

Answers to 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).

---Checking for rating values that don't fit in to uniformity---

Query Editor	Query History	Data Output	Explain	Messages	Notifications
1	SELECT DISTINCT rating	rating mpaa_rating			
2	FROM film				
3	GROUP BY rating				
		1 G			
		2 PG			
		3 PG-13			
		4 R			
		5 NC-17			

If values are not showing uniformity, then a constant value for each category should be chosen and then the values should be all changed to that value for each category. For future entries, each column could have constraints put on it to make sure that there are no odd values entered in the future.

---Duplicate search---

Query Editor	Query History	Data Output	Explain	Messages	Notifications
1	SELECT title, release_year,	title character varying (255)			
2	language_id, rental_duration,	release_year integer			
3	COUNT (*)	language_id smallint			
4	FROM film	rental_duration smallint			
5	GROUP BY title, release_year,	count bigint			
6	language_id, rental_duration				
7	HAVING COUNT (*) >1;				

Query Editor	Query History	Data Output	Explain	Messages	Notifications
1	SELECT first_name, last_name,	first_name character varying (45)			
2	address_id, email, active,	last_name character varying (45)			
3	COUNT (*)	address_id smallint			
4	FROM customer	email character varying (50)			
5	GROUP BY first_name, last_name,				
6	address_id, email, active				
7	HAVING COUNT (*) >1;				

Values coming up as duplicate can either be hidden with specific query types or if it is necessary, they can be deleted. Another option you can limit the query to distinct values as well.

---Missing value search---

Query Editor	Query History	Data Output	Explain	Messages	Notifications
1	SELECT	count_title bigint	count_rental_duration bigint	count_rental_rate bigint	count_len bigint
2	COUNT (title) AS count_title,	1	1000	1000	1000
3	COUNT (rental_duration) AS count_rental_duration,				
4	COUNT (rental_rate) AS count_rental_rate,				
5	COUNT (length) AS count_length,				
6	COUNT (replacement_cost) AS count_replacement_cost,				
7	COUNT (*) AS count_rows				
8	FROM film;				

Query Editor	Query History	Data Output	Explain	Messages	Notifications
1	SELECT	count_email bigint	count_last_name bigint	count_first_name bigint	count_custc bigint
2	COUNT (email) AS count_email,	1	599	599	599
3	COUNT (last_name) AS count_last_name,				
4	COUNT (first_name) AS count_first_name,				
5	COUNT (customer_id) AS count_customer_id,				
6	COUNT (store_id) AS count_store_id,				
7	COUNT (*) AS count_rows				
8	FROM customer;				

Any missing values can be filled in by finding column averages and inputting this value in to the blank spots. You could omit the entire column, however this is typically not a very desirable situation when it comes to data summarization.

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.

---descriptive statistics query film table---

Query Editor Query History

```

1  SELECT MIN(rental_rate) AS min_renatl_rate,
2  MAX(rental_rate) AS max_rental_rate,
3  AVG(rental_rate) AS avg_renatal_rate,
4  MIN(rental_duration) AS min_rental_duration,
5  MAX(rental_duration) AS max_rental_duration,
6  AVG(rental_duration) AS avg_rental_duration,
7  MIN(film_id) AS min_film,
8  MAX(film_id) AS max_film,
9  AVG(film_id) AS avg_film,
10 MIN(language_id) AS min_language,
11 MAX(language_id) AS max_language,
12 AVG(language_id) AS avg_language,
13 MIN(length) AS min_length,
14 MAX(length) AS max_length,
15 AVG(length) AS avg_length,
16 MIN(replacement_cost) AS min_replacement_cost,
17 MAX(replacement_cost) AS max_replacement_cost,
18 AVG(replacement_cost) AS avg_replacement_cost
19 FROM film

```

Data Output Explain Messages Notifications

	min_renatl_rate numeric	max_rental_rate numeric	avg_renatal_rate numeric	min_rental_duration smallint	max_rental_duration smallint	avg_rental_duration numeric
1	0.99	4.99	2.9800000000000000	3	7	4.9850000000000000

Data Output Explain Messages Notifications

min_film integer	max_film integer	avg_film numeric	min_language smallint	max_language smallint	avg_language numeric	min_length smallint
1	1000	500.5000000000000000	1	1	1.0000000000000000	46

lain Messages Notifications

max_length smallint	avg_length numeric	min_replacement_cost numeric	max_replacement_cost numeric	avg_replacement_cost numeric
185	115.2720000000000000	9.99	29.99	19.9840000000000000

---descriptive statistics query customer table---

Query Editor	Query History	Data Output	Explain	Messages	Notifications												
<pre>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;</pre>		<table><thead><tr><th></th><th>min_active integer</th><th>max_active integer</th><th>avg_active numeric</th><th>min_address smallint</th><th>max_address smallint</th></tr></thead><tbody><tr><td>1</td><td>0</td><td>1</td><td>0.97495826377295492487</td><td>5</td><td>605</td></tr></tbody></table>		min_active integer	max_active integer	avg_active numeric	min_address smallint	max_address smallint	1	0	1	0.97495826377295492487	5	605			
	min_active integer	max_active integer	avg_active numeric	min_address smallint	max_address smallint												
1	0	1	0.97495826377295492487	5	605												

Data Output	Explain	Messages	Notifications			
max_address smallint	avg_address numeric	min_customer integer	max_customer integer	avg_customer numeric	min_store smallint	max_store smallint
605	304.7245409015025042	1	599	300.0000000000000000	1	2

avg_store numeric
1.4557595993322204

---Finding the mode (most repeated) for non numeric columns film table---

Query Editor

Query History

```
1 SELECT mode() WITHIN GROUP (ORDER BY rating)
2 AS rating_value,
3 mode() WITHIN GROUP (ORDER BY special_features)
4 AS Feature_value,
5 mode() WITHIN GROUP (ORDER BY release_year)
6 AS year_value,
7 mode() WITHIN GROUP (ORDER BY title)
8 AS title_value
9 FROM film;
```

Data Output

Explain

Messages

Notifications

	rating_value mpaa_rating	feature_value text[]	year_value integer	title_value character varying
1	PG-13	(Trailers,Commentaries,"Behind the Scenes")	2006	Academy Dinosaur

---Finding the mode (most repeated) for non numeric columns customer table---

Query Editor		Query History			Data Output			Explain	Messages	Notifications
1	SELECT	mode()	WITHIN GROUP	(ORDER BY first_name)	first_name_value	last_name_value	email_value			
2	AS	first_name_value,			character varying	character varying	character varying			
3	mode()	WITHIN GROUP	(ORDER BY last_name)							
4	AS	last_name_value,								
5	mode()	WITHIN GROUP	(ORDER BY email)							
6	AS	email_value								
7	FROM	customer								

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.

Without a doubt SQL would be much quicker and more efficient for the above work done. Excel would not take a really long time, however building the pivot tables is just not as efficient with a lot more clicks on the screen to get the same result. For SQL you would just have to know how to ask the question properly, then your entire answer is there in milliseconds.