









3.6: Summarizing & Cleaning Data in SQL

Duplicates








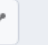
Film table

```
1 SELECT title, release_year, language_id,  
2 COUNT(*)  
3 FROM film  
4 GROUP BY title, release_year, language_id  
5 HAVING COUNT(*)>1;
```

| Data Output | Messages | Notifications |
|---|--------------|---------------|
|         | | |
| title | release_year | language_id |
| character varying (255) | integer | smallint |
| | | count |
| | | bigint |

Customer table

```
1 SELECT first_name, last_name, email, address_id,  
2 COUNT(*)  
3 FROM customer  
4 GROUP BY first_name, last_name, email, address_id  
5 HAVING COUNT(*)>1;
```

| Data Output | Messages | Notifications |
|---|------------------------|------------------------|
|         | | |
| first_name | last_name | email |
| character varying (45) | character varying (45) | character varying (50) |
| | | address_id |
| | | smallint |
| | | count |
| | | bigint |

In none of the tables could be duplicates identified. However, in case that yes, I could create a VIEW with only unique records. Another option would be to remove duplicate data with DELETE. If I don't have permission to do that, I'll need to write a query that returns only unique records. There are 2 ways to do this – GROUP BY and DISTINCT.

Non-Uniform Data

Film Table

Query Query History

```
1 SELECT release_year, language_id, rating
2 FROM film
3 GROUP BY release_year, language_id, rating;
```

Data Output Messages Notifications

| | release_year integer | language_id smallint | rating mpaa_rating |
|---|-------------------------|-------------------------|-----------------------|
| 1 | 2006 | 1 | PG-13 |
| 2 | 2006 | 1 | NC-17 |
| 3 | 2006 | 1 | G |
| 4 | 2006 | 1 | PG |
| 5 | 2006 | 1 | R |

Customer table

```
1 SELECT DISTINCT activebool, create_date
2 FROM customer;
```

Data Output Messages Notifications

| activebool boolean | create_date date |
|-----------------------|---------------------|
| true | 2006-02-14 |













No non-uniform could be identified. In case that yes, the UPDATE can be used to make values consistent.

Missing Values

Film Table

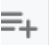






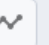



```
1 SELECT title, description, language_id, release_year
2 FROM film
3 WHERE release_year
4 IS null;
```

Data Output Messages Notifications

| | | | | | | | |
|--|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  |  |
| | | | | | | | |
| title character varying (255)  | | | | description text  | | language_id smallint  | |
| | | | | | | release_year integer  | |

```
1 SELECT first_name, last_name, email
2 FROM customer
3 WHERE email
4 IS null;
```

Data Output Messages Notifications

| | | | | | | | |
|--|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  |  |
| | | | | | | | |
| first_name character varying (45)  | | | last_name character varying (45)  | | | email character varying (50)  | |
| | | | | | | | |

No missing values could be identified. If a column contains a high number of missing values, the column should be excluded from the analyse by omitting it from the SELECT command. If a column contains a low amount of missing values, these missing values can be imputed with an estimate, such as column average.

2. Descriptive Statistics

Film table

```
1 SELECT
2     MIN(rental_duration) AS min_duration,
3     MAX(rental_duration) AS max_duration,
4     AVG(rental_duration) AS avg_duration,
5
6     MIN(rental_rate) AS min_rent,
7     MAX(rental_rate) AS max_rent,
8     AVG(rental_rate) AS avg_rent,
9
10    MIN(length) AS min_length,
11    MAX(length) AS max_length,
12    AVG(length) AS avg_length,
13
14    MIN(replacement_cost) AS min_replacement,
15    MAX(replacement_cost) AS max_replacement,
16    AVG(replacement_cost) AS avg_replacement,
17
18    MODE() WITHIN GROUP (ORDER BY release_year) AS mode_year,
19    MODE() WITHIN GROUP (ORDER BY language_id) AS mode_language_id,
20    MODE() WITHIN GROUP (ORDER BY rating) AS mode_rating
21
22 FROM film
```

Data Output Messages Notifications

| min_duration | max_duration | avg_duration | min_rent | max_rent | avg_rent | min_length | max_length | avg_length | min_replacement | max_replacement | avg_replacement | mode_year | mode_language_id | mode |
|--------------|--------------|--------------|----------|----------|----------|------------|------------|------------|-----------------|-----------------|-----------------|-----------|------------------|-------|
| smallint | smallint | numeric | numeric | numeric | numeric | smallint | smallint | numeric | numeric | numeric | numeric | integer | smallint | mpaa |
| 3 | 7 | 4.985 | 0.99 | 4.99 | 2.98 | 46 | 185 | 115.272 | 9.99 | 29.99 | 19.984 | 2006 | 1 | PG-13 |

Customer table

```
1 SELECT
2     MIN(customer_id) AS min_cust_id,
3     MAX(customer_id) AS customer_id,
4     AVG(customer_id) AS customer_id,
5
6     MIN(store_id) AS store_id,
7     MAX(store_id) AS store_id,
8     MODE() WITHIN GROUP (ORDER BY store_id) AS mode_store_id,
9
10    MIN(address_id) AS address_id,
11    MAX(address_id) AS address_id,
12
13    MIN(create_date) AS create_date,
14    MAX(create_date) AS create_date,
15    MODE() WITHIN GROUP (ORDER BY create_date) AS mode_create_date,
16
17    MIN(last_update) AS min_last_update,
18    MAX(last_update) AS last_update,
19    MODE() WITHIN GROUP (ORDER BY last_update) AS last_update,
20    mode() WITHIN GROUP (ORDER BY activebool) AS mode_activebool
21
22 FROM customer;
```

Data Output Messages Notifications

| min_cust_id | customer_id | customer_id | store_id | store_id | mode_store_id | address_id | address_id | create_date | create_date | mode_create_date | min_last_update | last_update | last |
|-------------|-------------|-------------|----------|----------|---------------|------------|------------|-------------|-------------|------------------|-----------------------------|-----------------------------|------|
| integer | integer | numeric | smallint | smallint | smallint | smallint | smallint | date | date | date | timestamp without time zone | timestamp without time zone | tim |
| 1 | 599 | 300 | 1 | 2 | 1 | 5 | 605 | 2006-02-14 | 2006-02-14 | 2006-02-14 | 2013-05-26 14:49:45.738 | 2013-05-26 14:49:45.738 | 20 |

3. Reflection

SQL makes data profiling very simple and quick (once you get used to the Syntax and the logic). Especially when data set is large, SQL provides a very efficient way to work with a data set, especially when it's about descriptive statistics. You can quickly identify the completeness and uniqueness of the dataset with DISTINCT and GROUP BY. However, since I am used to work with Excel already for several years, SQL is more complicated and slower for me. Especially the work with a smaller dataset would be faster with Excel. It's literally like learning a new language – first you need to set a steady base (the Syntax) and then it's practice and more practice to become fluent.