

## Achievement 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).

Non-uniform

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

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

Using DISTINCT to check if any data are entered differently. If so, for example,








```
1 UPDATE table
2 SET col1 = 'G'
3 WHERE col1 IN ('gen',
4               'g',
5               'General')
```

The above query will group all “gen” “g” and “General” and change them into “G” in col1.

## Duplicate Data

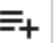






```
1 SELECT title,  
2         release_year,  
3         language_id,  
4         rental_duration,  
5         COUNT(*)  
6 FROM film  
7 GROUP BY title,  
8         release_year,  
9         language_id,  
10        rental_duration  
11 HAVING COUNT(*) >1
```

Data output Messages Notifications

						
title	release_year	language_id	rental_duration	count		
character varying (255)	integer	smallint	smallint	bigint		

```
1 SELECT customer_id,  
2         store_id,  
3         first_name,  
4         last_name,  
5         COUNT(*)  
6 FROM customer  
7 GROUP BY customer_id,  
8         store_id,  
9         first_name,  
10        last_name  
11 HAVING COUNT(*) >1
```

Data output Messages Notifications

						
customer_id	store_id	first_name	last_name	count		
[PK] integer	smallint	character varying (45)	character varying (45)	bigint		

If there are duplicated data, we can create a unique view table by choosing the right column ( CREATE VIEW viewname AS) or we can delete the data with query.

## Missing Values

- If there are many missing values in a column, we can just omit the column, and if there are a few missing data we can impute values such as using the AVG query.

```
1 UPDATE tablename
2 SET = AVG(col1)
3 WHERE col1 IS NULL
```

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.

```
1 SELECT MIN(rental_duration) AS min_duration,
2        MAX(rental_duration) AS max_duration,
3        AVG(rental_duration) AS avg_duration,
4        MIN(rental_rate) AS min_rent,
5        MAX(rental_rate) AS max_rent,
6        AVG(rental_rate) AS avg_rent,
7        MIN(length) AS min_length,
8        MAX(length) AS max_length,
9        AVG(length) AS avg_length,
10       MIN(replacement_cost) AS min_replacement_cost,
11       MAX(replacement_cost) AS max_replacement_cost,
12       AVG(replacement_cost) AS avg_replacement_cost
13 FROM film
```

Data output Messages Notifications

	min_duration smallint	max_duration smallint	avg_duration numeric	min_rent numeric	max_rent numeric	avg_rent numeric	min_length smallint	max_length smallint	avg_length numeric	min_replacement_cost numeric	max_replacement_cost numeric	avg_replacement_cost numeric
1	3	7	4.985	0.99	4.99	2.98	46	185	115.272	9.99	29.99	19.984

```
1 SELECT mode() WITHIN GROUP (ORDER BY title) AS modal_title,
2        mode() WITHIN GROUP (ORDER BY description) AS modal_description,
3        mode() WITHIN GROUP (ORDER BY release_year) AS modal_release_year,
4        mode() WITHIN GROUP (ORDER BY language_id) AS modal_language_id,
5        mode() WITHIN GROUP (ORDER BY rating) AS modal_rating
6 FROM film
```

Data output Messages Notifications

	modal_title character varying	modal_description text	modal_release_year integer	modal_language_id smallint	modal_rating mpaa_rating
1	Academy Dinosaur	A Action-Packed C...	2006	1	PG-13

```

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 FROM customer
14

```

	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_customer numeric	min_store smallint	max_store smallint	avg_store numeric
1	0	1	0.974958263	5	605	304.7245409	1	599	300	1	2	1.455759599

```

1 SELECT mode() WITHIN GROUP (ORDER BY first_name) AS modal_first_name,
2        mode() WITHIN GROUP (ORDER BY last_name) AS modal_last_name,
3        mode() WITHIN GROUP (ORDER BY email) AS modal_email
4 FROM customer

```

	modal_first_name character varying	modal_last_name character varying	modal_email character varying
1	Jamie	Abney	aaron.selby@sak...

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.

- It's more effective to use SQL for data profiling and getting information because SQL is able to collect all data from the database whereas in Excel you need to know what data to include first. Also, there can be many mistakes in Excel such as deleting or modifying data by accidents, whereas SQL will just show you the query is wrong, and you just need to adjust the query. Both programs need to have good understanding of the software to perform well. Overall, SQL is a better tool for company with lots of data and constantly pulling data from the database without making much mistakes.