1. Check for and clean dirty data:
   - Check for duplicates
     - Film Table

```
Query Editor

1   --Shows only those records that duplicate (based on columns selected)
2
3   SELECT film_id,
4          title,
5          description,
6          release_year,
7          language_id,
8          rental_duration,
9          rental_rate,
10         length,
11         replacement_cost,
12         rating,
13         last_update,
14         special_features,
15         fulltext,
16         COUNT(*)
17  FROM film
18  GROUP BY film_id,
19         title,
20         description,
21         release_year,
22         language_id,
23         rental_duration,
24         rental_rate,
25         length,
26         replacement_cost,
27         rating,
28         last_update,
29         special_features,
30         fulltext
31  HAVING COUNT(*) >1;
32  --no result set means we have no duplicates
```

     - Customer Table

```
Query Editor

1   --Shows only those records that duplicate (based on columns selected)
2
3   SELECT customer_id,
4          store_id,
5          first_name,
6          last_name,
7          email,
8          address_id,
9          activebool,
10         create_date,
11         last_update,
12         active,
13         COUNT(*)
14  FROM customer
15  GROUP BY customer_id,
16         store_id,
17         first_name,
18         last_name,
19         email,
20         address_id,
21         activebool,
22         create_date,
23         last_update,
24         active
25  HAVING COUNT(*) >1;
26  --no result set means we have no duplicates
```

- o  If duplicate data did exist, I would create a new view of the table to show only unique records. Then I would use that view to continue my analysis. If there were a lot of duplicates, I would also notify the data engineer so they could make sure data migration and everything on their end is correct.
- Check for non-uniform data
  - o  Film Table Queries
    - ▪ SELECT DISTINCT title
      FROM film
    - ▪ SELECT DISTINCT release_year
      FROM film
    - ▪ SELECT DISTINCT language_id
      FROM film
    - ▪ SELECT DISTINCT rental_rate
      FROM film
    - ▪ SELECT DISTINCT replacement_cost
      FROM film
    - ▪ SELECT DISTINCT rating
      FROM film
  - o  Customer Table Queries
    - ▪ SELECT DISTINCT email
      FROM customer
    - ▪ SELECT DISTINCT address_id
      FROM customer
    - ▪ SELECT DISTINCT active
      FROM customer
  - o  If there was inconsistent data, I would update it using the UPDATE and SET commands.
- Check for missing data
  - o  Film Table

```
Query Editor

1   SELECT film_id,
2          title,
3          description,
4          release_year,
5          language_id,
6          rental_duration,
7          rental_rate,
8          length,
9          replacement_cost,
10         rating,
11         last_update,
12         special_features,
13         fulltext
14  FROM film
15  WHERE film_id IS NULL
16         OR title IS NULL
17         OR description IS NULL
18         OR release_year IS NULL
19         OR language_id IS NULL
20         OR rental_duration IS NULL
21         OR rental_rate IS NULL
22         OR length IS NULL
23         OR replacement_cost IS NULL
24         OR rating IS NULL
25         OR last_update IS NULL
26         OR special_features IS NULL
27         OR fulltext IS NULL
```

- o Customer Table

```sql
Query Editor

1   SELECT 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
12  WHERE customer_id IS null
13        OR store_id IS null
14        OR first_name IS null
15        OR last_name IS null
16        OR email IS null
17        OR address_id IS null
18        OR activebool IS null
19        OR create_date IS null
20        OR last_update IS null
21        OR active IS null
```

- o If there was missing data in a column that I don't need for analysis, I would ignore the column entirely. If a small amount of data were missing from a column with numeric values that could be averaged, I would impute the average or mean in place of the NULL values and note it in my analysis.

2. Summarize data:
   - Film Table
     - o Rental Duration

```sql
Query Editor   Query History

1   SELECT MIN(rental_duration) AS min_rental_duration,
2          MAX(rental_duration) AS max_rental_duration,
3          AVG(rental_duration) AS avg_rental_duration,
4          COUNT(rental_duration) AS count_rental_duration,
5          mode() WITHIN GROUP (ORDER BY rental_duration)
6              AS mode_rental_duration,
7          COUNT(*) AS count_rows
8   FROM film
```

Data Output   Explain   Messages   Notifications

| min_rental_duration smallint | max_rental_duration smallint | avg_rental_duration numeric | count_rental_duration bigint | mode_rental_duration smallint | count_rows bigint |
|---|---|---|---|---|---|
| 1 | 3 | 7 | 4.9850000000000000 | 1000 | 6 | 1000 |

     - o Rental Rate

```sql
Query Editor   Query History

1   SELECT MIN(rental_rate) AS min_rental_rate,
2          MAX(rental_rate) AS max_rental_rate,
3          AVG(rental_rate) AS avg_rental_rate,
4          COUNT(rental_rate) AS count_rental_rate,
5          mode() WITHIN GROUP (ORDER BY rental_rate)
6              AS mode_rental_rate,
7          COUNT(*) AS count_rows
8   FROM film
```

Data Output   Explain   Messages   Notifications

| min_rental_rate numeric | max_rental_rate numeric | avg_rental_rate numeric | count_rental_rate bigint | mode_rental_rate numeric | count_rows bigint |
|---|---|---|---|---|---|
| 1 | 0.99 | 4.99 | 2.9800000000000000 | 1000 | 0.99 | 1000 |

- Length

```sql
1  SELECT MIN(length) AS min_film_length,
2         MAX(length) AS max_film_length,
3         AVG(length) AS avg_film_length,
4         COUNT(length) AS count_film_length,
5         mode() WITHIN GROUP (ORDER BY length)
6             AS mode_film_length,
7         COUNT(*) AS count_rows
8  FROM film
```

Data Output | Explain | Messages | Notifications

| min_film_length smallint | max_film_length smallint | avg_film_length numeric | count_film_length bigint | mode_film_length smallint | count_rows bigint |
|---|---|---|---|---|---|
| 1 | 46 | 185 | 15.2720000000000000 | 1000 | 85 | 1000 |

- Replacement Cost

```sql
1  SELECT MIN(replacement_cost) AS min_film_replacement_cost,
2         MAX(replacement_cost) AS max_replacement_cost,
3         AVG(replacement_cost) AS avg_replacement_cost,
4         COUNT(replacement_cost) AS count_replacement_cost,
5         mode() WITHIN GROUP (ORDER BY replacement_cost)
6             AS mode_replacement_cost,
7         COUNT(*) AS count_rows
8  FROM film
```

Data Output | Explain | Messages | Notifications

| min_film_replacement_cost numeric | max_replacement_cost numeric | avg_replacement_cost numeric | count_replacement_cost bigint | mode_replacement_cost numeric | count_rows bigint |
|---|---|---|---|---|---|
| 1 | 9.99 | 29.99 | 19.9840000000000000 | 1000 | 20.99 | 1000 |

- Rating

```sql
1  SELECT COUNT(rating) AS count_rating,
2         mode() WITHIN GROUP (ORDER BY rating)
3             AS mode_rating,
4         COUNT(*) AS count_rows
5  FROM film
```

Data Output | Explain | Messages | Notifications

| count_rating bigint | mode_rating mpaa_rating | count_rows bigint |
|---|---|---|
| 1 | 1000 | PG-13 | 1000 |

- Customer Table
  - Store ID

```sql
1  SELECT COUNT(store_id) AS count_store_id,
2         mode() WITHIN GROUP (ORDER BY store_id)
3             AS mode_store_id,
4         COUNT(*) AS count_rows
5  FROM customer
```
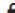
Data Output | Explain | Messages

| count_store_id bigint | mode_store_id smallint | count_rows bigint |
|---|---|---|
| 1 | 599 | 1 | 599 |

  - Active

```sql
1  SELECT COUNT(active) AS count_active,
2         mode() WITHIN GROUP (ORDER BY active)
3             AS mode_active,
4         COUNT(*) AS count_rows
5  FROM customer
```

Data Output | Explain | Messages | Notifications

| count_active bigint | mode_active integer | count_rows bigint |
|---|---|---|
| 1 | 599 | 1 | 599 |

3. SQL is more efficient than Excel for data profiling because SQL performs the summary calculations in one query and only returns the information requested. Excel in comparison, requires you to have an entire data set in hand, then perform additional steps and calculations to get data summary information.