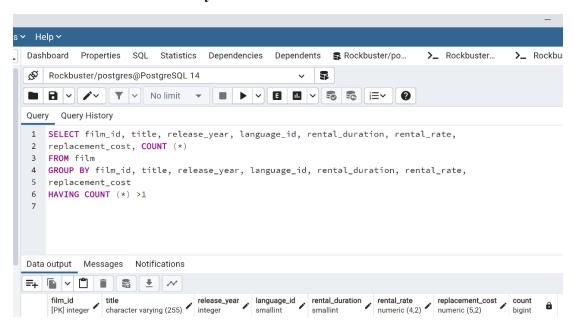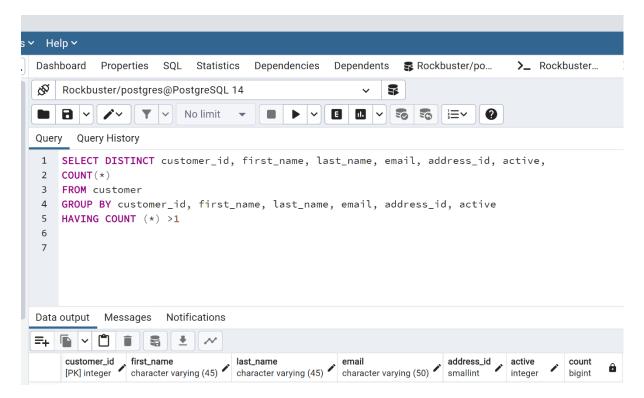# Summarizing & Cleaning Data in SQL
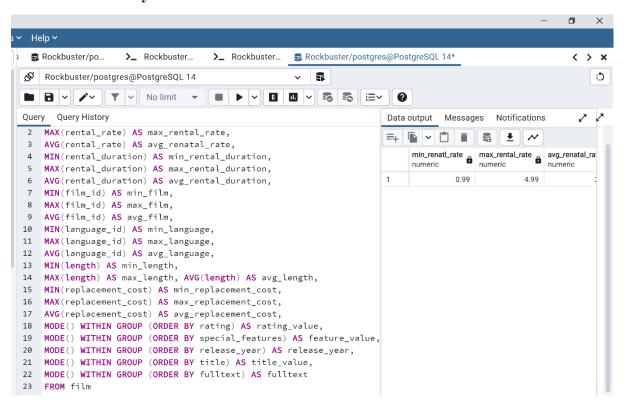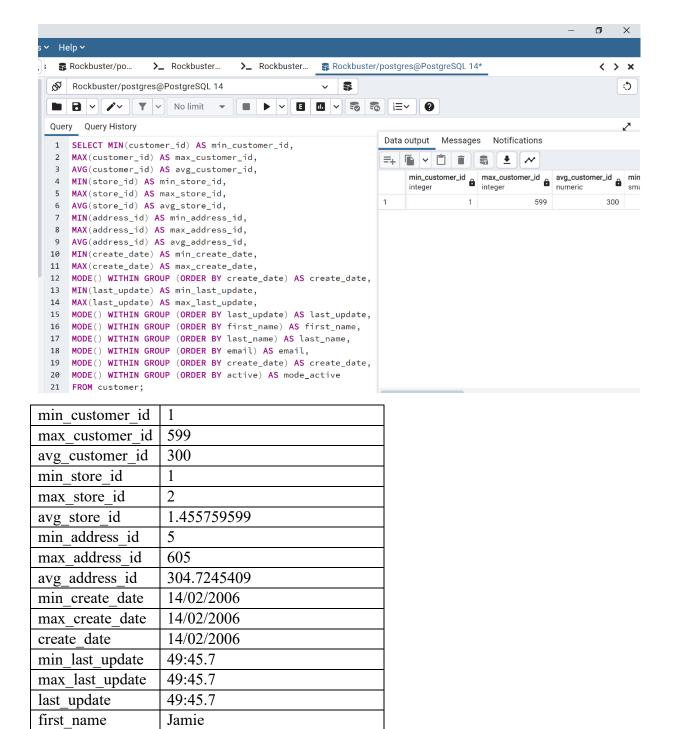
1. **Check for and clean dirty data:**





There is no returned duplicate value. Duplicate data can be solved by deleting the records if permissible or using a GROUP BY or DISTINCT to select required data.

## 2. Summarize your data:



```
 2   MAX(rental_rate) AS max_rental_rate,
 3   AVG(rental_rate) AS avg_renatal_rate,
 4   MIN(rental_duration) AS min_rental_duration,
 5   MAX(rental_duration) AS max_rental_duration,
 6   AVG(rental_duration) AS avg_rental_duration,
 7   MIN(film_id) AS min_film,
 8   MAX(film_id) AS max_film,
 9   AVG(film_id) AS avg_film,
10   MIN(language_id) AS min_language,
11   MAX(language_id) AS max_language,
12   AVG(language_id) AS avg_language,
13   MIN(length) AS min_length,
14   MAX(length) AS max_length, AVG(length) AS avg_length,
15   MIN(replacement_cost) AS min_replacement_cost,
16   MAX(replacement_cost) AS max_replacement_cost,
17   AVG(replacement_cost) AS avg_replacement_cost,
18   MODE() WITHIN GROUP (ORDER BY rating) AS rating_value,
19   MODE() WITHIN GROUP (ORDER BY special_features) AS feature_value,
20   MODE() WITHIN GROUP (ORDER BY release_year) AS release_year,
21   MODE() WITHIN GROUP (ORDER BY title) AS title_value,
22   MODE() WITHIN GROUP (ORDER BY fulltext) AS fulltext
23   FROM film
```

| | |
|---|---|
| min_renatl_rate | 0.99 |
| max_rental_rate | 4.99 |
| avg_renatal_rate | 2.98 |
| min_rental_duration | 3 |
| max_rental_duration | 7 |
| avg_rental_duration | 4.985 |
| min_film | 1 |
| max_film | 1000 |
| avg_film | 500.5 |
| min_language | 1 |
| max_language | 1 |
| avg_language | 1 |
| min_length | 46 |
| max_length | 185 |
| avg_length | 115.272 |
| min_replacement_cost | 9.99 |
| max_replacement_cost | 29.99 |
| avg_replacement_cost | 19.984 |
| rating_value | PG-13 |
| feature_value | {Trailers,Commentaries,"Behind the Scenes"} |
| release_year | 2006 |
| title_value | Academy Dinosaur |

| min_customer_id | 1 |
|---|---|
| max_customer_id | 599 |
| avg_customer_id | 300 |
| min_store_id | 1 |
| max_store_id | 2 |
| avg_store_id | 1.455759599 |
| min_address_id | 5 |
| max_address_id | 605 |
| avg_address_id | 304.7245409 |
| min_create_date | 14/02/2006 |
| max_create_date | 14/02/2006 |
| create_date | 14/02/2006 |
| min_last_update | 49:45.7 |
| max_last_update | 49:45.7 |
| last_update | 49:45.7 |
| first_name | Jamie |
| last_name | Abney |
| email | aaron.selby@sakilacustomer.org |
| create_date-2 | 14/02/2006 |
| mode_active | 1 |

3. **Reflect on your work:** Using Excel would be more effective when dealing with small amounts of data, or few columns within a data table to categorize and data profile. However, large amounts of data need SQL to be more effective.