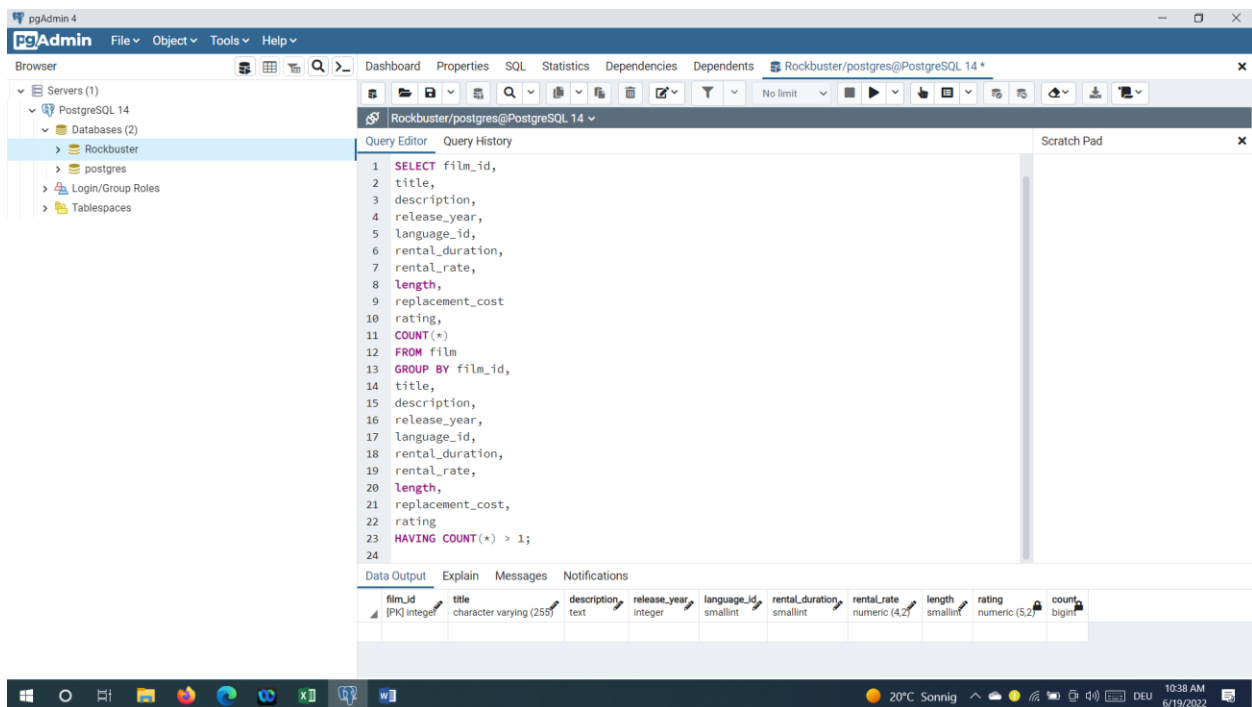# Task 3.6: Summarizing & Cleaning Data in SQL

**Directions**

Rockbuster's database engineers have loaded some new data into the database, and your manager has asked you to clean and profile it. Follow the instructions below to complete their request:

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)** Checking for duplicate value from the film table

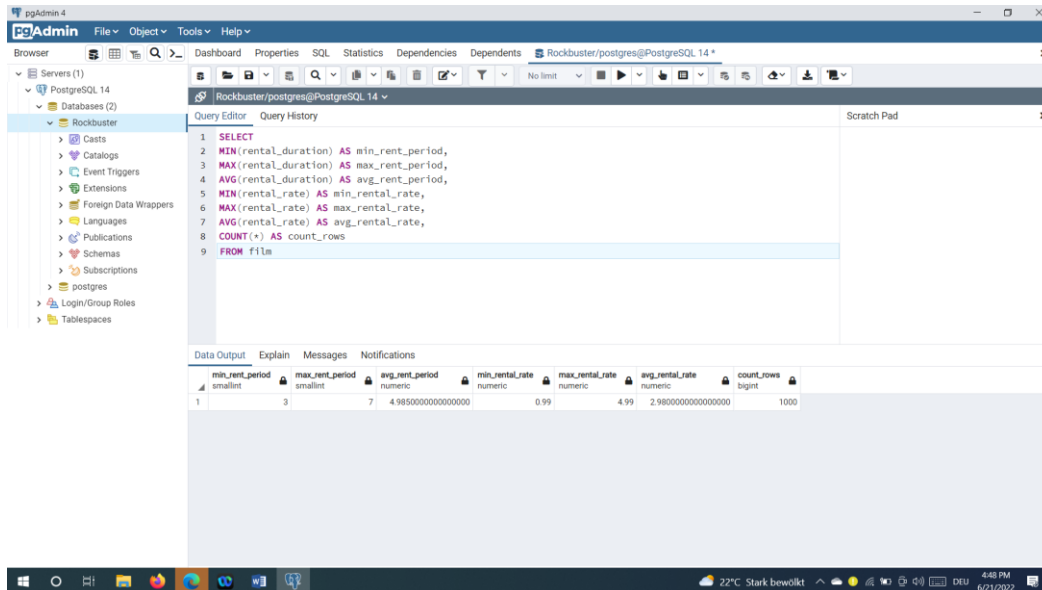b) Checking for duplicate value from the customer table



There are no duplicates in the above made tables. In case that there were any, I can:

• Create a virtual table, known as a "view," where you select only unique records.

• Delete the duplicate record from the table or view

**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.

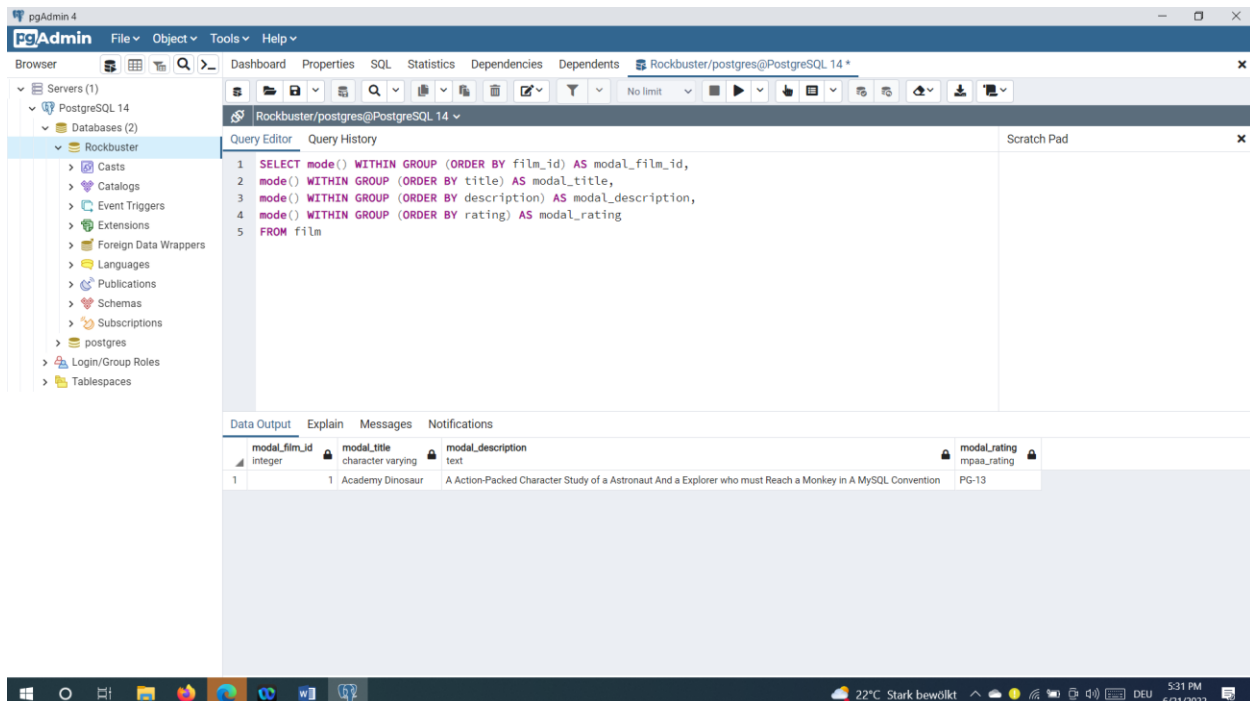a) Summary for numeric columns in film table



b) Summary for numeric columns in film table

## c) Summary for numeric columns in customer table



## d) 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 is a good tool in analyzing the smaller size of the data. On the other side SQL is more efficient and faster when it comes working with the bigger size of the data.