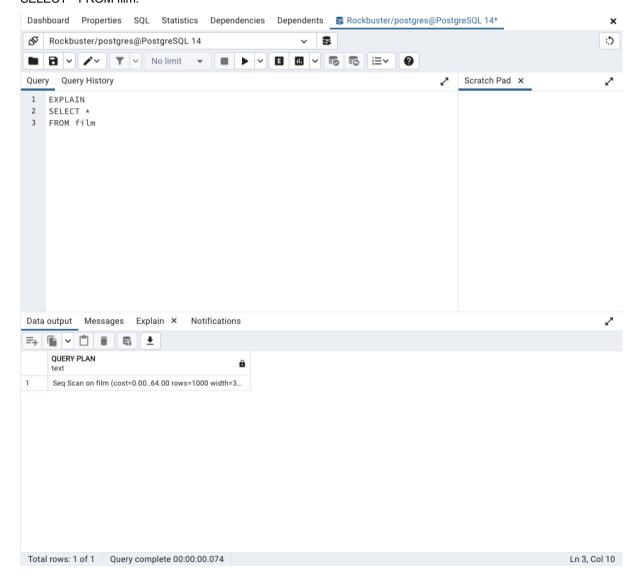# Task 3.4 Database Querying in SQL

By Lee Heng Chuah

It is time to put what you have learned into practice. In this task you will be optimising queries, sorting and grouping data, and reflecting on the database migration process outlined in the exercise.

**TIP:** You will need to use aliases to give appropriate names to the aggregate columns you calculate in step 3 and the bonus task. For a recap of how to assign aliases, see Rules and Best Practices in Exercise 3.3: SQL for Data Analysts
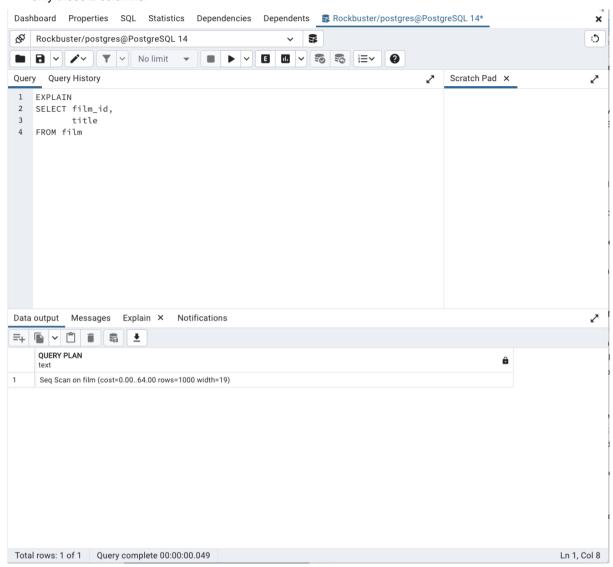
**Directions:**

As you have done for previous tasks, create a new text document for your answers and call it "Answers 3.4".

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

- You realise that only the "film_id" and the "title" columns are needed. Write a new query that selects only those 2 columns.
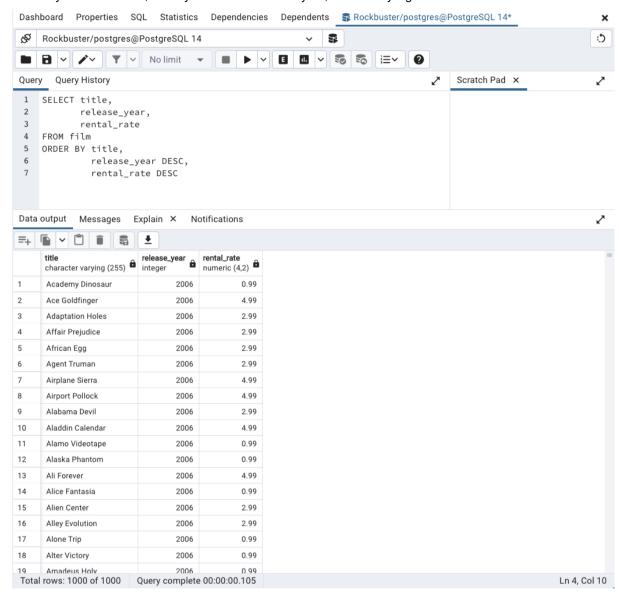


- Compare the cost of the original query and the revised query, and write a few sentences explaining the comparison. Can you suggest any ways to optimise this query?

  Both have same **cost** as we can see from the Data Output (cost=0.00…64.00). It means the "cost" or time, of retuning the first row is 0, but the cost of returning all the rows is 64.

  In this case, the query run time for the first query took 120 msec, while the query run time for the second query is 105 msec. The SQL becomes faster if we use the actual column names or more defined column in SELECT statement instead of than '*'.
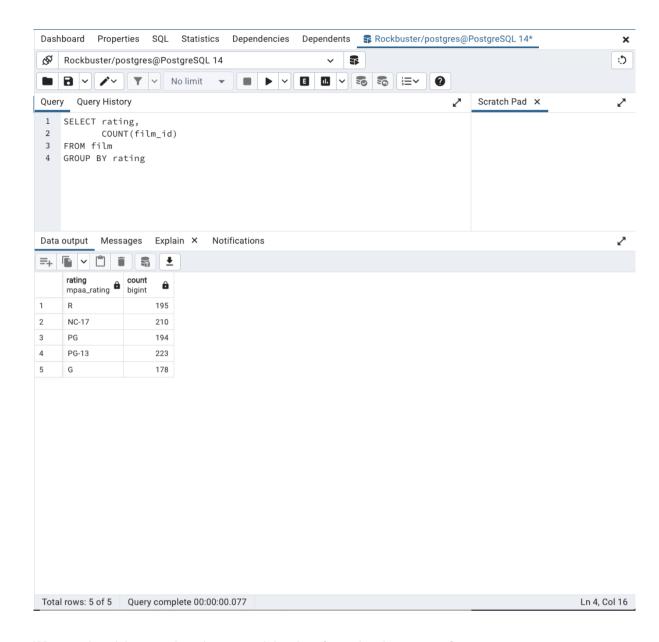
2. **Ordering the Data:**

- 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 analyse in Excel. (You may need to explore how to save your output as a csv file in the Query Tool)

  See separate csv file. - Completed

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

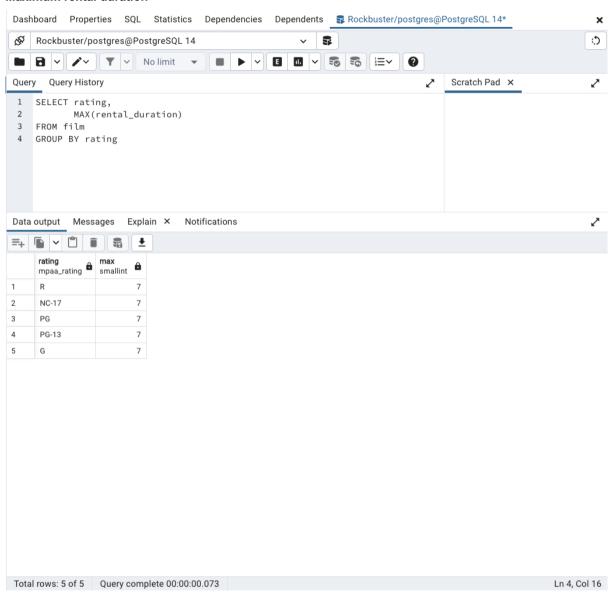- What is the average rental rate for each rating category?

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

  Maximum rental duration = 7

  Minimum rental duration = 3

**Maximum rental duration**

```
1   SELECT rating,
2          MAX(rental_duration)
3   FROM film
4   GROUP BY rating
```

| | rating mpaa_rating | max smallint |
|---|---|---|
| 1 | R | 7 |
| 2 | NC-17 | 7 |
| 3 | PG | 7 |
| 4 | PG-13 | 7 |
| 5 | G | 7 |

**Minimum rental duration**



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

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

     Data engineer would take the responsibility in migrating the data into data warehouse. Data engineer will start collecting the data from multiple data sources. Then, he or she will convert this data into another format which will be aligned with the new data warehouse format. The transformed data will then be inserted or loaded into the new database system. ETL process shall be followed in this case.

   - What problems do you foresee if you start analysing the data before it is been loaded into the data warehouse?

     The data from different systems typically doesn't play together very well. If we start analysing the data before it has been loaded into the data warehouse, the data could be full of inconsistencies and the key relationships across different data sources might be missed. The decision makers will lose

faith in its reliability e.g., not bringing in the information the users need the most, failing to support mission critical reporting workflow and fail to anticipating future data needs.

5. Combine your "Answers 3.4" document and csv files into a single PDF and upload it here for your tutor to review.