

# Kareemah Ashiru

## Exercise 3.6

### 1. Check for and clean dirty data

#### a. Duplicates: No duplicates found for film and customer fact tables.

i.

Dashboard × Processes × Rockbuster/postgres@PostgreSQL 16 ×

Rockbuster/postgres@PostgreSQL 16

Query 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 Messages Notifications

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

Dashboard × Processes × Rockbuster/postgres@PostgreSQL 16 ×

Rockbuster/postgres@PostgreSQL 16

Query Query History

```
1 SELECT
2     customer_id
3     ,store_id
4     ,first_name
5     ,last_name
6     ,email
7     ,address_id
8     ,activebool
9     ,create_date
10    ,last_update
11    ,active
12    ,COUNT(*)
13 FROM customer
14 GROUP BY
15     customer_id
16     ,store_id
17     ,first_name
18     ,last_name
19     ,email
20     ,address_id
21     ,activebool
22     ,create_date
23     ,last_update
24     ,active
25 HAVING COUNT(*) >1;
```

Data Output Messages Notifications

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

- b. Non Uniform data:** 1000 out of 1000 records were returned for the film table. This means that all records are unique. For the customer table; 599 out 599 records were returned. This means that all the records are unique.



	film_id	title character varying (255)	description	release_year integer	language_id smallint	rental_duration smallint	rental_rate numeric (4,2)	length smallint	replacement_cost numeric (5,2)	rating mpaa_rating	last_update timestamp without time zone	special_features
1	1	Academy Dinosaur	A Epic Drama of a Feminist And a Man who must Battle a Teacher in The Canadian Rockies	2006	1	6	6.99	86	20.99	PG	2013-05-26 14:50:58.951	(Trailers,Deleted)
2	2	Joe Goldberg	A Audacious Epistle of a Database Administrator And a Explorer who must Find a Car in Ancient China	2006	4	6	4.99	40	12.99	NC-17	2013-05-26 14:50:58.951	(Trailers,Deleted)
3	3	Adaptation Holes	A Astonishing Reflection of a Lumberjack And a Sci who must Sink a Leechman in A Russian Factory	2006	1	7	2.99	50	18.99	PG	2013-05-26 14:50:58.951	(Trailers,Deleted)
4	4	Alibi Prejudice	A Fanciful Documentary of a Fishbean And a Lumberjack who must Chase a Monkey in A Shark Tank	2006	1	5	2.99	110	22.99	G	2013-05-26 14:50:58.951	(Commentaries,)
5	5	African Egg	A Fast-Paced Documentary of a Patchy Chef And a Dentist who must Pursue a Forensic Psychologist in The Gulf of Mexico	2006	1	6	1.99	137	26.99	G	2013-05-26 14:50:58.951	(Deleted,Scenes)
6	6	Agent Truman	A Intrepid Panorama of a Robot And a Boy who must Escape a Sumo Wrestler in Ancient China	2006	3	2	4.99	169	17.99	PG	2013-05-26 14:50:58.951	(Deleted,Scenes)
7	7	Angrybe Sierra	A Touching Saga of a Hunter And a Butler who must Discover a Butler in A Jet Boat	2006	1	6	4.99	54	28.99	PG-13	2013-05-26 14:50:58.951	(Trailers,Deleted)
8	8	Arctic Pollock	A Trip of a Moose And a Girl who must Confront a Monkey in Ancient India	2006	1	6	4.98	24	15.99	PG	2013-05-26 14:50:58.951	(Trailers,Deleted)
9	9	Autismbe Devil	A Thoughtful Farcical of a Database Administrator And a Mast Scientist who must Outgun a Mast Scientist in A Jet Boat	2006	1	6	2.99	114	21.99	PG-13	2013-05-26 14:50:58.951	(Trailers,Deleted)
10	10	Aldriden Calendar	A Action-Packed Tale of a Man And a Lumberjack who must Reach a Feminist who must Outgun a Feminist in Ancient China	2006	1	6	4.99	63	24.99	PG-13	2013-05-26 14:50:58.951	(Trailers,Deleted)
11	11	Alamo Videotape	A Boring Epistle of a Butler And a Cat who must Fight a Patchy Chef in A MySQL Convention	2006	1	6	0.99	126	16.99	G	2013-05-26 14:50:58.951	(Commentaries,)
12	12	Alaska Phantom	A Fanciful Saga of a Hunter And a Patchy chef who must Vanquish a Boy in Australia	2006	1	6	0.99	136	22.99	PG	2013-05-26 14:50:58.951	(Commentaries,)
13	13	Ali Forever	A Mysterious Drama of a Dentist And a Crocodile who must Battle a Feminist in The Canadian Rockies	2006	1	4	4.99	150	21.99	PG	2013-05-26 14:50:58.951	(Trailers,Deleted)
14	14	Alce Fantasia	A Emotional Drama of a Shark And a Database Administrator who must Vanquish a Pioneer in Soviet Georgia	2006	1	6	0.99	94	23.99	NC-17	2013-05-26 14:50:58.951	(Trailers,Deleted)
15	15	Alibi Dreams of a Cat And a Mast Scientist who must Battle a Feminist in A MySQL Convent	2006	4	5	10.99	49	16.99	PG-13	2013-05-26 14:50:58.951	(Trailers,Deleted)	
16	16	Alley Evolution	A Fast-Paced Drama of a Robot And a Composer who must Chase a Butler a Borned in New Orleans	2006	1	5	2.99	180	23.99	NC-17	2013-05-26 14:50:58.951	(Trailers,Commentaries)
17	17	Alone Trip	A Fast-Paced Character Study of a Composer And a Dog who must Outgun a Boss in An Abandoned Fun House	2006	1	3	0.99	82	14.99	R	2013-05-26 14:50:58.951	(Trailers,Behind)
18	18	Alter Victory	A Thoughtful Drama of a Composer And a Feminist who must Meet a Secret Agent in The Canadian Rockies	2006	1	6	0.99	57	27.99	PG-13	2013-05-26 14:50:58.951	(Trailers,Behind)
19	19	Amadeus Holy	A Emotional Display of a Pioneer And a Technical Writer who must Battle a Mast Man in A Baloon	2006	1	6	0.99	113	20.99	PG	2013-05-26 14:50:58.951	(Commentaries,)
20	20	Amelie Heflighthers	A Boring Drama of a Woman And a Squirrel who must Conquer a Student in A Baloon	2006	1	4	4.99	79	23.99	R	2013-05-26 14:50:58.951	(Commentaries,)
21	21	American Circus	A Insightful Drama of a Girl And a Astronaut who must Chase a Database Administrator in A Shark Tank	2006	1	3	4.99	129	16.99	R	2013-05-26 14:50:58.951	(Commentaries,)

Dashboard
Processes
Rockbuster/postgres@PostgreSQL 16

Dashboard :buster/postgres@PostgreSQL 16

No limit

Query
Query History

```

1 SELECT DISTINCT
2   customer_id,
3   store_id,
4   first_name,
5   last_name,
6   email,
7   address_id,
8   activebool,
9   create_date,
10  last_update,
11  active
12 FROM customer;
13

```

Data Output
Messages
Notifications

	customer_id [PK] integer	store_id smallint	first_name character varying (45)	last_name character varying (45)	email character varying (50)	address_id smallint	activebool boolean	create_date date	last_update timestamp without time zone	active integer
1	357	1	Keith	Rico	keith.rico@sakilacustomer.org	362	true	2006-02-14	2013-05-26 14:49:45.738	1
2	171	2	Dolores	Wagner	dolores.wagner@sakilacustomer.org	175	true	2006-02-14	2013-05-26 14:49:45.738	1
3	139	1	Amber	Dixon	amber.dixon@sakilacustomer.org	143	true	2006-02-14	2013-05-26 14:49:45.738	1
4	471	1	Dean	Sauer	dean.sauer@sakilacustomer.org	476	true	2006-02-14	2013-05-26 14:49:45.738	1
5	594	1	Eduardo	Hiatt	eduardo.hiatt@sakilacustomer.org	600	true	2006-02-14	2013-05-26 14:49:45.738	1
6	401	2	Tony	Carranza	tony.carranza@sakilacustomer.org	406	true	2006-02-14	2013-05-26 14:49:45.738	1
7	157	2	Darlene	Rose	darlene.rose@sakilacustomer.org	161	true	2006-02-14	2013-05-26 14:49:45.738	1
8	154	2	Michele	Grant	michele.grant@sakilacustomer.org	158	true	2006-02-14	2013-05-26 14:49:45.738	1
9	530	2	Darryl	Ashcraft	darryl.ashcraft@sakilacustomer.org	536	true	2006-02-14	2013-05-26 14:49:45.738	1
10	493	1	Brent	Harkins	brent.harkins@sakilacustomer.org	498	true	2006-02-14	2013-05-26 14:49:45.738	1
11	542	2	Lonnie	Tirado	lonnie.tirado@sakilacustomer.org	548	true	2006-02-14	2013-05-26 14:49:45.738	1
12	566	1	Casey	Mena	casey.mena@sakilacustomer.org	572	true	2006-02-14	2013-05-26 14:49:45.738	1
13	186	2	Holly	Fox	holly.fox@sakilacustomer.org	190	true	2006-02-14	2013-05-26 14:49:45.738	1
14	128	1	Marjorie	Tucker	marjorie.tucker@sakilacustomer.org	132	true	2006-02-14	2013-05-26 14:49:45.738	1
15	466	1	Leo	Ebert	leo.ebert@sakilacustomer.org	471	true	2006-02-14	2013-05-26 14:49:45.738	1
16	494	2	Ramon	Choate	ramon.choate@sakilacustomer.org	499	true	2006-02-14	2013-05-26 14:49:45.738	1
17	178	2	Marion	Snyder	marion.snyder@sakilacustomer.org	182	true	2006-02-14	2013-05-26 14:49:45.738	1
18	65	2	Rose	Howard	rose.howard@sakilacustomer.org	69	true	2006-02-14	2013-05-26 14:49:45.738	1
19	450	1	Jay	Robb	jay.robb@sakilacustomer.org	455	true	2006-02-14	2013-05-26 14:49:45.738	1
20	234	1	Claudia	Fuller	claudia.fuller@sakilacustomer.org	238	true	2006-02-14	2013-05-26 14:49:45.738	1
21	343	1	Douglas	Graf	douglas.graf@sakilacustomer.org	348	true	2006-02-14	2013-05-26 14:49:45.738	1
22	288	1	Bobbie	Craig	bobbie.craig@sakilacustomer.org	293	true	2006-02-14	2013-05-26 14:49:45.738	1

c. **Missing data:** There is no missing data in both the Film table and the Customer table

Dashboard x Processes x Rockbuster/postgres@PostgreSQL 16 x

Rockbuster/postgres@PostgreSQL 16

Query Query History

```
1 SELECT *
2 FROM film
3 WHERE
4 film_id IS NULL OR
5 title IS NULL OR
6 description IS NULL OR
7 release_year IS NULL OR
8 language_id IS NULL OR
9 rental_duration IS NULL OR
10 rental_rate IS NULL OR
11 length IS NULL OR
12 replacement_cost IS NULL OR
13 rating IS NULL OR
14 last_update IS NULL OR
15 special_features IS NULL OR
16 fulltext IS NULL;
17
```

Data Output Messages Notifications

film_id	title	description	release_year	language_id	rental_duration	rental_rate	length	replacement_cost	rating	last_update	special_features	fulltext
[PK] integer	character varying (255)	text	integer	smallint	smallint	numeric (4,2)	smallint	numeric (5,2)	mpaa_rating	timestamp without time zone	text[]	tsvector

Dashboard x Processes x Rockbuster/postgres@PostgreSQL 16\* x

Rockbuster/postgres@PostgreSQL 16

Query Query History

```
1 SELECT *
2 FROM customer
3 WHERE
4 customer_id IS NULL OR
5 store_id IS NULL OR
6 first_name IS NULL OR
7 last_name IS NULL OR
8 Email IS NULL OR
9 address_id IS NULL OR
10 activebool IS NULL OR
11 create_date IS NULL OR
12 last_update IS NULL OR
13 active IS NULL;
14
```

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

## 2. Summarize your data using statistics

- a. **Film table numeric values:** release\_year, rental\_duration, rental\_rate, and replacement\_cost.

Dashboard × Processes × **Rockbuster/postgres@PostgreSQL 16\*** ×

Rockbuster/postgres@PostgreSQL 16

Query Query History

```
1 SELECT
2   MIN(release_year) AS release_year,
3   MAX(release_year) AS release_year,
4   AVG(release_year) AS release_year,
5   COUNT(release_year) AS release_year,
6   MIN(rental_duration) AS rental_duration,
7   MAX(rental_duration) AS rental_duration,
8   AVG(rental_duration) AS rental_duration,
9   COUNT(rental_duration) AS rental_duration,
10  MIN(rental_rate) AS min_rent,
11  MAX(rental_rate) AS max_rent,
12  AVG(rental_rate) AS avg_rent,
13  COUNT(rental_rate) AS rental_rate,
14  MIN(replacement_cost) AS replacement_cost,
15  MAX(replacement_cost) AS replacement_cost,
16  AVG(replacement_cost) AS replacement_cost,
17  COUNT(replacement_cost) AS replacement_cost
18 FROM film;
```

Data Output Messages Notifications

	release_year integer	release_year numeric	release_year bigint	release_year smallint	rental_duration smallint	rental_duration smallint	rental_duration numeric	rental_duration bigint	min_rent numeric	max_rent numeric	avg_rent numeric	rental_rate bigint	replacement_cost numeric	replacement_cost numeric	replacement_cost numeric	replacement_cost bigint
1	2006	2006	2006.0000000000000000	1000	3	7	4.9800000000000000	1000	0.99	4.99	2.9400000000000000	1000	9.99	25.99	19.9940000000000000	1000

- b. **Film table non-numeric values:** language\_id and rating

Dashboard × Processes × **Rockbuster/postgres@PostgreSQL 16\*** ×

Rockbuster/postgres@PostgreSQL 16

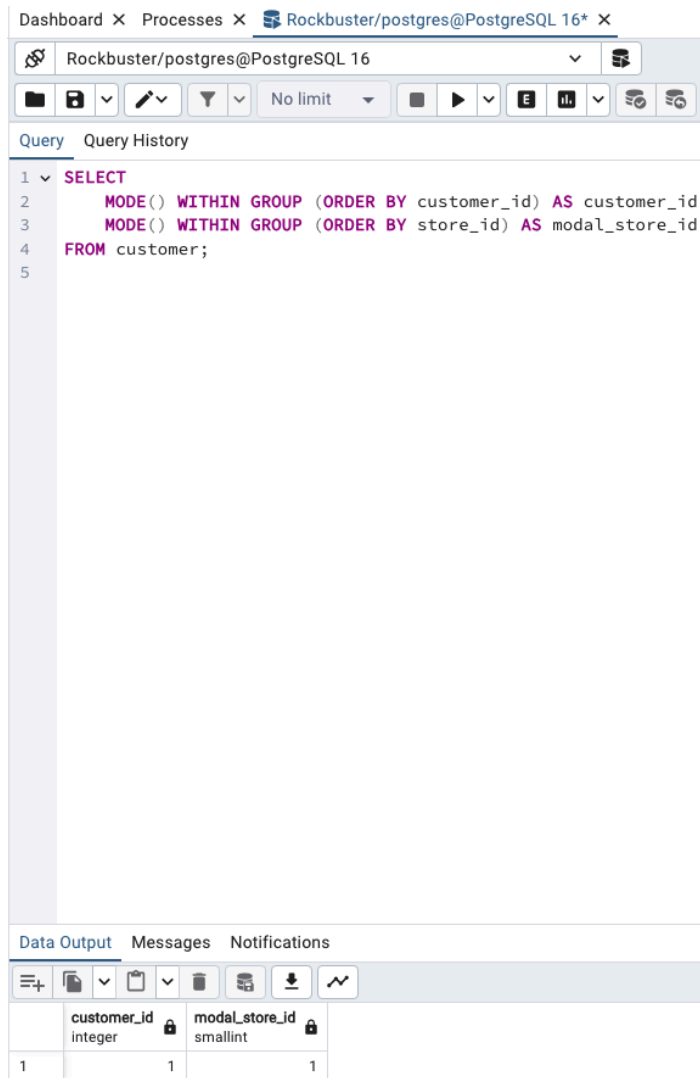
Query Query History

```
1 SELECT
2   MODE() WITHIN GROUP (ORDER BY rating) AS modal_rating,
3   MODE() WITHIN GROUP (ORDER BY language_id) AS modal_language_id
4 FROM film;
```

Data Output Messages Notifications

	modal_rating mpaa_rating	modal_language_id smallint
1	PG-13	1

c. **Customer table non-numeric values:** customer\_id and store\_id



The screenshot shows a PostgreSQL query editor interface. The top bar indicates the user is 'Rockbuster/postgres@PostgreSQL 16'. Below the toolbar, the query editor shows a SQL query:

```
1 SELECT
2     MODE() WITHIN GROUP (ORDER BY customer_id) AS customer_id,
3     MODE() WITHIN GROUP (ORDER BY store_id) AS modal_store_id
4 FROM customer;
5
```

Below the query editor, the 'Data Output' tab is active, showing the results of the query. The results are displayed in a table with two columns: 'customer\_id' (integer) and 'modal\_store\_id' (smallint). The table contains one row with the values 1 and 1.

	customer_id integer	modal_store_id smallint
1	1	1

3. **Reflection:** I find SQL to be faster and more efficient overall when it comes to cleaning dirty data and statistical data summary. What I prefer with SQL over Excel is the ability to immediately errors when performing a cleaning task. In excel the output of my results may not necessarily be what is correct. The only instance that I find Excel a bit easier is identifying missing, incorrect, or non uniform data. By clicking the small arrows on the filtered columns, I am able to quickly identify what values stand out from the crowd. However, with SQL, I have to manually sift through the data values for anomalies.