Data Immersion
Databases & SQL for Analysts
3.6: Summarizing & Cleaning Data in SQL
David Guillen Aroche

**1. Check for and clean dirty data for the film table and the customer table.**

- **Duplicates.**

*Film Table*

Query:
```sql
SELECT film_id,
       title,
       description,
       release_year,
       language_id,
       rental_duration,
       rental_rate,
       film.length,
       replacement_cost,
       rating,
       last_update,
       special_features,
       fulltext,
     COUNT(*)
FROM film
GROUP BY film_id,
       title,
       description,
       release_year,
       language_id,
       rental_duration,
       rental_rate,
       film.length,
       replacement_cost,
       rating,
       last_update,
       special_features,
       fulltext
HAVING COUNT(*) >1
```

Result:

No duplicates were found in film table, see screenshot below.

| ase_year<br>ger | language_id<br>smallint | rental_duration<br>smallint | rental_rate<br>numeric (4,2) | length<br>smallint | replacement_cost<br>numeric (5,2) | rating<br>mpaa_rating | last_update<br>timestamp without time zone | special_features<br>text[] | fulltext<br>tsvector | count<br>bigint |
|---|---|---|---|---|---|---|---|---|---|---|

Total rows: 0 of 0    Query complete 00:00:00.202                                                                Ln 254. Col 19

In case of finding duplicate values, I would have created a view with unique records with the following syntax to make each row unique:

```
CREATE VIEW viewname AS
SELECT col1,
       col2,
       col3 ...
FROM tablename
GROUP BY col1,
         col2,
         col3, ... ;
```

***Customer Table***

Query:
```
SELECT customer_id,
    store_id,
    first_name,
    email,
    address_id,
    activebool,
    last_update,
    active,
    COUNT(*)
FROM customer
GROUP BY customer_id,
    store_id,
    first_name,
    email,
    address_id,
    activebool,
```

last_update,
        active
HAVING COUNT(*) > 1

Result:
No duplicates were found in customers table, see screenshot below.



| ase_year<br>ger | language_id<br>smallint | rental_duration<br>smallint | rental_rate<br>numeric (4,2) | length<br>smallint | replacement_cost<br>numeric (5,2) | rating<br>mpaa_rating | last_update<br>timestamp without time zone | special_features<br>text[] | fulltext<br>tsvector | count<br>bigint |
|---|---|---|---|---|---|---|---|---|---|---|

Total rows: 0 of 0    Query complete 00:00:00.202                                                                Ln 254. Col 19

In case of finding duplicate values, I would have created a view with unique records with the following syntax to make each row unique:

```
CREATE VIEW viewname AS
SELECT col1,
       col2,
       col3 ...
FROM tablename
GROUP BY col1,
         col2,
         col3, ... ;
```

- **Non-uniform data.**

***Film Table***

Query (syntax):
SELECT DISTINCT column_name
FROM ctable_name
GROUP BY column_name

Results:
Non-unifrom values were searched in the table, but no irregularities were found. The total number of values for each column is as follows:
   o   *film_id, 1000 values*
   o   *title, 1000 values*
   o   *description, 1000 values*
   o   *release_year, 1 value*

- o  *language_id, 1 value*
- o  *rental_duration, 5 values*
- o  *rental_rate, 3 values*
- o  *film.length, 140 values*
- o  *replacement_cost, 21 values*
- o  *rating, 5 values*
- o  *last_update, 1 value*
- o  *special_features, 15 vlaues*
- o  *fulltext, 1000 values*

In case of finding a non-uniform value, first, I would have grouped the values to get familiar with them to choose a standard format, and checked on data types for each column. Then use the following syntax to update values:

```
UPDATE film

SET rating = 'G'

WHERE rating IN ('gen',

                 'g',

                 'General')
```

***Customer Table***

Query (syntax):
SELECT DISTINCT column_name
FROM ctable_name
GROUP BY column_name

Results:
No Non-uniform values were searched and have found no irregularities. The total number of values for each column as follows:
- o  *customer_id, 599 values*
- o  *store_id, 2 values*
- o  *first_name, 591 values*
- o  *last_name, 599 values*
- o  *email, 599 values*
- o  *address_id, 599 values*
- o  *activebool, 1 value*
- o  *create_date, 1 value*
- o  *active, 1 value*

e.g.:

```
304
305    SELECT DISTINCT active
306    FROM customer
307    GROUP BY active
```

**Data output**   **Messages**   **Notification**

| | active integer 🔒 |
|---|---|
| 1 | 0 |
| 2 | 1 |

Save results to file
F8

In case of finding a non-uniform value, first, I would have grouped the values to get familiar with them to choose a standard format, and checked on data types for each column. Then use the following syntax to update values:

```
UPDATE film
SET rating = 'G'
WHERE rating IN ('gen',
                 'g',
                 'General')
```
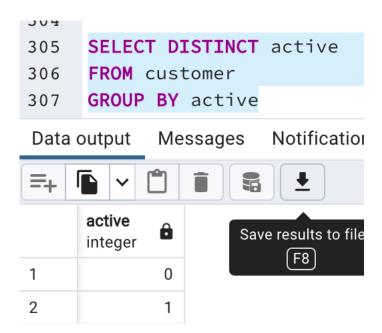
- **Missing values**

*Film Table*

While looking for **Non-uniform** values, no missing values were found.

In case of have found missing values, I would have checked on the proportion of the missing data for each column to determine if would have been valid to include them in the analysis or not, and adapt the analysis to mitigate the impact as much as possible.

*Customer Table*

While looking for **Non-uniform** values, no missing values were found.

In case of have found missing values, I would have checked on the proportion of the missing data for each column to determine if would have been valid to include them in the analysis or not, and adapt the analysis to mitigate the impact as much as possible.

## 2. Summarize Data.

***Film Table***
Query (Syntax) Numeric Columns:
*SELECT MIN(column_name) AS column_alias,*
    *MAX(column_name) AS column_alias,*
    *AVG(column_name) AS column_alias,*
    *COUNT(column_name) AS column_alias,*
    *COUNT(\*) AS column_alias*
*FROM table_name*

*--film_id.*
The average was not included for this query as it does not really say much.

```
296  SELECT MIN(film_id) AS lowest_id,
297         MAX(film_id) AS highest_id,
298         COUNT(film_id) AS films_count,
299         COUNT(*) AS tuples_count
300  FROM film
```

Data output    Messages    Notifications

| lowest_id<br>integer | highest_id<br>integer | films_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|
| 1 | 1000 | 1000 | 1000 |

*--release_year*
The average **years count** was not included since there is only one value for this column.

```
296  SELECT MIN(release_year) AS minimum_year,
297         MAX(release_year) AS maximum_year,
298         COUNT(release_year) AS years_count,
299         COUNT(*) AS tuples_count
300  FROM film
```

Data output    Messages    Notifications

| minimum_year<br>integer | maximum_year<br>integer | years_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|
| 2006 | 2006 | 1000 | 1000 |

## --language_id

There is only one language value for the total tuples in the table.

```
296  SELECT language_id,
297         MIN(language_id) AS minimum_language_id,
298         MAX(language_id) AS maximum_language_id,
299         AVG (language_id) AS average_languages_count,
300         COUNT(language_id) AS languages_count,
301         COUNT(*) AS tuples_count
302  FROM film
303  GROUP BY language_id
```

Data output  Messages  Notifications

| | language_id<br>smallint | minimum_language_id<br>smallint | maximum_language_id<br>smallint | average_languages_count<br>numeric | languages_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | 1000 | 1000 |

## --rental_duration

```
296  SELECT MIN(rental_duration) AS minimum_rental_duration,
297         MAX(rental_duration) AS maximum_rental_duration,
298         AVG (rental_duration) AS average_rental_duration,
299         COUNT(rental_duration) AS rental_duration_values_count,
300         COUNT(*) AS tuples_count
301  FROM film
```

Data output  Messages  Notifications

| | minimum_rental_duration<br>smallint | maximum_rental_duration<br>smallint | average_rental_duration<br>numeric | rental_duration_values_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|---|
| 1 | 3 | 7 | 4.985 | 1000 | 1000 |

## --rental_rate

```
296  SELECT MIN(rental_rate) AS minimum_rental_rate,
297         MAX(rental_rate) AS maximum_rental_rate,
298         AVG (rental_rate) AS average_rental_rate,
299         COUNT(rental_rate) AS rental_rate_values_count,
300         COUNT(*) AS tuples_count
301  FROM film
```

Data output  Messages  Notifications

| | minimum_rental_rate<br>numeric | maximum_rental_rate<br>numeric | average_rental_rate<br>numeric | rental_rate_values_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|---|
| 1 | 0.99 | 4.99 | 2.98 | 1000 | 1000 |

*--film.length*

```
296   SELECT MIN(film.length) AS minimum_film_length,
297          MAX(film.length) AS maximum_film_length,
298          AVG (film.length) AS average_film_length,
299          COUNT(film.length) AS film_length_count,
300          COUNT(*) AS tuples_count
301   FROM film
```

Data output    Messages    Notifications

| | minimum_film_length smallint | maximum_film_length smallint | average_film_length numeric | film_length_count bigint | tuples_count bigint |
|---|---|---|---|---|---|
| 1 | 46 | 185 | 115.272 | 1000 | 1000 |

*--replacement_cost*

```
296   SELECT MIN(replacement_cost) AS minimum_replacement_cost,
297          MAX(replacement_cost) AS maximum_replacement_cost,
298          AVG (replacement_cost) AS average_replacement_cost,
299          COUNT(replacement_cost) AS replacement_cost_count,
300          COUNT(*) AS tuples_count
301   FROM film
```

Data output    Messages    Notifications

| | minimum_replacement_cost numeric | maximum_replacement_cost numeric | average_replacement_cost numeric | replacement_cost_count bigint | tuples_count bigint |
|---|---|---|---|---|---|
| 1 | 9.99 | 29.99 | 19.984 | 1000 | 1000 |

Query (Syntax) Non-Numeric Columns:
SELECT mode() WITHIN GROUP (ORDER BY column_name)
     AS modal_value
FROM tablename

*--title*

```
305   SELECT mode() WITHIN GROUP (ORDER BY title)
306          AS modal_value
307   FROM film
```

Data output    Messages    Notifications

| | modal_value character varying |
|---|---|
| 1 | Academy Dinosaur |

*--description*

```
305   SELECT mode() WITHIN GROUP (ORDER BY description)
306           AS modal_value
307   FROM film
```

Data output | Messages | Notifications

| | modal_value<br>text |
|---|---|
| 1 | A Action-Packed Character Study of a Astronaut And a Explorer who must Reach a Monkey in A MySQL Convention |

*--rating*

```
305   SELECT mode() WITHIN GROUP (ORDER BY rating)
306           AS modal_value
307   FROM film
```

Data output | Messages | Notifications

| | modal_value<br>mpaa_rating |
|---|---|
| 1 | PG-13 |

*--special_features*

```
305   SELECT mode() WITHIN GROUP (ORDER BY special_features)
306           AS modal_value
307   FROM film
```

Data output | Messages | Notifications

| | modal_value<br>text[] |
|---|---|
| 1 | {Trailers,Commentaries,"Behind the Scenes"} |

*--fulltext*

```
305   SELECT mode() WITHIN GROUP (ORDER BY fulltext)
306           AS modal_value
307   FROM film
```

Data output | Messages | Notifications

| | modal_value<br>tsvector |
|---|---|
| 1 | 'baloon':19 'confront':14 'documentari':5 'feminist':8,11,16 'mile':2 'must':13 'spi':1 'thrill':4 |

***Customer Table***

Query (Syntax) Numeric Columns:
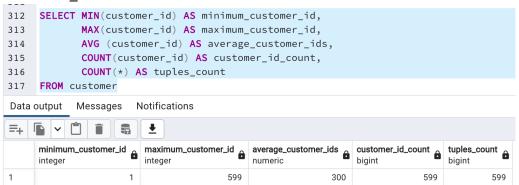*SELECT MIN(column_name) AS column_alias,*
    *MAX(column_name) AS column_alias,*
    *AVG(column_name) AS column_alias,*
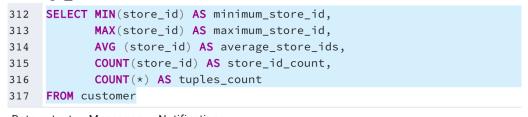    *COUNT(column_name) AS column_alias,*
    *COUNT(\*) AS column_alias*
*FROM table_name*

*--customer_id*
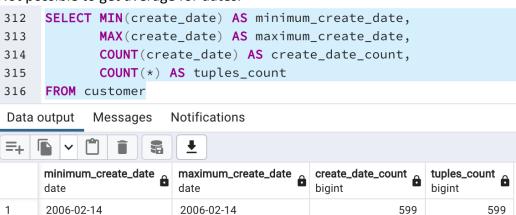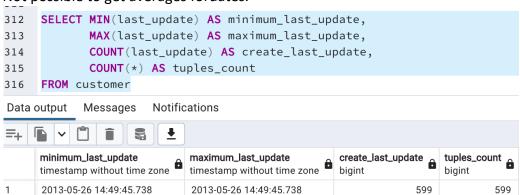
```
312   SELECT MIN(customer_id) AS minimum_customer_id,
313         MAX(customer_id) AS maximum_customer_id,
314         AVG (customer_id) AS average_customer_ids,
315         COUNT(customer_id) AS customer_id_count,
316         COUNT(*) AS tuples_count
317   FROM customer
```

Data output   Messages   Notifications

| | minimum_customer_id<br>integer | maximum_customer_id<br>integer | average_customer_ids<br>numeric | customer_id_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|---|
| 1 | 1 | 599 | 300 | 599 | 599 |

*--storage_id*

```
312   SELECT MIN(store_id) AS minimum_store_id,
313         MAX(store_id) AS maximum_store_id,
314         AVG (store_id) AS average_store_ids,
315         COUNT(store_id) AS store_id_count,
316         COUNT(*) AS tuples_count
317   FROM customer
```

Data output   Messages   Notifications

| | minimum_store_id<br>smallint | maximum_store_id<br>smallint | average_store_ids<br>numeric | store_id_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|---|
| 1 | 1 | 2 | 1.4557595993322203 | 599 | 599 |

*--address_id*

```
312   SELECT MIN(address_id) AS minimum_address_id,
313          MAX(address_id) AS maximum_address_id,
314          AVG (address_id) AS average_address_ids,
315          COUNT(address_id) AS address_id_count,
316          COUNT(*) AS tuples_count
317   FROM customer
```

Data output    Messages    Notifications

| | minimum_address_id<br>smallint | maximum_address_id<br>smallint | average_address_ids<br>numeric | address_id_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|---|
| 1 | 5 | 605 | 304.7245409015025 | 599 | 599 |

*--create_date*

Not possible to get average for dates.

```
312   SELECT MIN(create_date) AS minimum_create_date,
313          MAX(create_date) AS maximum_create_date,
314          COUNT(create_date) AS create_date_count,
315          COUNT(*) AS tuples_count
316   FROM customer
```

Data output    Messages    Notifications

| | minimum_create_date<br>date | maximum_create_date<br>date | create_date_count<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|
| 1 | 2006-02-14 | 2006-02-14 | 599 | 599 |

*--last_update*

Not possible to get averages fordates.

```
312   SELECT MIN(last_update) AS minimum_last_update,
313          MAX(last_update) AS maximum_last_update,
314          COUNT(last_update) AS create_last_update,
315          COUNT(*) AS tuples_count
316   FROM customer
```

Data output    Messages    Notifications

| | minimum_last_update<br>timestamp without time zone | maximum_last_update<br>timestamp without time zone | create_last_update<br>bigint | tuples_count<br>bigint |
|---|---|---|---|---|
| 1 | 2013-05-26 14:49:45.738 | 2013-05-26 14:49:45.738 | 599 | 599 |

*--active*

Averages was not calculated for this columns as does not make sense to get it; only to values.

```
312    SELECT MIN(active) AS minimum_active,
313           MAX(active) AS maximum_active,
314           COUNT(active) AS create_active,
315           COUNT(*) AS tuples_count
316    FROM customer
```

Data output    Messages    Notifications

| | minimum_active integer | maximum_active integer | create_active bigint | tuples_count bigint |
|---|---|---|---|---|
| 1 | 0 | 1 | 599 | 599 |

Query (Syntax) Non-Numeric Columns:
SELECT mode() WITHIN GROUP (ORDER BY column_name)
     AS modal_value
FROM tablename

*--first_name*

```
320    SELECT mode() WITHIN GROUP (ORDER BY first_name)
321           AS modal_value
322    FROM customer
```

Data output    Messages    Notifications

| | modal_value character varying |
|---|---|
| 1 | Jamie |

*--last_name*

```
320    SELECT mode() WITHIN GROUP (ORDER BY last_name)
321           AS modal_value
322    FROM customer
```

Data output    Messages    Notifications

| | modal_value character varying |
|---|---|
| 1 | Abney |

*--email*
**Please consider there are 599, and 599 different emails. This results is for one of those 599 different emails.**

```
320    SELECT mode() WITHIN GROUP (ORDER BY email)
321           AS modal_value
322    FROM customer
```

Data output    Messages    Notifications

| modal_value<br>character varying 🔒 |
|---|
| 1    aaron.selby@sakilacustomer.org |

3. **Which tool (Excel or SQL) is more efficient for data profiling and why? Consider functions, ease of use, and speed.**

In general, SQL is way more effective at providing data types, data integrity, and calculating descriptive statistics of the data. Personally, I find the difference essentially in the capacity that both have in terms of processing data, and how they show the data; SQL's results are succinct, and is considerably quicker coding a query rather than making the data profile from scratch in Excel—though making a data profile in Excel is slower and has less data processing potence, it also is more intuitive once you know what you are doing, and there are more resources around to troubleshoot in excel than in SQL.

*Functions.*
I believe SQL is more effective at providing results; the used functions are basically the same as in Excel, but the one difference would be knowing how to code queries with everything at once in SQL, rather than making a longer process in Excel for each one of the functions—leaving aside the data processing capacity. The downside of SQL is that you must know how to query, you have to have the technical knowledge to do so, otherwise it might be difficult to impossible to retrieve the data the way you want it.

*Ease of Use.*
It's easier to use Excel even though the data integrity analysis represents more time. The issue with Excel is that it has less capacity to process data. On the other hand, once you have (at least the basic) technical knowledge to code queries, it's easier to use SQL, again, because of how succinct the process and the results are

*Speed.*
SQL is by far the fastest way to make a data integrity analysis. It's faster to code a query, even if it's a long and/or complext one (once you have the knowledge).

In conclusion, SQL is better to make a data integrity analysis, but not the easiest.