

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).

Checking for Duplicate Data

Film table

Query Query History Scratch Pad x

```
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
```

Data Output Messages Notifications

film_id	title	description	release_year	language_id	rental_duration	rental_rate	length	replacement_cost	rating	last_update	special_features	fulltext	count
[PK] integer	character varying (255)	text	integer	smallint	smallint	numeric (4,2)	smallint	numeric (5,2)	mpaa_rating	timestamp without time zone	text[]	tsvector	bigint

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

customer_id	store_id	first_name	last_name	email	address_id	activebool	create_date	last_update	active	count
[PK] integer	smallint	character varying (45)	character varying (45)	character varying (50)	smallint	boolean	date	timestamp without time zone	integer	bigint

Duplicate Data

If you find duplicate records in your database, there are two ways to fix them:

1. Create a virtual table, known as a “view,” where you select only unique records.

Option 1: Creating a View with Unique Records

```
CREATE VIEW viewname AS
SELECT col1,
       col2,
       col3 ...
FROM tablename
GROUP BY col1,
         col2,
         col3, ... --Group by will make each row
unique
```

2. Delete the duplicate record from the table or view.

```
/* This will delete all duplicate records, keepin
g only one of its versions with the lowest key/un
ique ID value. There must be a unique column in t
he table for this to work.
unique_id is any column acting as a primary key u
nique for each row
*/
DELETE
FROM tablename
WHERE unique_id NOT IN
      (SELECT MIN(unique_id)
       FROM tablename
       GROUP BY col1,
                col2,
                col3, ...)
```

Check for Non-Uniform Data

Film table

Query

Query History

1

SELECT DISTINCT rating

2

FROM film

Data Output

Messages

Notifications

rating

mpaa_rating

1

PG

2

R

3

NC-17

4








PG-13

5

G

Customer Table

Query		Query History	
1	SELECT DISTINCT	first_name,	
2		last_name	
3	FROM	customer	
4	ORDER BY	first_name ASC,	
5		last_name DESC;	

Data Output		Messages	Notifications
			
			
	first_name	last_name	
	character varying (45) 🔒	character varying (45) 🔒	
152	Diane	Collins	
153	Dianne	Shelton	
154	Dolores	Wagner	
155	Don	Bone	
156	Donald	Mahon	
157	Donna	Thompson	
Total rows: 599 of 599		Query complete 00:00:00.101	

Non-Uniform Data

To check for inconsistencies, you would use the GROUP BY and DISTINCT keywords. To fix and update the values in a column for consistency, you would need to use the SQL command below:

```
UPDATE film

SET rating = 'G'

WHERE rating IN ('gen',

                'g',

                'General')
```

Check for Missing Data

Film table

The screenshot shows a database query editor with a query window and a scratch pad. The query is as follows:

```
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 WHERE rating IS null
```

Below the query window, there is a table structure for the 'film' table:

film_id	title	description	release_year	language_id	rental_duration	rental_rate	length	replacement_cost	rating	last_update	special_features	fulltext
[PK] integer	character varying (255)	text	integer	smallint	smallint	numeric (4,2)	smallint	numeric (5,2)	mpaa_rating	timestamp without time zone	text[]	tsvector

You can run the same query for each individual column to determine any missing values. However, there are no missing values in the film table.

Customer Table

The screenshot shows a SQL IDE interface. At the top, there are tabs for 'Query' and 'Query History'. The 'Query' tab is active, displaying a SQL query. To the right of the query editor is a 'Scratch Pad' tab. Below the query editor, there are tabs for 'Data Output', 'Messages', and 'Notifications'. The 'Data Output' tab is active, showing a table with 11 columns: customer_id, store_id, first_name, last_name, email, address_id, activebool, create_date, last_update, and active. The table has a light blue header and a light gray body. The data is as follows:

customer_id	store_id	first_name	last_name	email	address_id	activebool	create_date	last_update	active
1	1	John	Deere	john.deere@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
2	1	Jane	Smith	jane.smith@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
3	1	Mike	Johnson	mike.johnson@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
4	1	Sarah	Williams	sarah.williams@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
5	1	David	Brown	david.brown@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
6	1	Emily	White	emily.white@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
7	1	Chris	Black	chris.black@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
8	1	Alex	Green	alex.green@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
9	1	Olivia	Gray	olivia.gray@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
10	1	Noah	Blue	noah.blue@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1
11	1	Liam	Orange	liam.orange@example.com	1	TRUE	2023-01-01 10:00:00	2023-01-01 10:00:00	1

Missing Data

To avoid missing values or minimize their impact on an analysis. The first approach is to simply ignore columns with a high percentage of missing values. To do so in SQL, you simply omit whichever column you want to ignore from your SELECT statement.

```
SELECT col1,  
  
    col2,  
  
    col4... --col3 ignored in select because it has a lot of missing values  
  
FROM tablename
```

The second approach is to impute the missing values using statistical methods.

```
--imputing missing values with the AVG value
```

```
UPDATE tablename
```

```
SET = AVG(col1)
```

WHERE col1 IS NULL

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.

Minimum, maximum, and average values for film table

Numerical columns (film_id, release_year, language_id, rental_duration, rental_rate, length, replacement_cost, last_update)

Query

Query History

```
1 SELECT MIN(film_id) AS min_film_id,
2        MAX(film_id) AS max_film_id,
3        AVG(film_id) AS avg_film_id,
4        COUNT(film_id) AS count_film_id_values,
5        MIN(release_year) AS min_release_year,
6        MAX(release_year) AS max_release_year,
7        AVG(release_year) AS avg_release_year,
8        COUNT(release_year) AS count_release_year_values,
9        MIN(language_id) AS min_language_id,
10       MAX(language_id) AS max_language_id,
11       AVG(language_id) AS avg_language_id,
12       COUNT(language_id) AS count_language_id_values,
13       MIN(rental_duration) AS min_rental_duration,
14       MAX(rental_duration) AS max_rental_duration,
15       AVG(rental_duration) AS avg_rental_duration,
16       COUNT(rental_duration) AS count_rental_duration_values,
17       MIN(rental_rate) AS min_rental_rate,
18       MAX(rental_rate) AS max_rental_rate,
19       AVG(rental_rate) AS avg_rental_rate,
20       COUNT(rental_rate) AS count_rental_rate_values,
21       MIN(length) AS min_length,
22       MAX(length) AS max_length,
23       AVG(length) AS avg_length,
24       COUNT(length) AS count_length_values,
25       MIN(replacement_cost) AS min_replacement_cost,
26       MAX(replacement_cost) AS max_replacement_cost,
27       AVG(replacement_cost) AS avg_replacement_cost,
28       COUNT(replacement_cost) AS count_replacement_cost_values,
29       MIN(last_update) AS min_last_update,
30       MAX(last_update) AS max_last_update,
31       COUNT(last_update) AS count_last_update_values,
32       COUNT(*) AS count_rows
33 FROM film;
```

Data Output

Messages

Notifications

	min_film_id integer	max_film_id integer	avg_film_id numeric	count_film_id_values bigint	min_release_year integer	max_release_year integer	avg_release_year numeric	count_release_year_values bigint	min_language_id smallint	max_language_id smallint	count_language_id_values smallint
1	1	1000	500.50000000000000000000	1000	2006	2006	2006.00000000000000000000	1000	1	1	1

Mode values for film table

Non-numerical columns (title, description, rating, special_features, fulltext)

Query

Query History

Scratch Pad

✕

1

SELECT mode() WITHIN GROUP (ORDER BY title)

2

AS modal_value,

3

mode() WITHIN GROUP (ORDER BY description)

4

AS modal_value,

5

mode() WITHIN GROUP (ORDER BY rating)

6

AS modal_value,

7

mode() WITHIN GROUP (ORDER BY special_features)

8

AS modal_value,

9

mode() WITHIN GROUP (ORDER BY fulltext)

10

AS modal_value

11

FROM film;

Data Output

Messages

Notifications

modal_value

modal_value

modal_value

modal_value

1

Academy Dinosaur

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

PG-13

{Trailers,Commentaries,'Behind the Scenes'}

'balloon':19 'confront':14 'documentari':5 'feminist':8,11,16 'mile':2 'must':13 'sp':1 'thri...

Minimum, maximum, and average values for customer table

Numerical columns (customer_id, store_id, address_id, create_date, last_update, active)

Query

Query History

Scratch Pad

```

1 SELECT MIN(customer_id) AS min_customer_id,
2        MAX(customer_id) AS max_customer_id,
3        AVG(customer_id) AS avg_customer_id,
4        COUNT(customer_id) AS count_customer_id_values,
5        MIN(store_id) AS min_store_id,
6        MAX(store_id) AS max_store_id,
7        AVG(store_id) AS avg_store_id,
8        COUNT(store_id) AS count_store_id_values,
9        MIN(address_id) AS min_address_id,
10       MAX(address_id) AS max_address_id,
11       AVG(address_id) AS avg_address_id,
12       COUNT(address_id) AS count_address_id_values,
13       MIN(create_date) AS min_create_date,
14       MAX(create_date) AS max_create_date,
15       COUNT(create_date) AS count_create_date_values,
16       MIN(last_update) AS min_last_update,
17       MAX(last_update) AS max_last_update,
18       COUNT(last_update) AS avg_last_update_values,
19       MIN(active) AS min_active,
20       MAX(active) AS max_active,
21       AVG(active) AS avg_active,
22       COUNT(active) AS count_active_values,
23       COUNT(*) AS count_rows
24 FROM customer;

```

Data Output

Messages

Notifications

	min_customer_id integer	max_customer_id integer	avg_customer_id numeric	count_customer_id_values bigint	min_store_id smallint	max_store_id smallint	avg_store_id numeric	count_store_id_values bigint	min_address_id smallint	max_address_id smallint	avg_address_id numeric
1	1	599	300.000000000000000000	599	1	2	1.4557595993322204	599	5	605	304.724540901502504

Query

Query History

```

1 SELECT mode() WITHIN GROUP (ORDER BY first_name)
2     AS modal_value,
3     mode() WITHIN GROUP (ORDER BY last_name)
4     AS modal_value,
5     mode() WITHIN GROUP (ORDER BY email)
6     AS modal_value,
7     mode() WITHIN GROUP (ORDER BY activebool)
8     AS modal_value
9 FROM customer;

```

Data Output

Messages

Notifications

	modal_value character varying	modal_value character varying	modal_value character varying	modal_value boolean
1	Jamie	Abney	aaron.selby@sakilacustomer.org	true

- 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.

SQL is the language used to interact with relational databases. It can work with large volumes of data and operates fast to get results. A well-written SQL query can fetch results from a few million rows within a minute. By using SQL, you can write a query once and reuse it again. As the data increases in your database, you are not required to change much in your query to accommodate similar results.