

3.6: Summarizing & Cleaning Data in SQL

1. Check for and clean dirty data: Find out if the film table and the customer table contain any dirty data, specifically non-uniform or duplicate data, or missing values. Create a new “Answers 3.6” document and copy-paste your queries into it. Next to each query write 2 to 3 sentences explaining how you would clean the data (even if the data is not dirty).

Duplicate Data from Film

The screenshot shows a SQL query editor with a query tab selected. The query is as follows:

```
1 --looking for duplicate values
2 SELECT film_id,
3 title,
4 release_year,
5 rental_duration,
6 rental_rate,
7 count(*)
8 FROM film
9 GROUP BY film_id,
10 title,
11 release_year,
12 rental_duration,
13 rental_rate
14 HAVING COUNT(*) > 1;
```

Below the query editor, there is a toolbar with icons for various actions. At the bottom, a table structure is displayed with the following columns and data types:

film_id	title	release_year	rental_duration	rental_rate	count
[PK] integer	character varying (255)	integer	smallint	numeric (4,2)	bigint

Duplicate Data from Customer

The screenshot shows a SQL query editor with a query tab selected. The query is as follows:

```
1 --looking for duplicate values
2 SELECT customer_id
3 first_name,
4 last_name,
5 email,
6 address_id,
7 active,
8 count(*)
9 FROM customer
10 GROUP BY customer_id,
11 first_name,
12 last_name,
13 email,
14 address_id,
15 active
16 HAVING COUNT(*) > 1;
```

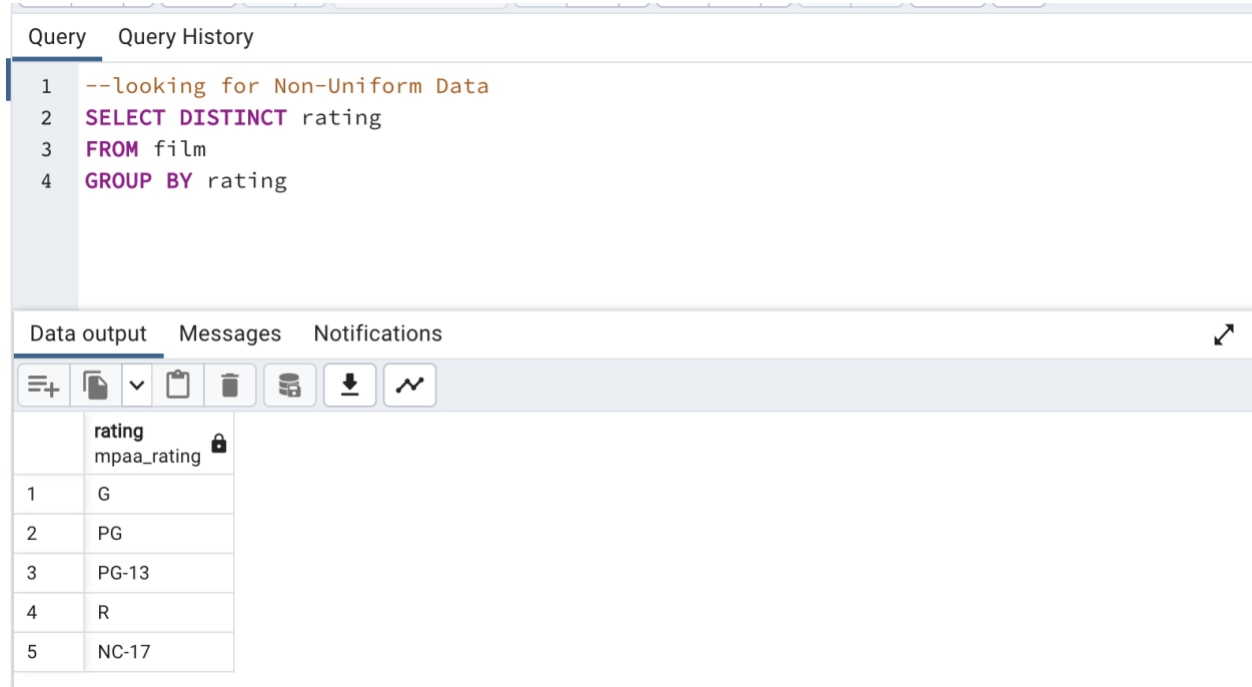
Below the query editor, there is a toolbar with icons for various actions. At the bottom, a table structure is displayed with the following columns and data types:

first_name	last_name	email	address_id	active	count
integer	character varying (45)	character varying (50)	smallint	integer	bigint

There is no returned duplicate value. There are two ways of dealing with duplicate value if you have permission to alter the database. • Create a virtual table “View” where unique records can be

selected • Delete duplicate record from the table or View However, if altering table is not permitted, we can use GROUP BY or DISTINCT to select unique records

Non-Uniform Data



The screenshot shows a SQL query editor with a query window and a data output window. The query window contains the following SQL code:

```
1 --looking for Non-Uniform Data
2 SELECT DISTINCT rating
3 FROM film
4 GROUP BY rating
```

The data output window shows the results of the query, which are the distinct ratings from the film table. The output is displayed in a table with two columns: 'rating' and 'mpaa_rating'.

	rating mpaa_rating
1	G
2	PG
3	PG-13
4	R
5	NC-17

Values are homogeneous, but in order to fix it on the table if there was an issue, I would use the UPDATE command combined with SET and WHERE, to replace the values that should be differently represented.

Incorrect Data

For incorrect data, logic and critical thinking can be used instead of commands, as you can analyse the tables and data for outliers or things that are wrongly displayed.

A mix of the commands show before can also help, by grouping and organizing the information in a more visible way in order for the analyst to review it.

Missing Data

This can also be fixed with logic, as you can use commands to better visualize the data and find missing or null information.

In some cases where a column has too much information missing, it might even be valid to remove it from queries.

A solution for it after finding the missing data, is to replace it with the average of the remaining informed data (if appropriate). The following command can assist with that:

```
UPDATE tablename
SET = AVG(col1)
WHERE col1 IS NULL
```

2. Summarize your data: Use SQL to calculate descriptive statistics for both the film table and the customer table. For numerical columns, this means finding the minimum, maximum, and average values. For non-numerical columns, calculate the mode value. Copy-paste your SQL queries and their outputs into your answers document.

Summary for customer

--descriptive statistics for customer table

```
SELECT MIN(customer_id) AS min_customer_id,  
       MAX(customer_id) AS max_customer_id,  
       AVG(customer_id) AS avg_customer_id,  
       MIN(store_id) AS min_store_id,  
       MAX(store_id) AS max_store_id,  
       AVG(store_id) AS avg_store_id,  
       MIN(address_id) AS min_address_id,  
       MAX(address_id) AS max_address_id,  
       AVG(address_id) AS avg_address_id,  
       MIN(create_date) AS min_create_date,  
       MAX(create_date) AS max_create_date,  
       MODE() WITHIN GROUP (ORDER BY create_date) AS create_date,  
       MIN(last_update) AS min_last_update,  
       MAX(last_update) AS max_last_update,  
       MODE() WITHIN GROUP (ORDER BY last_update) AS last_update,  
       MODE() WITHIN GROUP (ORDER BY first_name) AS first_name,  
       MODE() WITHIN GROUP (ORDER BY last_name) AS last_name,  
       MODE() WITHIN GROUP (ORDER BY email) AS email,  
       MODE() WITHIN GROUP (ORDER BY create_date) AS create_date,  
       MODE() WITHIN GROUP (ORDER BY active) AS mode_active  
FROM customer;
```

Rockbuster/postgres@posgre

Query Query History

```

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_address_id,
10        AVG(address_id) AS avg_address_id,
11        MIN(create_date) AS min_create_date,
12        MAX(create_date) AS max_create_date,
13        MODE() WITHIN GROUP (ORDER BY create_date) AS create_date,
14        MIN(last_update) AS min_last_update,
15        MAX(last_update) AS max_last_update,
16        MODE() WITHIN GROUP (ORDER BY last_update) AS last_update,
17        MODE() WITHIN GROUP (ORDER BY first_name) AS first_name,
18        MODE() WITHIN GROUP (ORDER BY last_name) AS last_name,
19        MODE() WITHIN GROUP (ORDER BY email) AS email,
20        MODE() WITHIN GROUP (ORDER BY create_date) AS create_date,
21        MODE() WITHIN GROUP (ORDER BY active) AS mode_active
22  FROM customer;
23

```

Data output 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
1	1	599	300	1	2	1.4557595993322203	

For some information instead of MAX, MIN and AVG I used MAX, MIN and MODE, it makes more sense to see when something happened more often for dates rather than the average date.

Summary for film

```

--descriptive statistics for film table
SELECT MIN(rental_rate) AS min_rental_rate,
MAX(rental_rate) AS max_rental_rate,
AVG(rental_rate) AS avg_rental_rate,
MIN(rental_duration) AS min_rental_duration,
MAX(rental_duration) AS max_rental_duration,
AVG(rental_duration) AS avg_rental_duration,
MIN(film_id) AS min_film,
MAX(film_id) AS max_film,
AVG(film_id) AS avg_film,
MIN(language_id) AS min_language,
MAX(language_id) AS max_language,
AVG(language_id) AS avg_language,
MIN(length) AS min_length,

```

```

MAX(length) AS max_length,
AVG(length) AS avg_length,
MIN(replacement_cost) AS min_replacement_cost,
MAX(replacement_cost) AS max_replacement_cost,
AVG(replacement_cost) AS avg_replacement_cost,
MODE() WITHIN GROUP (ORDER BY rating) AS rating_value,
MODE() WITHIN GROUP (ORDER BY special_features) AS feature_value,
MODE() WITHIN GROUP (ORDER BY release_year) AS release_year,
MODE() WITHIN GROUP (ORDER BY title) AS title_value,
MODE() WITHIN GROUP (ORDER BY fulltext) AS fulltext
FROM film

```

The screenshot shows a PostgreSQL query editor interface. The query is a SELECT statement that calculates various descriptive statistics for the 'film' table. The results are displayed in a table below the query.

Query:

```

1  --descriptive statistics for film table
2  SELECT MIN(rental_rate) AS min_rental_rate,
3  MAX(rental_rate) AS max_rental_rate,
4  AVG(rental_rate) AS avg_rental_rate,
5  MIN(rental_duration) AS min_rental_duration,
6  MAX(rental_duration) AS max_rental_duration,
7  AVG(rental_duration) AS avg_rental_duration,
8  MIN(film_id) AS min_film,
9  MAX(film_id) AS max_film,
10 AVG(film_id) AS avg_film,
11 MIN(language_id) AS min_language,
12 MAX(language_id) AS max_language,
13 AVG(language_id) AS avg_language,
14 MIN(length) AS min_length,
15 MAX(length) AS max_length,
16 AVG(length) AS avg_length,
17 MIN(replacement_cost) AS min_replacement_cost,
18 MAX(replacement_cost) AS max_replacement_cost,
19 AVG(replacement_cost) AS avg_replacement_cost,
20 MODE() WITHIN GROUP (ORDER BY rating) AS rating_value,
21 MODE() WITHIN GROUP (ORDER BY special_features) AS feature_value,
22 MODE() WITHIN GROUP (ORDER BY release_year) AS release_year,
23 MODE() WITHIN GROUP (ORDER BY title) AS title_value,
24 MODE() WITHIN GROUP (ORDER BY fulltext) AS fulltext
25 FROM film

```

Data output:

	min_rental_rate numeric	max_rental_rate numeric	avg_rental_rate numeric	min_rental_duration smallint	max_rental_duration smallint	avg_rental_duration numeric
1	0.99	4.99	2.9800000000000000	3	7	4.9850000000000000

For some information instead of MAX, MIN and AVG I used MAX, MIN and MODE, it makes more sense to see when something happened more often for dates rather than the average date.

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.

Excel works great with a small amount of data. At the same time, it will be easy to view data using pivot tables. However, renaming the output (alias) for an aggregation column will take longer, it's easy to work with huge data in SQL. Using SQL data profiling is getting easier and faster. The specific result/request for details will be returned immediately with the correct syntax.

It might be nice to use Excel for data cleansing and processing small datasets, but then for working with it or when processing huge amounts of data, I would choose SQL.