#### Answers 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).

#### Non-uniform data

SELECT DISTINCT rating FROM film GROUP BY rating

Non-uniform values can be updated to uniform values. Use UPDATE, SET, WHERE commands to replace non-uniform values.

# **Duplicate data**

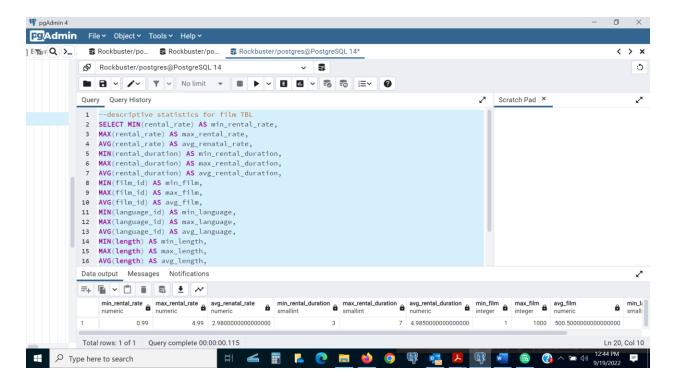
```
--looking for duplicate data in Film TBL
SELECT film id,
title,
release year,
rental_duration,
rental rate,
       COUNT(*)
FROM film
GROUP BY film_id,
title,
release_year,
rental_duration,
rental rate
HAVING COUNT(*) > 1
---- looking for duplicate values customer TBL
SELECT first name,
         last name
         email,
         address_id
         active,
         COUNT(*)
FROM customer
GROUP BY first name,
    last name,
                email,
                address id,
                active
HAVING COUNT(*) > 1
```

For Duplicate values a VIEW function can be used to create a view table to store duplicate values to be later cleaned or the DELETE command can be used to delete the duplicate values by appropriate data engineers with permissions.

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.

### --descriptive statistics for film TBL

SELECT MIN(rental\_rate) AS min\_rental\_rate, MAX(rental rate) AS max rental rate, AVG(rental rate) AS avg renatal rate, MIN(rental\_duration) AS min\_rental\_duration, MAX(rental\_duration) AS max\_rental\_duration, AVG(rental duration) AS avg rental duration, MIN(film id) AS min film, MAX(film\_id) AS max\_film, AVG(film id) AS avg film, MIN(language\_id) AS min\_language, MAX(language id) AS max language, AVG(language\_id) AS avg\_language, MIN(length) AS min length, MAX(length) AS max length, AVG(length) AS avg length, MIN(replacement\_cost) AS min\_replacement\_cost, MAX(replacement\_cost) AS max\_replacement\_cost, AVG(replacement\_cost) AS avg\_replacement\_cost FROM film



# --descriptive statistics for customer TBL

SELECT MIN(customer\_id) AS min\_customer\_id,

MAX(customer\_id) AS max\_customer\_id,

AVG(customer\_id) AS avg\_customer\_id,

MIN(store id) AS min store id,

MAX(store\_id) AS max\_store\_id,

AVG(store\_id) AS avg\_store\_id,

MIN(address\_id) AS min\_address\_id,

MAX(address id) AS max address id,

AVG(address\_id) AS avg\_address\_id,

MIN(create date) AS min create date,

MAX(create\_date) AS max\_create\_date,

MODE() WITHIN GROUP (ORDER BY create\_date) AS create\_date,

MIN(last\_update) AS min\_last\_update,

MAX(last update) AS max last update,

MODE() WITHIN GROUP (ORDER BY last update) AS last update,

MODE() WITHIN GROUP (ORDER BY first\_name) AS first\_name,

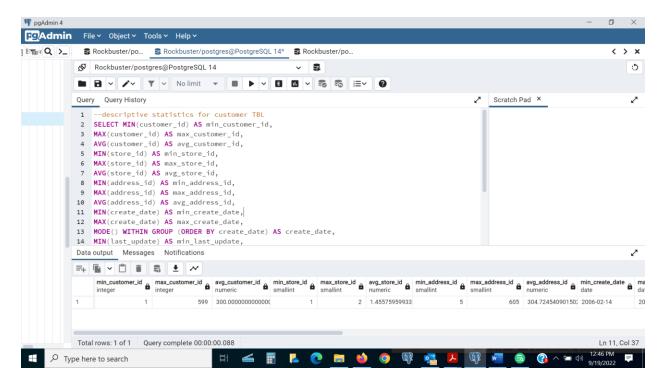
MODE() WITHIN GROUP (ORDER BY last name) AS last name,

MODE() WITHIN GROUP (ORDER BY email) AS email,

MODE() WITHIN GROUP (ORDER BY create date) AS create date,

MODE() WITHIN GROUP (ORDER BY active) AS mode active

FROM customer



- 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.
  - Based on my prior experience with data profiling in Excel and now SQL, I feel that SQL if
    more effective for data profiling. SQL works well with large amounts of data where
    Excel does not. SQL results / output faster if using correct functions/syntax, results are
    returned faster and easier to modify, update and store.