Peter Antonaros
Math 290 Homework 3

## Exercise 1: (Complete)

**Filled out the google spreadsheet**

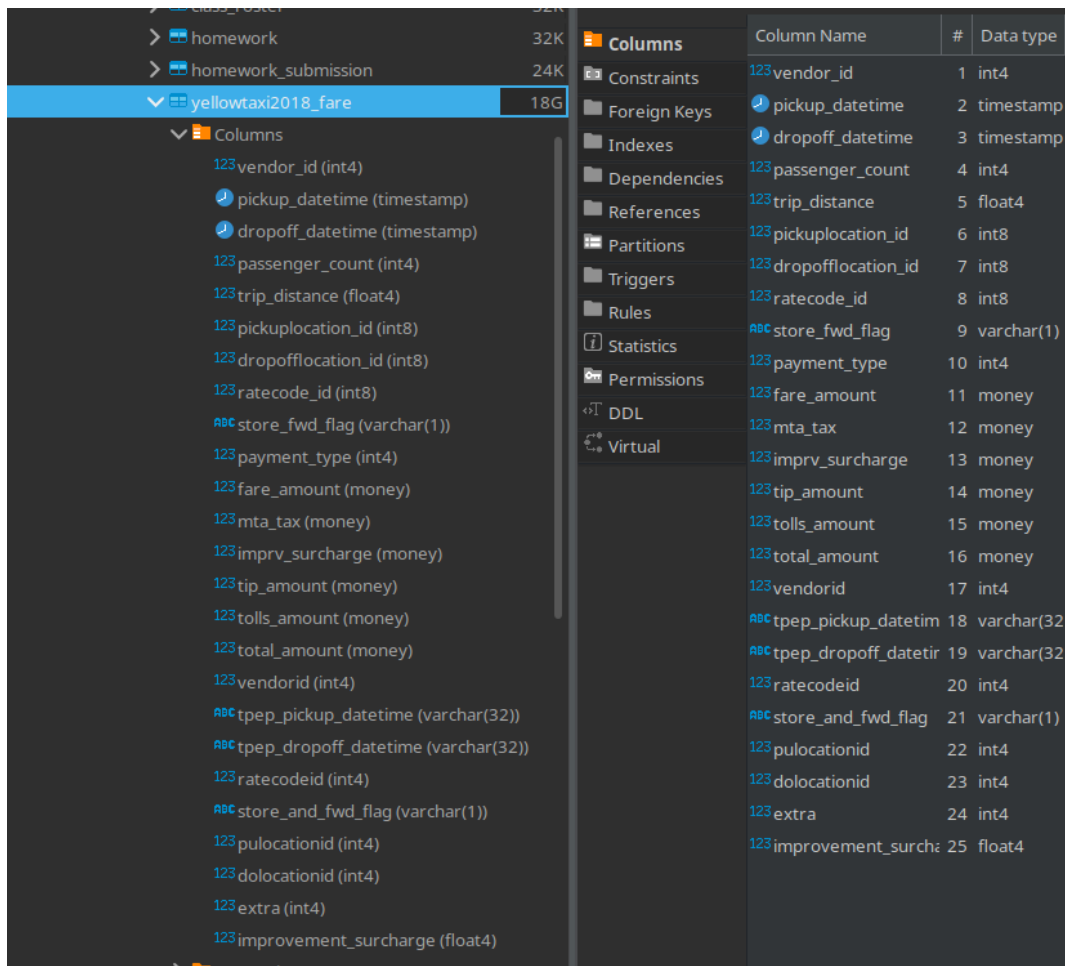ERROR: could not open file "" for reading: Permission denied

**This was solved by changing the linux file permissions to the CSV taxi data set.  I did this with**

**chmod 764 taxiDataset.csv**

**Here the 7 represents (READ+WRITE+EXECUTE) for the OWNER**
**6 (READ+WRITE) for the USERGROUP**
**4(READ)  for the WORLD**

**My copy statement in Dbeaver then worked how I expected it to**
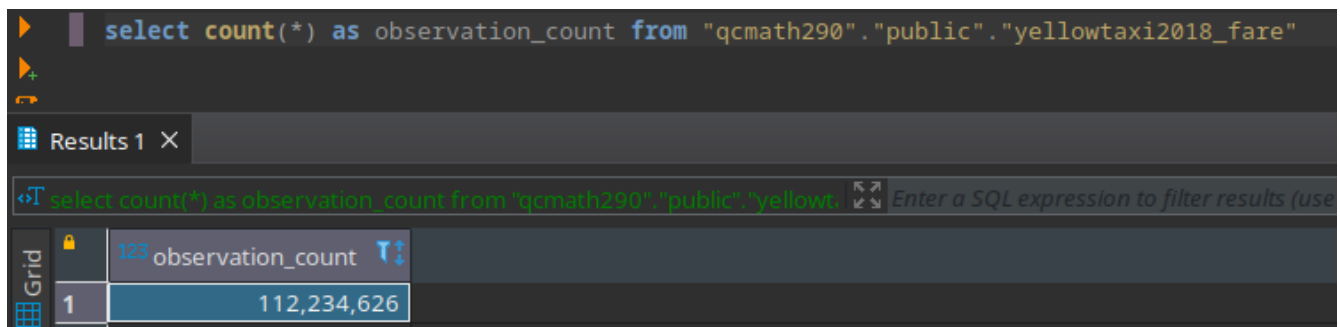
## Exercise 2: (Complete)

Peter Antonaros
Math 290 Homework 3


**Exercise 3: (Complete)**

**select count(*) as observation_count from "qcmath290"."public"."yellowtaxi2018_fare"**
This statement took 59s to execute



**This matches the number of rows provided by the website from where the dateset originated from**

**Exercise 4: (Complete)**

**select count(distinct "vendor_id") from "qcmath290"."public"."yellowtaxi2018_fare.vendor_id"**
This statement took 57s to execute

**There were only 3 distinct numbers.**

**Exercise 5: (Complete)**

**Given the fact that there are only 3 distinct numbers in vendor_id, it would be a bad idea to use this as a primary key. For a primary key we want something that is distinct to that we can efficiently retrieve data. With only 3 distinct numbers, our 112 million rows would be broken into 3 parts each of which would then need a "WHERE" statement to narrow down the 3 row groups into a single row. This would be extremely slow since each "row group" has about 37m rows in it.**

**Exercise 6: (Complete)**

**My first impression was the perhaps the datetime would be a reasonable guess for a primary key since it has the greatest level of granularity, thus the highest chance of being unique.**

select count(distinct "pickup_datetime") from "qcmath290"."public"."yellowtaxi2018_fare.vendor_id"

select count(distinct "dropoff_datetime") from "qcmath290"."public"."yellowtaxi2018_fare.vendor_id"

Peter Antonaros
Math 290 Homework 3

**Neither of these columns had ALL distinct values which in hindsight makes sense, since with so many rides per day I guess our time measurement is bound to have conflict.**

**At this point rather than trying more columns which were obviously not going to be distinct it would be best to add our own primary key that is unique to each row!**

**Exercise 7: (Complete)**

**None of the rows in the vendor_id column are null, which is good from the perspective that we don't need to drop rows without a vendor_id.**

Select * from  "qcmath290"."public"."yellowtaxi2018_fare" where (COLUMS HERE) is null

**I used this statement which is essentially saying get all the rows where there is a null value in any of the columns.  I excluded all the column names in this document as it is a wall of text, but once that is inserted we get our result of NO NULL VALUES IN THE COLUMNS**