Katie Goyal

Data Analytics Immersion 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).

## Duplicate data:

### Film Table
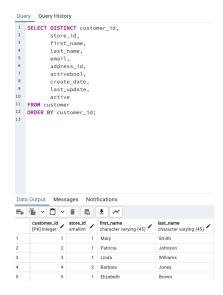


### Customer Table



- *There were no duplicates found when running the query. If duplicates where found one would just need to run a query to delete the duplicates*

## Non_uniform:

### Film Table



### Customer Table



- *There were no non_uniform data in the tables. If there were to be non_uniform data one would the UPDATE Statement instead of select to filter and uniform the data.*

## Missing:

### Film Table



### Customer Table



- *There was no missing data found in the tables above. If there was missing data I would filter clean and update the data to fill in the missing values.*

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.

### Numerical

```
Query   Query History
1   SELECT
2       MIN(release_year) AS min_release_year,
3       MIN(rental_duration) AS min_rentdur,
4       MIN(rental_rate) AS min_rate,
5       MIN(length) AS min_leng,
6       MIN(replacement_cost) AS min_replac,
7       MAX(release_year) AS max_release_year,
8       MAX(rental_duration) AS max_rentdur,
9       MAX(rental_rate) AS max_rate,
10      MAX(length) AS max_leng,
11      MAX(replacement_cost) AS max_replac,
12      AVG(release_year) AS avg_release_year,
13      AVG(rental_duration) AS avg_rentdur,
14      AVG(rental_rate) AS avg_rate,
15      AVG(length) AS avg_leng,
16      AVG(replacement_cost) AS avg_replac
17  FROM film;
18
```

Data Output   Messages   Notifications

| | min_release_year<br>integer | min_rentdur<br>smallint | min_rate<br>numeric | min_leng<br>smallint |
|---|---|---|---|---|
| 1 | 2006 | 3 | 0.99 | 46 |

- *There is no numerical data that would help in understanding this data. Knowing the number of customers who rented would have been more beneficial*

### Nonnumerical

```
Query   Query History
1  SELECT MODE() WITHIN GROUP (ORDER BY language_id) AS most_lang,
2         MODE() WITHIN GROUP (ORDER BY rating) AS most_rat
3  FROM film;
```

Data Output   Messages   Notifications

| | most_leng<br>smallint | most_rat<br>mpaa_rating |
|---|---|---|
| 1 | 1 | PG-13 |

- *When looking at this data we are given a numerical number for most languages rented and not a language which does not help us in*

*understanding what was the most rented movie by language. We can understand that the most rented movies are PG 13.*

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.

   - *When working with data in SQL versus Excel we can control, manage, edit, and manipulate the data much faster and easier than if we were to do it in Excel. The amount of time it takes to write and run a query I would have only finished one column in Excel.*