

## ANSWERS 3.6 - Summarizing & Cleaning Data in SQL

### 1a. Finding Duplicate values

The image shows two screenshots of the SQL Developer interface. The top screenshot displays a query to find duplicate values in the 'customer' table based on customer\_id, store\_id, first\_name, last\_name, and address\_id. The bottom screenshot displays a query to find duplicate values in the 'film' table based on film\_id, title, release\_year, language\_id, and rental\_duration. Both queries use the HAVING COUNT(\*) > 1 clause to identify duplicates.

**Query 1: Finding duplicates in customer table**

```
1 --Shows duplicate values based on columns selected
2 SELECT customer_id,store_id,first_name,last_name,
3 address_id,
4 Count(*)
5 FROM customer
6 GROUP BY customer_id,store_id,first_name,last_name,
7 address_id
8 HAVING COUNT(*)>1
```

**Data output table:**

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

**Query 2: Finding duplicates in film table**

```
1 --Shows duplicate values based on columns selected
2 SELECT film_id,title,release_year,language_id,
3 rental_duration,
4 Count(*)
5 FROM film
6 GROUP BY film_id, title,release_year,
7 language_id,rental_duration
8 HAVING COUNT(*)>1
```

**Data output table:**

film_id	title	release_year	language_id	rental_duration	count
[PK] integer	character varying (255)	integer	smallint	smallint	bigint

There are no duplicates, but to remove a duplicate, one would either create a view with unique records which is also known as a virtual table or the other is to delete the duplicates. But I will rather create a virtual table by selecting unique records and if I have to delete any duplicates, I will better use a DSP model.

## 1b. Finding non-uniform values

Query Query History

```
1 --show non-uniform value in film table base on rating
2 SELECT DISTINCT rating FROM film
3 GROUP BY rating
4
5
```

Data output Messages Notifications

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

There are not any non-uniform values but to fix this if there are any, I will need to update any values that represent the real values in different formats to be in the exact format as the real value format e.g GEN, g, would have to be updated to be G.

## 1c. Finding Missing Values

Query Query History

```
1 --Finding missing values in film
2 SELECT
3 COUNT(film_id) AS count_of_film_id,
4 COUNT(title) AS count_of_title,
5 COUNT(rental_duration) AS count_of_rental_duration,
6 COUNT(rental_rate) AS count_of_rental_rate,
7 COUNT(replacement_cost) AS count_of_replacement_cost,
8 COUNT(*) AS All_row_count
9 FROM film
```

Data output Messages Notifications

	count_of_film_id bigint	count_of_title bigint	count_of_rental_duration bigint	count_of_rental_rate bigint	count_of_replacement_cost bigint	all_row_count bigint
1	1000	1000	1000	1000	1000	1000

Query

Query History

```

1  --Finding missing values from customer table
2  SELECT
3  COUNT(customer_id)AS count_of_customer_id,
4  COUNT(first_name)AS count_of_first_name,
5  COUNT(last_name)AS count_of_last_name,
6  COUNT(*) AS count_of_row
7  FROM customer

```

Data output

Messages

Notifications

+

📄

▼

📋

🗑️

🔄

⬇️

	count_of_customer_id bigint	count_of_first_name bigint	count_of_last_name bigint	count_of_row bigint
1	599	599	599	599

If there are a lot of missing values in a column, I would better ignore that column to minimize the impact on the analysis. If there are just a few missing values and I know what the values are I can fill it in or just leave it as it is.

## 2. Descriptive statistics for numerical value in the film table

Query

Query History

Scratch Pad

1

--Calculating descriptive statistics for numerical columns

2

SELECT

3

MAX(rental\_duration) AS Maximum\_rental\_duration,

4

MIN(rental\_duration) AS Minimum\_rental\_duration,

5

AVG(rental\_duration) AS Average\_rental\_duration,

6

MAX(rental\_rate) AS Maximum\_rental\_rate,

7

MIN(rental\_rate) AS Minimum\_rental\_rate,

8

AVG(rental\_rate) AS Average\_rental\_rate,

9

MAX(length) AS Maximum\_length,

10

MIN(length) AS Minimum\_length,

11

AVG(length) AS Average\_length,

12

MAX(replacement\_cost) AS Maximum\_replacement\_cost,

13

MIN(replacement\_cost) AS Minimum\_replacement\_cost,

14

AVG(replacement\_cost) AS Average\_replacement\_cost

15

FROM film

Data output

Messages

Notifications

	maximum_rei smallint	minimum_ren smallint	average_rental_dura numeric	maximum_rental_rate numeric	minimum_ren numeric	average_renti numeric	maximum_leng smallint	minimum_len smallint	average_leng numeric	maximum_rei numeric	minimum_rep numeric	average_reple numeric
1	7	3	4.9850000000000000	4.99	0.99	2.9800000000	185	46	115.2720000	29.99	9.99	19.98400000

## 2b. Descriptive statistics for the numerical column in the customer table

I didn't calculate the Average for these numerical values as we are with "id" so I stuck with Max and Minimum to know the highest and lowest value for the ids.

QueryQuery History

Scratch Pad x

```
1 SELECT
2 MAX(customer_id)AS Maximum_customer_id,
3 MIN(customer_id)AS Minimum_customer_id,
4 MAX(store_id)AS Maximum_store_id,
5 MIN(store_id)AS Minimum_store_id,
6 MAX(address_id)AS Maximum_address_id,
7 MIN(address_id) AS Minimum_address_id,
8 MAX(active)AS Maximum_active,
9 MIN(active)AS Minimum_active
10 FROM customer
```

Data outputMessagesNotifications

	maximum_customer_id integer	minimum_customer_id integer	maximum_store_id smallint	minimum_store_id smallint	maximum_address_id smallint	minimum_address_id smallint	maximum_active integer	minimum_active integer
1	599	1	2	1	605	5	1	0

## 2c. Descriptive statistics for non-numerical value in the film table

QueryQuery History

```
1 --Calculating the mode value for non-numerical columns.
2 SELECT 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
```

Data outputMessagesNotifications

	modal_title character varying	modal_description text	modal_rating mpaa_rating
1	Academy Dinosaur	A Action-Packed C...	PG-13

## 2d. Descriptive statistics for non-numerical value in the customer table

Query

Query History

1

--Calculating the mode value for non-numerical columns.

2

SELECT mode()WITHIN GROUP(ORDER BY first\_name)AS modal\_first\_name,

3

mode()WITHIN GROUP(ORDER BY last\_name)AS modal\_last\_name,

4

mode()WITHIN GROUP(ORDER BY activebool) AS modal\_activebool

5

FROM customer



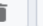



6

7

Data output

Messages

Notifications



	modal_first_name character varying	modal_last_name character varying	modal_activebool boolean
1	Jamie	Abney	true

- Using a pivot table for profiling in excel is quite easy when it comes to profiling big data, I think SQL is much easier and faster as much as one knows how to write the syntax and command properly. With Excel, you must go through various stages, like creating the table first, then following some procedures if you want to aggregate the data but with SQL one just needs to know the SQL syntax, write, and click execute and the answer is right there. Overall, using SQL is much easy and faster when it comes to large data and excel is also great when the data isn't large.