

## 1. Film Table

Check for duplicates:

```
SELECT title, film_id,  
       COUNT(*)  
FROM film  
GROUP BY title, film_id  
HAVING COUNT (*)>1
```

To clean the data, I would use CREATE VIEW -AS. Then I could use DELETE to get rid of any duplicates.

Check for unique data

```
SELECT DISTINCT *  
FROM film
```

To clean this data, I would CREATE VIEW first. Then I would use UPDATE to change any misnamed or data that is not unique. To fix any missing values I would use SELECT and possibly exclude any columns with a few missing values. I can also UPDATE the table by imputing missing values using an average or other aggregate.

## Customer Table

```
SELECT email, customer_id,address_id,  
       COUNT(*)  
FROM customer  
GROUP BY email, customer_id, address_id  
HAVING COUNT (*)>1
```

Check for unique data:

```
SELECT DISTINCT *  
FROM customer
```

For both these, I would follow the same procedures as above.

## 2. Film Table

```
SELECT  
    MIN(rental_duration) AS minimum_rental_duration,  
    MAX(rental_duration) AS maximum_rental_duration,  
    AVG(rental_duration) AS average_rental_duration,  
    MIN(rental_rate) AS minimum_rental_rate,  
    MAX(rental_rate) AS maximum_rental_rate,  
    AVG(rental_rate) AS average_rental_rate,  
    MIN(length) AS minimum_length,  
    MAX(length) AS maximum_length,
```

AVG(length) AS average\_length,  
 MIN(replacement\_cost) AS minimum\_replacement\_cost,  
 MAX(replacement\_cost) AS maximum\_replacement\_cost,  
 AVG(replacement\_cost) AS average\_replacement\_cost,  
 MODE () WITHIN GROUP (ORDER BY rating) AS rating\_mode,  
 MODE() WITHIN GROUP (ORDER BY language\_id) AS language\_mode,  
 MODE() WITHIN GROUP (ORDER BY release\_year) AS year\_mode

FROM film

Results:

<b>minimum_rental_duration</b> smallint	<b>maximum_rental_duration</b> smallint	<b>average_rental_duration</b> numeric
3	7	4.9850000000000000
<b>minimum_rental_rate</b> numeric	<b>maximum_rental_rate</b> numeric	<b>average_rental_rate</b> numeric
0.99	4.99	2.9800000000000000
<b>minimum_length</b> smallint	<b>maximum_length</b> smallint	<b>average_length</b> numeric
46	185	115.2720000000000000
<b>minimum_replacement_cost</b> numeric	<b>maximum_replacement_cost</b> numeric	<b>average_replacement_cost</b> numeric
9.99	29.99	19.9840000000000000
<b>rating_mode</b> mpaa_rating	<b>language_mode</b> smallint	<b>year_mode</b> integer
PG-13	1	2006

### Customer Table

SELECT

MIN(customer\_id) AS minimum\_customer\_id,  
 MAX(customer\_id) AS maximum\_rental\_customer\_id,  
 MIN(store\_id) AS minimum\_store\_id,  
 MAX(store\_id) AS maximum\_store\_id,  
 MODE() WITHIN GROUP (ORDER BY store\_id) AS store\_id\_mode,  
 MIN(address\_id) AS minimum\_address\_id,  
 MAX(address\_id) AS maximum\_address\_id,  
 MIN(create\_date) AS minimum\_create\_date,

```

MAX(create_date) AS maxium_create_date,
MODE () WITHIN GROUP (ORDER BY create_date) AS rcreate_date_mode,
MIN(last_update) AS minimum_last_update,
MAX(last_update) AS maxium_last_update,
MODE () WITHIN GROUP (ORDER BY last_update) AS last_update_mode,
MODE() WITHIN GROUP (ORDER BY activebool) AS active_mode

```

FROM customer

<b>minimum_customer_id</b> integer	<b>maximum_rental_customer_id</b> integer	<b>minimum_store_id</b> smallint	<b>maximum_store_id</b> smallint
1	599	1	2
<b>store_id_mode</b> smallint	<b>minimum_address_id</b> smallint	<b>maxium_address_id</b> smallint	<b>minimum_create_date</b> date
1	5	605	2006-02-14
<b>maxium_create_date</b> date	<b>rcreate_date_mode</b> date	<b>minimum_last_update</b> timestamp without time zone	<b>maxium_last_update</b> timestamp without time zone
2006-02-14	2006-02-14	2013-05-26 14:49:45.738	2013-05-26 14:49:45.738
<b>maxium_last_update</b> timestamp without time zone	<b>last_update_mode</b> timestamp without time zone	<b>active_mode</b> boolean	
2013-05-26 14:49:45.738	2013-05-26 14:49:45.738	true	

3.

SQL makes the data cleaning and profiling pretty easy. Once I get the 'grammar' it is much more intuitive than Excel. I do think Excel works well with smaller data sets and it does have some great graphic features that SQL does not have. Excel will also automatically find a lot of these descriptive statistics which is nice to have at your fingertips too.