# Summarizing & 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).

**Film table / Duplicate Data**

Query   Query History

```
 1    SELECT film_id,
 2           title,
 3           description,
 4           release_year,
 5           language_id,
 6           rental_duration,
 7           rental_rate,
 8           replacement_cost,
 9           rating,
10           count(*)
11    FROM film
12    GROUP BY film_id,
13           title,
14           description,
15           release_year,
16           language_id,
17           rental_duration,
18           rental_rate,
19           replacement_cost,
20           rating
21    HAVING COUNT(*) > 1
```
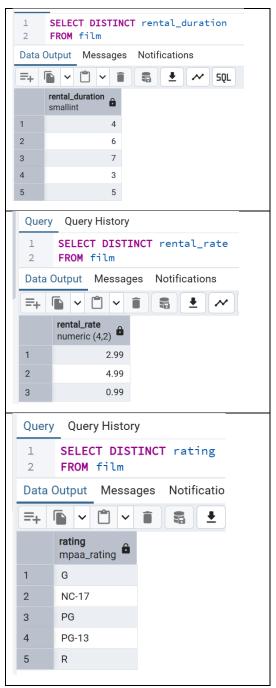
Total rows: 0    Query complete 00:00:00.061

No duplicate data was found in the film table.

If there is duplicate data,

(1) use DISTINCT or GROUP BY in queries to return unique records for analysis.

(2) create a view to create unique records if permitted by supervisor.

**Film table / Non-Uniform Data**

```
1    SELECT DISTINCT rental_duration
2    FROM film
```

Data Output    Messages    Notifications

| | rental_duration 🔒 smallint |
|---|---|
| 1 | 4 |
| 2 | 6 |
| 3 | 7 |
| 4 | 3 |
| 5 | 5 |

Query    Query History

```
1    SELECT DISTINCT rental_rate
2    FROM film
```

Data Output    Messages    Notifications

| | rental_rate 🔒 numeric (4,2) |
|---|---|
| 1 | 2.99 |
| 2 | 4.99 |
| 3 | 0.99 |

Query    Query History

```
1    SELECT DISTINCT rating
2    FROM film
```

Data Output    Messages    Notificatio

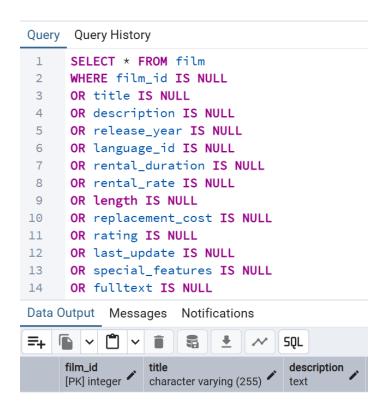| | rating 🔒 mpaa_rating |
|---|---|
| 1 | G |
| 2 | NC-17 |
| 3 | PG |
| 4 | PG-13 |
| 5 | R |

No non-uniform data was found.

If there is non-uniform data, I will use UPDATE the table using the commands as follows:

UPDATE film
SET column_name = 'value'
Where column_name IN ('value1', 'value2', 'value3')

**Film table/Missing Values**

```
Query    Query History
1     SELECT * FROM film
2     WHERE film_id IS NULL
3     OR title IS NULL
4     OR description IS NULL
5     OR release_year IS NULL
6     OR language_id IS NULL
7     OR rental_duration IS NULL
8     OR rental_rate IS NULL
9     OR length IS NULL
10    OR replacement_cost IS NULL
11    OR rating IS NULL
12    OR last_update IS NULL
13    OR special_features IS NULL
14    OR fulltext IS NULL
```

Data Output    Messages    Notifications

| film_id<br>[PK] integer | title<br>character varying (255) | description<br>text |
|---|---|---|

No missing values were identified in the table. If missing values are present in future analyses and account for less than 5% of the data, they will be addressed using either listwise deletion or mean imputation. If missing values exceed 5%, I will ignore them and clearly explained why I do that.

The SQL commands for mean imputations are as follows:

UPDATE film
SET = AVG(column_name)
WHERE column_name IS NULL

**Customer table / Duplicate Data**

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

Data Output   Messages   Notifications

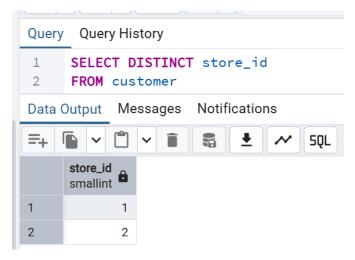| customer_id [PK] integer | store_id smallint | first_name character varying (45) | last_name character varying (45) | email character varying (50) | address_id smallint | count bigint |
|---|---|---|---|---|---|---|

No duplicate data was found in the customer table.

If there is duplicate data,

    (1)  use DISTINCT or GROUP BY in queries to return unique records for analysis.

    (2)  create a view  to create unique records if permitted by supervisor.

**Customer table / Non-Uniform Data**

Query   Query History

```
1    SELECT DISTINCT store_id
2    FROM customer
```

Data Output   Messages   Notifications

| store_id smallint |
|---|
| 1 |
| 2 |

No non-uniform data was found.
If there is non-uniform data, I will use UPDATE the table using the commands as follows:

UPDATE customer
SET column_name = 'value'
Where column_name IN ('value1', 'value2', 'value3')

**Customer table/Missing Values**

```
1    SELECT customer_id,
2           store_id,
3           first_name,
4           last_name,
5           email,
6           address_id
7    FROM customer
8    WHERE customer_id IS NULL
9    OR    store_id IS NULL
10   OR    first_name IS NULL
11   OR    last_name IS NULL
12   OR    email IS NULL
13   OR    address_id IS NULL
```

Data Output   Messages   Notifications

| customer_id [PK] integer | store_id smallint | first_name character varying (45) | last_name character varying (45) | email character varying (50) | address_id smallint |
|---|---|---|---|---|---|

No missing values were identified in the table. If missing values are present in future analyses and account for less than 5% of the data, they will be addressed using either listwise deletion or mean imputation. If missing values exceed 5%, I will ignore them and clearly explained why I do that.

The SQL commands for mean imputations are as follows:

UPDATE customer
SET = AVG(column_name)
WHERE column_name 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. Copy-paste your SQL queries and their outputs into your answers document.**
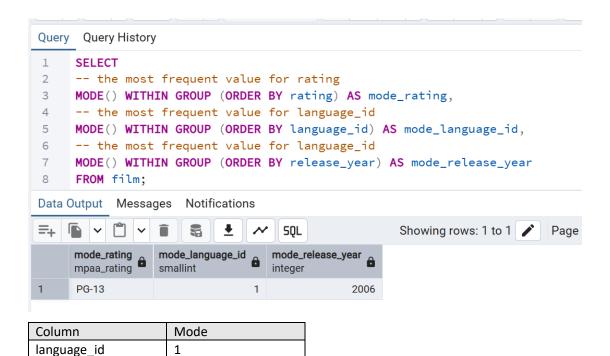
**Film Table**

```sql
1   SELECT
2   -- relase_duration
3   AVG(rental_duration) AS avg_rental_duration,
4   MAX(rental_duration) AS max_rental_duration,
5   MIN(rental_duration) AS min_rental_duration,
6   -- rental_rate
7   AVG(rental_rate) AS avg_rental_rate,
8   MAX(rental_rate) AS max_rental_rate,
9   MIN(rental_rate) AS min_rate,
10  -- length
11  AVG(length) AS avg_length,
12  MAX(length) AS max_length,
13  MIN(length) AS min_length,
14  -- replacement_cost
15  AVG(replacement_cost) AS avg_replacement_cost,
16  MAX(replacement_cost) AS max_replacement_cost,
17  MIN(replacement_cost) AS min_replacement_cost
18  FROM film
19
```

| Column | Average | Maximum | Minimum |
|---|---|---|---|
| rental_duration | 4.99 | 7 | 3 |
| rental_rate | 2.98 | 4.99 | 0.99 |
| length | 115.272 | 185 | 46 |
| replacement_cost | 19.984 | 29.99 | 9.99 |

```sql
1   SELECT
2   -- the most frequent value for rating
3   MODE() WITHIN GROUP (ORDER BY rating) AS mode_rating,
4   -- the most frequent value for language_id
5   MODE() WITHIN GROUP (ORDER BY language_id) AS mode_language_id,
6   -- the most frequent value for language_id
7   MODE() WITHIN GROUP (ORDER BY release_year) AS mode_release_year
8   FROM film;
```

Data Output   Messages   Notifications

Showing rows: 1 to 1   Page

| | mode_rating mpaa_rating | mode_language_id smallint | mode_release_year integer |
|---|---|---|---|
| 1 | PG-13 | 1 | 2006 |

| Column | Mode |
|---|---|
| language_id | 1 |
| rating | PG- 13 |
| release_year | 2006 |

## Customer Table

```sql
1   SELECT
2   -- create_date
3   COUNT(create_date) AS count_create_date,
4   MAX(create_date) AS max_create_date,
5   MIN(create_date) AS min_create_date,
6   -- last_update
7   COUNT(last_update) AS count_last_update,
8   MAX(last_update) AS max_last_update,
9   MIN(last_update) AS min_last_update,
10  -- the most frequent value for activebool
11  MODE() WITHIN GROUP (ORDER BY activebool) AS mode_activebool,
12  -- the most frequent value for active
13  MODE() WITHIN GROUP (ORDER BY active) AS mode_active,
14  -- the most frequent value for first_name
15  MODE() WITHIN GROUP (ORDER BY first_name) AS mode_first_name,
16  -- the most frequent value for last_name
17  MODE() WITHIN GROUP (ORDER BY last_name) AS mode_last_name
18  FROM customer
```

| Column | Count | Maximum | Minimum |
|---|---|---|---|
| creation_date | 599 | 2006-02-14 | 2006-02-14 |
| last_update | 599 | 2013-05-26 14:49:45.738 | 2013-05-26 14:49:45.738 |

| Column | Mode |
|---|---|
| activebool | true |
| active | 1 |
| first_name | Jamie |
| Last_name | Abney |

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. Write a short paragraph in the running document that you have started.**

Ease of use: Excel is easier for starters. SQL might get easier after you get familiar with the commands and rules.

Functions: SQL is more powerful because it can deal with complicated tasks with several lines of commands. In comparison, Excel is much less automatic.

Speed: SQL is faster than Excel especially when dealing with large datasets.