3.6

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).

Checking for duplicate value from the film table



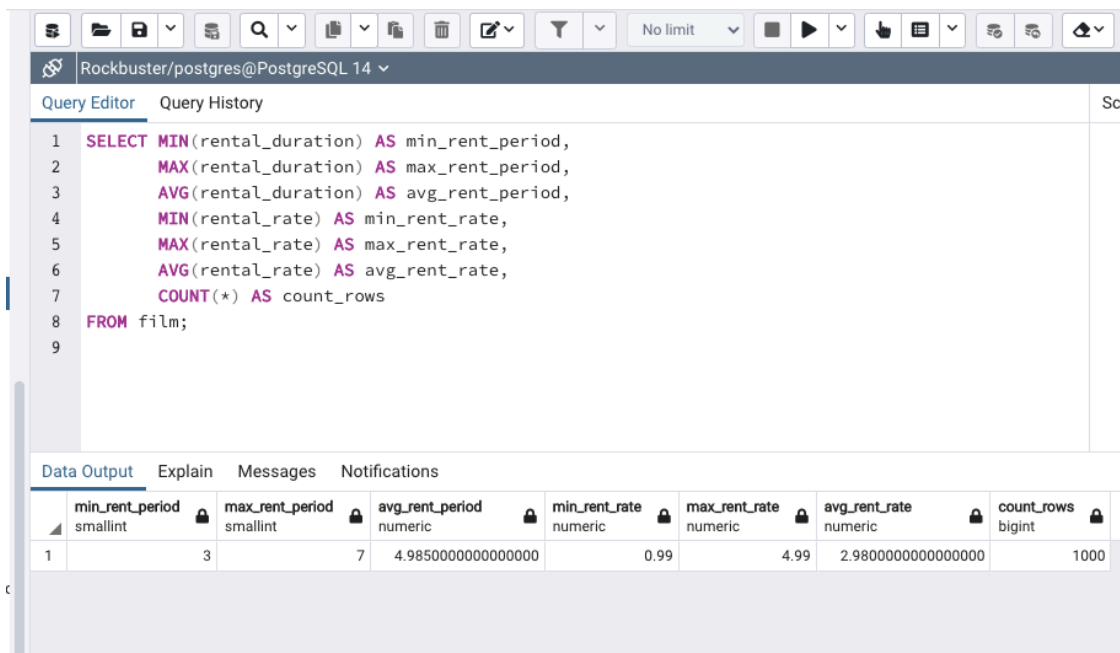Checking for duplicate value from the customer table

There is no returned duplicate value. Essentially, there are two ways of dealing with duplicate value if you have permission to alter the database.

- Create a virtual table "View" where unique records can be selected
- Delete duplicate record from the table or View

However, if altering table is not permitted, we can use GROUP BY or DISTINCT to select unique records.

2. **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.

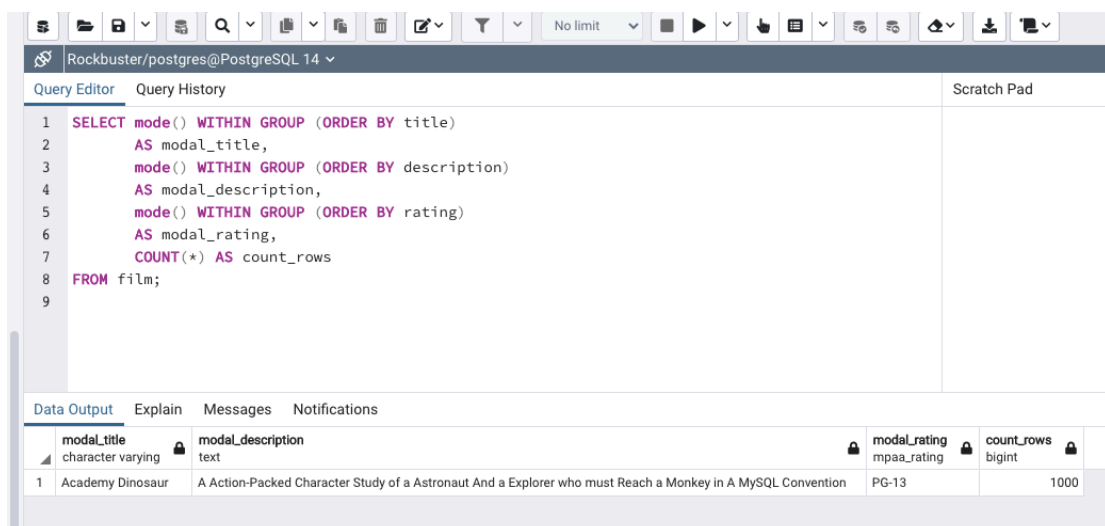Summary for numeric columns in film table



Summary for non-numeric columns in film table

Summary for numeric columns in customer table



Summary for non-numeric columns in customer table



3. **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.

Excel works perfectly with small data size. With this it will be easy to view the data using the pivot table. However, renaming the output (Aliasing) for aggregate column would take more time in excel. Essentially, it is easy to work with huge data in SQL. Using SQL data profiling becomes much easier and faster. Specific result/query to details would be returned in a glance with the right

syntax. Conclusively, SQL works perfectly for data profiling based on the speed and ability to work with huge data.