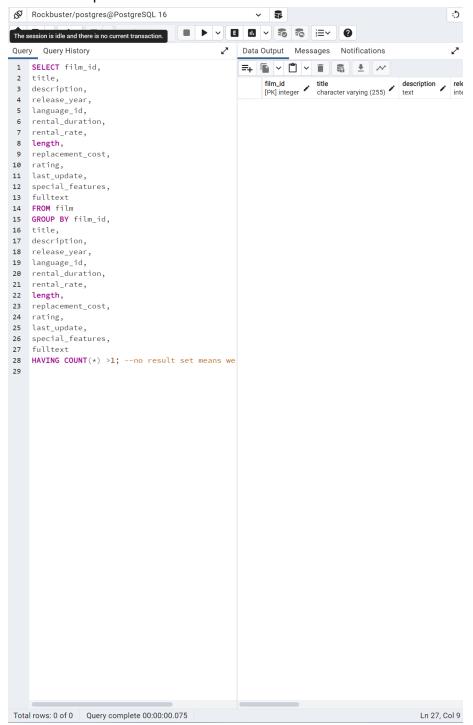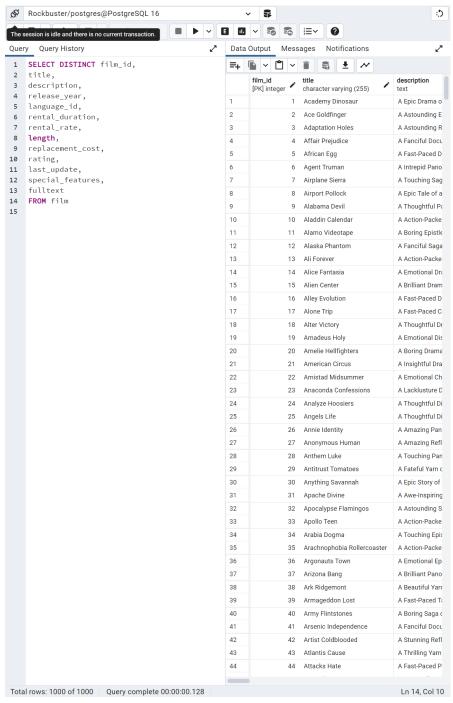Alexander Coley
11-19-23
Exercise 3.6
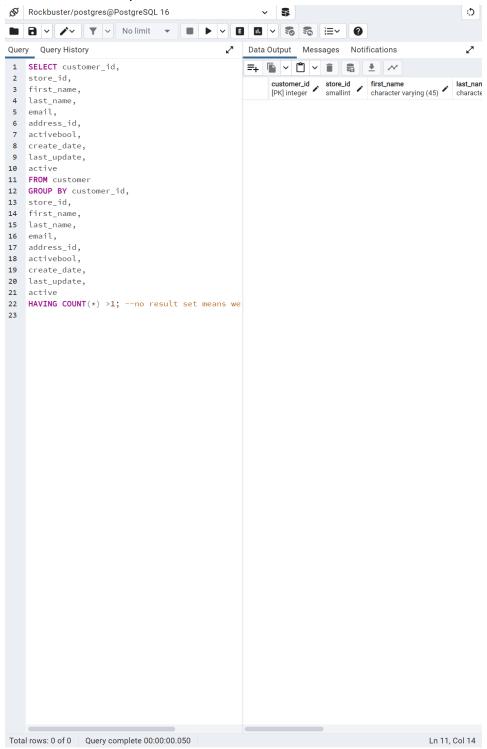
## Question 1:
Film table-Duplicate Data

I didn't find any duplicates in this data set. Although, if there were any duplicates, I would delete the duplicate columns to keep the data set consistent. Deleting these duplicates would also help with visualizations later on.
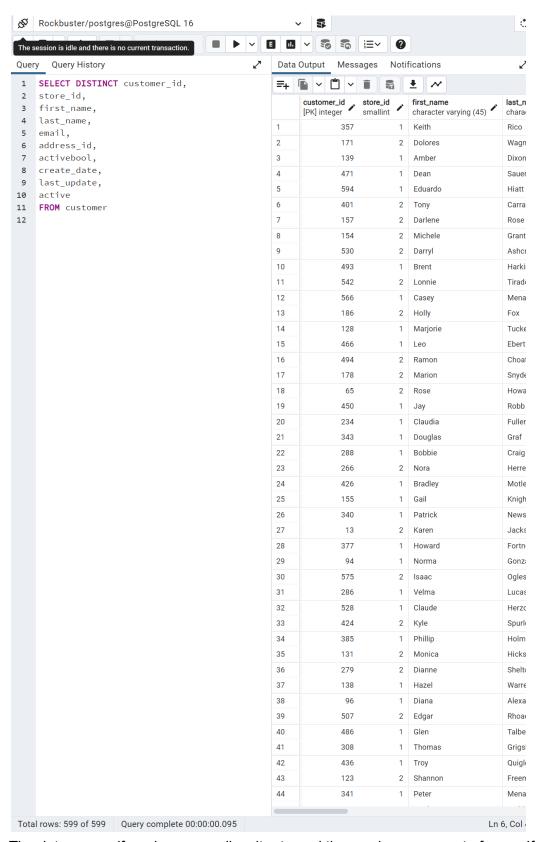
Non-uniform data



There isn't any non-uniform data. The data sent back 1000 out of 1000 rows so everything is consistent.

## Customer Table-Duplicate Data

```
Rockbuster/postgres@PostgreSQL 16

Query   Query History                              Data Output   Messages   Notifications

 1  SELECT customer_id,
 2  store_id,
 3  first_name,
 4  last_name,
 5  email,
 6  address_id,
 7  activebool,
 8  create_date,
 9  last_update,
10  active
11  FROM customer
12  GROUP BY customer_id,
13  store_id,
14  first_name,
15  last_name,
16  email,
17  address_id,
18  activebool,
19  create_date,
20  last_update,
21  active
22  HAVING COUNT(*) >1; --no result set means we
23
```

customer_id [PK] integer | store_id smallint | first_name character varying (45) | last_nan characte

Total rows: 0 of 0    Query complete 00:00:00.050                         Ln 11, Col 14

There weren't any duplicates for the customer table. If there were any duplicates I would do the same as before and delete the duplicates as to not hurt my progress in making the visualizations.

Non-Uniform Data

The data was uniform here as well as it returned the maximum amount of rows. If there were any missing values, if I didn't have access to what was missing, I would make the values zero.

Question 2:
Numerical Values for Film Table

Release year
SELECT MIN (release_year) AS min_release_year ,
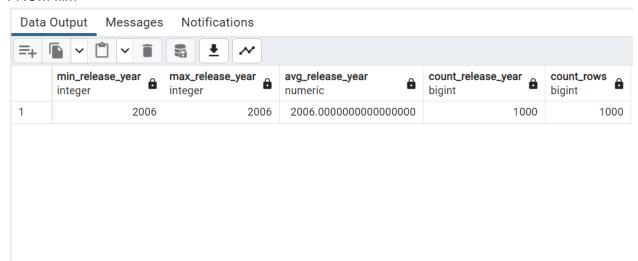MAX (release_year) AS max_release_year,
AVG (release_year) AS avg_release_year,
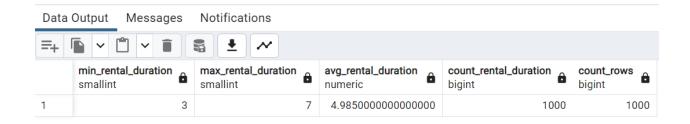COUNT (release_year) AS count_release_year,
COUNT (*) AS count_rows
FROM film

Data Output    Messages    Notifications

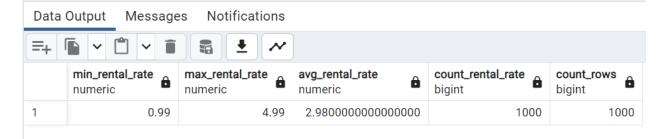| | min_release_year 🔒<br>integer | max_release_year 🔒<br>integer | avg_release_year 🔒<br>numeric | count_release_year 🔒<br>bigint | count_rows 🔒<br>bigint |
|---|---|---|---|---|---|
| 1 | 2006 | 2006 | 2006.0000000000000000 | 1000 | 1000 |

Rental Duration
SELECT MIN (rental_duration) AS min_rental_duration ,
MAX (rental_duration) AS max_rental_duration,
AVG (rental_duration) AS avg_rental_duration,
COUNT (rental_duration) AS count_rental_duration,
COUNT (*) AS count_rows
FROM film

## Data Output

| | min_rental_duration<br>smallint 🔒 | max_rental_duration<br>smallint 🔒 | avg_rental_duration<br>numeric 🔒 | count_rental_duration<br>bigint 🔒 | count_rows<br>bigint 🔒 |
|---|---|---|---|---|---|
| 1 | 3 | 7 | 4.9850000000000000 | 1000 | 1000 |

Rental Rate
SELECT MIN (rental_rate) AS min_rental_rate ,
MAX (rental_rate) AS max_rental_rate,
AVG (rental_rate) AS avg_rental_rate,
COUNT (rental_rate) AS count_rental_rate,
COUNT (*) AS count_rows
FROM film

## Data Output   Messages   Notifications

| | min_rental_rate<br>numeric 🔒 | max_rental_rate<br>numeric 🔒 | avg_rental_rate<br>numeric 🔒 | count_rental_rate<br>bigint 🔒 | count_rows<br>bigint 🔒 |
|---|---|---|---|---|---|
| 1 | 0.99 | 4.99 | 2.9800000000000000 | 1000 | 1000 |

Length
SELECT MIN (length) AS min_length ,
MAX (length) AS max_length,
AVG (length) AS avg_length,

COUNT (length) AS count_length,
COUNT (*) AS count_rows
FROM film

Data Output    Messages    Notifications

| | min_length smallint | max_length smallint | avg_length numeric | count_length bigint | count_rows bigint |
|---|---|---|---|---|---|
| 1 | 46 | 185 | 115.2720000000000 | 1000 | 1000 |

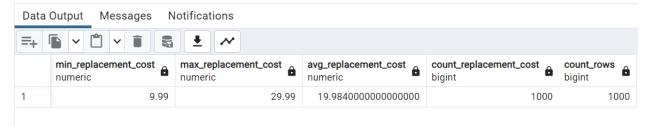Replacement Cost
SELECT MIN (replacement_cost) AS min_replacement_cost ,
MAX (replacement_cost) AS max_replacement_cost,
AVG (replacement_cost) AS avg_replacement_cost,
COUNT (replacement_cost) AS count_replacement_cost,
COUNT (*) AS count_rows
FROM film

Data Output    Messages    Notifications

| | min_replacement_cost numeric | max_replacement_cost numeric | avg_replacement_cost numeric | count_replacement_cost bigint | count_rows bigint |
|---|---|---|---|---|---|
| 1 | 9.99 | 29.99 | 19.9840000000000 | 1000 | 1000 |

Non numerical values

SELECT MODE() WITHIN GROUP (ORDER BY film_id) AS modal_film_id,
MODE() WITHIN GROUP (ORDER BY title) AS modal_title,
MODE() WITHIN GROUP (ORDER BY description) AS modal_description,
MODE() WITHIN GROUP (ORDER BY language_id) AS modal_language_id,
MODE() WITHIN GROUP (ORDER BY rating) AS modal_rating,
MODE() WITHIN GROUP (ORDER BY last_update) AS modal_last_update,
MODE() WITHIN GROUP (ORDER BY special_features) AS modal_special_features,
MODE() WITHIN GROUP (ORDER BY fulltext) AS modal_fulltext
FROM film

Data Output    Messages    Notifications

| min_replacement_cost numeric | max_replacement_cost numeric | avg_replacement_cost numeric | count_replacement_cost bigint | count_rows bigint |
|---|---|---|---|---|
| 1 | 9.99 | 29.99 | 19.9840000000000000 | 1000 | 1000 |

Customer Table *only contains non numerical values
SELECT MODE() WITHIN GROUP (ORDER BY customer_id) AS modal_customer_id,
MODE() WITHIN GROUP (ORDER BY store_id) AS modal_store_id,
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,
MODE() WITHIN GROUP (ORDER BY address_id) AS modal_address_id,
MODE() WITHIN GROUP (ORDER BY activebool) AS modal_activebool,
MODE() WITHIN GROUP (ORDER BY create_date) AS modal_create_date,
MODE() WITHIN GROUP (ORDER BY last_update) AS modal_last_update,
MODE() WITHIN GROUP (ORDER BY active) AS modal_active
FROM customer

| modal_customer_id integer | modal_store_id smallint | modal_first_name character varying | modal_last_name character varying | modal_email character varying | modal_address_id smallint | modal_activebool boolean | modal_create_date date | modal_last_update timestamp without time zone | modal_active integer |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | Jamie | Abney | aaron.selby@sakilacustomer.org | 5 | true | 2006-02-14 | 2013-05-26 14:49:45.738 | 1 |

Question 3:
As much as I've learned from SQL so far I definitely prefer it over excel. Excel is still very useful but I feel as though the functions in SQL are easier to use and understand. It's also more visually appealing to me personally. I hope to learn more about SQL functions in the future as I grow as an analyst.