**3.6: Summarizing & Cleaning Data in SQL**

**1. Check for and clean dirty data:**

Query Editor    Query History

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

Data Output    Explain    Messages    Notifications

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

Query Editor    Query History

```
1   --Checking for duplicate data
2   SELECT first_name,
3          last_name,
4          email,
5          address_id,
6          COUNT(*)
7   FROM customer
8   GROUP BY first_name,
9          last_name,
10         email,
11         address_id
12  HAVING COUNT(*) >1;
```

Data Output    Explain    Messages    Notifications

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

If the data needed to be cleaned, I could do so by creating a view table showing only unique records or I could delete the duplicate records. If I don't have permissions to update, delete, or create views then I could write a query that only returns unique records by using group by or distinct.

## 2. Summarize your data:

### Film Table

```
1   --descriptive statistics for film table
2   SELECT  MIN(rental_duration) AS min_rental_duration,
3           MAX(rental_duration) AS max_rental_duration,
4           AVG(rental_duration) AS avg_rental_duration,
5           MIN(rental_rate) AS min_rental_rate,
6           MAX(rental_rate) AS max_rental_rate,
7           AVG(rental_rate) AS avg_rental_rate,
8           MIN(length) AS min_length,
9           MAX(length) AS max_length,
10          AVG(length) AS avg_length,
11          MIN(replacement_cost) AS min_replacement_cost,
12          MAX(replacement_cost) AS max_replacement_cost,
13          AVG(replacement_cost) AS avg_replacement_cost,
14          COUNT(*) AS count_rows
15  FROM film
16
```

Data Output | Explain | Messages | Notifications

| min_rental_duration smallint | max_rental_duration smallint | avg_rental_duration numeric | min_rental_rate numeric | max_rental_rate numeric | avg_rental_rate numeric | min_length smallint | max_length smallint | avg_length numeric | min_replacement_cost numeric |
|---|---|---|---|---|---|---|---|---|---|
| 3 | 7 | 4.985 | 0.99 | 4.99 | 2.98 | 46 | 185 | 115.272 | 9.99 |

```
1   --descriptive statistics for film table
2   SELECT mode() WITHIN GROUP (ORDER BY title) AS modal_title,
3          mode() WITHIN GROUP (ORDER BY description) AS modal_description,
4          mode() WITHIN GROUP (ORDER BY release_year) AS modal_release_year,
5          mode() WITHIN GROUP (ORDER BY language_id) AS modal_language_id,
6          mode() WITHIN GROUP (ORDER BY rating) AS modal_rating
7   FROM film
```

Data Output | Explain | Messages | Notifications

| modal_title character varying | modal_description text | modal_release_year integer | modal_language_id smallint | modal_rating mpaa_rating |
|---|---|---|---|---|
| Academy Dinosaur | A Action-Packed Character Study of a Astronaut And a Explorer who must Reach a Monkey in A MySQL Convention | 2006 | 1 | PG-13 |

### Customer Table

```
1   --descriptive statistics for 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_addredd_id,
10         AVG(address_id) AS avg_address_id
11  FROM customer
```

Data Output | Explain | Messages | Notifications

| min_customer_id integer | max_customer_id integer | avg_customer_id numeric | min_store_id smallint | max_store_id smallint | avg_store_id numeric | min_address_id smallint | max_addredd_id smallint | avg_address_id numeric |
|---|---|---|---|---|---|---|---|---|
| 1 | 599 | 300 | 1 | 2 | 1.4557595993322203 | 5 | 605 | 304.7245409015025 |

```
Query Editor    Query History

  1   --descriptive statistics for customer table
  2   SELECT mode() WITHIN GROUP (ORDER BY first_name) AS modal_first_name,
  3          mode() WITHIN GROUP (ORDER BY last_name) AS modal_last_name,
  4          mode() WITHIN GROUP (ORDER BY email) AS modal_email
  5   FROM customer
```

Data Output    Explain    Messages    Notifications

| modal_first_name character varying | modal_last_name character varying | modal_email character varying |
|---|---|---|
| 1 Jamie | Abney | aaron.selby@sakilacustomer.org |

**3. Reflect on your work:**

When working with smaller data sets, I feel excel is easier for me as typing the query out correctly in SQL takes more time than getting the info in excel. When working with larger data sets the ability to write the query, add comments of what your doing, always having the query there in history if you need to go back to it, and how quickly it returns results makes SQL more effective in data profiling.