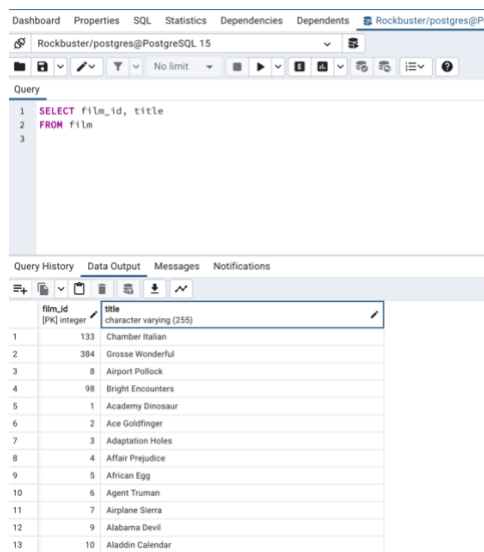# 3.4- Database Querying in SQL

## Answers 3.4.

1. **Refining Your Query:** You need to get data from the "film" table and decide to use the query SELECT * FROM film.

   o You realize that only the "film_id" and "title" columns are needed. Write a new query that selects only those two columns.
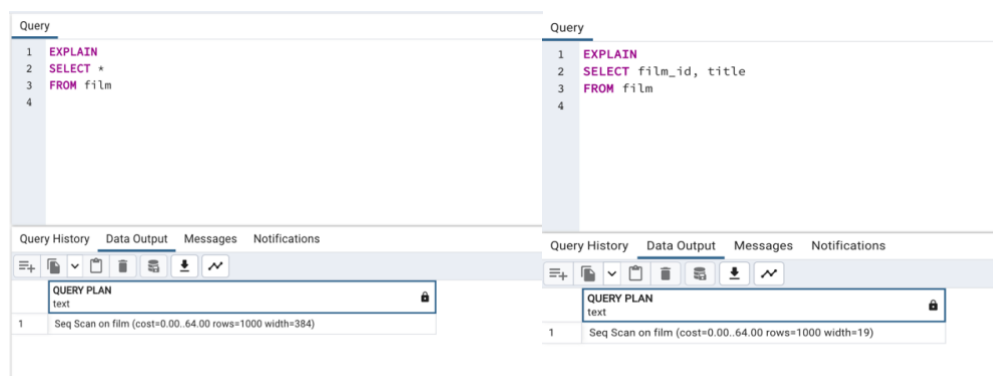


   o **Compare the cost of the original and revised query, and write a few sentences explaining the comparison. Can you suggest any ways to optimize this query?**



   The costs are the same in both queries (cost=0.00..64.00). The important thing about this example is that creating a more specific query can save time and costs when selecting the desired information at once. That is why to optimize it as much as possible before execution. After all, faster scripts mean lower prices.

## 2. Ordering the Data:

- **In the pgAdmin Query Tool, run a query that selects every film from the "film" table, with the movies sorted by title from A to Z, then by most recent release year, and then by highest to lowest rental rate.**



- **Extract the data output of your query into a CSV file for the film collection department to analyze in Excel. To do this, click the button "Save results to file":**

  **DONE**

3. **Grouping Data: The strategy department has asked you the questions below. Write a SQL query to retrieve the correct answers, then extract your results as a CSV file.**

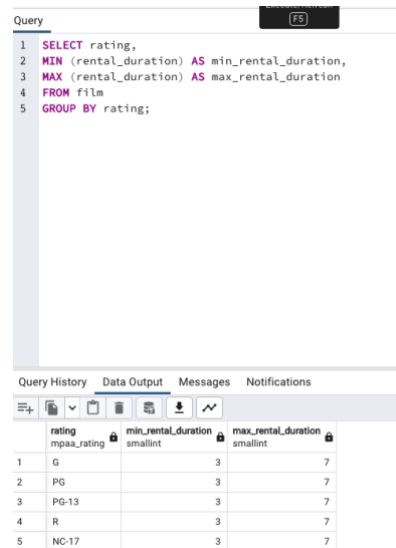   o **What is the average rental rate for each rating category?**

**What are the minimum and maximum rental durations for each rating category?**



4. **Database Migration: Your team has decided to use an external tool to collect data on user behavior in the new Rockbuster Android app. Data collected from this new source will need to be loaded into the data warehouse before you can analyze it.**

   o **Can you outline the procedure for migrating the data and who will be responsible for it?**

   **The migration will be done via ETL (Extract, Transform, Load). Data engineers carry this out.**

   **The steps; Extract: This is the first step, and it involves the collection of data from various data sources Transform: In this step, the extracted data is converted into another format. This could mean calculating ages from dates of birth or combining multiple data points like area codes and telephone numbers to get a contact number, for example. Load: In this step, the transformed data is inserted or loaded into the new database**

   o **What problems do you foresee if you start analyzing the data before it's been loaded into the data warehouse?**

Getting data from different sources could always be problematic. They might be a problem with the data, such as formatting issues or unrelated data, to mention a few.