


Answers 3.6


Find out if the film table and the customer table contain any dirty data, specifically non-uniform or duplicate data, or missing values. 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 Duplicates in “film”

Query Query History 

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

Data output Messages Notifications



title	release_year	language_id	rental_duration	length	count
character varying (255)	integer	smallint	smallint	smallint	bigint

Checking for Duplicates in “customer”

Query Query History

```
1 SELECT first_name, last_name, email, address_id, COUNT(*)
2 FROM customer
3 GROUP BY first_name, last_name, email, address_id
4 HAVING COUNT(*) > 1
```

Data output Messages Notifications



first_name	last_name	email	address_id	count
character varying (45)	character varying (45)	character varying (50)	smallint	bigint

- There are no duplicates in either.
- If there were duplicate records, I would run a view command that will make each row unique based on the desired above command.

Checking for Non-Uniform values

Query

Query History

1

SELECT DISTINCT

rating

2

FROM

film

3

GROUP BY

rating

Data output

Messages

Notification

+

▼

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

- There are no non-uniform values.
- If there were, a command using “UPDATE film” followed by SET rating = ‘[desired rating]’ followed by WHERE listing all of the non-uniform values that need to be altered. This could be repeated for multiple ratings if there were multiple errors.

Checking for Missing Values

- If there are missing values, we must determine whether they need to be replaced.
- If there are missing numeric values, a possible solution, if there are only very few, is to run the following command:
UPDATE name_of_table
SET = AVG(desired_column)
WHERE desired_column IS NULL

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.

For “film” table

Query

Query History

1

SELECT MIN(release_year) AS oldest_release_year,

2

MAX(release_year) AS newest_release_year,

3

mode() WITHIN GROUP (ORDER BY release_year) AS most_common_release_year,

4

MIN(rental_duration) AS min_rental_duration,

5

MAX(rental_duration) AS max_rental_duration,

6

mode() WITHIN GROUP (ORDER BY rental_duration) AS most_common_rental_duration,

7

MIN(rental_rate) AS min_rental_rate,

8

MAX(rental_rate) AS max_rental_rate,

9

AVG(rental_rate) AS avg_rental_rate,

10

MIN(length) AS shortest_film_length,

11

MAX(length) AS longest_film_length,

12

AVG(length) AS avg_film_length,

13

mode() WITHIN GROUP (ORDER BY rating) AS most_common_rating

14

FROM film

Scratch Pad

X

Data output

Messages

Notifications

oldest_release_year

integer

newest_release_year

integer

most_common_release_year

integer

min_rental_duration

smallest

max_rental_duration

smallest

most_common_rental_duration

smallest

min_rental_rate

numeric

max_rental_rate

numeric

avg_rental_rate

numeric

shortest_film_length

smallest

longest_film_length

smallest

avg_film_length

numeric

most_common_rating

mpaa_rating

1

2006

2006

2006

3

7

6

0.99

4.99

2.98

46

185

115.272

PG-13

For “customer” table

Query

Query History

1

SELECT mode() WITHIN GROUP (ORDER BY first_name) AS most_common_first_name,

2

mode() WITHIN GROUP (ORDER BY last_name) AS most_common_larst_name,

3

MIN(create_date) AS oldest_create_date,

4

MAX(create_date) AS newest_create_date,

5

AVG(active) AS customer_activity_rate

6

FROM customer

Data output

Messages

Notifications

	most_common_first_name character varying	most_common_larst_name character varying	oldest_create_date date	newest_create_date date	customer_activity_rate numeric
1	Jamie	Abney	2006-02-14	2006-02-14	0.9749582637729549

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 a significantly better data profiling tool. In SQL you can simply type a command, and you will find all the information you want, whereas in Excel you must fetch it from the table. Looking for errors is even easier in SQL, as the commands save you the trouble of sifting through mass amounts of data, where mistakes are easy to miss. The difference in time efficiency too, gives SQL a boost, which can be crucial in the corporate world.