3.6 Answer Sheet
Katherine Lecce

**1.**
<mark>**Non uniformed Data**</mark>
**→ Film Table:**

```
Query    Query History

1    SELECT DISTINCT rating,
2    description,
3    release_year,
4    language_id,
5    rental_duration,
6    rental_rate,length,
7    replacement_cost,
8    rating,
9    last_update,
10   special_features,
11   fulltext
12   FROM film;
```

Data Output    Messages    Notifications

| | rating mpaa_rating 🔒 | description text |
|---|---|---|
| 1 | G | A Action-Packed Character Study of a Dog And a Lumberjack who |
| 2 | G | A Action-Packed Display of a Mad Cow And a Astronaut who must |
| 3 | G | A Action-Packed Panorama of a Husband And a Feminist who mus |
| 4 | G | A Action-Packed Story of a Pioneer And a Technical Writer who mu |
| 5 | G | A Action-Packed Yarn of a Boat And a Crocodile who must Build a |
| 6 | G | A Amazing Character Study of a Robot And a Student who must Ch |
| 7 | G | A Amazing Display of a Mad Cow And a Pioneer who must Redeem |

**→ Customer Table:**

```
Query    Query History

1    SELECT DISTINCT customer_id,
2                store_id,
3                first_name,
4                last_name,
5                email,
6                address_id,
7                activebool,
8                create_date,
9                last_update,
10               active
11   From Customer
```

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) |
|---|---|---|---|---|---|
| 1 | 357 | 1 | Keith | Rico | keith.rico@sakilacustomer.org |
| 2 | 171 | 2 | Dolores | Wagner | dolores.wagner@sakilacustomer.org |
| 3 | 139 | 1 | Amber | Dixon | amber.dixon@sakilacustomer.org |
| 4 | 471 | 1 | Dean | Sauer | dean.sauer@sakilacustomer.org |
| 5 | 594 | 1 | Eduardo | Hiatt | eduardo.hiatt@sakilacustomer.org |
| 6 | 401 | 2 | Tony | Carranza | tony.carranza@sakilacustomer.org |
| 7 | 157 | 2 | Darlene | Rose | darlene.rose@sakilacustomer.org |

## Missing Data:

**Film Data →**



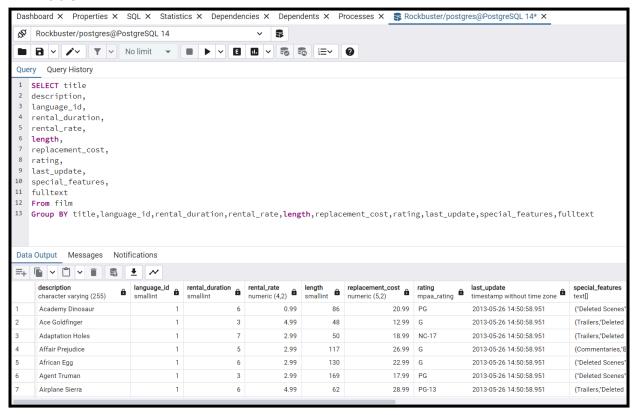**Customer Table →**
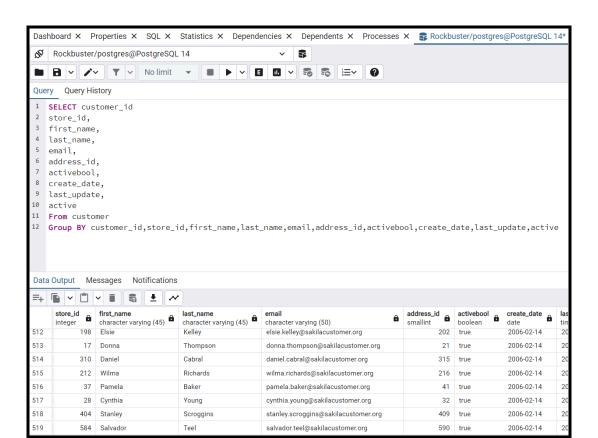
## Duplicate Data:

### Film Table:



### Customer Table:

## 2. Summarize your data:

### Film Table: Numeric

```
Dashboard ×   Properties ×   SQL ×   Statistics ×   Dependencies ×   Dependents ×   Processes ×   🗄 Rockbuster/postgres@PostgreSQL 14* ×

Rockbuster/postgres@PostgreSQL 14                    ▼   🗄

📁  💾 ▼   🖉 ▼   ▼ ▼   No limit  ▼   ■ ▶ ▼   E �𝄜 ▼   🗄 🗄   ☰ ▼   ❓

Query   Query History

1   SELECT MIN (release_year) AS min_release_year,
2   MIN (rental_duration) AS min_rentdur,
3   MIN (rental_rate) AS min_rate,
4   Min (length) AS min_leng,
5   MIN (replacement_cost) AS min_replac,
6   MAX (rental_duration) AS max_rate,
7   MAX (length) AS max_leng,
8   MAX (replacement_cost) AS max_replac,
9   AVG (release_year) AS avg_release_year,
10  AVG (rental_duration) AS avg_rentdur,
11  AVG (rental_rate) AS avg_rate,
12  AVG (length) AS avg_leng,
13  AVG (replacement_cost) AS avg_replac
14  FROM film
```
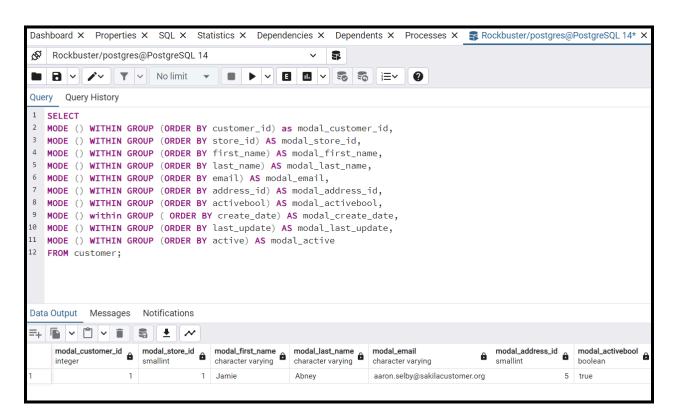
Data Output   Messages   Notifications

| | min_release_year<br>integer | min_rentdur<br>smallint | min_rate<br>numeric | min_leng<br>smallint | min_replac<br>numeric | max_rate<br>smallint | max_leng<br>smallint | max_replac<br>numeric | avg_release_year<br>numeric | avg_rentdur<br>numeric | avg_rate<br>numeric |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2006 | 3 | 0.99 | 46 | 9.99 | 7 | 185 | 29.99 | 2006.0000000000000000 | 4.9850000000000000 | 2.9800000000000 |

### Film Table - Non Numeric
### <mark>Movie Rating</mark>

```
Dashboard ×   Properties ×   SQL ×   Statistics ×   Dependencies ×   Dependents ×   Processes ×

Rockbuster/postgres@PostgreSQL 14                    ▼   🗄

📁  💾 ▼   🖉 ▼   ▼ ▼   No limit  ▼   ■ ▶ ▼   E ⬛ ▼   🗄 🗄   ☰ ▼   ❓

Query   Query History

1   SELECT rating AS rating_lang
2   FROM (
3   SELECT rating, COUNT(*) AS rating_count
4   FROM film
5   GROUP BY rating
6   ORDER BY rating DESC
7   LIMIT 1
8   ) AS mode_subquery;
```

Data Output   Messages   Notifications

| | rating_lang<br>mpaa_rating |
|---|---|
| 1 | NC-17 |

**Language**

```
Rockbuster/postgres@PostgreSQL 14

1   SELECT language_id AS most_lang
2   FROM (
3   SELECT language_id, COUNT(*) AS lang_count
4   FROM film
5   GROUP BY language_id
6   ORDER BY lang_count DESC
7   LIMIT 1
8   ) AS mode_subquery;
```

Data Output   Messages   Notifications

| most_lang 🔒 smallint |
| --- |
| 1 | 1 |

**Customer Table:**

Dashboard ✕  Properties ✕  SQL ✕  Statistics ✕  Dependencies ✕  Dependents ✕  Processes ✕  Rockbuster/postgres@PostgreSQL 14* ✕

```
Rockbuster/postgres@PostgreSQL 14

1    SELECT
2    MODE () WITHIN GROUP (ORDER BY customer_id) as modal_customer_id,
3    MODE () WITHIN GROUP (ORDER BY store_id) AS modal_store_id,
4    MODE () WITHIN GROUP (ORDER BY first_name) AS modal_first_name,
5    MODE () WITHIN GROUP (ORDER BY last_name) AS modal_last_name,
6    MODE () WITHIN GROUP (ORDER BY email) AS modal_email,
7    MODE () WITHIN GROUP (ORDER BY address_id) AS modal_address_id,
8    MODE () WITHIN GROUP (ORDER BY activebool) AS modal_activebool,
9    MODE () within GROUP ( ORDER BY create_date) AS modal_create_date,
10   MODE () WITHIN GROUP (ORDER BY last_update) AS modal_last_update,
11   MODE () WITHIN GROUP (ORDER BY active) AS modal_active
12   FROM customer;
```

Data Output   Messages   Notifications

| modal_customer_id 🔒 integer | modal_store_id 🔒 smallint | modal_first_name 🔒 character varying | modal_last_name 🔒 character varying | modal_email 🔒 character varying | modal_address_id 🔒 smallint | modal_activebool 🔒 boolean |
| --- | --- | --- | --- | --- | --- | --- |
| 1 | 1 | Jamie | Abney | aaron.selby@sakilacustomer.org | 5 | true |

1. **How would you clean each set of data?**

A. **Missing Data:** Save Data, then you'll want to use the 'SELECT' statement to omit any columns you'd like to skip over when checking this data (make note in query). You can also use the 'Impute' or fill in where there are missing sections of the data.

B. **Duplicate Data:** First and foremost, backup your data before starting to edit the query. Then use the 'Group By' function to identify the records that you want to check in a specific column of the data. Then, if duplicates are found, use the 'Delete' statement to remove any information in the set that is invalid or 'Merge' multiple statements into one single record by using the 'update' function to join the multiple sections you'd like to combine.

C. **Non-Uniformed Data:** After saving your data, use the GROUP By and DISTINCT functions to check a few areas in your data to see if your outputs have any inconsistencies. If you you find areas that need to be uniformed use the

**UPDATE** (name of data set)
**SET** (section you'd like the format to change FROM) **=** ( section name you'd like the area to be formatted TO)
**WHERE** (column) **IN** ( section you'd like the format to change FROM)


**3 - 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.**

Because both systems have different functions it's important to first consider how large the data sets you're looking at are. If this is a task that only one or two people are looking at altering any type of data within the spreadsheet. While SQL is more practical for large amounts of data across a range of multiple operations and avenues within a business, company, or program. SQL is better for tasks that require the data to alter and change into migrating into one system format rather than keeping it in a spreadsheet format. Personally, because I am more familiar with spreadsheets, it is more comfortable and easier for me to profile data when Excel and more familiar with using functions such as filter, pivot tables, and highlighting to more clearly see where information might need to be altered.