)) to ensure your output matches the formatting (in addition to being correct): Enter your answer into `queries/q1a.sql` and run `make`; the `diff` script should show no difference.

Easy Packaging. For your convenience, we automated the packaging of the submission: When you are done, type `make submission`, and this should automatically generate the tar-file you need to submit. However, it still is **your responsibility** to make sure that all the sql query files are included.

FAQs

- *Q: May we use views?*
- *A: Yes - you may use any feature of SQL that is supported by postgres on the andrew/ghc linux machines. But, make sure that no extra output is generated - remember, we'll do a diff, against our answers.*
- *Q: What if a question is unclear?*
- *A: Our apologies - please post on blackboard; or write down your assumptions, and solve your interpretation of the query. We will accept all reasonable interpretations.*
- *Q: What if my assigned machine is not responding?*
- *A: Our apologies again - as we said earlier, please use another machine, in the range ghc25..86 but with your assigned port number, YYYYYY.*

Question 1: Yelp Reviews: Part 1 - Statistics, etc. . . [55 points]

For this question you will look into Yelp reviews collected from multiple cities across the world, including Pittsburgh. Figure 1 gives the schema of the tables you will use. For example, the tuple in the `review` table

(“1202”, “1001”, 3.0, “2011-08-10”, 2, 1, 0)

means that the #1202 business was reviewed by #1001 user on August 10th, 2010 and got 3 stars. 2 other users think this review is useful, 1 user thinks this review is funny and 0 users think this review is cool.

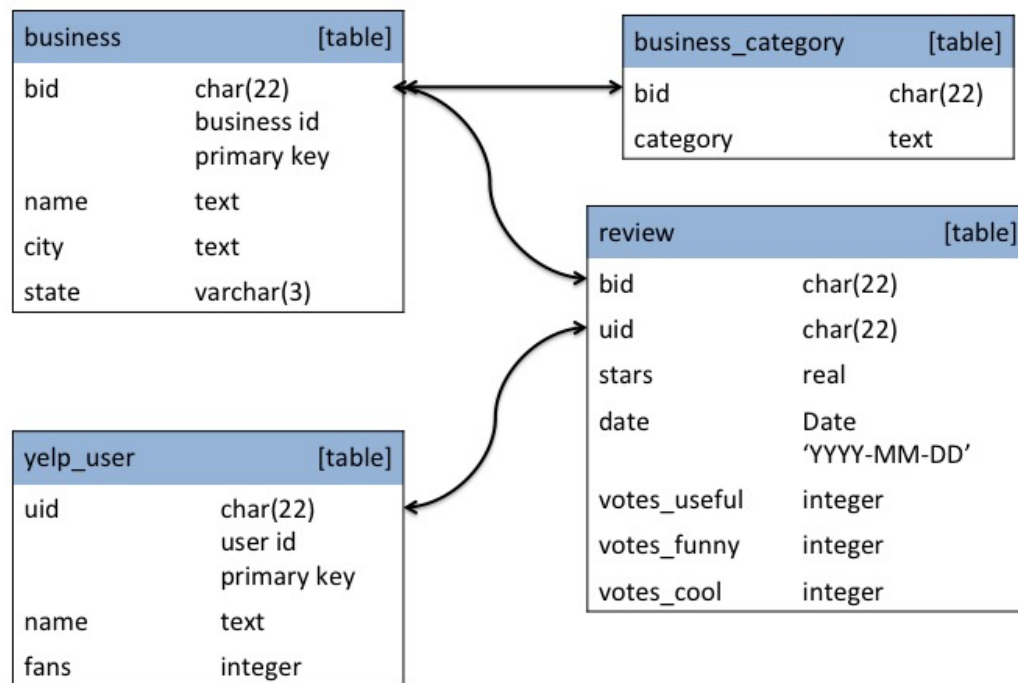


Figure 1: Tables for Yelp Reviews

- [0 points]** (“Sample”): Count the distinct categories of businesses. The purpose of this query is to make sure that the formatting of your output matches exactly the formatting of our auto-grading script.
- [5 points]** (“Warm-up”): Count the number of businesses in Pennsylvania (“PA”).
- [5 points]** (“Misclassified”): Find every business in Pennsylvania that has the word “Coffee” (case-sensitive) in its `name` but is not classified as a coffee place (i.e. has no word “Coffee” in any of its categories). List the business’ `bid` (i.e. business id) and `name` in ascending `bid` order.
- [10 points]** (“Most talked-about bar, per state”): For each state, give the most popular bar in that state. “Popular” means that it has the highest *count* of reviews, among all establishments that have one category as “Bars”, in said state. Sort by state name (increasing). In the remote case of a tie in first place, list all bars, in

ascending **bid** order. More specifically, for each bar, print its **bid** (i.e. business id), **name**, number of reviews and **state**.

- (e) **[10 points]** (“Best breakfast in Pittsburgh, PA”): We define *score* as average number of stars. List the top 10 businesses in Pittsburgh, Pennsylvania¹ for breakfast or brunch. “Top” means that it has the highest score among all businesses that have one category as “Breakfast & Brunch”. Sort by their scores and break ties with number of reviews received (decreasing), then **bid** (increasing). For each business, print its **bid** (business id), **name**, average # of stars, and number of reviews.
- (f) **[5 points]** (“Frequent travelers”): List users that have been to more than 5 **distinct** states. Order by number of states traveled to (decreasing), break ties with **uid** (user id, increasing). For each user, list his/her **uid**, **name** and number of states traveled to.
- (g) **[10 points]** (“Burgers for travelers”): Consider only reviews by travelers (users who have been to more than 2 states). We define *score* as the average number of stars received, considering only reviews by travelers (i.e., *ignore* reviews by everybody else). List the top 5 burger restaurants (i.e. have the highest score) in Pittsburgh, Pennsylvania. “Top” means that it has the highest score among all businesses in Pittsburgh, Pennsylvania that have one category as “Burgers”. Sort by score (highest, first) and break ties with number of reviews by travelers (decreasing), then **bid** (increasing). For each restaurant, print its **bid**, **name**, score, and number of traveler reviews received.
- (h) **[10 points]** (“The most helpful user”): We define as “usefulness count” of a user, say “Smith”, as the sum of the ‘useful’ votes that his/her reviews have attracted. Find the user with the highest such count. If there is a tie in first place, list all such users, ordered by **uid** (i.e. user id). For each user, print his/her **uid** (i.e. user id), **name**, and usefulness count.

¹There is Pittsburgh, California, too - sigh!

Question 2: Yelp Reviews - Part 2 - Fraud Detection [45 points]

Here we look for a few suspicious patterns in the reviews, which may signal fraud.

- (a) **[15 points]** (“Too good to be true”): List Pennsylvania businesses that have more than 10 reviews and all of them are “5 stars”. (This is probably a sign of fraud). Order by the number of reviews received (decreasing), and `bid` (i.e. business id, increasing). For each business, print `bid`, `name` and number of reviews.
- (b) **[15 points]** (“Too many one-hit wonders”): We refer to users that have only voted once as one-hit wonders and such users are likely to be fake accounts. In this query, find the businesses that got more than 150 5-star reviews from such one-hit wonders. For each business, print its `bid` (i.e. business id), `name` and the number of reviews from one-hit wonders. Order by the number of reviews received from one-hit wonders (decreasing), then business name (increasing) then `bid` (increasing).
- (c) **[15 points]** (“Big jump”): Find the businesses whose average rating was raised by more than 1 stars from May 2011 to June 2011. (Such sudden jumps are also suspicious, often due to fabricated reviews). Order your results by the magnitude of the jump (largest, first), then business names (increasing) and `id` (increasing). Print the first 10 if there are more than 10. For each business, print its `bid` (i.e. business id), `name` and jump magnitude.