# *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).

*   Manual check one by one for non-uniform

```
45    --/*
46 ∨  SELECT
47    --COUNT(rating),
48    --rating
49    --COUNT(title),
50    --title
51    --COUNT(description),
52    --description
53    --COUNT(release_year),
54    --release_year
55    --COUNT(language_id),
56    --language_id
57    --COUNT(rental_duration),
58    --rental_duration
59    --COUNT(rental_rate),
60    --rental_rate
61    --COUNT(length),
62    --length
63    --COUNT(replacement_cost),
64    --replacement_cost
65    --COUNT(last_update),
66    --last_update
67    --COUNT(special_features),
68    --special_features
69    --COUNT(fulltext),
70    --fulltext
71    --FROM film
72    --COUNT(store_id),
73    --store_id
74    --COUNT(first_name),
75    --first_name
76    --COUNT(last_name),
77    --last_name
78    --COUNT(email),
79    --email
80    --COUNT(address_id),
81    --address_id
82    --COUNT(activebool),
83    --activebool
84    --COUNT(create_date),
85    --create_date
86    --COUNT(last_update),
87    --last_update
88    COUNT(active),
89    active
90    FROM customer
91    GROUP BY 2
92    --HAVING COUNT(*) > 1
93    ORDER BY COUNT(*) DESC
94    --*/
```

Data Output   Messages   Notifications

| | count bigint 🔒 | active integer 🔒 |
|---|---|---|
| 1 | 584 | 1 |
| 2 | 15 | 0 |

- duplicate data, or missing values

```
1  /*
2  SELECT *
3  FROM film
4  WHERE title IS NULL
5  OR description IS NULL
6  OR release_year IS NULL
7  OR language_id IS NULL
8  */
9  /*
10  SELECT *
11  FROM customer
12  WHERE first_name IS NULL
13  OR last_name IS NULL
14  OR email IS NULL
15  OR active IS NULL
16  */
17  /*
18  SELECT
19  COUNT(*),
20  c.customer_id,
21  c.store_id,
22  c.first_name,
23  c.last_name,
24  c.email,
25  c.address_id,
26  c.activebool,
27  c.create_date,
28  c.last_update,
29  c.active
30  FROM customer c
31  GROUP BY 2,3,4,5,6,7,8,9,10,11
32  HAVING COUNT(*) >1
33  */
34  /*
35  SELECT
36  COUNT(*),
37  title,
38  description,
39  release_year,
40  language_id
41  FROM film
42  GROUP BY title, description, release_year, language_id
43  HAVING COUNT(*) > 1
44  */
```

- Query for the potential fix:

```
 96    UPDATE film
 97    SET title = 'Unknown'
 98    WHERE title IS NULL
 99    */
100  ∨ /*
101    UPDATE film
102    SET release_year = 2000 -- or any default year
103    WHERE release_year IS NULL
104    */
105  ∨ /*
106    DELETE FROM film
107    WHERE language_id IS NULL; -- If no valid language_id exists
108    */
109  ∨ /*
110    UPDATE customer
111    SET email = 'unknown@example.com'
112    WHERE email IS NULL
113    */
114  ∨ /*
115    UPDATE customer
116    SET first_name = 'Anonymous', last_name = 'Anonymous'
117    WHERE first_name IS NULL OR last_name IS NULL;
118    */
119  ∨ /*
120    UPDATE customer
121    SET active = 0
122    WHERE active IS NULL
123    */
124  ∨ /*
125    DELETE FROM customer
126    WHERE customer_id NOT IN (
127        SELECT MIN(customer_id)
128        FROM customer
129        GROUP BY store_id, first_name, last_name, email, address_id)
130    */
131  ∨ /*
132    DELETE FROM film
133    WHERE film_id NOT IN (
134        SELECT MIN(film_id)
135        FROM film
136        GROUP BY title, description, release_year, language_id)
137    */
```

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

```
138    --/*
139  ∨ SELECT
140    COUNT(*) AS count_rows,
141    MIN(release_year) AS min_release_year,
142    MAX(release_year) AS max_release_year,
143    ROUND(AVG(release_year)) AS avg_release_year,
144    COUNT(release_year) AS count_release_year_values,
145    MIN(rental_duration) AS min_rental_duration,
146    MAX(rental_duration) AS max_rental_duration,
147    ROUND(AVG(rental_duration)) AS avg_rental_duration,
148    COUNT(rental_duration) AS count_rental_duration_values,
149    MIN(length) AS min_length,
150    MAX(length) AS max_length,
151    ROUND(AVG(length)) AS avg_length,
152    COUNT(length) AS count_length_values,
153    MIN(replacement_cost) AS min_replacement_cost,
154    MAX(replacement_cost) AS max_replacement_cost,
155    ROUND(AVG(replacement_cost)) AS avg_replacement_cost,
156    COUNT(replacement_cost) AS count_replacement_cost_values,
157    MIN(rental_rate) AS min_rent,
158    MAX(rental_rate) AS max_rent,
159    ROUND(AVG(rental_rate)) AS avg_rent,
160    COUNT(rental_rate) AS count_rent_values,
161    MODE() WITHIN GROUP (ORDER BY rating) AS modal_rating_value,
162    MODE() WITHIN GROUP (ORDER BY special_features) AS modal_special_features_value,
163    MODE() WITHIN GROUP (ORDER BY fulltext) AS modal_fulltext_value
164    FROM film
165    --*/
```

| count_rows bigint | min_release_year integer | max_release_year integer | avg_release_year numeric | count_release_year_values bigint | min_rental_duration smallint | max_rental_duration smallint | avg_rental_duration numeric | count_rental_duration_values bigint | min_len smallin |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1004 | 1987 | 2006 | 2006 | 1004 | 3 | 7 | 5 | 1004 |

- Customer table

```
166    --/*
167  ∨ SELECT
168    COUNT(*) AS count_rows,
169    MODE() WITHIN GROUP (ORDER BY store_id) AS modal_store_id_value,
170    MODE() WITHIN GROUP (ORDER BY first_name) AS modal_first_name_value,
171    MODE() WITHIN GROUP (ORDER BY last_name) AS modal_last_name_value,
172    MODE() WITHIN GROUP (ORDER BY email) AS modal_email_value,
173    MODE() WITHIN GROUP (ORDER BY activebool) AS modal_activebool_value,
174    MODE() WITHIN GROUP (ORDER BY create_date) AS modal_create_date_value,
175    MODE() WITHIN GROUP (ORDER BY last_update) AS modal_last_update_value
176    FROM customer
177    --*/
```

| count_rows bigint | modal_store_id_value smallint | modal_first_name_value character varying | modal_last_name_value character varying | modal_email_value character varying | modal_activebool_value boolean | modal_create_date_value date | modal_last_update_value timestamp without time zone |
|---|---|---|---|---|---|---|---|
| 1 | 599 | 1 | Jamie | Abney | aaron.selby@sakilacustomer.org | true | 2006-02-14 | 2013-05-26 14:49:45.738 |

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.

- **Reflection: Excel vs. SQL for Data Profiling**

From my experience, both Excel and SQL are valuable tools for data profiling, but they serve different purposes depending on the dataset size and the complexity of the task. Excel is more intuitive and user-friendly, especially for smaller datasets or visualizing

trends and patterns quickly through features like pivot tables and conditional formatting. However, it becomes cumbersome and slow for large datasets or complex operations.

SQL, on the other hand, is far more efficient for handling large-scale data and complex queries. It allows for precise and automated operations like filtering, grouping, and joining multiple tables with ease. The ability to write repeatable queries in SQL ensures consistency and reduces the chances of manual errors. While SQL requires a steeper learning curve compared to Excel, it excels in speed and scalability, making it the better choice for large datasets or data stored in relational databases.

In conclusion, for quick, small-scale profiling tasks, Excel might be more practical, but for robust, large-scale analysis, SQL is the more effective and reliable tool.