# Task 3.6 Summarising & Cleaning Data in SQL
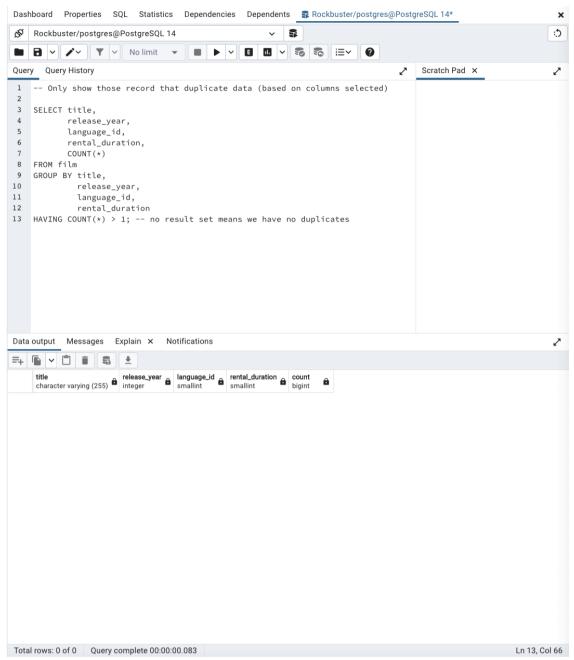
By Lee Heng Chuah

In this task you will calculate some descriptive statistics using the MIN, MAX, AVG, COUNT, SUM and MODE() aggregates discussed in this Exercise, and you will reflect on what you learned about data profiling back in Exercise 1.5: Data Profiling & Integrity.

**Directions:**

Rockbuster's database engineers have loaded some new data into the database, and your manager has asked you to clean and profile it. Follow the instructions below to complete their request:

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 "Answer 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)

**Find the duplicate data (film):**

Query   Query History                                                          Scratch Pad  ✕

```
1   -- Only show those record that duplicate data (based on columns selected)
2
3   SELECT title,
4          release_year,
5          language_id,
6          rental_duration,
7          COUNT(*)
8   FROM film
9   GROUP BY title,
10           release_year,
11           language_id,
12           rental_duration
13   HAVING COUNT(*) > 1; -- no result set means we have no duplicates
```

Data output   Messages   Explain  ✕   Notifications

| title character varying (255) | release_year integer | language_id smallint | rental_duration smallint | count bigint |
|---|---|---|---|---|

Total rows: 0 of 0     Query complete 00:00:00.083                                      Ln 13, Col 66

**Find the duplicate data (customer):**

Query    Query History                                                    ↙    Scratch Pad  ✕                ↗

```
 1   -- Only show those record that duplicate data (based on column selected)
 2
 3   SELECT store_id,
 4          first_name,
 5          last_name,
 6          email,
 7          address_id,
 8          active,
 9          COUNT(*)
10   FROM customer
11   GROUP BY store_id,
12            first_name,
13            last_name,
14            email,
15            address_id,
16            active
17   HAVING COUNT (*) > 1; -- no result set means we have no duplicates
18   |
```

Data output    Messages    Explain  ✕    Notifications                                                    ↗

≣+  🗐 ⌄  📋  🗑  🗄  ⬇

| store_id<br>smallint 🔒 | first_name<br>character varying (45) 🔒 | last_name<br>character varying (45) 🔒 | email<br>character varying (50) 🔒 | active<br>smallint 🔒 | count<br>bigint 🔒 |
|---|---|---|---|---|---|

Total rows: 0 of 0    Query complete 00:00:00.050                                                    Ln 18, Col 1
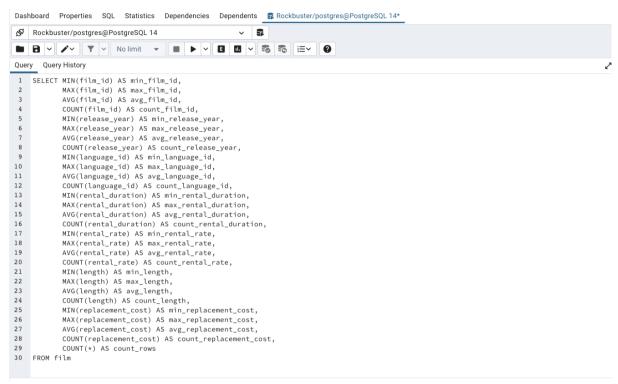
There are no duplications on both film and customer data, as we can see from the Data Output above.

As described in Task 3.6, there are two ways that we can fix them if we find duplicate records in our database:
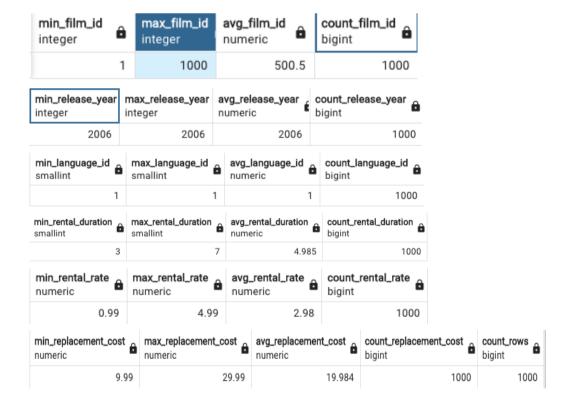
- Create a virtual table, known as a "view", where you select only unique records – standard way
- Delete the duplicate record from the table or view

2. **Summarise 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.
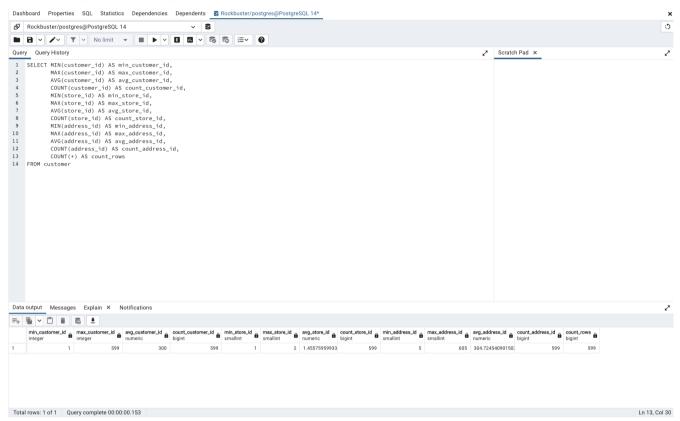
**Film – numerical columns**

```
1   SELECT MIN(film_id) AS min_film_id,
2          MAX(film_id) AS max_film_id,
3          AVG(film_id) AS avg_film_id,
4          COUNT(film_id) AS count_film_id,
5          MIN(release_year) AS min_release_year,
6          MAX(release_year) AS max_release_year,
7          AVG(release_year) AS avg_release_year,
8          COUNT(release_year) AS count_release_year,
9          MIN(language_id) AS min_language_id,
10         MAX(language_id) AS max_language_id,
11         AVG(language_id) AS avg_language_id,
12         COUNT(language_id) AS count_language_id,
13         MIN(rental_duration) AS min_rental_duration,
14         MAX(rental_duration) AS max_rental_duration,
15         AVG(rental_duration) AS avg_rental_duration,
16         COUNT(rental_duration) AS count_rental_duration,
17         MIN(rental_rate) AS min_rental_rate,
18         MAX(rental_rate) AS max_rental_rate,
19         AVG(rental_rate) AS avg_rental_rate,
20         COUNT(rental_rate) AS count_rental_rate,
21         MIN(length) AS min_length,
22         MAX(length) AS max_length,
23         AVG(length) AS avg_length,
24         COUNT(length) AS count_length,
25         MIN(replacement_cost) AS min_replacement_cost,
26         MAX(replacement_cost) AS max_replacement_cost,
27         AVG(replacement_cost) AS avg_replacement_cost,
28         COUNT(replacement_cost) AS count_replacement_cost,
29         COUNT(*) AS count_rows
30  FROM film
```

**Film – numerical columns** – Output (Note: The output has been split into a few screens due to the length of information)

| min_film_id integer | max_film_id integer | avg_film_id numeric | count_film_id bigint |
|---|---|---|---|
| 1 | 1000 | 500.5 | 1000 |

| min_release_year integer | max_release_year integer | avg_release_year numeric | count_release_year bigint |
|---|---|---|---|
| 2006 | 2006 | 2006 | 1000 |

| min_language_id smallint | max_language_id smallint | avg_language_id numeric | count_language_id bigint |
|---|---|---|---|
| 1 | 1 | 1 | 1000 |

| min_rental_duration smallint | max_rental_duration smallint | avg_rental_duration numeric | count_rental_duration bigint |
|---|---|---|---|
| 3 | 7 | 4.985 | 1000 |

| min_rental_rate numeric | max_rental_rate numeric | avg_rental_rate numeric | count_rental_rate bigint |
|---|---|---|---|
| 0.99 | 4.99 | 2.98 | 1000 |

| min_replacement_cost numeric | max_replacement_cost numeric | avg_replacement_cost numeric | count_replacement_cost bigint | count_rows bigint |
|---|---|---|---|---|
| 9.99 | 29.99 | 19.984 | 1000 | 1000 |

| min_length smallint | max_length smallint | avg_length numeric | count_length bigint |
|---|---|---|---|
| 46 | 185 | 115.272 | 1000 |

## Customer – numerical columns

Rockbuster/postgres@PostgreSQL 14

Query   Query History                                                                                                                                   Scratch Pad ✕

```
1   SELECT MIN(customer_id) AS min_customer_id,
2          MAX(customer_id) AS max_customer_id,
3          AVG(customer_id) AS avg_customer_id,
4          COUNT(customer_id) AS count_customer_id,
5          MIN(store_id) AS min_store_id,
6          MAX(store_id) AS max_store_id,
7          AVG(store_id) AS avg_store_id,
8          COUNT(store_id) AS count_store_id,
9          MIN(address_id) AS min_address_id,
10         MAX(address_id) AS max_address_id,
11         AVG(address_id) AS avg_address_id,
12         COUNT(address_id) AS count_address_id,
13         COUNT(*) AS count_rows
14  FROM customer
```

Data output   Messages   Explain ✕   Notifications

| | min_customer_id integer | max_customer_id integer | avg_customer_id numeric | count_customer_id bigint | min_store_id smallint | max_store_id smallint | avg_store_id numeric | count_store_id bigint | min_address_id smallint | max_address_id smallint | avg_address_id numeric | count_address_id bigint | count_rows bigint |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 599 | 300 | 599 | 1 | 2 | 1.45575959933 | 599 | 5 | 605 | 304.72454090150: | 599 | 599 |

Total rows: 1 of 1   Query complete 00:00:00.153                                                                                                                                            Ln 13, Col 30

# Film – non-numerical columns

Rockbuster/postgres@PostgreSQL 14

Query    Query History

```
1  SELECT mode() WITHIN GROUP (ORDER BY title)
2      AS mode_title_value,
3          mode() WITHIN GROUP (ORDER BY rating)
4      AS mode_rating_value,
5          mode() WITHIN GROUP (ORDER BY special_features)
6      AS mode_special_feature_value
7  FROM film;
8
```

Scratch Pad

Data output    Messages    Explain    Notifications

| | mode_title_value<br>character varying | mode_rating_value<br>mpaa_rating | mode_special_feature_value<br>text[] |
|---|---|---|---|
| 1 | Academy Dinosaur | PG-13 | {Trailers,Commentaries,"Behind the Scenes"} |

Total rows: 1 of 1    Query complete 00:00:00.079    Ln 8, Col 1

**Customer – non-numerical columns**

| Dashboard | Properties | SQL | Statistics | Dependencies | Dependents | Rockbuster/postgres@PostgreSQL 14* | ✛ | ✕ |
|---|---|---|---|---|---|---|---|---|

Rockbuster/postgres@PostgreSQL 14 ▾

No limit ▾

Query   Query History                                                 Scratch Pad ✕

```
 1  SELECT mode() WITHIN GROUP (ORDER BY first_name)
 2      AS mode_first_name_value,
 3          mode() WITHIN GROUP (ORDER BY last_name)
 4      AS mode_last_name_value,
 5          mode() WITHIN GROUP (ORDER BY email)
 6      AS mode_email_value,
 7      mode() WITHIN GROUP (ORDER by activebool)
 8      AS mode_activebool_value
 9  FROM customer;
10
```

Data output   Messages   Explain ✕   Notifications

| | mode_first_name_value<br>character varying | mode_last_name_value<br>character varying | mode_email_value<br>character varying | mode_activebool_value<br>boolean |
|---|---|---|---|---|
| 1 | Jamie | Abney | aaron.selby@sakilacustomer.org | true |

Total rows: 1 of 1     Query complete 00:00:00.117                                                        Ln 6, Col 25

---

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.

Back in Achievement 1, it took me approximate 2 weeks to find the duplication and missing data with Excel. In this achievement, it took me less than 3 hours to find the repetition and missing data using the SQL function. Excel took me a lengthy time to average duplicated data and eliminate the duplication of data, while SQL could average the data by using the function. Clearly, the SQL is faster, easier and safer than Excel.

4. Save your "Answer 3.6" document as a PDF and upload it here for your tutor to review.