**Task 6. Clean and summarize data**

**Answers 3.6**

## Contents

### 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**. Next to each query write 2 to 3 sentences explaining how you would clean the data (even if the data is not dirty).
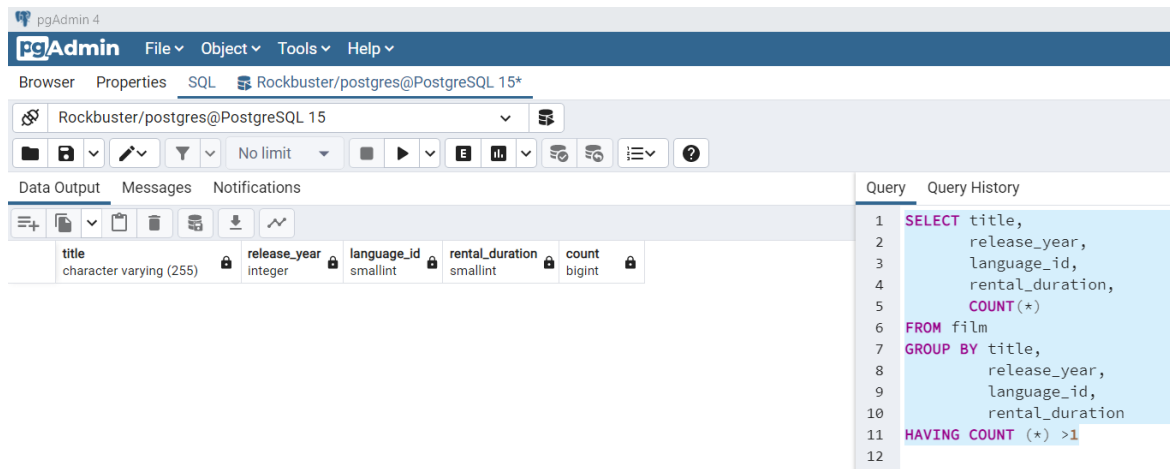
Film table
Looking for Duplicates:

```
SELECT title,
     release_year,
     language_id,
     rental_duration,
     COUNT(*)
FROM film
GROUP BY title,
      release_year,
      language_id,
      rental_duration
```

HAVING COUNT (*) >1


Printscreen:



There are no duplicates returned in the film table. If there were any, I would keep them but use other queries, such as GROUP BY or DISTINCT so unique records could be analysed.

Looking for missing values and non-uniform data:
SELECT title,
    release_year,
    language_id,
    rental_duration,
        replacement_cost,
        rating
FROM film
GROUP BY title,
    release_year,
    language_id,
    rental_duration,
            replacement_cost,
            rating;

Or using the distinct statement instead:



There are no non-uniform and missing values. Also, the data type has been set for all the columns, using constraints, which prevents the entry of non-uniform data.

If they existed, to fix non-uniform values we could use the UPDATE statement to update the values, like this:

UPDATE table_name

SET column_name = 'here list correct value'

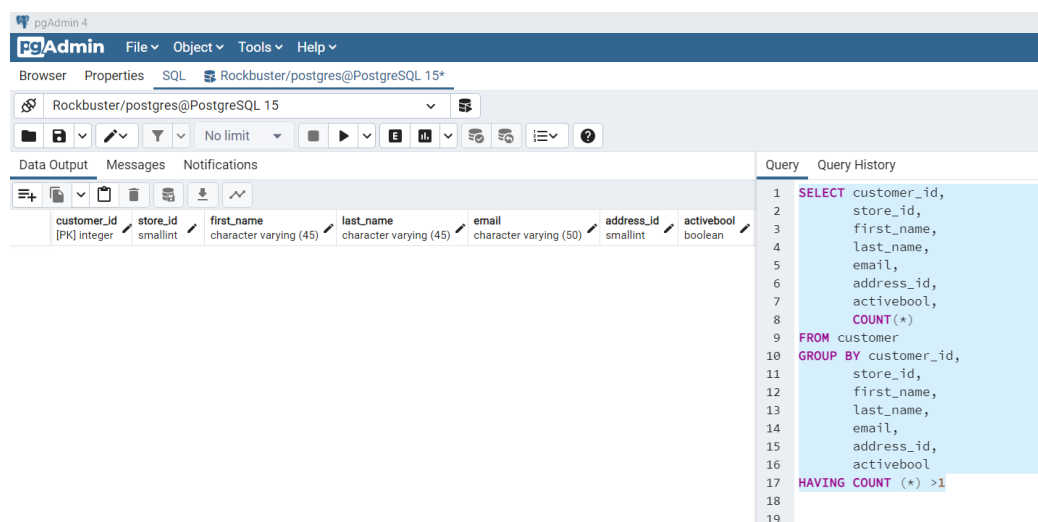WHERE column_name IN (…, here list the non-uniform values that we found)

If there were missing values, we could either ignore the column if there were a lot of values missing or we could impute values, if there were only a few values missing. For instance, if the value is numeric, we could input the mean.

## Customer table
## Looking for duplicates

```
SELECT customer_id,
     store_id,
     first_name,
     last_name,
         email,
         address_id,
         activebool,
     COUNT(*)
FROM customer
GROUP BY customer_id,
     store_id,
     first_name,
     last_name,
         email,
         address_id,
         activebool
HAVING COUNT (*) >1
```
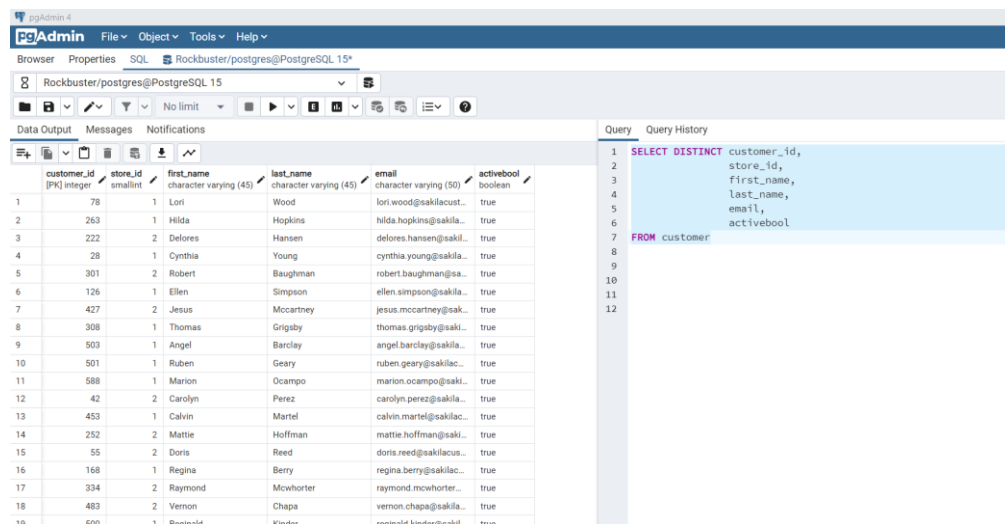
Screenshot:



There are also no duplicate values. If there were, they should be handled very carefully, yet because the table includes unique ID values, it would be possible to use a delete statement to remove the duplicate value. For instance:

DELETE
FROM costumer
WHERE customer_id NOT IN
(SELECT MIN (customer_id)
FROM customer
GROUP BY customer_id,
        store_id)

Look for non-uniform data and missing values:

SELECT DISTINCT customer_id,
        store_id,
        first_name,
        last_name,
            email,
        activebool
FROM customer



There were also no missing values and non-uniform data in the customer table. If there were, the solutions would be similar those in the film table. We could update data, to ensure all values were uniform, assuming we knew the right form of those values. We would either ignore columns with missing values if these were extensive, or input values to columns, if there were only a few missing values.

## 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. Copy-paste your SQL queries and their outputs into your answers document.

Descriptive statistics of film table and customer table

Film Table

Numeric values: film_id, release_year, language_id, rental_duration, rental_rate, length, replacement_cost

film_id:

SELECT MIN (film_id),
    MAX (film_id),
        AVG (film_id)
        FROM film



release_year:



language_id:



rental_duration:

rental_rate:



length:



replacement_cost:



**Non-numeric values**: title, rating, special_features, fulltext

Title:

Statement for MODE:

SELECT mode() WITHIN GROUP (ORDER BY title)
    AS modal_value
FROM film



Rating:



special_features:



Fulltext:

## Customer table
Numeric values: customer_id, store_id, address_id
(note: the activebool and active, are both bolean values, - i.e., true or false)

customer_id:

store_id:

address_id:

## Non-numeric values: first_name, last_name, email

first_name:



last_name



email:



### 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. Write a short paragraph in the running document that you have started.

I think it is very easy to clean data in SQL, once you become comfortable with the statements and how to use them. In terms of data profiling, I also find it easier in SQL, we can check the data type in the columns of the table, while in excel the value could be 'general' or 'numeric' without specifications, so it is not always obvious exactly what the

range of the numbers are for instance. It is fast to perform a quick summary of the data (using descriptive statistics) in SQL, compared to the time it would take to the same calculations in excel.

Yet, I think excel is far superior when it comes to types of statistical analysis that can be done. In this aspect SQL is very basic. That said, using the MODE in SQL is interesting, I don't think we can retrieve this type of result with all text data in excel.