# ANSWERS 3.6

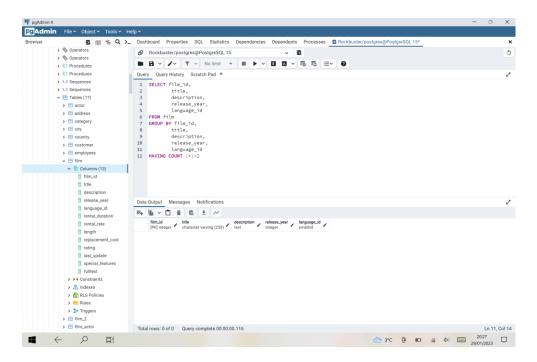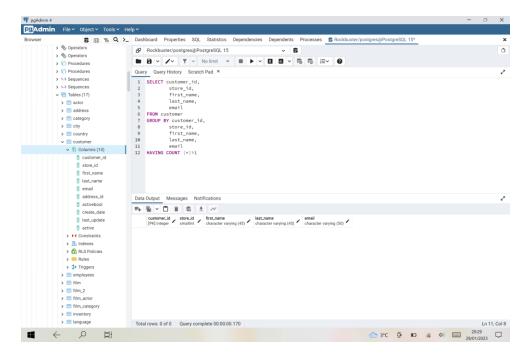1. **Check for and clean dirty data:** Find out if the film table and the customer table contain any dirty data, specifically non-uniform or duplicate data, or missing values.

## FOR FINDING DUPLICATE DATA
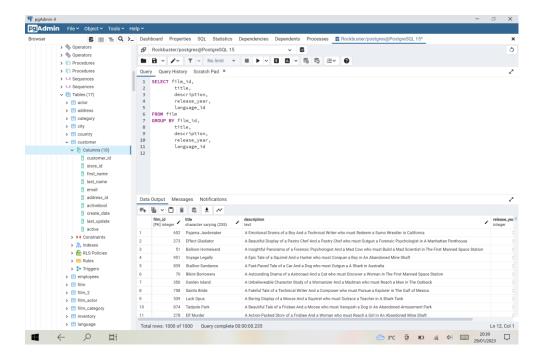
### FILM TABLE



### CUSTOMER TABLE

# ANSWERS 3.6

Both the film and customer tables have no duplicate records, but it they did we could use 2 methods to clean the data:-

a) We could create a virtual VIEW table, where we could use only unique values

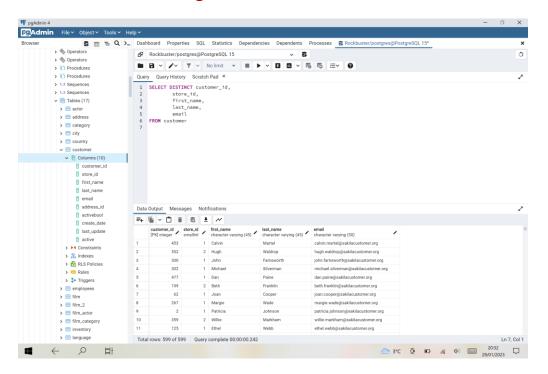b) We could delete all the duplicate values and then use the existing data.

## FOR FINDING NON-UNIFORM DATA

We can either use DISTINCT or GROUP BY to get unique records. Both the table have unique data, but if there was non-uniform data in any column, then we could use UPDATE to get uniform values.

## FILM TABLE USING GROUP BY

# ANSWERS 3.6

## CUSTOMER TABLE using DISTINCT



## FOR FINDING MISSING DATA

This data has no missing data.  Missing values are often null, empty, or replaced with a dummy or default value instead. Like incorrect data, there's no quick fix for missing data. If there is a lot of missing data then we can either ignore the data completely or we can impute average values .

*-Ignore columns with high percentage of missing values*

 SELECT column1, column2, column4

 --column3 ignored in select because it has a lot of missing values

FROM tablename

 *--Imputing missing values with the average value*

 UPDATE tablename

SET=AVG (column1)

WHERE column1 IS NULL

# ANSWERS 3.6

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.

**FILM TABLE (NUMERICAL VALUES)**

SELECT MAX (film_id )AS max_film_id,

        MIN (film_id )AS min_film_id,

        AVG (film_id) AS avg_film_id,

        AVG (film_id) AS avg_film_id,

        MAX (release_year) AS max_release_year,

        MIN (release_year) AS min_release_year,

        AVG (release_year) AS avg_release_year,

        MAX (language_id) AS max_language_id,

        MIN (language_id) AS min_language_id,

        AVG (language_id) AS avg_language_id,

        MAX (rental_duration) AS max_rental_duration,

        MIN (rental_duration) AS min_rental_duration,

        AVG (rental_duration) AS avg_rental_duration,

        MAX (rental_rate) AS max_rental_rate,

        MIN (rental_rate) AS min_rental_rate,

        AVG (rental_rate) AS avg_rental_rate,

        MAX (length) AS max_length,

        MIN (length) AS min_length,

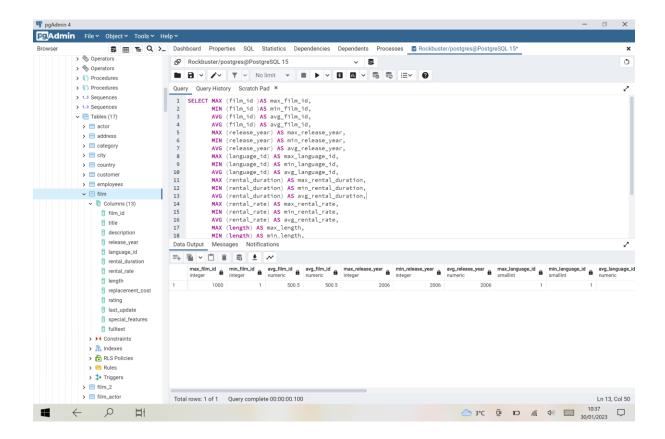        AVG (length) AS avg_length,

        MAX (replacement_cost) AS max_replacement_cost,

        MIN (replacement_cost) AS min_replacement_cost,

        AVG (replacement_cost) AS avg_replacement_cost

FROM film

# ANSWERS 3.6



**FILM TABLE (NON-NUMERICAL VALUES)**

SELECT mode() WITHIN GROUP (ORDER BY title)AS mode_title,
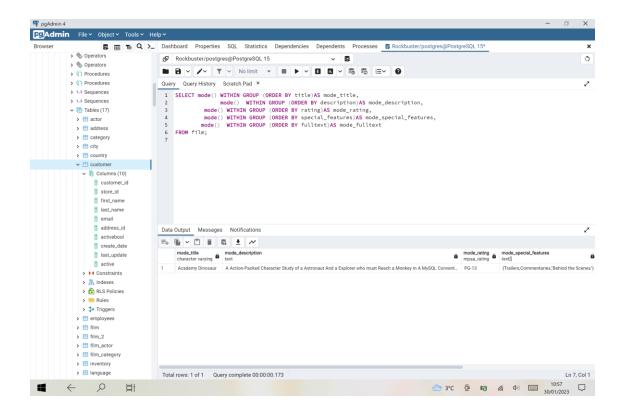
      WITHIN GROUP (ORDER BY description)AS mode_description,
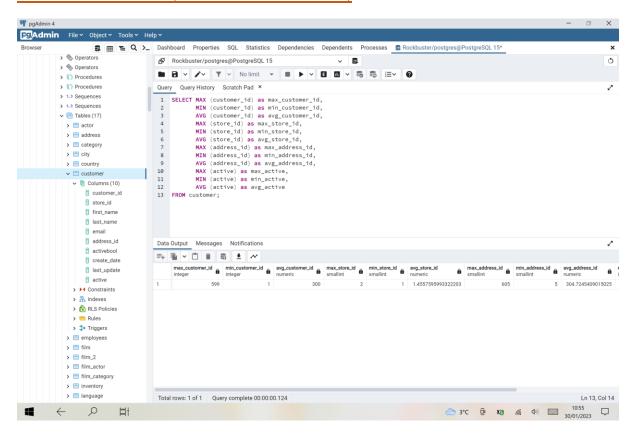
      WITHIN GROUP (ORDER BY rating)AS mode_rating,

      WITHIN GROUP (ORDER BY special_features)AS mode_special_features,
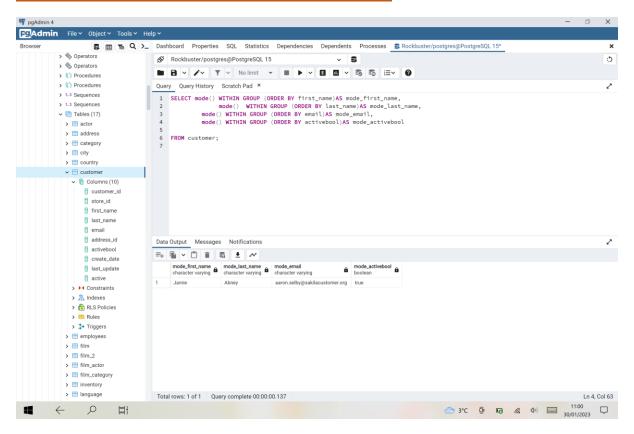
      WITHIN GROUP (ORDER BY fulltext)AS mode_fulltext

FROM film;

# ANSWERS 3.6



## CUSTOMER TABLE (NUMERICAL VALUES)

# ANSWERS 3.6

## CUSTOMER TABLE (NON-NUMERICAL VALUES)



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.

   - In my view, I feel SQL is a lot better as a tool for data profiling. It is irrespective of the data size; we get quicker results without scanning through the complete data.

   - We can just write the relevant queries and find out the dirty data and clean it.

   - With appropriate queries we can find out the unique data.