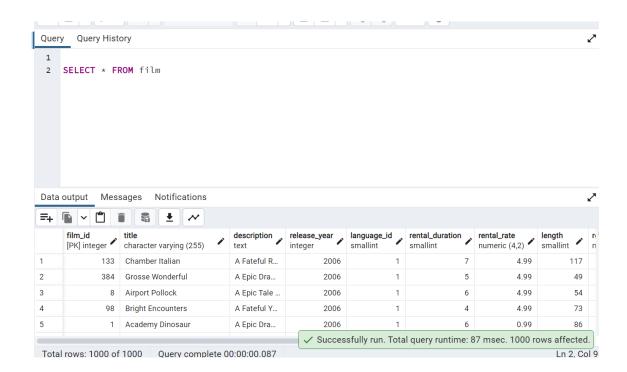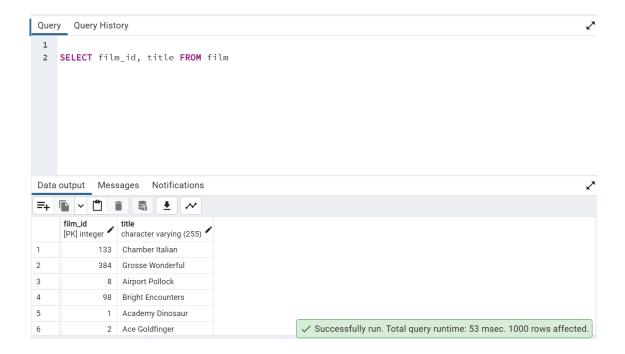# 3.4: Database Querying in SQL

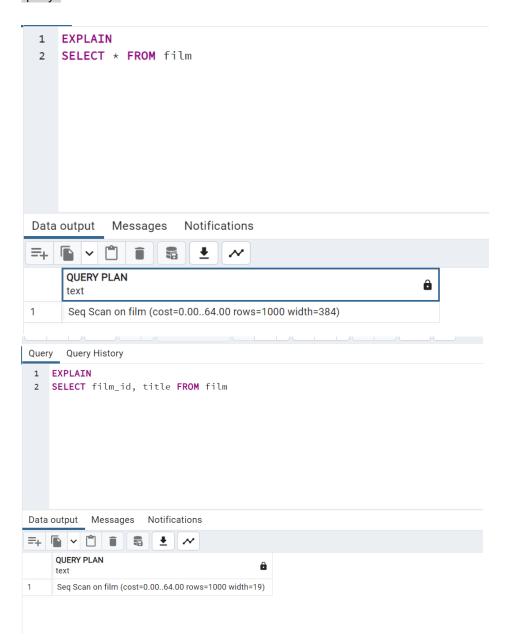1.

Refining Your Query: You need to get some 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 2 columns.

o 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?

```
1  EXPLAIN
2  SELECT * FROM film
```

Data output    Messages    Notifications

| | QUERY PLAN text | |
|---|---|---|
| 1 | Seq Scan on film (cost=0.00..64.00 rows=1000 width=384) | |

Query    Query History

```
1  EXPLAIN
2  SELECT film_id, title FROM film
```

Data output    Messages    Notifications

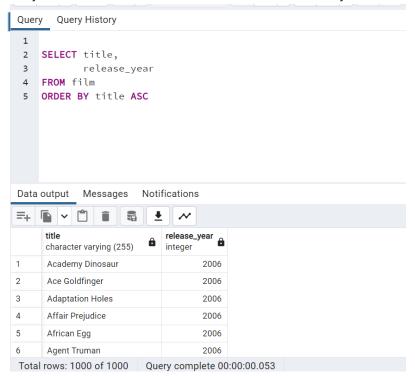| | QUERY PLAN text | |
|---|---|---|
| 1 | Seq Scan on film (cost=0.00..64.00 rows=1000 width=19) | |

Of course if we reduce the information we are asking for the query should be optimized and answer faster. We can see as the width has reduced from 384 to 19. The more accurate we build the query, the more optimized will be the answer
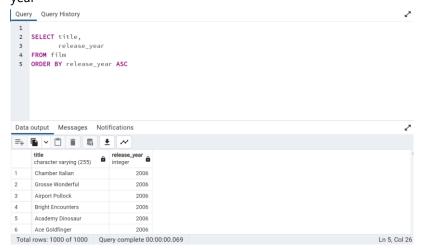
2. Ordering the Data:
   o In the pgAdmin Query Tool, run a query that selects every film from the "film" table, with the movies sorted by title from A to Z, then by most recent release year, and then by highest to lowest rental rate.
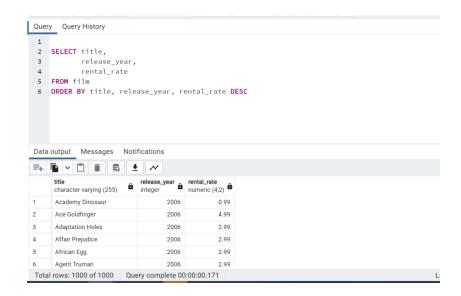
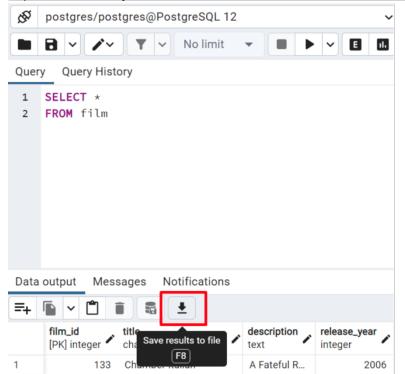o every film from the "film" table, with the movies sorted by title from A to Z

Query    Query History

```
1
2  SELECT title,
3          release_year
4  FROM film
5  ORDER BY title ASC
```

Data output    Messages    Notifications

| | title<br>character varying (255) 🔒 | release_year 🔒<br>integer |
|---|---|---|
| 1 | Academy Dinosaur | 2006 |
| 2 | Ace Goldfinger | 2006 |
| 3 | Adaptation Holes | 2006 |
| 4 | Affair Prejudice | 2006 |
| 5 | African Egg | 2006 |
| 6 | Agent Truman | 2006 |

Total rows: 1000 of 1000    Query complete 00:00:00.053

o every film from the "film" table, with the movies sorted by most recent release year

Query    Query History

```
1
2  SELECT title,
3          release_year
4  FROM film
5  ORDER BY release_year ASC
```

Data output    Messages    Notifications

| | title<br>character varying (255) 🔒 | release_year 🔒<br>integer |
|---|---|---|
| 1 | Chamber Italian | 2006 |
| 2 | Grosse Wonderful | 2006 |
| 3 | Airport Pollock | 2006 |
| 4 | Bright Encounters | 2006 |
| 5 | Academy Dinosaur | 2006 |
| 6 | Ace Goldfinger | 2006 |

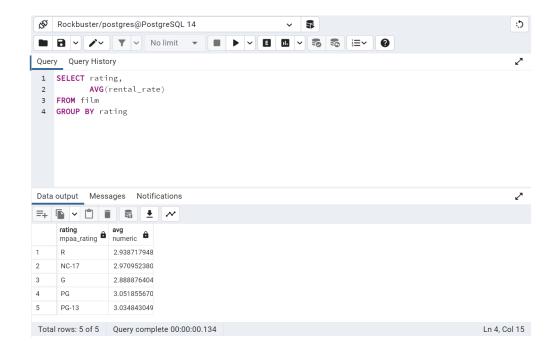Total rows: 1000 of 1000    Query complete 00:00:00.069    Ln 5, Col 26
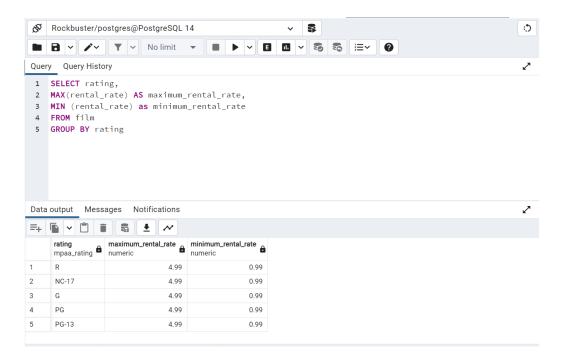
o everything together:

o    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":



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?

o   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?

      As it´s a new data base to add to the previous information we already have on the warehouse, we should analize and normalize so we can after that use and mix

the information to give an answer and correlate to the insights we may extract also from this new new source

This migration should follow the classic procedure of ETL:
- **Extract:** The first step involves collecting the data from multiple data sources.
- **Transform:** During 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:** At this point 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?
  Precisely if we haven´t done the previous analysis and normalization, we can find incoherences on the information so we can´t extract conclusion from it

5. Save your "Answers 3.4" document as a pdf (with screenshots) and your csv files as a single .xlsx Excel file and upload it here for your tutor to review.