3.6 Summarizing & Cleaning Data in SQL

1. Check for dirty data:

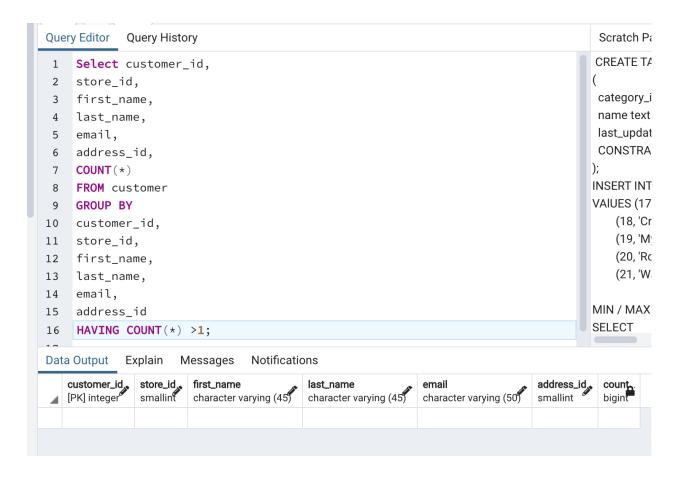Film table:

```
13    GROUP BY film_id,
14        title,
15        description,
16        release_year,
17        language_id,
18        rental_duration,
19        rental_rate,
20        length,
21        replacement_cost,
22        rating
23    HAVING COUNT(*) > 1;
24
```

```
 1   Select film_id,
 2       title,
 3       description,
 4       release_year,
 5       language_id,
 6       rental_duration,
 7       Rental_rate,
 8       length,
 9       replacement_cost,
10       rating,
11       count(*)
12   FROM film
13   GROUP BY film_id,
```

Customer table:

```
 1   Select customer_id,
 2   store_id,
 3   first_name,
 4   last_name,
 5   email,
 6   address_id,
 7   COUNT(*)
 8   FROM customer
 9   GROUP BY
10   customer_id,
11   store_id,
12   first_name,
13   last_name,
14   email,
15   address_id
16   HAVING COUNT(*) >1;
```

CREATE TA
(
  category_i
  name text
  last_updat
  CONSTRA
);
INSERT INT
VAIUES (17
        (18, 'Cr
        (19, 'M
        (20, 'Rc
        (21, 'W

MIN / MAX
SELECT

Data Output    Explain    Messages    Notifications

| customer_id [PK] integer | store_id smallint | first_name character varying (45) | last_name character varying (45) | email character varying (50) | address_id smallint | count bigint |
|---|---|---|---|---|---|---|
| | | | | | | |

To clean the tables, I would create a new table in the view format and delete any duplicate records found using the having count above.
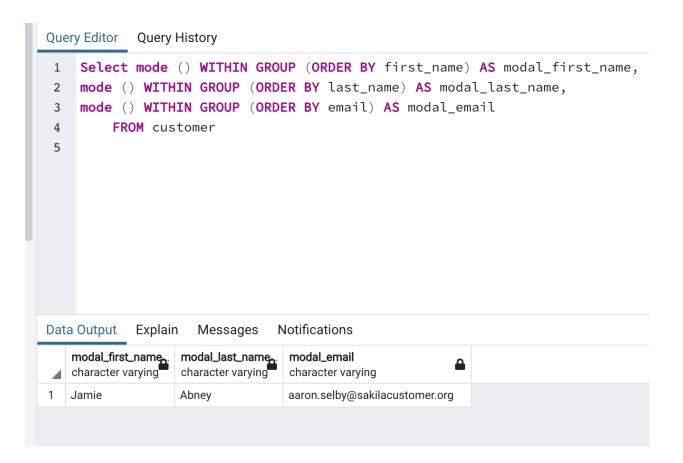
2. Summarize your data:

```
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
9
```

CREAT
(
categ
name
last_u
CONS
);
INSERT
VAIUES
    (1:
    (1:
    (2(

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 | 1 | 599 | 300 | 1 | 2 | 1.4557595993322203 |

```
1    Select
2    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 (length) AS min_length,
6    MAX (length) AS max_length,
7    AVG (length) AS avg_length,
8    MIN (replacement_cost) AS min_replacement_cost,
9    MAX (replacement_cost) AS max_replacement_cost,
10   AVG (replacement_cost) AS avg_replacement_cost
11   FROM film
12
```

CREATE
(
categor
name t
last_up
CONST
);
INSERT I
VAIUES (
    (18,
    (19,
    (20,

Data Output    Explain    Messages    Notifications

| | min_rental_duration smallint | max_rental_duration smallint | avg_rental_duration numeric | min_length smallint | max_length smallint | avg_length numeric | min_replac numeric |
|---|---|---|---|---|---|---|---|
| 1 | 3 | 7 | 4.985 | 46 | 185 | 115.272 | |

```
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
6
```

```
CREATE TABLE category
(
  category_id integer NOT NULL
  name text COLLATE pg_catalo
  last_update timestamp with tir
  CONSTRAINT category_pkey P
);
INSERT INTO category(category
VAlUES (17, 'Thriller')
     (18, 'Crime' ),
     (19, 'Mystery' ),
     (20, 'Romance' ),
```

Data Output   Explain   Messages   Notifications

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

---

```
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
5
```

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: 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.

Larger databases are best suited for SQL and smaller data sets are best suited for excel. I think the speed depends on which program you are more comfortable with and how large your data set is. SQL takes me a little longer due to my lack of experience with it, although I can see how fast it could be once you do have to background in it.