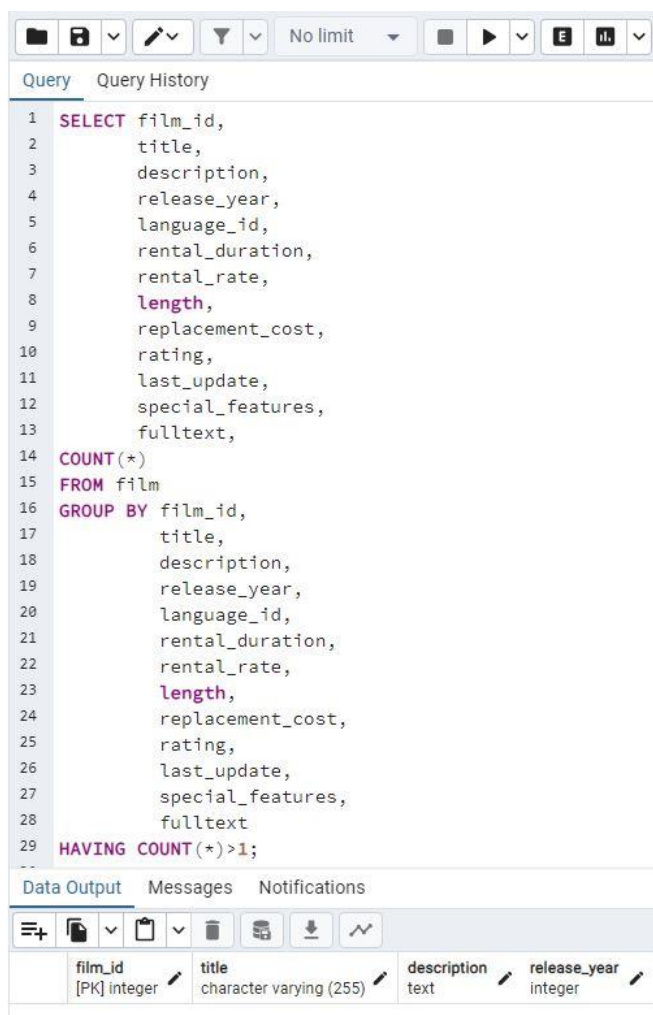


Data Immersion 3.6

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).

Duplicate Data

Film Table



The screenshot shows a SQL query editor with a toolbar at the top containing icons for file operations, query execution, and settings. Below the toolbar, there are tabs for "Query" and "Query History". The main area displays a SQL query designed to identify duplicate rows in the "film" table based on all columns except "film_id".

```
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 COUNT(*)
15 FROM film
16 GROUP BY film_id,
17          title,
18          description,
19          release_year,
20          language_id,
21          rental_duration,
22          rental_rate,
23          length,
24          replacement_cost,
25          rating,
26          last_update,
27          special_features,
28          fulltext
29 HAVING COUNT(*) > 1;
```

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: "film_id" (PK) integer, "title" character varying (255), "description" text, and "release_year" integer.

I found no duplicates in the film_table and customer_table; if duplicates were present, I would have removed them using GROUP BY or DISTINCT, or created a view table to manage the duplicates.

Customer Table

Query

Query History

```
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 COUNT(*)
12 FROM customer
13 GROUP BY customer_id,
14          store_id,
15          first_name,
16          last_name,
17          email,
18          address_id,
19          activebool,
20          create_date,
21          last_update,
22          active
23 HAVING COUNT(*) < 1;
```

Data Output

Messages

Notifications

+

📄

▼

📋

▼

🗑️

🔍

⬇️

📶

store_id	last_name	email
integer	character varying (45)	character varying (50)

Non-uniform Data

Film Table

Query		Query History	
1	SELECT DISTINCT	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	ORDER BY	film_id;	
16			
Data Output		Messages	
		Notifications	
	film_id [PK] integer	title character varying (255)	description text
1	1	Academy Dinosaur	A Epic Drama of a Femir
2	2	Ace Goldfinger	A Astounding Epistle of
3	3	Adaptation Holes	A Astounding Reflection
4	4	Affair Prejudice	A Fanciful Documentary
5	5	African Egg	A Fast-Paced Document
6	6	Agent Truman	A Intrepid Panorama of
7	7	Airplane Sierra	A Touching Saga of a Hu
8	8	Airport Pollock	A Epic Tale of a Moose
9	9	Alabama Devil	A Thoughtful Panorama
10	10	Aladdin Calendar	A Action-Packed Tale of
11	11	Alamo Videotape	A Boring Epistle of a But
12	12	Alaska Phantom	A Fanciful Saga of a Hur
13	13	Ali Forever	A Action-Packed Drama
14	14	Alice Fantasia	A Emotional Drama of a
15	15	Alien Center	A Brilliant Drama of a Ce
16	16	Alley Evolution	A Fast-Paced Drama of










There were no uniformity issues in the mentioned tables; however, if encountered, I would correct them using SQL commands like UPDATE, GROUP BY, or DISTINCT to ensure data consistency.

Customer Table

Query		Query History	
1	SELECT DISTINCT	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	ORDER BY film_id;		
16			
Data Output		Messages	
film_id	title	description	
[PK] integer	character varying (255)	text	
1	Academy Dinosaur	A Epic Drama of a Femir	
2	Ace Goldfinger	A Astounding Epistle of	
3	Adaptation Holes	A Astounding Reflection	
4	Affair Prejudice	A Fanciful Documentary	
5	African Egg	A Fast-Paced Document	
6	Agent Truman	A Intrepid Panorama of	
7	Airplane Sierra	A Touching Saga of a Hu	
8	Airport Pollock	A Epic Tale of a Moose	
9	Alabama Devil	A Thoughtful Panorama	
10	Aladdin Calendar	A Action-Packed Tale of	
11	Alamo Videotape	A Boring Epistle of a But	
12	Alaska Phantom	A Fanciful Saga of a Hur	
13	Ali Forever	A Action-Packed Drama	
14	Alice Fantasia	A Emotional Drama of a	
15	Alien Center	A Brilliant Drama of a Ce	
16	Alley Evolution	A Fast-Paced Drama of	












Missing Values

Film Table

Query	Query History	
3	WHERE (film_id,	
4	title,	
5	description,	
6	release_year,	
7	language_id,	
8	rental_duration,	
9	rental_rate,	
10	length,	
11	replacement_cost,	
12	rating,	
13	last_update,	
14	special_features,	
15	fulltext) IS NULL	
16	ORDER BY film_id;	
17		
Data Output	Messages	Notifications
        		
film_id [PK] integer	title character varying (255)	description text

I encountered no missing values in the tables, but if necessary, I would impute missing data using SQL syntax to fill values with column averages or other relevant statistics, applying the WHERE clause to specify conditions for imputing data.

Customer Table

Query	Query History	
1	SELECT *	
2	FROM customer	
3	WHERE (store_id,	
4	first_name,	
5	last_name,	
6	email,	
7	address_id,	
8	activebool,	
9	create_date,	
10	last_update,	
11	active) IS NULL;	
12		
Data Output	Messages	Notifications
        		
customer_id [PK] integer 	store_id smallint 	first_name character varyi

2. Summarise 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.

Numeric

Query

Query History

```
1  SELECT
2    MIN(release_year) AS min_release_year,
3    MIN(rental_duration) AS min_rentdur,
4    MIN(rental_rate) AS min_rate,
5    MIN(length) AS min_leng,
6    MIN(replacement_cost) AS min_replac,
7    MAX(release_year) AS max_release_year,
8    MAX(rental_duration) AS max_rentdur,
9    MAX(rental_rate) AS max_rate,
10   MAX(length) AS max_leng,
11   MAX(replacement_cost) AS max_replac,
12   AVG(release_year) AS avg_release_year,
13   AVG(rental_duration) AS avg_rentdur,
14   AVG(rental_rate) AS avg_rate,
15   AVG(length) AS avg_leng,
16   AVG(replacement_cost) AS avg_replac
17  FROM film;
```

Data Output

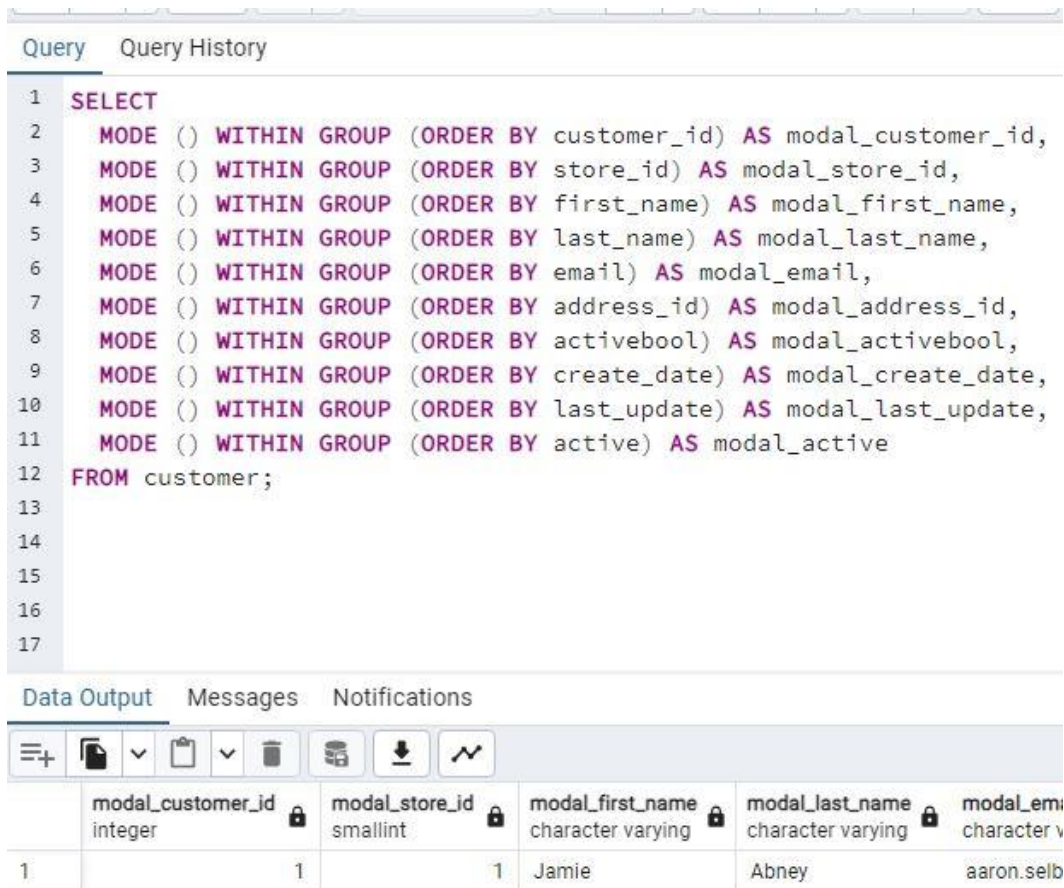
Messages

Notifications

	min_release_year integer	min_rentdur smallint	min_rate numeric	min_leng smallint
1	2006	3	0.99	46

The data could reveal insights such as customer preferences and rental behaviours, using SQL functions like MIN, MAX, AVG for numeric data and MODE for non-numeric data to summarize customer interactions and movie popularity, although the customer_table lacks numerical data for detailed statistical analysis.

Non-numeric



The screenshot shows a SQL query editor with a query window and a data output window. The query is as follows:

```
1 SELECT
2     MODE () WITHIN GROUP (ORDER BY customer_id) AS modal_customer_id,
3     MODE () WITHIN GROUP (ORDER BY store_id) AS modal_store_id,
4     MODE () WITHIN GROUP (ORDER BY first_name) AS modal_first_name,
5     MODE () WITHIN GROUP (ORDER BY last_name) AS modal_last_name,
6     MODE () WITHIN GROUP (ORDER BY email) AS modal_email,
7     MODE () WITHIN GROUP (ORDER BY address_id) AS modal_address_id,
8     MODE () WITHIN GROUP (ORDER BY activebool) AS modal_activebool,
9     MODE () WITHIN GROUP (ORDER BY create_date) AS modal_create_date,
10    MODE () WITHIN GROUP (ORDER BY last_update) AS modal_last_update,
11    MODE () WITHIN GROUP (ORDER BY active) AS modal_active
12 FROM customer;
13
14
15
16
17
```

The data output window shows the following table:

	modal_customer_id integer	modal_store_id smallint	modal_first_name character varying	modal_last_name character varying	modal_email character varying
1	1	1	Jamie	Abney	aaron.selb

3. Reflect on your work: Back in Achievement 1 you learned about data profiling in Excel. 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. Write a short paragraph in the running document that you have started.

From my perspective, if the task at hand involves data profiling for smaller datasets and a need for quick, visual insights, I would lean towards Excel. Its robust features for visualisation allow for immediate graphical representation of data, which can be very intuitive for spotting trends and outliers. However, when dealing with large datasets or requiring scalability, SQL is a more

powerful choice. Despite having a steeper learning curve and being less user-friendly, SQL can handle complex queries and large volumes of data efficiently. It's a tool designed with performance in mind, capable of profiling data directly within the database, reducing the need for data transfer. In essence, Excel shines for accessibility and visualisation, while SQL excels in data handling capacity and is well-suited for in-depth, large-scale data profiling tasks. I definitely find it logistically easier to use Excel though and its user interface more appealing.