

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.

Checking for duplicate records (this could also be done with the DISTINCT queries shown for the non-uniform answer below):

Query

Query History

Scratch Pad

1

2

3

4

5

6

7

8

--Show only unique records in film table

SELECT DISTINCT film_id,

title,

description

FROM film

Data Output

Messages

Notifications

+

📄

▼

📋

🗑️

📦

📥

📈

	film_id [PK] integer	title character varying (255)	description text
1	503	Kramer Chocolate	A Amazing Yarn of a Robot And a Pastry Chef who must Redeem a Mad Scientist in The Outback
2	182	Control Anthem	A Fateful Documentary of a Robot And a Student who must Battle a Cat in A Monastery
3	158	Clones Pinocchio	A Amazing Drama of a Car And a Robot who must Pursue a Dentist in New Orleans
4	643	Orient Closer	A Astounding Epistle of a Technical Writer And a Teacher who must Fight a Squirrel in The Sahara Desert
5	375	Grail Frankenstein	A Unbelievable Saga of a Teacher And a Monkey who must Fight a Girl in An Abandoned Mine Shaft
6	35	Arachnophobia Rollercoaster	A Action-Packed Reflection of a Pastry Chef And a Composer who must Discover a Mad Scientist in The First M


```

1  --Show only unique records in customer table
2
3  CREATE VIEW customer_unique AS
4  SELECT customer_id,
5         first_name,
6         last_name,
7         email
8  FROM customer
9  GROUP BY customer_id,
10         first_name,
11         last_name,
12         email
13

```

Checking for non-uniform records: I show the DISTINCT queries here to differentiate them from the duplicate queries above. However, quickly identifying differences in the ways data had been entered would be more easily accomplished with the GROUP BY queries shown above.

The screenshot shows a database query editor interface. The top toolbar includes icons for file operations, query execution, and settings. The main query area contains the following SQL code:

```

1  --Show only unique records in film table
2
3  SELECT DISTINCT film_id,
4                 title,
5                 description
6  FROM film
7
8

```

Below the query editor, the 'Data Output' tab is active, displaying the results of the query in a table format. The table has four columns: an index, film_id, title, and description. The data is as follows:

	film_id [PK] integer	title character varying (255)	description text
29	792	Shrek License	A Fateful Yarn of a Secret Agent And a Feminist who must Find a Feminist in A Jet Boat
30	289	Eve Resurrection	A Awe-Inspiring Yarn of a Pastry Chef And a Database Administrator who must Challenge a Teacher in A Baloon
31	50	Baked Cleopatra	A Stunning Drama of a Forensic Psychologist And a Husband who must Overcome a Waitress in A Monastery
32	208	Dares Pluto	A Fateful Story of a Robot And a Dentist who must Defeat a Astronaut in New Orleans
33	897	Torque Bound	A Emotional Display of a Crocodile And a Husband who must Reach a Man in Ancient Japan
34	40	Army Flintstones	A Boring Saga of a Database Administrator And a Womanizer who must Battle a Waitress in Nigeria

Query Query History Scratch Pad x

```

1  --Show only unique records in customer table
2
3  SELECT DISTINCT customer_id,
4     first_name,
5     last_name,
6     email
7  FROM customer
8
9

```

Data Output Messages Notifications

	customer_id [PK] integer	first_name character varying (45)	last_name character varying (45)	email character varying (50)
1	114	Grace	Ellis	grace.ellis@sakilacustomer.org
2	374	Jeremy	Hurtado	jeremy.hurtado@sakilacustomer.org
3	172	Bernice	Willis	bernice.willis@sakilacustomer.org
4	556	Armando	Gruber	armando.gruber@sakilacustomer.org
5	420	Jacob	Lance	jacob.lance@sakilacustomer.org
6	354	Justin	Ngo	justin.ngo@sakilacustomer.org

If I identified non-uniform entries with a GROUP BY query, I could use an UPDATE query to SET one uniform value.

Checking for missing records: I could use COUNT queries to identify columns that have missing values.

Query Query History Scratch Pad x

```

1  --Identifying columns with missing values
2
3  SELECT COUNT(film_id) AS film_id_count,
4     COUNT(title) AS title_count,
5     COUNT(description) AS description_count,
6     COUNT(release_year) AS release_year_count,
7     COUNT(language_id) AS language_id_count,
8     COUNT(rental_duration) AS rental_duration_count,
9     COUNT(rental_rate) AS rental_rate_count
10 FROM film
11
12
13
14
15
16

```

Data Output Messages Notifications

	film_id_count bigint	title_count bigint	description_count bigint	release_year_count bigint	language_id_count bigint	rental_duration_count bigint	rental_rate_count bigint
1	1000	1000	1000	1000	1000	1000	1000

If there were many missing values in a column, I would eliminate that column from my analysis. If there were only a few missing values, I would impute values for the missing records.

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.

Query
Query History
Scratch Pad

```

1  --Summarizing the data from film table
2
3  SELECT MIN(release_year) AS min_release_year,
4         MAX(release_year) AS max_release_year,
5         MODE() WITHIN GROUP (ORDER BY release_year) AS modal_release_year,
6         MODE() WITHIN GROUP (ORDER BY language_id) AS modal_language_id,
7         MIN(rental_duration) AS min_rental_duration,
8         MAX(rental_duration) AS max_rental_duration,
9         AVG(rental_duration) AS avg_rental_duration,
10        MIN(rental_rate) AS min_rental_rate,
11        MAX(rental_rate) AS max_rental_rate,
12        AVG(rental_rate) AS avg_rental_rate,
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        MODE() WITHIN GROUP (ORDER BY rating) AS modal_rating,
20        MIN(last_update) AS min_last_update,
21        MAX(last_update) AS max_last_update,
22        MODE() WITHIN GROUP (ORDER BY last_update) AS modal_last_update
23  FROM film
24

```

Data Output
Messages
Notifications

Query
Query History
Scratch Pad

```

1  --Summarizing the data from customer table
2
3  SELECT MODE() WITHIN GROUP (ORDER BY store_id) AS modal_store_id,
4         MODE() WITHIN GROUP (ORDER BY activebool) AS modal_activebool,
5         MIN(create_date) AS min_create_date,
6         MAX(create_date) AS max_create_date,
7         MODE() WITHIN GROUP (ORDER BY create_date) AS modal_create_date,
8         MIN(last_update) AS min_last_update,
9         MAX(last_update) AS max_last_update,
10        MODE() WITHIN GROUP (ORDER BY last_update) AS modal_last_update
11  FROM customer
12
13
14
15

```

Data Output
Messages
Notifications

	modal_store_id smallint	modal_activebool boolean	min_create_date date	max_create_date date	modal_create_date date	min_last_update timestamp without time zone	max_last_update timestamp without time zone	modal_last_update timestamp without time zone
1	1	true	2006-02-14	2006-02-14	2006-02-14	2013-05-26 14:49:45.738	2013-05-26 14:49:45.738	2013-05-26 14:49:45.738