

## Exercise 3.6 - Summarizing & Cleaning Data in SQL

### Step 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. 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

Query

Query History

```
1 SELECT film_id,
2 title,
3 release_year,
4 rental_duration,
5 rental_rate,
6 COUNT(*)
7 FROM film
8 GROUP BY film_id,
9 title, release_year,
10 rental_duration,
11 rental_rate
12 HAVING COUNT(*) > 1
```

Data output

Messages

Notifications

≡

📄

▼

📋

🗑️

🗄️

⬇️

📈

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

### Duplicate Data from Customer

Query

Query History

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

Data output

Messages

Notifications

≡

📄

▼

📋

🗑️

🗄️

⬇️

📈

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

→ Here there is no returned duplicate value. In case we found some and we need to clean it, the process would have been to delete them (rows) by using the DELETE function.  
To do so we would need to use the VIEW function to create another view of the table in order to clean it in a more practical way.

### Non-Uniform Data

Query	Query History
1	SELECT DISTINCT rating
2	FROM film
3	GROUP BY rating

Data output	Messages	Notifications
<div><div>rating</div><div>mpaa_rating</div><div>1</div><div>2</div><div>3</div><div>4</div><div>5</div><div>G</div><div>PG</div><div>PG-13</div><div>R</div><div>NC-17</div></div>		

→ Here 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:

Query	Query History
1	UPDATE tablename
2	SET = AVG(col1)
3	WHERE col1 IS NULL

## Step 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.

### Film Table

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

Data output	Messages	Notifications
	min_renatl_rate numeric	max_rental_rate numeric
	avg_renatal_rate numeric	min_rental_duration smallint
		max_rental_duration smallint
		avg_rental_duration numeric
		min_film integer
		max_film integer
		avg_film numeric
1	0.99	4.99
	2.98	3
	7	4.985
	1	1000
		500.5

NB: In order to makes more sense to see when something happened more often for dates rather than the average date, MAX, MIN and MODE were used instead of MAX, MIN and AVG for Film and Customer Table.

### Customer Table

```
Query Query History
```

```
1 SELECT MIN(customer_id) AS min_customer_id,
2 MAX(customer_id) AS max_customer_id,
3 AVG(customer_id) AS avg_customer_id,
4 MIN(store_id) AS min_store_id,
5 MAX(store_id) AS max_store_id,
6 AVG(store_id) AS avg_store_id,
7 MIN(address_id) AS min_address_id,
8 MAX(address_id) AS max_address_id,
9 AVG(address_id) AS avg_address_id,
10 MIN(create_date) AS min_create_date,
11 MAX(create_date) AS max_create_date,
12 MODE() WITHIN GROUP (ORDER BY create_date) AS create_date,
13 MIN(last_update) AS min_last_update,
14 MAX(last_update) AS max_last_update,
15 MODE() WITHIN GROUP (ORDER BY last_update) AS last_update,
16 MODE() WITHIN GROUP (ORDER BY first_name) AS first_name,
17 MODE() WITHIN GROUP (ORDER BY last_name) AS last_name,
18 MODE() WITHIN GROUP (ORDER BY email) AS email,
19 MODE() WITHIN GROUP (ORDER BY create_date) AS create_date,
20 MODE() WITHIN GROUP (ORDER BY active) AS mode_active
21 FROM customer;
```

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	max_address_id smallint	avg_address_id numeric	min_create_date timestamp	max_create_date timestamp	mode_create_date timestamp	min_last_update timestamp	max_last_update timestamp	mode_last_update timestamp	min_first_name text	max_first_name text	mode_first_name text	min_last_name text	max_last_name text	mode_last_name text	min_email text	max_email text	mode_email text	min_create_date timestamp	max_create_date timestamp	mode_create_date timestamp	min_active boolean	max_active boolean	mode_active boolean
1	1	599	300	1	2	1.4557595993322203	5	5	3.3333333333333335	2008-09-08 09:09:00	2008-09-08 09:09:00	2008-09-08 09:09:00	2008-09-08 09:09:00	2008-09-08 09:09:00	2008-09-08 09:09:00	John	John	John	Deo	Deo	Deo	john.doe@freesoft.com	john.doe@freesoft.com	john.doe@freesoft.com	2008-09-08 09:09:00	2008-09-08 09:09:00	2008-09-08 09:09:00	1	1	1

### Step 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.

Excel is a great tool to quickly visualize all your data in spread sheet or on pivot table view.

However, it is limited in term of data quantity and SQL give you the possibility to easily reorganize your dataset, create alias, correct syntax etc...

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.