

### 3.4: DATABASE QUERYING IN SQL

#### STEP 1

You need to get some data from the “film” table and decide to use the query `SELECT * FROM film`.

- You realize that only the “film\_id” and “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 optimize this query?

The screenshot shows two instances of the pgAdmin Query Tool. The left instance shows a query with three lines: `1 EXPLAIN`, `2 SELECT*`, and `3 FROM film`. The right instance shows a query with seven lines: `1 EXPLAIN`, `2 SELECT*`, `3 FROM film`, `4`, `5 EXPLAIN`, `6 SELECT film_id,title`, and `7 FROM film`. Both instances show a 'Data Output' tab with a 'QUERY PLAN' section. The plan for the first query is 'Seq Scan on film (cost=0.00..64.00 rows=1000 width=384)'. The plan for the second query is 'Seq Scan on film (cost=0.00..64.00 rows=1000 width=19)'. The width of the second query is significantly smaller than the first.

The cost of both queries is the same, however the width is much smaller with the updated query.

The updated query helps us narrow down the columns that we are interested in and therefore show less chaotic table.

A way to optimize the query further would be to either to order the list in ascending or descending order. Or, put a LIMIT on the query to search for the TOP or BOTTOM values that interest us.

#### STEP 2

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.

The screenshot shows a query in the pgAdmin Query Tool: `9 SELECT title,release_year,rental_rate`, `10 FROM film`, `11 ORDER BY title ASC,`, `12 release_year DESC,`, and `13 rental_rate DESC;`. The 'Data Output' tab shows the results of the query. The results are sorted by title (A to Z), then by release\_year (most recent to oldest), and then by rental\_rate (highest to lowest). The results are displayed in a table with columns: title, release\_year, and rental\_rate. The first row is 'Academy Dinosaur' with release\_year 2006 and rental\_rate 0.99. The last row is 'Ali Forever' with release\_year 2006 and rental\_rate 4.99. The status bar at the bottom indicates 'Total rows: 1000 of 1000' and 'Query complete 00:00:00.579'.

---

### STEP 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.

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

| Query       |                                              | Query History       |   |   |  |
|-------------|----------------------------------------------|---------------------|---|---|--|
| 15          | SELECT rating,                               |                     |   |   |  |
| 16          | AVG(rental_rate) AS avg_rental_rate,         |                     |   |   |  |
| 17          | MIN(rental_duration) AS min_rental_duration, |                     |   |   |  |
| 18          | MAX(rental_duration) AS max_rental_duration  |                     |   |   |  |
| 19          | FROM film                                    |                     |   |   |  |
| 20          | GROUP BY rating                              |                     |   |   |  |
| Data Output |                                              | Messages            |   |   |  |
| rating      |                                              | avg_rental_rate     |   |   |  |
| mpaa_rating |                                              | numeric             |   |   |  |
|             |                                              | min_rental_duration |   |   |  |
|             |                                              | smallint            |   |   |  |
|             |                                              | max_rental_duration |   |   |  |
|             |                                              | smallint            |   |   |  |
| 1           | PG                                           | 3.0518556701030928  | 3 | 7 |  |
| 2           | R                                            | 2.9387179487179487  | 3 | 7 |  |
| 3           | NC-17                                        | 2.970952380952381   | 3 | 7 |  |
| 4           | PG-13                                        | 3.034843049327354   | 3 | 7 |  |
| 5           | G                                            | 2.888876404494382   | 3 | 7 |  |

---

### STEP 4

---

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.

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

A procedure called ETL should be followed. Data engineers should be responsible for this task with support from data analysts.

- EXTRACT - Gather the data from new Rockbuster Android APP.
- TRANSFORM - Convert the data into suitable format.
- LOAD - Transformed data is inserted or loaded into the new database.
- What problems do you foresee if you start analyzing the data before it's been loaded into the data warehouse?

If we do not follow the ETL procedure, we will have to clean and analyze the data sources separately and have them in different format which will make our work less efficient.

Having all the data collected, transformed and loaded into one singular data warehouse makes the data more accessible and allows us to work faster and more efficient.