# Answers 3.6

3.6: Summarizing & Cleaning Data in SQL

---

- 1) No duplicates found.
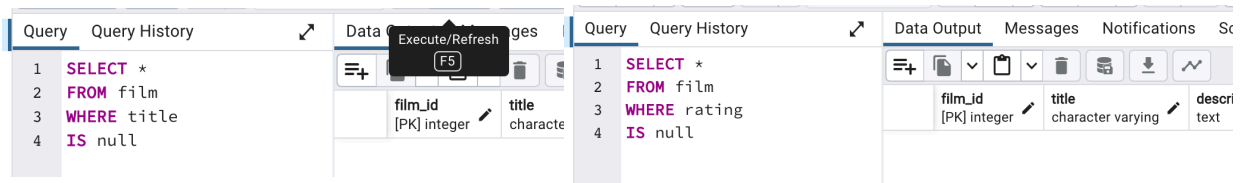
```
1   SELECT title,
2          release_year,
3          language_id,
4          rental_duration,
5          COUNT(*)
6   FROM film
7   GROUP BY title,
8          release_year,
9          language_id,
10         rental_duration
11  HAVING COUNT(*) >1; --no result set means we have no duplicates
```

Data Output

title
character v

```
1   SELECT store_id,
2       last_name,
3       first_name,
4       COUNT(*)
5   FROM customer
6   GROUP BY store_id, last_name, first_name
7   HAVING COUNT(*) >1; --no result set means we have no duplicates
```
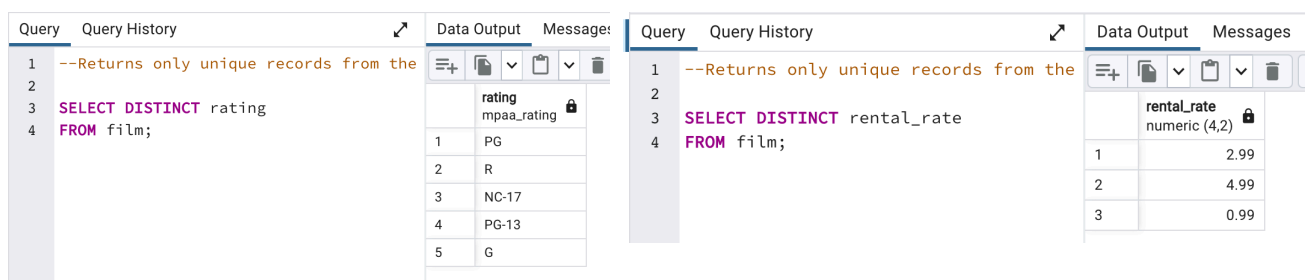
Data Output    Messag

store_id    last
smallint    cha

- No missing information found with 2 searches in both Customer and Film.

```
1   SELECT *
2   FROM customer
3   WHERE first_name
4   IS null
```

Data Output    Messages

customer_id    store_id
[PK] integer   smallint

```
1   SELECT *
2   FROM customer
3   WHERE last_name
4   IS null
```

Data Output    Messages

customer_id    store_id
[PK] integer   smallint

- Using DISTINCT I looked at ratings and rental_rate. There are 5 variables for rating and 3 for rental_rate.



— Cleaning the affected units would require permissions firstly. To change the duplicates, we could GROUP the duplicates.

—As we would likely not be able to DELETE data, we may be able to add an ACTIVE or INACTIVE filter to the accounts with duplicate names (for Customers).

—For the missing data, we can use an educated guess for that. (*New releases are 4.99; less than 5 years old 2.99; over 5 years old .99¢*)

- 2) Film Data summaries.

```
1  SELECT MIN(rental_duration)AS min_duration,
2      MAX(rental_duration)AS max_duration,
3      AVG(rental_duration)AS avg_duration,
4      COUNT(rental_duration)AS count_rental_duration_values,
5      MIN(rental_rate)AS min_rate,
6      MAX(rental_rate)AS max_rate,
7      AVG(rental_rate)AS avg_rate,
8      COUNT(rental_rate)AS count_rental_rate_values,
9      COUNT(length)AS count_length_values,
10     MIN(replacement_cost)AS min_replacement_cost,
11     MAX(replacement_cost)AS max_replacement_cost,
12     AVG(replacement_cost)AS avg_replacement_cost,
13     COUNT(replacement_cost)AS count_replacement_cost_values,
14     MIN(length)AS min_length,
15     MAX(length)AS max_length,
16     AVG(length)AS avg_length
17 FROM film
```

| min_duration | max_duration | avg_duration | count_rental_duration_values |
|---|---|---|---|
| 3 | 7 | 4,985 | 1000 |
| min_rate | max_rate | avg_rate | count_rental_rate_values |
| 0.99 | Apr 99 | 2,98 | 1000 |
| min_replacement_cost | max_replacement_cost | avg_replacement_cost | count_replacement_cost_values |
| 9,99 | 29.99 | 19,98 | 1000 |
| min_length | max_length | avg_length | count_length_values |
| 46 | 185 | 115,27 | 1000 |

| title_modal_value | rating_modal_value | description_modal_value |
|---|---|---|
| Academy Dinosaur | PG-13 | A Action-Packed Character Study of a Astronaut And a Explorer who must Reach a Monkey in A MySQL Convention |

| release_year_modal_value | special_features_modal_value |
|---|---|
| 2006 | {Trailers,Commentaries,"Behind the Scenes"} |

Query    Query History                    Data Output    Messages

```
1  SELECT MODE() WITHIN GROUP (ORDER BY title)
2      AS title_modal_value,
3      MODE() WITHIN GROUP (ORDER BY rating)
4      AS rating_modal_value,
5      MODE() WITHIN GROUP (ORDER BY description)
6      AS description_modal_value,
7      MODE() WITHIN GROUP (ORDER BY release_year)
8      AS release_year_modal_value,
9      MODE() WITHIN GROUP (ORDER BY special_features)
10     AS special_features_modal_value
11 FROM film
```

| | title_modal_value<br>character varying |
|---|---|
| 1 | Academy Dinosaur |

- 3). I am slowly becoming more comfortable with SQL. Thinking about how to best clean and organise data, Excel seems slightly ahead with its ability to filter the data within each column. If you want PG-13 films from 1994 for example. SQL is easier because I just had to copy-paste the queries one after the other, changing what it was I was looking for, so I could create larger/longer strings. (*See the MIN/MAX example).* SQL may end up being more useful for me (in the sense of I could use it more easily) than Excel, but I assume it's all about using both daily for long enough to really grasp the strengths and weaknesses of each.