# Task 3.6 – Summarizing & Cleaning Data in SQL

1.

```
Query   Query History

 1  SELECT title,
 2          release_year,
 3          language_id,
 4          rental_duration,
 5          COUNT(*)
 6  FROM film
 7  GROUP BY title,
 8          release_year,
 9          language_id,
10          rental_duration
11  HAVING COUNT(*) > 1
```

```
Query   Query History

 1  SELECT  first_name,
 2          last_name,
 3          email,
 4          address_id,
 5          COUNT(*)
 6  FROM customer
 7  GROUP BY first_name,
 8          last_name,
 9          email,
10          address_id
11  HAVING COUNT(*) > 1
```
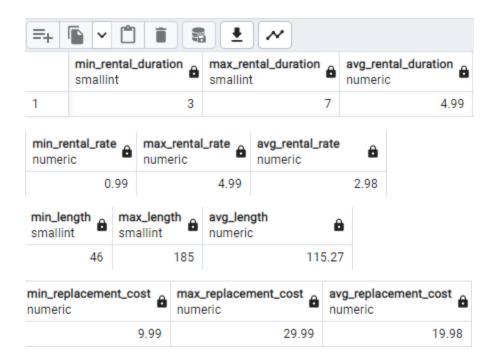
This query will check for any duplicate data, and in this case, no result set was returned meaning there were no duplicates. If there were duplicate records, there are two ways to fix them.

   a.  I could create a virtual table (VIEW) and select only the unique records.
   b.  I could delete the duplicate record from the table or view.

We could also use GROUP BY or DISTINCT to find any non-uniform data.

2.

```
 1  -- descriptive statistics for the film table
 2  SELECT MIN(rental_duration) AS min_rental_duration,
 3          MAX(rental_duration) AS max_rental_duration,
 4          ROUND(AVG(rental_duration), 2) AS avg_rental_duration,
 5          MIN(rental_rate) AS min_rental_rate,
 6          MAX(rental_rate) AS max_rental_rate,
 7          ROUND(AVG(rental_rate), 2) AS avg_rental_rate,
 8          MIN(length) AS min_length,
 9          MAX(length) AS max_length,
10          ROUND(AVG(length), 2) AS avg_length,
11          MIN(replacement_cost) AS min_replacement_cost,
12          MAX(replacement_cost) AS max_replacement_cost,
13          ROUND(AVG(replacement_cost), 2) AS avg_replacement_cost
14  FROM film
```

| min_rental_duration<br>smallint | max_rental_duration<br>smallint | avg_rental_duration<br>numeric |
|---|---|---|
| 3 | 7 | 4.99 |

| min_rental_rate<br>numeric | max_rental_rate<br>numeric | avg_rental_rate<br>numeric |
|---|---|---|
| 0.99 | 4.99 | 2.98 |

| min_length<br>smallint | max_length<br>smallint | avg_length<br>numeric |
|---|---|---|
| 46 | 185 | 115.27 |

| min_replacement_cost<br>numeric | max_replacement_cost<br>numeric | avg_replacement_cost<br>numeric |
|---|---|---|
| 9.99 | 29.99 | 19.98 |

```sql
--modal values for non-numerical columns for the film table
SELECT mode() WITHIN GROUP (ORDER BY title)
    AS modal_title,
mode() WITHIN GROUP (ORDER BY rating)
    AS modal_rating,
mode() WITHIN GROUP (ORDER BY special_features)
    AS modal_special_features,
mode() WITHIN GROUP (ORDER BY fulltext)
    AS modal_fulltext
FROM film
```

| modal_title<br>character varying | modal_rating<br>mpaa_rating | modal_special_features<br>text[] | modal_fulltext<br>tsvector |
|---|---|---|---|
| Academy Dinosaur | PG-13 | {Trailers,Commentari... | 'baloon':19 'con... |

```sql
1   -- descriptive statistics for the customer table
2   SELECT MIN(customer_id) AS min_customer_id,
3          MAX(customer_id) AS max_customer_id,
4          AVG(customer_id) AS avg_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          MIN(address_id) AS min_address_id,
9          MAX(address_id) AS max_address_id,
10         AVG(address_id) AS avg_address_id
11  FROM customer
```

| min_customer_id integer | max_customer_id integer | avg_customer_id numeric |
|---|---|---|
| 1 | 599 | 300.00 |

| min_store_id smallint | max_store_id smallint | avg_store_id numeric |
|---|---|---|
| 1 | 2 | 1.46 |

| min_address_id smallint | max_address_id smallint | avg_address_id numeric |
|---|---|---|
| 5 | 605 | 304.72 |

```sql
-- modal values for non-numerical columns for the customer table
SELECT mode() WITHIN GROUP (ORDER BY first_name)
    AS modal_first_name,
mode() WITHIN GROUP (ORDER BY last_name)
    AS modal_last_name,
mode() WITHIN GROUP (ORDER BY email)
    AS modal_email
FROM customer
```

| modal_first_name character varying | modal_last_name character varying | modal_email character varying |
|---|---|---|
| Jamie | Abney | aaron.selby@sakilacustomer.org |

3. I lean more towards SQL just because I find it quicker to query or to filter data. There are a ton of commands that can be used, or a combination of commands that can be used to extract exactly what you need. In addition, you can easily save it somewhere and type a new query. Although Excel is more user-friendly at first, I think SQL is more efficient in the long run.