



DECEMBER 18, 2025

## 3.6 SUMMARIZING & CLEANING DATA IN SQL

JANELLA VALENCIA



## Table of Contents

<b>STEP 1- CHECK FOR AND CLEAN DIRTY DATA:</b> .....	<b>2</b>
<b>FILM DUPLICATES</b> .....	<b>2</b>
<b>FILM NON-UNIFORM DATA</b> .....	<b>2</b>
<b>FILM MISSING DATA</b> .....	<b>3</b>
<b>CUSTOMER DUPLICATES</b> .....	<b>3</b>
<b>CUSTOMER NON-UNIFORM DATA</b> .....	<b>4</b>
<b>CUSTOMER MISSING DATA</b> .....	<b>4</b>
<b>STEP 2- SUMMARIZE YOUR DATA:</b> .....	<b>5</b>
<b>NUMERIC FILM COLUMNS</b> .....	<b>5</b>
<b>NON-NUMERIC FILM COLUMN</b> .....	<b>5</b>
<b>NUMERIC CUSTOMER COLUMNS</b> .....	<b>5</b>
<b>NON-NUMERIC CUSTOMER COLUMN</b> .....	<b>5</b>
<b>STEP 3- REFLECT ON YOUR WORK:</b> .....	<b>6</b>

## Step 1- Check for and clean dirty data:

### Film Duplicates




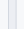








```
1 SELECT title, release_year, language_id, rental_duration,
2 COUNT(*)
3 FROM film
4 GROUP BY title, release_year, language_id, rental_duration
5 HAVING COUNT(*) >1
```

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

- No records showed as duplicates under the film table. In order to delete duplicate records I would create a view in order to visualize the unique records. I would then delete the records in question. Depending on my position in the company I would either be responsible for the deletion OR I wouldn't be, in that case I would use GROUP BY or DISTINCT to show unique records for analyzing purposes.

### Film Non-Uniform Data

```
33 SELECT release_year,rating,
34 COUNT(*)
35 FROM film
36 GROUP BY release_year,rating
37 ORDER BY release_year,rating;
38
```

Data Output	Messages	Notifications	
<div></div>			
	release_year 	rating 	count 
	integer	mpaa_rating	bigint
1	2006	G	178
2	2006	PG	194
3	2006	PG-13	223
4	2006	R	195
5	2006	NC-17	210

- No records showed as non-uniform data. I checked for release year and rating because I believed in the rating is where we may see the most non-uniform data. If there would have been non-uniformed data I would have used the UPDATE command to make the values uniformed.



## Customer Non-Uniform Data

```
1 SELECT customer_id, store_id, first_name, last_name, email, create_date, activebool
2 FROM customer
3 GROUP BY customer_id
4 ORDER BY customer_id
```

Data Output Messages Notifications

Showing rows: 1 to 599 Page No: 1

	customer_id [PK] integer	store_id smallint	first_name character varying (45)	last_name character varying (45)	email character varying (50)	create_date date	activebool boolean
1	1	1	Mary	Smith	mary.smith@sakilacustomer.org	2006-02-14	true
2	2	1	Patricia	Johnson	patricia.johnson@sakilacustomer.org	2006-02-14	true
3	3	1	Linda	Williams	linda.williams@sakilacustomer.org	2006-02-14	true
4	4	2	Barbara	Jones	barbara.jones@sakilacustomer.org	2006-02-14	true
5	5	1	Elizabeth	Brown	elizabeth.brown@sakilacustomer.org	2006-02-14	true
6	6	2	Jennifer	Davis	jennifer.davis@sakilacustomer.org	2006-02-14	true
7	7	1	Maria	Miller	maria.miller@sakilacustomer.org	2006-02-14	true
8	8	2	Susan	Wilson	susan.wilson@sakilacustomer.org	2006-02-14	true
9	9	2	Margaret	Moore	margaret.moore@sakilacustomer.org	2006-02-14	true
10	10	1	Dorothy	Taylor	dorothy.taylor@sakilacustomer.org	2006-02-14	true
11	11	2	Lisa	Anderson	lisa.anderson@sakilacustomer.org	2006-02-14	true
12	12	1	Nancy	Thomas	nancy.thomas@sakilacustomer.org	2006-02-14	true
13	13	2	Karen	Jackson	karen.jackson@sakilacustomer.org	2006-02-14	true
14	14	2	Betty	White	betty.white@sakilacustomer.org	2006-02-14	true
15	15	1	Helen	Harris	helen.harris@sakilacustomer.org	2006-02-14	true
16	16	2	Sandra	Martin	sandra.martin@sakilacustomer.org	2006-02-14	true
17	17	1	Donna	Thompson	donna.thompson@sakilacustomer.org	2006-02-14	true
18	18	2	Carol	Garcia	carol.garcia@sakilacustomer.org	2006-02-14	true
19	19	1	Ruth	Martinez	ruth.martinez@sakilacustomer.org	2006-02-14	true
20	20	2	Sharon	Robinson	sharon.robinson@sakilacustomer.org	2006-02-14	true

- There is no non-uniform data for customers. If there were non-uniformed data I would UPDATE command to make the consistent

## Customer Missing Data

```
6 SELECT *
7 FROM customer
8 WHERE first_name IS NULL
9       OR last_name IS NULL
10      OR email IS NULL
11      OR store_id IS NULL
12      OR address_id IS NULL;
13
```

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	activebool boolean	create_date date	last_update timestamp without time zone	active integer
-----------------------------	----------------------	--------------------------------------	-------------------------------------	---------------------------------	------------------------	-----------------------	---------------------	--	-------------------

- There is no missing data for customer table. I checked for the main columns. If there was missing data I would check to see within what percentage it fell. If it was 5% I would go ahead and use the AVG of the data if it was more than 30% I wouldn't use that information.

## Step 2- Summarize your data:

- Use SQL to calculate descriptive statistics for both the film table and the customer table.

### Numeric Film Columns

```
45 SELECT
46 MIN(rental_duration) AS min_rental,
47 MAX(rental_duration) AS max_rental,
48 AVG(rental_duration) AS avg_rental,
49 MIN(length) AS min_length,
50 MAX(length) AS max_length,
51 AVG(length) AS avg_length,
52 MIN(rental_rate) AS min_rate,
53 MAX(rental_rate) AS max_rate,
54 AVG(rental_rate) AS avg_rate,
55 MIN(replacement_cost) AS min_cost,
56 MAX(replacement_cost) AS max_cost,
57 AVG(replacement_cost) AS avg_cost
58 FROM film;
```

	min_rental smallint	max_rental smallint	avg_rental numeric	min_length smallint	max_length smallint	avg_length numeric	min_rate numeric	max_rate numeric	avg_rate numeric	min_cost numeric	max_cost numeric	avg_cost numeric
1	3	7	4.9850000000000000	46	185	115.2720000000000000	0.99	4.99	2.9800000000000000	9.99	29.99	19.9840000000000000

### Non-Numeric Film Column

60	SELECT MODE() WITHIN GROUP
61	(ORDER BY language_id)
62	AS modal_value
63	FROM film;

	modal_value smallint
1	1

65	SELECT MODE() WITHIN GROUP
66	(ORDER BY rating)
67	AS modal_value
68	FROM film;

	modal_value mpaa_rating
1	PG-13

65	SELECT MODE() WITHIN GROUP
66	(ORDER BY title)
67	AS modal_value
68	FROM film;

	modal_value character varying
1	Academy Dinosaur

### Numeric Customer Columns

- No non numeric columns in customer table

### Non-Numeric Customer column

86	SELECT MODE() WITHIN GROUP
87	(ORDER BY first_name)
88	AS modal_value
89	FROM customer;

	modal_value character varying
1	Jamie

86	SELECT MODE() WITHIN GROUP
87	(ORDER BY address_id)
88	AS modal_value
89	FROM customer;

	modal_value smallint
1	5

86	SELECT MODE() WITHIN GROUP
87	(ORDER BY email)
88	AS modal_value
89	FROM customer;

	modal_value character varying
1	aaron.selby@sakilacustomer...

86	SELECT MODE() WITHIN GROUP
87	(ORDER BY activebool)
88	AS modal_value
89	FROM customer;

Data Output	Messages	Notifications
-------------	----------	---------------

≡+	📄	▼	📋	▼	🗑️	🗄️	⬇️	📈
----	---	---	---	---	----	----	----	---

	modal_value boolean	🔒
1	true	

Step 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? Consider their respective functions, ease of use, and speed.
  - o In excel we were able to go a step above and find stand deviations and correlation coefficient ; not to say we can't do that in SQL, but I haven't searched or learned it yet in the program. In excel we had to manually find the MIN,MAX,AVG for each section we needed. The advantage was that we were able to work with different sheets at once.
  - o For SQL, I found it much quicker to get information like MIN,MAX or AVG in one SQL command for multiple columns. It was all arranged in a chart and easy to read.