

In this task, you'll calculate some descriptive statistics using the *MIN*, *MAX*, *AVG*, *COUNT*, *SUM*, and *MODE()* aggregates discussed in this Exercise, and you'll reflect on what you learned about data profiling back in [Exercise 1.5: Data Profiling & Integrity](#).

Directions

Rockbuster's database engineers have loaded some new data into the database, and your manager has asked you to clean and profile it. Follow the instructions below to complete their request:

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).

Film table

Checking duplicate data

```

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

```

film_id	title	release_year	language_id	rental_duration	count
[PK] integer	character varying (255)	integer	smallint	smallint	bigint

No duplicates found. For cleaning I could either create a virtual table "VIEW" with unique records or delete the duplicate records. As creating a "View" is safest and preferred way I would use query to return unique records with **GROUP BY** or **DISTINCT**.

Checking missing values

Whether values are missing WHERE and OR commands can be combined.

```

1 SELECT * FROM film
2 WHERE film_id IS NULL
3 OR title IS NULL
4 OR description IS NULL
5 OR release_year IS NULL
6 OR release_year IS NULL
7 OR language_id IS NULL
8 OR rental_duration IS NULL
9 OR rental_rate IS NULL
10 OR length IS NULL
11 OR replacement_cost IS NULL
12 OR rating IS NULL
13 OR last_update IS NULL
14 OR special_features IS NULL
15 OR fulltext IS NULL;

```

film_id	title	description	release_year	language_id	rental_duration	rental_rate	length	replacement_cost	rating
[PK] integer	character varying (255)	text	integer	smallint	smallint	numeric (4,2)	smallint	numeric (5,2)	tinyint

Or, checking for NULLs for some random important variables can be done, e.g., as shown below.

Query

Query History

1

2

3

4

5

6

SELECT

SUM(CASE WHEN rental_duration IS NULL THEN 1 ELSE 0 END) AS rental_duration_missing,

SUM(CASE WHEN rental_rate IS NULL THEN 1 ELSE 0 END) AS rental_rate_missing,

SUM(CASE WHEN replacement_cost IS NULL THEN 1 ELSE 0 END) AS replacement_cost_missing,

SUM(CASE WHEN rating IS NULL THEN 1 ELSE 0 END) AS rating_missing

FROM film;

Data Output

Messages

Notifications

SQL

	rental_duration_missing bigint	rental_rate_missing bigint	replacement_cost_missing bigint	rating_missing bigint
1	0	0	0	0

In this case no NULLs are found and we assume no value is missing. If any NULL is found, then further investigations can be done, e.g., identifying the rows with missing values to determine if missing value is less or more than 5%. If it is less than 5%, imputing averages can be done, else the column should be omitted from analysis and reporting for the output transparency.

Checking Non-Uniform data

To check the Non-Uniform data I would use **GROUP BY** and **DISTINCT** query functions to find inconsistency randomly for any column, such as checking 'rating' or release year column as example.

Query		Query History
1	SELECT DISTINCT rating	
2	FROM film	
3	ORDER BY rating	

Data Output		Messages	Notifications
	rating		
	mpaa_rating		
1	G		
2	PG		
3	PG-13		
4	R		
5	NC-17		

Query		Query History
1	-- For example checking release year	
2	SELECT release_year	
3	FROM film	
4	GROUP BY release_year	

Data Output		Messages	Notifications
	release_year		
	integer		
1	2006		

No non-uniform data is found. If non-uniform data exist in the table, values in the table can be updated using **UPDATE** command.

Customer table

Checking duplicate data

Query Query History

```

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

```

Data Output Messages Notifications

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

If duplicates are found, similar ways of cleaning can be approached.

Checking missing values

Query Query History

```

1 SELECT * FROM customer
2 WHERE customer_id IS NULL
3       OR store_id IS NULL
4       OR first_name IS NULL
5       OR last_name IS NULL
6       OR email IS NULL
7       OR address_id IS NULL
8       OR activebool IS NULL
9       OR create_date IS NULL
10      OR last_update IS NULL
11      OR active IS NULL;

```

Data Output Messages Notifications

customer_id	store_id	first_name	last_name	email	address_id	activebool	create_date	last_update	active
[PK] integer	smallint	character varying (45)	character varying (45)	character varying (50)	smallint	boolean	date	timestamp without time zone	integer

Similar approach as above for film table.

Checking Non-Uniform Data

Query Query History

```

1 SELECT DISTINCT store_id
2 FROM customer;

```

Data Output Messages Notifications

Showing rows: 1 to 2 Page No: 1

store_id
1
2

Query Query History

```

1 SELECT create_date
2 FROM customer
3 WHERE create_date < '1900-01-01'
4       OR create_date > CURRENT_DATE;

```

Data Output Messages Notifications

create_date
date

In the second query last_update has been checked whether the date lies within the usual range as example. No non-uniform data is found. In case of non-uniform data I would have use UPDATE command to correct the data to make it consistent.

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.

Film Table:

For numerical columns

Query	Query History
1	<code>SELECT MIN(film_id),</code>
2	<code>MAX(film_id),</code>
3	<code>ROUND(AVG(film_id), 2),</code>
4	<code>COUNT(film_id)</code>
5	<code>FROM film;</code>
Data Output	Messages
min integer	max integer
1	1000
round numeric	count bigint
500.50	1000

Query	Query History
1	<code>SELECT MIN(rental_rate) AS min_rent,</code>
2	<code>MAX(rental_rate) AS max_rent,</code>
3	<code>ROUND(AVG(rental_rate), 2) AS ave_rent,</code>
4	<code>COUNT(rental_rate) AS count_rent_value</code>
5	<code>FROM film;</code>
Data Output	Messages
min_rent numeric	max_rent numeric
0.99	4.99
ave_rent numeric	count_rent_value bigint
2.98	1000

For non-numerical columns

Query	Query History
1	<code>SELECT MODE() WITHIN GROUP (ORDER BY fulltext)</code>
2	<code>AS modal_value</code>
3	<code>FROM film;</code>
Data Output	Messages
modal_value	tsvector
1	'balloon':19 'confront':14 'documentari':5 'feminist':8,11,16 'mile':2 'must':13 'spi':1 'thrill':4

Summary:

Columns	MIN	MAX	AVG	COUNT	COUNT ROWS	MODE
film id	1	1000	500.50	1000	1000	
title				1000	1000	Academy Dinosaur
release year	2006	2006	2006	1000	1000	
language id	1	1	1	1000	1000	
rental duration	3	7	4.99	1000	1000	
rental rate	0.99	4.99	2.98	1000	1000	
length	46	185	115.27	1000	1000	
replacement cost	9.99	29.99	19.98	1000	1000	
rating				1000	1000	PG-13

last_update				1000	1000	2013-05-26 14:50:58.951
special_features				1000	1000	{Trailers,Commentaries, ""Behind the Scenes""}
full_text				1000	1000	'baloon':19 'confront':14 'documentari':5 'feminist':8,11,16 'mile':2 'must':13 'spi':1 'thrill':4

Customer Table:

For numerical columns

Query	Query History
<pre> 1 SELECT MIN(active), 2 MAX(active), 3 ROUND(AVG(active), 2), 4 COUNT(active), 5 COUNT(*) 6 FROM customer;</pre>	
Data Output	Messages Notifications
<div> <div>min integer</div> <div>max integer</div> <div>round numeric</div> <div>count bigint</div> <div>count bigint</div> </div>	
1	0 1 0.97 599 599

For non-numerical columns

Query	Query History
<pre> 1 SELECT MODE() WITHIN GROUP (ORDER BY first_name) 2 AS modal_value 3 FROM customer;</pre>	
Data Output	Messages Notifications
<div>modal_value character varying</div>	
1	Jamie

Summary:

Columns	MIN	MAX	AVG	COUNT	COUNT ROWS	MODE
customer_id	1	599	300.00	599	599	
store_id	1	2	1.46	599	599	
first_name				599	599	Jamie
last_name				599	599	Abney
email				599	599	aaron.selby@sakilacustomer.org
address_id	5	605	304.72	599	599	
activebool				599	599	true
create_date				599	599	2006-02-14
last_update				599	599	2013-05-26 14:49:45.738
active	0	1	0.97	599	599	

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.

As it is known that SQL is more effective for larger data set. Various functions such as MIN, MAX, AVG, COUNT etc. can be executed at once in SQL, which makes SQL faster than Excel. Furthermore, SQL provide powerful querying capabilities to aggregate, filter, and identify data quality issue quickly. Excel is user-friendly and great for quick, visual exploration of smaller datasets, offering built-in functions and filters that are easy to use without deep technical knowledge.

4. Save your “Answers 3.6” document as a PDF and upload it here for your tutor to review.