

## Task 3.6 – Summarizing and 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).

Query Query History

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

Data Output Messages Notifications

title	release_year	language_id	rental_duration	count		
character varying (255)	integer	smallint	smallint	bigint		

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

There are no duplicates in either data table. In the case that there was “dirty” data, I would clean the data by either creating a secondary virtual table, or a view, that selects only the unique data, or I would delete the duplicate records from the table.

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

Data Output		Messages		Notifications	
rating	mpaa_rating				
1	G				
2	PG				
3	PG-13				
4	R				
5	NC-17				

There was no non-uniform data in the film table. If there was, this could be corrected by using the UPDATE, SET, WHERE commands.

In the case of missing data, it would depend on how much of the data is missing, and whether or not the data can be substituted. For example, if a large amount of data values is missing, the column could be ignored when performing queries. If the data is numerical, then an average can be imputed instead.

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.

```
SELECT MIN(film_id) AS minimum_film_id,  
       MAX(film_id) AS maximum_film_id,  
       AVG(film_id) AS average_film_id,  
       MIN(release_year) AS minimum_release_year,  
       MAX(release_year) AS maximum_release_year,  
       AVG(release_year) AS average_release_year,  
       MIN(language_id) AS minimum_language_id,  
       MAX(language_id) AS maximum_language_id,  
       AVG(language_id) AS average_language_id,  
       MIN(rental_duration) AS min_rental_duration,  
       MAX(rental_duration) AS max_rental_duration,  
       AVG(rental_duration) AS avg_rental_duration,  
       MIN(rental_rate) AS min_rental_rate,  
       MAX(rental_rate) AS max_rental_rate,  
       AVG(rental_rate) AS avg_rental_rate,  
       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 title) AS modal_title,  
       MODE() WITHIN GROUP (ORDER BY description) AS modal_description,  
       MODE() WITHIN GROUP (ORDER BY rating) AS modal_rating,  
       MODE() WITHIN GROUP (ORDER BY special_features) AS modal_specials,  
       MODE() WITHIN GROUP (ORDER BY fulltext) AS modal_fulltext  
FROM film
```

Scratch Pad ✕

Scratch Pad X

[illegible]

```
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(active) AS min_active,  
       MAX(active) AS max_active,  
       AVG(active) AS avg_active,  
       MIN(create_date) AS min_create_date,  
       MAX(create_date) AS max_create_date,  
       MODE() WITHIN GROUP (ORDER BY create_date) AS modal_create_date,  
       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 activebool) AS modal_activebool  
FROM customer
```

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(active) AS min\_active,

11

MAX(active) AS max\_active,

12

AVG(active) AS avg\_active,

13

MIN(create\_date) AS min\_create\_date,

14

MAX(create\_date) AS max\_create\_date,

15

MODE() WITHIN GROUP (ORDER BY create\_date) AS modal\_create\_date,

16

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

17

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

18

MODE() WITHIN GROUP (ORDER BY email) AS modal\_email,

19

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

20

FROM customer

21

Scratch Pad

X

Data Output

Messages

Notifications

	min_customer_id	max_customer_id	avg_customer_id	min_store_id	max_store_id	avg_store_id	min_address_id	max_address_id	avg_address_id	min_active	max_active	avg_active	min_create_date	max_create_date	modal_create_date
	integer	integer	numeric	smallint	smallint	numeric	smallint	smallint	numeric	integer	integer	numeric	date	date	date
1	1	599	300.00000000000000	1	2	1.45575959933	5	605	304.7245409015021	0	1	0.9749582637	2006-02-14	2006-02-14	2006-02-14

- 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 seems to have felt faster while trying to come up with data summarizing. While working on problem 2, it felt very time-consuming typing up the same query for a different column over and over again, as well as including an alias for all of them. I do understand, however, that it does take some time and practice to be able to use SQL quickly and efficiently, and it is much faster than Excel when working with large amounts of data.