# Summarizing & Cleaning Data in SQL

*Question 1*

Film Table

```
SELECT title,
release_year,
language_id,
rental_duration,
COUNT(*)
FROM film
GROUP BY title,
release_year,
language_id,
rental_duration
HAVING COUNT (*) >1;
```

Customer Table



**There are no duplicates in the above.**

## Question 2

## Film table with non-numerical columns



```sql
1   SELECT mode() WITHIN GROUP (ORDER BY film_id) AS modal_film_id,
2   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 rating) AS modal_rating
5   FROM film
```

Data Output   Explain   Messages   Notifications

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

Customer table with numerical columns

Query Editor   Query History

```
1  SELECT
2  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  FROM customer
```
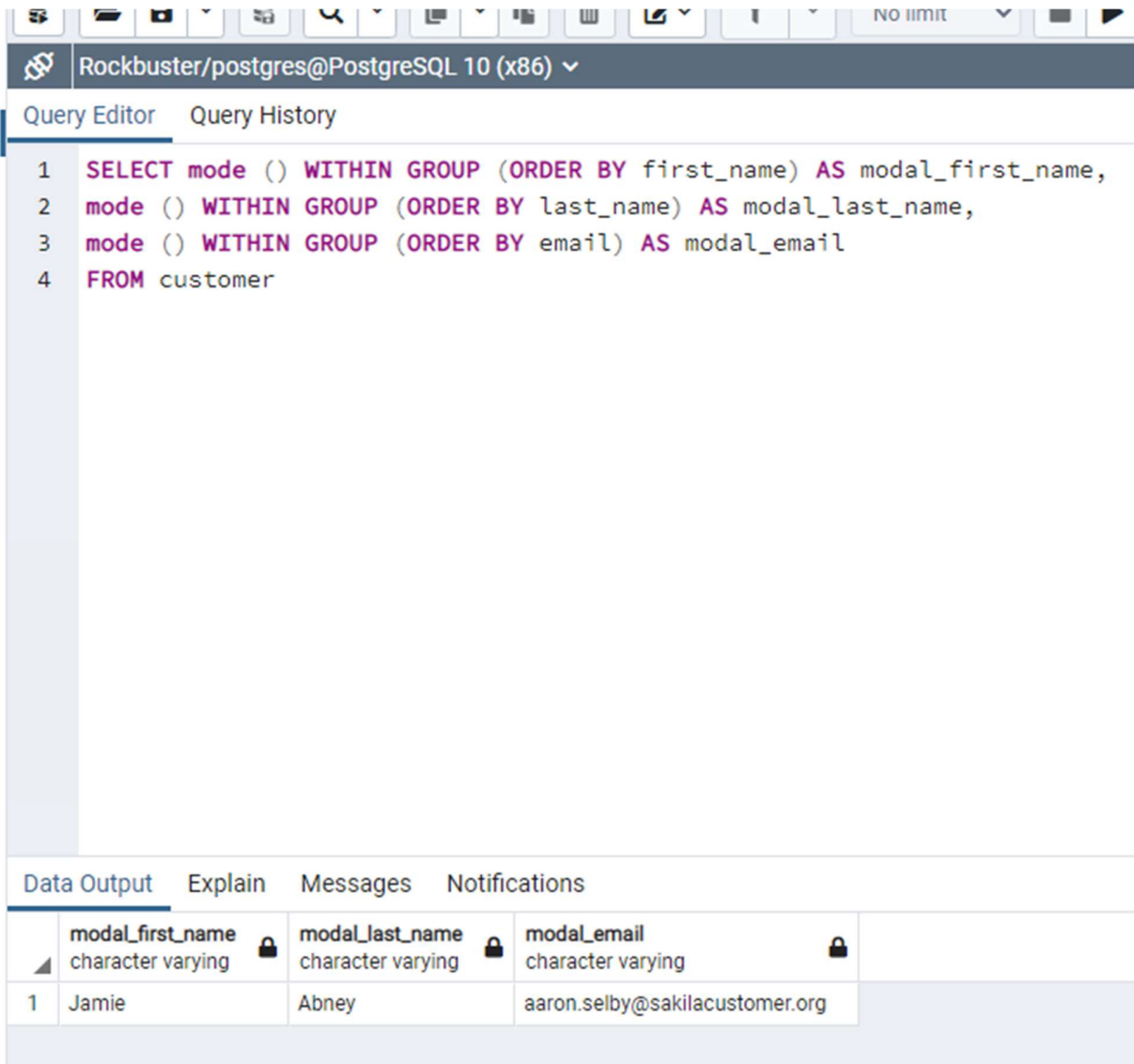
Data Output   Explain   Messages   Notifications

| min_customer_id integer | max_customer_id integer | avg_customer_id numeric | min_store_id smallint | max_store_id smallint | avg_store_id numeric |
|---|---|---|---|---|---|
| 1 | 599 | 300.0000000000000000 | 1 | 2 | 1.4557595993322204 |

Customer table with non-numerical columns



```
1   SELECT mode () WITHIN GROUP (ORDER BY first_name) AS modal_first_name,
2   mode () WITHIN GROUP (ORDER BY last_name) AS modal_last_name,
3   mode () WITHIN GROUP (ORDER BY email) AS modal_email
4   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 |

Question 3

When it comes to Data profiling I think it is much each to use SQL dependent on the size of the Database.