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).

- Duplicates – None found, if found duplicate data can be cleaned by creating a view with unique records, or duplicate records can be deleted.
  - o Film Table

| Query | Query History | | | | | | | | | | | | Scratch Pad X |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

```
1  SELECT film_id, title, description, release_year, language_id, rental_duration, rental_rate, length,
2  replacement_cost, rating, last_update, special_features, fulltext, COUNT(*)
3  FROM film
4  GROUP By film_id, title, description, release_year, language_id, rental_duration, rental_rate, length,
5  replacement_cost, rating, last_update, special_features, fulltext
6  HAVING COUNT(*)>1
```

Data output   Messages   Notifications

| 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) | rating mpaa_rating | last_update timestamp without time zone | special_features text[] | fulltext tsvector | count bigint |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

  - o Customer Table

Query   Query History

```
1  SELECT customer_id, store_id, first_name, last_name, email, address_id, activebool, create_date, last_update, active, COUNT(*)
2  FROM customer
3  GROUP By customer_id, store_id, first_name, last_name, email, address_id, activebool, create_date, last_update, active
4  HAVING COUNT(*)>1
```

Data output   Messages   Notifications

| 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 | last_update timestamp without time zone | active integer | count bigint |
|---|---|---|---|---|---|---|---|---|---|---|

- Non-Uniform Values – None found, if found while searching through a few random values to check for inconsistencies then the record can be verified and updated to match similar records

- Film Table

```
1  SELECT DISTINCT rating
2  FROM film
3  GROUP BY rating
```

Data output   Messages   Notifications

| | rating<br>mpaa_rating 🔒 |
|---|---|
| 1 | G |
| 2 | PG |
| 3 | PG-13 |
| 4 | R |
| 5 | NC-17 |

Query   Query History

```
1  SELECT DISTINCT rental_rate
2  FROM film
3  GROUP BY rental_rate
```

Data output   Messages   Notifications

| | rental_rate<br>numeric (4,2) 🔒 |
|---|---|
| 1 | 0.99 |
| 2 | 2.99 |
| 3 | 4.99 |

- Customer Table

Query   Query History

```
1  SELECT DISTINCT store_id
2  FROM customer
3  GROUP BY store_id
```

Data output   Messages   Notifications

| | store_id<br>smallint 🔒 |
|---|---|
| 1 | 1 |
| 2 | 2 |

Query   Query History

```
1  SELECT DISTINCT activebool
2  FROM customer
3  GROUP BY activebool
```

Data output   Messages   Notifications

| | activebool<br>boolean 🔒 |
|---|---|
| 1 | true |

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.
   - Film Table

```
 1   SELECT MIN(rental_rate) AS min_rental_rate,
 2   MAX(rental_rate) AS max_rental_rate,
 3   AVG(rental_rate) AS avg_renatal_rate,
 4   MIN(rental_duration) AS min_rental_duration,
 5   MAX(rental_duration) AS max_rental_duration,
 6   AVG(rental_duration) AS avg_rental_duration,
 7   MIN(film_id) AS min_film,
 8   MAX(film_id) AS max_film,
 9   AVG(film_id) AS avg_film,
10   MIN(language_id) AS min_language,
11   MAX(language_id) AS max_language,
12   AVG(language_id) AS avg_language,
13   MIN(length) AS min_length,
14   MAX(length) AS max_length,
15   AVG(length) AS avg_length,
16   MIN(replacement_cost) AS min_replacement_cost,
17   MAX(replacement_cost) AS max_replacement_cost,
18   AVG(replacement_cost) AS avg_replacement_cost,
19   MODE() WITHIN GROUP (ORDER BY rating) AS rating_value,
20   MODE() WITHIN GROUP (ORDER BY special_features) AS feature_value,
21   MODE() WITHIN GROUP (ORDER BY release_year) AS release_year,
22   MODE() WITHIN GROUP (ORDER BY title) AS title_value,
23   MODE() WITHIN GROUP (ORDER BY fulltext) AS fulltext
24   FROM film
```

Data output    Messages    Notifications

| | min_rental_rate numeric | max_rental_rate numeric | avg_renatal_rate numeric | min_rental_duration smallint | max_rental_duration smallint | avg_r nume |
|---|---|---|---|---|---|---|
| 1 | 0.99 | 4.99 | 2.98000000000000 | 3 | 7 | 4.985 |

- Customer Table

Query    Query History

```
 1   SELECT MIN(active) AS min_active,
 2   MAX(active) AS max_active,
 3   AVG(active) AS avg_active,
 4   MIN(address_id) AS min_address,
 5   MAX(address_id) AS max_address,
 6   AVG(address_id) AS avg_address,
 7   MIN(customer_id) AS min_customer,
 8   MAX(customer_id) AS max_customer,
 9   AVG(customer_id) AS avg_customer,
10   MIN(store_id) AS min_store,
11   MAX(store_id) AS max_store,
12   AVG(store_id) AS avg_store,
13   MODE() WITHIN GROUP (ORDER BY last_update) AS last_update,
14   MODE() WITHIN GROUP (ORDER BY first_name) AS first_name,
15   MODE() WITHIN GROUP (ORDER BY last_name) AS last_name,
16   MODE() WITHIN GROUP (ORDER BY email) AS email,
17   MODE() WITHIN GROUP (ORDER BY create_date) AS create_date,
18   MODE() WITHIN GROUP (ORDER BY active) AS mode_active
19   FROM customer;
```

Data output    Messages    Notifications

| min_active integer | max_active integer | avg_active numeric | min_address smallint | max_address smallint | avg_address numeric | min_customer integer | max_customer integer | avg_ num |
|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0.9749582637 | 5 | 605 | 304.724540901 | 1 | 599 | 300 |

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 more effective than Excel at data profiling when the data set is very large, and when the data is housed in a shared storage method. With smaller data sets, Excel and SQL are both proficient in data profiling but sharing and collaboration would be hindered with excel. SQL's language allows analysts to efficiently process, clean, and analyze data in a streamlined wat making it overall more effective than Excel.

4. Save your "Answers 3.6" document as a PDF and upload it here for your tutor to review.