

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.

DUPLICATE DATA

The screenshot shows the PostgreSQL IDE interface. On the left, the Object Explorer displays the database schema, with the 'film' table selected under 'Columns (13)'. The main query editor contains the following SQL query:

```
1 SELECT film_id,
2 title,
3 description,
4 release_year,
5 language_id,
6 rental_duration,
7 rental_rate,
8 length,
9 replacement_cost,
10 rating,
11 last_update,
12 special_features,
13 fulltext
14 FROM film
15 GROUP BY film_id,
16 title,
17 description,
18 release_year,
19 language_id,
20 rental_duration,
21 rental_rate,
22 length,
23 replacement_cost,
24 rating,
25 last_update,
26 special_features,
27 fulltext
28 HAVING COUNT(*) > 1
29
```

Below the query editor, the Data Output tab shows the results of the query. The table structure is displayed with columns: film_id [PK] integer, title character varying (255), description text, release_year integer, language_id smallint, rental_duration smallint, rental_rate numeric (4,2), length smallint, replacement_cost numeric (5,2), and rating numeric (2,1). The status bar at the bottom indicates 'Total rows: 0 of 0' and 'Query complete 00:00:00.351'.

The screenshot shows the PostgreSQL IDE interface. On the left, the Object Explorer displays the database schema, with the 'customer' table selected under 'Columns (10)'. The main query editor contains the following SQL query:

```
1 SELECT customer_id,
2 store_id,
3 first_name,
4 last_name,
5 email,
6 address_id,
7 activebool,
8 create_date,
9 last_update,
10 active
11 FROM customer
12 GROUP BY customer_id,
13 store_id,
14 first_name,
15 last_name,
16 email,
17 address_id,
18 activebool,
19 create_date,
20 last_update,
21 active
22 HAVING COUNT(*) > 1
23
```

Below the query editor, the Data Output tab shows the results of the query. The table structure is displayed with columns: 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, and last_update timestamp with time zone. The status bar at the bottom indicates 'Total rows: 0 of 0' and 'Query complete 00:00:00.248'.

No duplicates were found in either table, however, if there had been duplicated data, there are two ways to fix them:

1. Create a virtual table, known as a “view,” where you select only unique records.
2. Delete the duplicate record from the table or view.

According to the lesson of this exercise, creating a new view is the standard way for a data analyst to handle duplicate records. On the other hand, deleting data is not recommended and sometimes it is even unauthorised.

NON-UNIFORM

Rockbuster/postgres@PostgreSQL 16

Query Query History

```

1  --Returns only unique records from the film table
2
3  SELECT DISTINCT
4  film_id,
5  title,
6  description,
7  release_year,
8  rental_duration,
9  rental_rate,
10 length,
11 replacement_cost,
12 rating,
13 last_update,
14 special_features,
15 fulltext
16 FROM film;

```

Data Output Messages Notifications

	film_id [PK] integer	title character varying (255)	description text
1	1	Academy Dinosaur	A Epic Drama of a Feminist And a Man

Total rows: 1000 of 1000 Query complete 00:00:00.527

Query Query History

```

1  --Returns only unique records from the customer table
2
3  SELECT DISTINCT
4  customer_id,
5  store_id,
6  first_name,
7  last_name,
8  email,
9  address_id,
10 activebool,
11 create_date,
12 last_update,
13 active
14 FROM customer;

```

Data Output Messages Notifications

	customer_id [PK] integer	store_id smallint	first_name character varying (45)	last_name character varying (45)
1	257	1	Keith	Diana

Total rows: 599 of 599 Query complete 00:00:00.323

There were no non-uniform data found in either film/customer tables; there is no standard approach to deal with non-uniform data, but in case it had been found, what I would have done to fix that is try to update the values using the following command:

```

UPDATE x
SET x = 'x'
WHERE x IN ('x')...

```

MISSING DATA

Query	Query History
1	SELECT *
2	FROM customer
3	WHERE (customer_id,
4	store_id,
5	first_name,
6	last_name,
7	email,
8	address_id,
9	activebool,
10	create_date,
11	last_update,
12	active)
13	IS NULL;
14	

Data Output	Messages	Notifications
<div> <div>+</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> </div>		
customer_id [PK] integer	store_id smallint	first_name character varying (45)
Total rows: 0 of 0 Query complete 00:00:00.678		

Query	Query History
1	SELECT *
2	FROM film
3	WHERE (title,
4	description,
5	release_year,
6	language_id,
7	rental_duration,
8	rental_rate,
9	length,
10	replacement_cost,
11	rating,
12	last_update,
13	special_features,
14	fulltext)
15	IS NULL;
16	

Data Output	Messages	Notifications
<div> <div>+</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> <div>▼</div> </div>		
film_id [PK] integer	title character varying (255)	description text
Total rows: 0 of 0 Query complete 00:00:00.565		

There were no missing values found in both tables; but in any case the suggested technique to avoid that is to ignore columns with a high percent of missing values by omitting whichever column you want to ignore, with the following SELECT statement: *ignored in select because it has a lot of missing values.*

The second suggested approach is to impute the missing values using statistical methods such as AVG, MIN, MAX. Using the UPDATE command.

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.

NUMERICAL VALUES FILM COLUMN

- rental_duration
- rental_rate
- Length
- replacement_cost

```

Query  Query History
1  SELECT MAX (rental_duration) AS max_rental_duration,
2      MIN (rental_duration) AS min_rental_duration,
3      AVG (rental_duration) AS avg_rental_duration,
4      MAX (rental_rate) AS max_rental_rate,
5      MIN (rental_rate) AS min_rental_rate,
6      AVG (rental_rate) AS avg_rental_rate,
7      MAX (length) AS max_length,
8      MIN (length) AS min_length,
9      AVG (length) AS avg_length,
10     MAX (replacement_cost) AS max_replacement_cost,
11     MIN (replacement_cost) AS min_replacement_cost,
12     AVG (replacement_cost) AS avg_replacement_cost
13 FROM film;
14

```

	max_rental_duration smallint	min_rental_duration smallint	avg_rental_duration numeric
1	7	3	4.9850000000000000

max_rental_rate numeric	min_rental_rate numeric	avg_rental_rate numeric
4.99	0.99	2.9800000000000000

max_length smallint	min_length smallint	avg_length numeric
185	46	115.27200000000000

max_replacement_cost numeric	min_replacement_cost numeric	avg_replacement_cost numeric
29.99	9.99	19.984000000000000

NON-NUMERICAL VALUES FILM TABLE

- film_id
- Title
- language_id
- Rating
- last_update
- special_features
- Fulltext

```
1 SELECT MODE () WITHIN GROUP (ORDER BY film_id) AS film_id_modal_value,
2 MODE () WITHIN GROUP (ORDER BY title) AS title_modal_value,
3 MODE () WITHIN GROUP (ORDER BY language_id) AS language_id_modal_value,
4 MODE () WITHIN GROUP (ORDER BY rating) AS rating_modal_value,
5 MODE () WITHIN GROUP (ORDER BY last_update) AS last_update_modal_value,
6 MODE () WITHIN GROUP (ORDER BY special_features) AS special_features_modal_value,
7 MODE () WITHIN GROUP (ORDER BY fulltext) AS fulltext_modal_value
8 FROM film;
9
10
11
```

	film_id_modal_value integer	title_modal_value character varying	language_id_modal_value smallint	rating_modal_value mpaa_rating	last_update_modal_value timestamp without time zone
1	1	Academy Dinosaur	1	PG-13	2013-05-26 14:50:58.951

special_features_modal_value text[]	fulltext_modal_value tsvector
{Trailers,Commentaries,"Behind the Scenes"}	'baloon':19 'confront':14 'documentari':5 'feminist':8,11,16 'mile':2 'must':13 'spi':1 'thri...

NON-NUMERICAL VALUES CUSTOMER TABLE

*(no numerical values for the customer table)

- store_id
- Activebool
- create_date
- last_update

1

2

3

4

5

6

7

8

```
SELECT MODE () WITHIN GROUP (ORDER BY store_id) AS store_id_modal_value,
MODE () WITHIN GROUP (ORDER BY activebool) AS activebool_modal_value,
MODE () WITHIN GROUP (ORDER BY create_date) AS create_date_modal_value,
MODE () WITHIN GROUP (ORDER BY last_update) AS last_update_modal_value
FROM customer;
```

Data Output

Messages

Notifications

<

3. 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.

Neither is absolutely better than the other because each is designed and intended to be used in different records keeping and manipulation situations, but in this case, in my opinion, once you learn how to structure queries on SQL, it might be faster than the long step-by-step guide for doing this in Excel: Import data, sort data, filter certain values, delete unwanted rows, and so on.