

### 3.6: Summarizing & Cleaning Data in SQL

Prepared by: Rob Rowland

#### Film Table Duplicate Data:

Query Editor

Query History

```
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

Explain

Messages

Notifications

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

#### Customer Table Duplicate Data:

Query Editor

Query History

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

Data Output

Explain

Messages

Notifications

customer_id	store_id	first_name	last_name	email	address_id	count
[PK] integer	smallint	character varying (45)	character varying (45)	character varying (50)	smallint	bigint

If either of these tables had duplicate data, I could either a) create a view with unique records using GROUP BY or DISTINCT or b) delete the duplicate record using DELETE. It would likely be most wise to create the view with unique records in case any other users needed to access the duplicate rows for any reason. Permanently deleting a record should be done with care.

Film Table, Numerical Columns (columns continue to right, not seen on picture):

Query Editor

Query History

1

SELECT AVG(release\_year) AS average\_release\_year,

2

MIN(release\_year) AS minimum\_release\_year,

3

MAX(release\_year) AS maximum\_release\_year,

4

AVG(rental\_duration) AS average\_rental\_duration,

5

MIN(rental\_duration) AS minimum\_rental\_duration,

6

MAX(rental\_duration) AS maximum\_rental\_duration,

7

AVG(rental\_rate) AS average\_rental\_rate,

8

MIN(rental\_rate) AS minimum\_rental\_rate,

9

MAX(rental\_rate) AS maximum\_rental\_rate,

10

AVG(length) AS average\_length,

11

MIN(length) AS minimum\_length,

12

MAX(length) AS maximum\_length,

13

AVG(replacement\_cost) AS average\_replacement\_cost,

14

MIN(replacement\_cost) AS minimum\_replacement\_cost,

15

MAX(replacement\_cost) AS maximum\_replacement\_cost

16

FROM film

Data Output

Explain

Messages

Notifications

	average_release_year numeric	minimum_release_year integer	maximum_release_year integer	average_rental_duration numeric	minimum_rental_duration smallint	maximum_rental_duration smallint	average_rental_rate numeric	minimum_rental_rate numeric	maximum_rental_rate numeric
1	2006	2006	2006	4.985	3	7	2.98	0.99	4.99

Film Table, Non-Numerical Columns:

1

SELECT mode() WITHIN GROUP (ORDER BY rating) AS modal\_rating,

2

mode() WITHIN GROUP (ORDER BY title) AS modal\_title,

3

mode() WITHIN GROUP (ORDER BY language\_id) AS modal\_language

4

FROM film;

Data Output

Explain

Messages

Notifications

	modal_rating mpaa_rating	modal_title character varying	modal_language smallint
1	PG-13	Academy Dinosaur	1

**Customer Table, Numerical and Non=Numerical Columns (columns continue to right, not seen on picture):**

Query Editor

Query History

```
1 SELECT MIN(customer_id) AS min_customer_id,
2 MAX(customer_id) AS max_customer_id,
3 AVG(customer_id) AS avg_customer_id,
4 MIN(address_id) AS min_address_id,
5 MAX(address_id) AS max_address_id,
6 AVG(address_id) AS avg_address_id,
7 MIN(store_id) AS min_store_id,
8 MAX(store_id) AS max_store_id,
9 AVG(store_id) AS avg_store_id,
10 MODE() WITHIN GROUP (ORDER BY first_name) AS modal_first_name,
11 MODE() WITHIN GROUP (ORDER BY last_name) AS modal_last_name,
12 MODE() WITHIN GROUP (ORDER BY email) AS modal_email
13 FROM customer;
```

Data Output

Explain

Messages

Notifications

	min_customer_id integer	max_customer_id integer	avg_customer_id numeric	min_address_id smallint	max_address_id smallint	avg_address_id numeric
1	1	599	300.0000000000000000	5	605	304.72454

I remain more familiar with Excel currently, so I can conduct the data profiling much quicker that way. As data sets get larger, I can certainly see the merit for using SQL. I think that, as I learn and become fluent in SQL, it will be the preferred option for quickly referencing data points from large data sets. However, when trying to build a larger profile involving multiple columns, Excel seems easier to use at this time.