

TASK 3.6

Summarizing & Cleaning Data in SQL

1: DUPLICATE

The screenshot shows a PostgreSQL query editor interface. The top bar indicates the connection is 'Rockbuster/postgres@PostgreSQL 16'. The 'Query' tab is active, displaying the following SQL query:

```
1 SELECT film_id, title, description, release_year, language_id, rental_rate, rental_duration,
2 length, replacement_cost, rating, last_update, special_features, fulltext,
3 COUNT (*)
4 from film
5 Group By film_id, title, description, release_year, language_id, rental_rate, rental_duration
6 length, replacement_cost, rating, last_update, special_features, fulltext
7 HAVING COUNT (*)>1;
```

Below the query editor, the 'Data Output' tab is active, showing a table with the following columns and data types:

film_id	title	description	release_year	language_id	rental_rate	rental_duration	length	replacement_cost	rating	mpaa_rating	last_update
---------	-------	-------------	--------------	-------------	-------------	-----------------	--------	------------------	--------	-------------	-------------

The bottom screenshot shows the same PostgreSQL query editor interface. The 'Query' tab is active, displaying the following SQL query:

```
1 SELECT customer_id, store_id, first_name, last_name, email, address_id, activebool, create_date,
2 COUNT(*)
3 FROM customer
4 Group By customer_id, store_id, first_name, last_name, email, address_id, activebool, create_date
5 HAVING COUNT (*)>1;
```

Below the query editor, the 'Data Output' tab is active, showing a table with the following columns and data types:

customer_id	store_id	first_name	last_name	email	address_id	activebool	create_date	last_update
-------------	----------	------------	-----------	-------	------------	------------	-------------	-------------

I didn't come across any duplicate entries in the film_table or customer_table. But if there were any, I would address them by either removing the duplicate entry or creating a "view table" that selects only one instance of the duplicates.

MISSING VALUES

Rockbuster/postgres@PostgreSQL 16

Query Query History Scratch Pad X

```
1 SELECT *
2 FROM customer
3 Where (customer_id,store_id, first_name,last_name, email, address_id, activebool, create_date
4 is NULL;
```

Data Output Messages Notifications

customer_id	store_id	first_name	last_name	email	address_id	activebool	create_date	last_update
[PK] integer	smallint	character varying (45)	character varying (45)	character varying (50)	smallint	boolean	date	timestamp without time zone

Rockbuster/postgres@PostgreSQL 16

Query Query History Scratch Pad X

```
1 SELECT *
2 FROM film
3 Where (film_id,title, description, release_year,
4 language_id, rental_rate, rental_duration, length,
5 replacement_cost, last_update, rating, special_features, fulltext)
6 IS NULL;
```

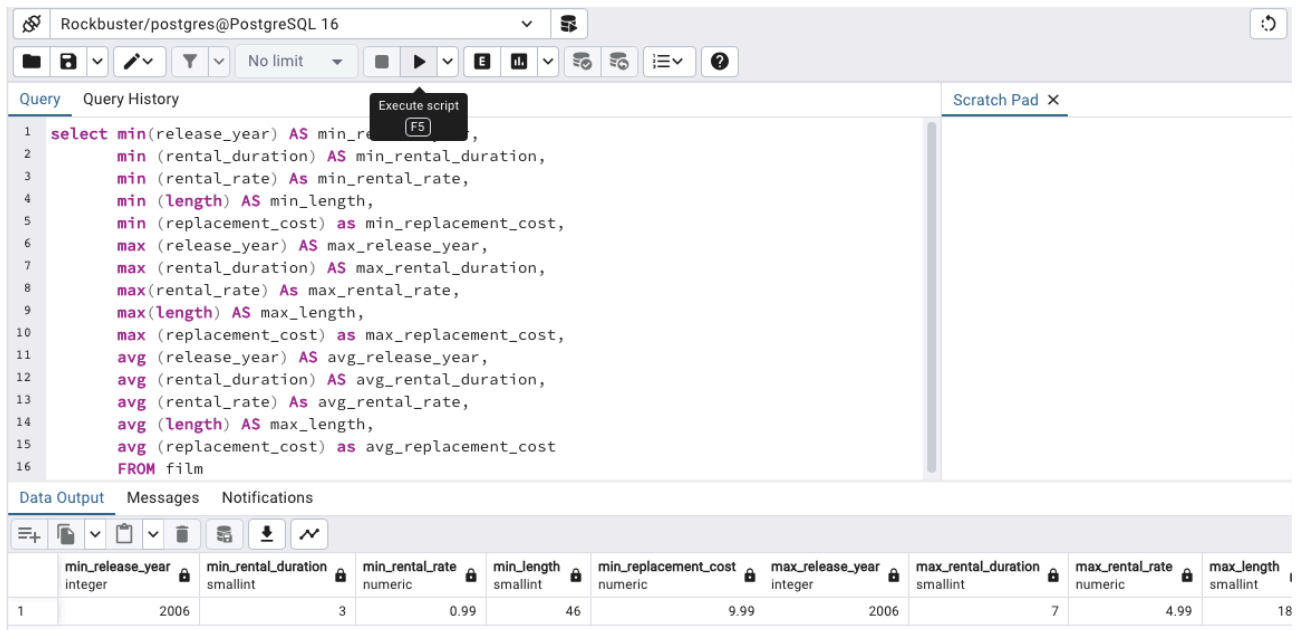
Data Output Messages Notifications

film_id	title	description	release_year	language_id	rental_duration	rental_rate	length	replacement_cost	rating	last_update
[PK] integer	character varying (255)	text	integer	smallint	smallint	numeric (4,2)	smallint	numeric (5,2)	mpaa_rating	timestamp without time zone

I haven't encountered any missing values in the tables so far. However, if any were present, I would fill them in as needed. For instance, if the last_name field in the customer_table was blank, I would skip it. But if the release_year of a film was missing, I would make sure to find the information and populate the record accordingly.

2:

NUMERICAL



The screenshot shows a PostgreSQL query editor interface. The top bar indicates the connection is 'Rockbuster/postgres@PostgreSQL 16'. The query editor contains a SQL query that calculates various statistical values from the 'film' table. The query is as follows:

```
1 select min(release_year) AS min_release_year,
2        min (rental_duration) AS min_rental_duration,
3        min (rental_rate) AS min_rental_rate,
4        min (length) AS min_length,
5        min (replacement_cost) AS min_replacement_cost,
6        max (release_year) AS max_release_year,
7        max (rental_duration) AS max_rental_duration,
8        max (rental_rate) AS max_rental_rate,
9        max (length) AS max_length,
10       max (replacement_cost) AS max_replacement_cost,
11       avg (release_year) AS avg_release_year,
12       avg (rental_duration) AS avg_rental_duration,
13       avg (rental_rate) AS avg_rental_rate,
14       avg (length) AS max_length,
15       avg (replacement_cost) AS avg_replacement_cost
16 FROM film
```

The 'Data Output' tab at the bottom shows the results of the query. The results are displayed in a table with 10 columns, each corresponding to a field in the query. The data for the first row is as follows:

	min_release_year integer	min_rental_duration smallint	min_rental_rate numeric	min_length smallint	min_replacement_cost numeric	max_release_year integer	max_rental_duration smallint	max_rental_rate numeric	max_length smallint
1	2006	3	0.99	46	9.99	2006	7	4.99	18

These data records offer crucial insights into several aspects:

- The overall satisfaction level of customers with the rented movies.
- The average duration customers typically keep a rented movie.
- The potential cost associated with replacing a specific movie.
- Additionally, analyzing the popularity of different movie lengths.

In the Customer_table, there isn't pertinent numerical data for retrieving minimum, maximum, and average values. For instance, customer_id, being a unique identifier for each customer, wouldn't yield meaningful insights through minimum, maximum, or average calculations.

NON NUMERICAL

Rockbuster/postgres@PostgreSQL 16

Query Query History

```

1 SELECT MODE () WITHIN GROUP (ORDER BY language_id) AS longest,
2     MODE () WITHIN GROUP (ORDER BY rating) AS most_rated
3 from film

```

Data Output Messages Notifications

	longest smallint	most_rated mpaa_rating
1	1	PG-13

In this analysis, it's feasible to determine the most prevalent language among the movie collection and identify which audience demographic could derive the most value from the current selection. In my perspective, querying data from text fields like description wouldn't be pertinent as the information contained within those columns wouldn't offer substantial insights.

CUSTOMER TABLE

Rockbuster/postgres@PostgreSQL 16

Query Query History

```

1 SELECT MODE () WITHIN GROUP (ORDER BY store_id) AS most_store,
2     MODE () WITHIN GROUP (ORDER BY activebool) AS most_activebool
3 from customer

```

Data Output Messages Notifications

	most_store smallint	most_activebool boolean
1	1	true

In this case, the analysis shows that Store 1 has the most customers signed up. It also asks if there are more active customers than inactive ones. But checking things like first names or last names wouldn't really help find useful information, like the most common names in the database.

3:

SQL is really handy because it gives you exactly what you ask for when you make a query. It's great for connecting data from different places and makes sure your data stays safe from accidental changes. But since I'm still learning, I sometimes find it hard to use the right words. With practice, though, I think I'll get better at it.

Excel is good too because it helps automate tasks with lots of functions. I'm more used to Excel, so I can do similar things by moving tables around and using functions. But when it comes to big amounts of data, Excel can struggle a bit compared to SQL. Even though Excel is easier to use and understand, I'm enjoying learning SQL and getting better at it. Just like learning any new language, it takes time to get the hang of it.