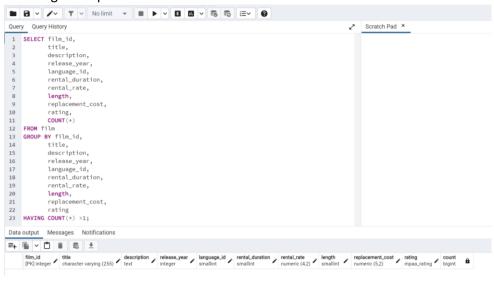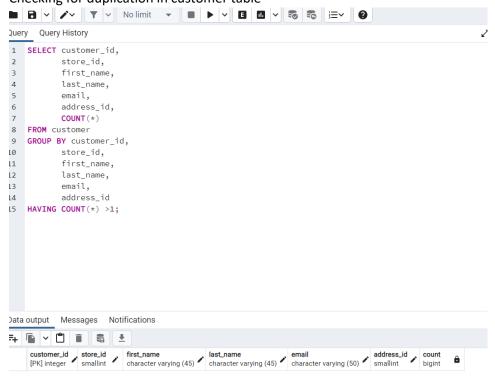Answers 3.6

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. Create a new "Answers 3.6" document and copy-paste your queries into it. Next to each query write 2 to 3 sentences explaining how you would clean the data (even if the data is not dirty).
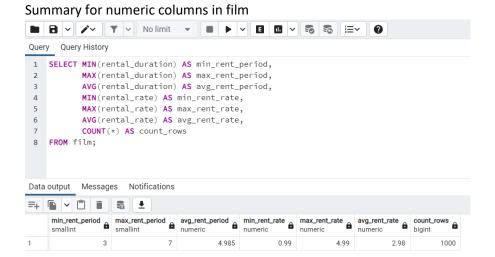
Checking for duplication in film



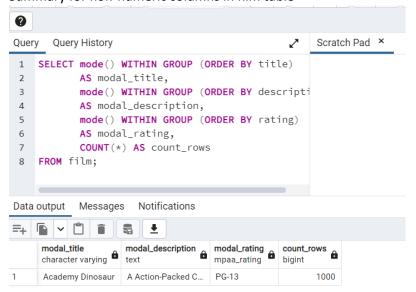Checking for duplication in customer table

There is no returned duplicate value. If there were duplication, there would have been two ways to deal with it. One, create a virtual table "view" where unique records can be selected. Two, Delete duplicate record from the table or view. If neither of those are permitted by the company, we can use GROUP BY or DISTINCT to select unique records.

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.
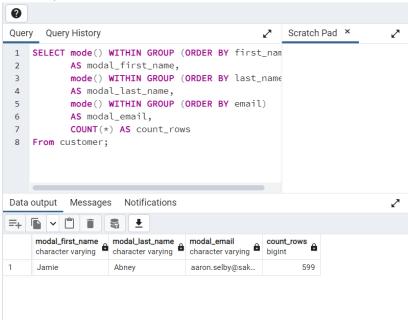
## Summary for numeric columns in film

```
1  SELECT MIN(rental_duration) AS min_rent_period,
2         MAX(rental_duration) AS max_rent_period,
3         AVG(rental_duration) AS avg_rent_period,
4         MIN(rental_rate) AS min_rent_rate,
5         MAX(rental_rate) AS max_rent_rate,
6         AVG(rental_rate) AS avg_rent_rate,
7         COUNT(*) AS count_rows
8  FROM film;
```

Data output  Messages  Notifications

| | min_rent_period smallint | max_rent_period smallint | avg_rent_period numeric | min_rent_rate numeric | max_rent_rate numeric | avg_rent_rate numeric | count_rows bigint |
|---|---|---|---|---|---|---|---|
| 1 | 3 | 7 | 4.985 | 0.99 | 4.99 | 2.98 | 1000 |

## Summary for non-numeric columns in film table

Query  Query History  ↗  Scratch Pad  ✕

```
1  SELECT mode() WITHIN GROUP (ORDER BY title)
2         AS modal_title,
3         mode() WITHIN GROUP (ORDER BY descripti
4         AS modal_description,
5         mode() WITHIN GROUP (ORDER BY rating)
6         AS modal_rating,
7         COUNT(*) AS count_rows
8  FROM film;
```

Data output  Messages  Notifications

| | modal_title character varying | modal_description text | modal_rating mpaa_rating | count_rows bigint |
|---|---|---|---|---|
| 1 | Academy Dinosaur | A Action-Packed C… | PG-13 | 1000 |

Summary for numeric columns in customer table



Summary of non-numeric columns in customer table



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.

Excel works best with smaller data, while using pivot tables are easy, it is harder where there is a massive about of data. SQL is easier to manipulate with large data, which is also faster. You can also answer specific questions, more detailed questions and with the right query answers pop up much quicker.