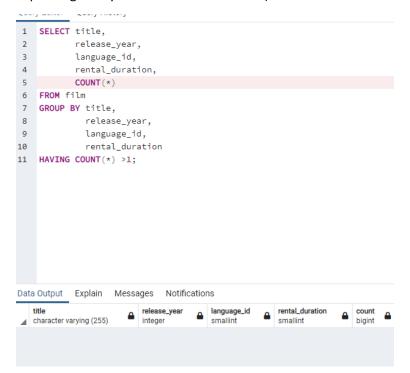
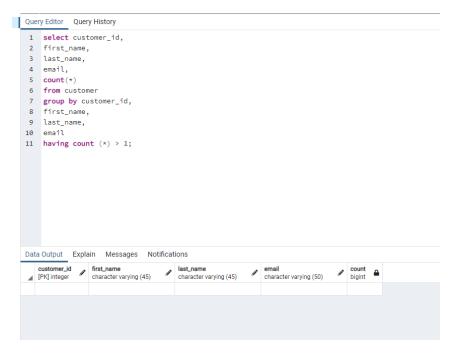
3.6 Summarizing and Cleaning Data

1. **Check for and clean dirty data:** Find out if the **film** table and the **customer** table contain any dirty data, specifically non-uniform or duplicate data, or missing values. Create a new "Answers 3.6" document and copy-paste your queries into it. Next to each query write 2 to 3 sentences explaining how you would clean the data (even if the data is not dirty).



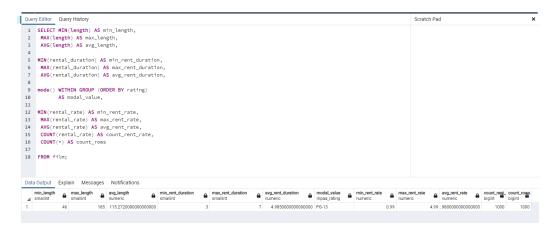


- a. The first step is to use a general select * from [table_name] so the analyst can see all of the factors they have to check for. After, the analyst should select a few column names to check for duplicates amongst them followed by the count(*) and then the table name [select columns, count(*)]. Next, group those same columns and follow it with the 'having count (*) >1;' command. This counts any duplicates in the listed columns.
- b. The above screen shots show how I did the described steps, yielding nothing in my results output. This means that, within these two tables, there are no duplicates.
- 2. Were there duplicates, I would utilize a command to delete the duplicates:

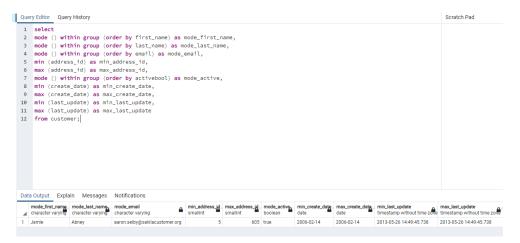
```
3. DELETE
4. FROM tablename
5. WHERE unique_id NOT IN
6. (SELECT MIN(unique_id)
7. FROM tablename
8. GROUP BY col1,
9. col2,
10. col3, ...)
```

11. **Summarize your data:** Use SQL to calculate descriptive statistics for both the film table and the customer table. For numerical columns, this means finding the minimum, maximum, and average values. For non-numerical columns, calculate the mode value. Copy-paste your SQL queries and their outputs into your answers document.

Film:



Customer:



I don't see the relevance of the above statistics for a group related to something like customers. In reality, I'd just spell out the columns and compare it to a [select distinct] or utilize the ['having count (*) >1] to find duplicates. These sort of data don't really need a modal representation. Namely, there's no first name, last name, or email that should occur more than once.

- 12. **Reflect on your work:** Back in Achievement 1 you learned about data profiling in Excel. Based on your previous experience, which tool (Excel or SQL) do you think is more effective for data profiling, and why? Consider their respective functions, ease of use, and speed. Write a short paragraph in the running document that you have started.
 - a. Excel seems faster when one is looking at things like averages, min, and max utilizing the pivot chart. One could simply print [select * from table] the whole category as a .csv. Using the min and max prompts for each category is tedious. The advantage is having a copy/paste in the sql syntax that can be modified in small ways to represent different categories; that is, if you're looking at min(), max(), avg(), in one category, simply copy and paste the syntax and edit it to target the new category/column in the table. Don't forget proper comma placement/removal when doing this.
- 13. Save your "Answers 3.6" document as a PDF and upload it here for your tutor to review.