

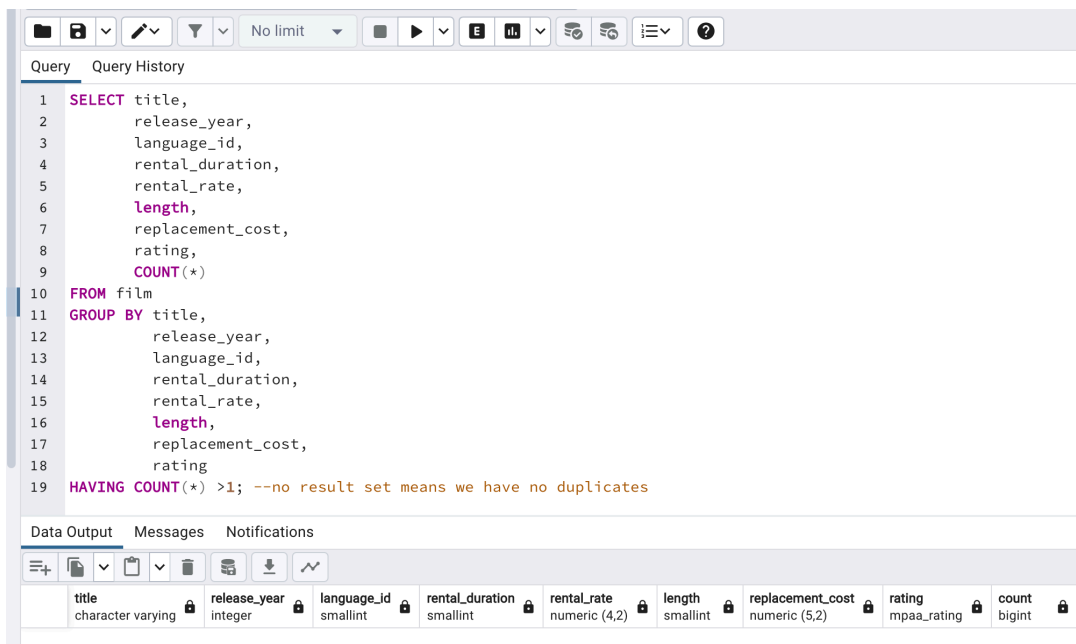
## 3.6: Summarizing & Cleaning Data in SQL

- 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).

Answer: [checking for duplicate data](#)

For Film Table:

```
SELECT title, release_year, language_id, rental_duration, rental_rate, length,  
       replacement_cost, rating,  
       COUNT(*) FROM film  
GROUP BY title, release_year, language_id, rental_duration, rental_rate, length,  
       replacement_cost, rating  
HAVING COUNT(*) >1; --no result set means we have no duplicates
```



The screenshot shows a SQL query editor interface. The query is as follows:

```
1 SELECT title,  
2       release_year,  
3       language_id,  
4       rental_duration,  
5       rental_rate,  
6       length,  
7       replacement_cost,  
8       rating,  
9       COUNT(*)  
10 FROM film  
11 GROUP BY title,  
12          release_year,  
13          language_id,  
14          rental_duration,  
15          rental_rate,  
16          length,  
17          replacement_cost,  
18          rating  
19 HAVING COUNT(*) >1; --no result set means we have no duplicates
```

Below the query editor, there is a 'Data Output' section with a table showing the schema of the 'film' table:

title	release_year	language_id	rental_duration	rental_rate	length	replacement_cost	rating	count
character varying	integer	smallint	smallint	numeric (4,2)	smallint	numeric (5,2)	mpaa_rating	bigint

For Customer Table:

```
SELECT store_id, first_name, last_name, email,  
       COUNT(*)  
FROM customer  
GROUP BY store_id,  
       first_name, last_name, email  
HAVING COUNT(*) >1; --no result set means we have no duplicates
```

The screenshot shows a SQL query editor with a toolbar at the top. The query is as follows:

```

1 SELECT store_id,
2       first_name,
3       last_name,
4       email,
5       COUNT(*)
6 FROM customer
7 GROUP BY store_id,
8         first_name,
9         last_name,
10        email
11 HAVING COUNT(*) > 1; --no result set means we have no duplicates

```

Below the query editor, there are tabs for "Data Output", "Messages", and "Notifications". The "Data Output" tab is active, showing a table with the following columns and data types:

store_id	first_name	last_name	email	count
smallint	character varying	character varying	character varying	bigint

At the bottom, a status bar indicates "Total rows: 0 of 0" and "Query complete 00:00:00.225".

I would create View and select only unique records and delete the duplicate records. As a junior data analyst, I wouldn't be responsible for deleting the records.

### Checking for non-uniform data:

For Film Table:

The screenshot shows a SQL query editor with a toolbar at the top. The query is as follows:

```

1 SELECT DISTINCT
2       release_year
3 FROM film;

```

Below the query editor, there are tabs for "Data Output", "Messages", and "Notifications". The "Data Output" tab is active, showing a table with the following columns and data types:

release_year
integer

At the bottom, a status bar indicates "Total rows: 0 of 0" and "Query complete 00:00:00.225".

**SELECT DISTINCT  
release\_year FROM film;**

Rockbuster/postgres@PostgreSQL 15

Query Query History

```

1 SELECT DISTINCT
2   rating
3 FROM film;

```

Data Output Messages Notifications

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

SELECT DISTINCT  
rating  
FROM film;

Rockbuster/postgres@Po:

Query Query History

```

1 SELECT DISTINCT
2   rental_rate
3 FROM film;

```

Data Output Messages Noti

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

SELECT DISTINCT  
rental\_rate  
FROM film;

For customer Table:

Query Query History

```

1 SELECT DISTINCT
2   store_id
3 FROM customer;

```

Data Output Messages Notif

	store_id smallint
1	1
2	2

SELECT DISTINCT  
store\_id  
FROM customer;

Query		Query History
1	SELECT DISTINCT	
2	active	
3	FROM customer;	

Data Output		Messages
	active	
	integer	
1	0	
2	1	

SELECT DISTINCT  
active  
FROM customer;

After checking with these queries, we can say that there is no non-uniform data. If there is any non-uniform data then we could perform UPDATE, SET; WHERE functions to adjust entries.

## Checking for missing data:

For Film Table:

Query		Query History
1	SELECT *	
2	FROM film	
3	WHERE rating IS NULL	

Data Output		Messages	Notifications
	film_id	title	description
	[PK] integer	character varying	text
			release_year
			integer
			language_id
			smallint
			rental_duration
			smallint
			rental_rate
			numeric (4,2)
			length
			smallint
			replacement_cost
			numeric (5,2)
			rating
			mpaa_rating
			last_upda
			timestamp

Query		Query History
1	SELECT *	
2	FROM film	
3	WHERE rental_rate IS NULL	

Data Output		Messages	Notifications
	film_id	title	description
	[PK] integer	character varying	text
			release_year
			integer
			language_id
			smallint
			rental_duration
			smallint
			rental_rate
			numeric (4,2)
			length
			smallint
			replacement_cost
			numeric (5,2)
			rating
			mpaa_rating
			last_upda
			timestamp

Query		Query History
1	SELECT *	
2	FROM film	
3	WHERE film_id IS NULL	

Data Output		Messages	Notifications
	film_id	title	description
	[PK] integer	character varying	text
			release_year
			integer
			language_id
			smallint
			rental_duration
			smallint
			rental_rate
			numeric (4,2)
			length
			smallint
			replacement_cost
			numeric (5,2)
			rating
			mpaa_rating
			last_upda
			timestamp

Query	Query History
1 <b>SELECT</b> *	
2 <b>FROM</b> film	
3 <b>WHERE</b> title <b>IS</b> NULL	

Data Output	Messages	Notifications
<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>		
<div> <div>film_id</div> <div>[PK] integer</div> </div> <div> <div>title</div> <div>character varying</div> </div> <div> <div>description</div> <div>text</div> </div> <div> <div>release_year</div> <div>integer</div> </div> <div> <div>language_id</div> <div>smallint</div> </div> <div> <div>rental_duration</div> <div>smallint</div> </div> <div> <div>rental_rate</div> <div>numeric (4,2)</div> </div> <div> <div>length</div> <div>smallint</div> </div> <div> <div>replacement_cost</div> <div>numeric (5,2)</div> </div> <div> <div>rating</div> <div>mpaa_rating</div> </div> <div> <div>last_upd</div> <div>timestamp</div> </div>		

For Customer Table:

Query	Query History
1 <b>SELECT</b> *	
2 <b>FROM</b> customer	
3 <b>WHERE</b> store_id <b>IS</b> NULL	

Data Output	Messages	Notifications
<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>		
<div> <div>customer_id</div> <div>[PK] integer</div> </div> <div> <div>store_id</div> <div>smallint</div> </div> <div> <div>first_name</div> <div>character varying</div> </div> <div> <div>last_name</div> <div>character varying</div> </div> <div> <div>email</div> <div>character varying</div> </div> <div> <div>address_id</div> <div>smallint</div> </div> <div> <div>activebool</div> <div>boolean</div> </div> <div> <div>create_date</div> <div>date</div> </div> <div> <div>last_update</div> <div>timestamp without time zone</div> </div> <div> <div>active</div> <div>integer</div> </div>		

Query	Query History
1 <b>SELECT</b> *	
2 <b>FROM</b> customer	
3 <b>WHERE</b> first_name <b>IS</b> NULL	

Data Output	Messages	Notifications
<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>		
<div> <div>customer_id</div> <div>[PK] integer</div> </div> <div> <div>store_id</div> <div>smallint</div> </div> <div> <div>first_name</div> <div>character varying</div> </div> <div> <div>last_name</div> <div>character varying</div> </div> <div> <div>email</div> <div>character varying</div> </div> <div> <div>address_id</div> <div>smallint</div> </div> <div> <div>activebool</div> <div>boolean</div> </div> <div> <div>create_date</div> <div>date</div> </div> <div> <div>last_update</div> <div>timestamp without time zone</div> </div> <div> <div>active</div> <div>integer</div> </div>		

Query	Query History
1 <b>SELECT</b> *	
2 <b>FROM</b> customer	
3 <b>WHERE</b> email <b>IS</b> NULL	

Data Output	Messages	Notifications
<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>		
<div> <div>customer_id</div> <div>[PK] integer</div> </div> <div> <div>store_id</div> <div>smallint</div> </div> <div> <div>first_name</div> <div>character varying</div> </div> <div> <div>last_name</div> <div>character varying</div> </div> <div> <div>email</div> <div>character varying</div> </div> <div> <div>address_id</div> <div>smallint</div> </div> <div> <div>activebool</div> <div>boolean</div> </div> <div> <div>create_date</div> <div>date</div> </div> <div> <div>last_update</div> <div>timestamp without time zone</div> </div> <div> <div>active</div> <div>integer</div> </div>		

I did not find any missing data in any column. If I would find any missing data then I would perform these two options:

1. Impute missing values with average value by calculating average function and then would perform UPDATE function to set the average value.
2. Ignore columns with high number of missing values.

**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.

```
SELECT film_id, MAX(rental_duration) AS
maximum_rental_duration, MIN(rental_duration) AS
minimum_rental_duration, AVG(rental_duration) AS
average_rental_duration, MAX(rental_rate) AS
maximum_rental_rate, MIN(rental_rate) AS minimum_rental_rate,
AVG(rental_rate) AS average_rental_rate, MAX(length) AS
maximum_movie_length, MIN(length) AS minimum_movie_length,
AVG(length) AS average_movie_length, MAX(replacement_cost)
AS maximum_replacement_cost, MIN(replacement_cost) AS
minimum_replacement_cost, AVG(replacement_cost) AS
average_replacement_cost
```

```
FROM film GROUP BY film_id, rental_duration, rental_rate, length,
replacement_cost
```

Data Output Messages Notifications							
	maximum_rental_duration smallint	minimum_rental_duration smallint	average_rental_duration numeric	count_rental_duration bigint	maximum_rental_rate numeric	minimum_rental_rate numeric	average_rental_rate numeric
1	7	3	4.9850000000000000	1000	4.99	0.99	2.9800000000000000

Data Output Messages Notifications							
	count_rental_rate bigint	maximum_movie_length smallint	minimum_movie_length smallint	average_movie_length numeric	count_length bigint	maximum_replacement_cost numeric	minimum_replacement_cost numeric
1	1000	185	46	115.2720000000000000	1000	29.99	9.99

Data Output Messages Notifications							
	average_movie_length numeric	count_length bigint	maximum_replacement_cost numeric	minimum_replacement_cost numeric	average_replacement_cost numeric	count_replacement_cost bigint	count_rows bigint
1	115.2720000000000000	1000	29.99	9.99	19.9840000000000000	1000	1000

Query

Query History

1

-- Mode for non-numeric columns.

2

SELECT MODE() WITHIN GROUP (ORDER BY title) AS title\_mode\_value,

3

MODE() WITHIN GROUP (ORDER BY description) AS description\_mode\_value,

4

MODE() WITHIN GROUP (ORDER BY last\_update) AS last\_update\_mode\_value,

5

MODE() WITHIN GROUP (ORDER BY release\_year) AS release\_year\_mode\_value,

6

MODE() WITHIN GROUP (ORDER BY rating) AS rating\_mode\_value,

7

MODE() WITHIN GROUP (ORDER BY fulltext) AS full\_text\_mode\_value

8

FROM film

Data Output

Messages

Notifications

title\_mode\_value

character varying

description\_mode\_value

text

last\_update\_mode\_value

timestamp without time zone

release\_year\_mode\_value

integer

1

Academy Dinosaur

A Action-Packed Character Study of a Astronaut And a Explorer who must Reach a Monkey in A MySQL Convent...

2013-05-26 14:50:58.951

2006

Data Output

Messages

Notifications

		last_update_mode_value timestamp without time zone	release_year_mode_value integer	rating_mode_value mpaa_rating	full_text_mode_value tsvector
1	in MySQL Convent...	2013-05-26 14:50:58.951	2006	PG-13	'baloon':19 'confront':14 'documentari':5 'feminist':8,11,16 'mile':2 'must':13 'spi':1 'thri...

Query

Query History

```
1  -- Mode for non-numeric columns.
2  SELECT  MODE() WITHIN GROUP (ORDER BY first_name) AS firstname_mode_value,
3          MODE() WITHIN GROUP (ORDER BY last_name) AS lastname_mode_value,
4          MODE() WITHIN GROUP (ORDER BY last_update) AS last_update_mode_value,
5          MODE() WITHIN GROUP (ORDER BY create_date) AS createdate_mode_value,
6          MODE() WITHIN GROUP (ORDER BY active) AS active_mode_value
7
8          FROM customer
```

Data Output

Messages

Notifications

	firstname_mode_value character varying	lastname_mode_value character varying	last_update_mode_value timestamp without time zone	createdate_mode_value date	active_mode_value integer
1	Jamie	Abney	2013-05-26 14:49:45.738	2006-02-14	1

**3. Reflect on your work:** At this stage, we can not say much more about both the languages. They both have their own pros and cons. Some Excel to SQL points are there :

- **SQL** is not actually harder than **Excel**, it's just a little different. Once we learned the language, acts much more like “talking” as opposed to simply managing or manipulating data given its syntax.
- It can take minutes in **SQL** to do what it takes nearly an hour to do in **Excel**.
- **Excel** is a great program for simplicity and flexibility. **SQL** databases are excellent choices for storage, manipulation and analysis of large amount of data.

- **Excel** is extremely intuitive. We can see the data and calculations at a glance. The formulas are simple for anyone to understand. It's inexpensive and there's a free version.
- **Excel** is prone to human error. It's easy to type over a critical formula.
- Data is much safer in **SQL** data storage than in **Excel**, since it's more difficult for a user to delete data by mistake.
- Collaboration isn't a strong point for **Excel**, since it's easy for other users to delete critical information or formulas.
- **Excel** can technically handle one million rows, but that's before the pivot tables, multiple tabs, and functions you're probably using.
- **SQL** is much faster than **Excel**. **SQL** also separates analysis from data.