

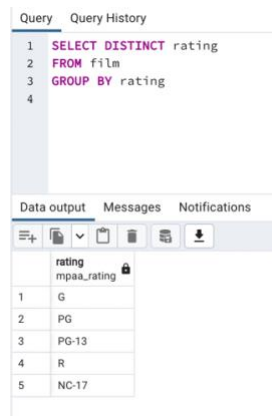
## 3.6 Answers- Summarizing & Cleaning

---

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.

### Finding non-uniform data

```
SELECT DISTINCT rating
FROM film
GROUP BY rating
```



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

```
1 SELECT DISTINCT rating
2 FROM film
3 GROUP BY rating
4
```

The results window shows the output of the query, which is a list of distinct ratings from the film table. The results are as follows:

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

- Review a few random values to check for inconsistencies using the `GROUP BY` and `DISTINCT` keywords.

### Finding Duplicates in film table

```
SELECT title,
release_year,
language_id,
rental_duration,
COUNT(*)
FROM film
GROUP BY title,
release_year,
language_id,
rental_duration
HAVING COUNT(*) >1
```

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
12

```

Data output Messages Notifications

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

### Finding Duplicates in Customer table

```

SELECT first_name,
last_name,
address_id,
email,
active,
COUNT(*)
FROM customer
GROUP BY first_name,
last_name,
address_id,
email,
active
HAVING COUNT(*) >1;

```

Query Query History Scratch Pad x

```

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

```

Data output Messages Notifications

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

- **Duplicate - There are 2 options.** Create a virtual table, known as a "view," where you select only unique records or delete the duplicate record from the table or view (not advised).

### Finding Missing Values in film table

```

SELECT
COUNT(title) AS count_title,
COUNT(rental_duration) AS count_rental_duration,
COUNT(rental_rate) AS count_rental_rate,

```

COUNT(length) AS count\_length,  
COUNT(replacement\_cost) AS count\_replacement\_cost,  
COUNT(\*) AS count\_rows  
FROM film;

Query

Query History

Scratch Pad

X

```

1 SELECT
2 COUNT(title) AS count_title,
3 COUNT(rental_duration) AS count_rental_duration,
4 COUNT(rental_rate) AS count_rental_rate,
5 COUNT(length) AS count_length,
6 COUNT(replacement_cost) AS count_replacement_cost,
7 COUNT(*) AS count_rows
8 FROM film;
9

```

Data output

Messages

Notifications

+

+

+

+

+

+

	count_title bigint	count_rental_duration bigint	count_rental_rate bigint	count_length bigint	count_replacement_cost bigint	count_rows bigint
1	1000	1000	1000	1000	1000	1000

### Finding Missing Values in customer table

SELECT  
COUNT(email) AS count\_email,  
COUNT(last\_name) AS count\_last\_name,  
COUNT(first\_name) AS count\_first\_name,  
COUNT(customer\_id) AS count\_customer\_id,  
COUNT(store\_id) AS count\_store\_id,  
Answers 3.6 2  
COUNT(\*) AS count\_rows  
FROM customer;

Query

Query History

Scratch Pad

1

SELECT

2

COUNT(email) AS count\_email,

3

COUNT(last\_name) AS count\_last\_name,

4

COUNT(first\_name) AS count\_first\_name,

5

COUNT(customer\_id) AS count\_customer\_id,

6

COUNT(store\_id) AS count\_store\_id,

7

COUNT(\*) AS count\_rows

8

FROM customer;

Data output

Messages

Notifications

count\_email

count\_last\_name

count\_first\_name

count\_customer\_id

count\_store\_id

count\_rows

bigint

bigint

bigint

bigint

bigint

bigint

1

599

599

599

599

599

- Missing Values - We can either ignore columns with a high percentage of missing values or *impute* missing values with a column average: Ex.  
 imputing missing values with the AVG value  
 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.

### Descriptive Stats for film table (numerical columns)

```
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
FROM film
```



The screenshot shows a SQL query editor with a query history pane on the left and a scratch pad on the right. The query is the same as the one above. Below the editor, there is a data output table with 14 columns and 1 row of data.

	min_rental_rate	max_rental_rate	avg_rental_rate	min_rental_duration	max_rental_duration	avg_rental_duration	min_film	max_film	avg_film	min_language	max_language	avg_language	min_length	max_length	avg_length	min_replacement_cost	max_replacement_cost	avg_replacement_cost
1	0.99	4.99	2.98	5	7	4.985	1	1000	500.5	1	1000	500.5	45	185	115.272	9.99	29.99	19.984

### Mode value for film table (non-numerical)

```
SELECT 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 year_value,
mode() WITHIN GROUP (ORDER BY title)
AS title_value
FROM film;
```



Query Query History

### 3. Reflect on your work

Based on your previous experience, which tool (Excel or SQL) do you think is more effective for data profiling, and why?

SQL is more efficient at computing data profiles. Once the code is written, it can be applied to other data sets, it just takes time to get use to and lot's of trial and error.