1.1 DUPLICATES

Film

```
1 SELECT film_1d,
2 title,
3 title,
4 crites_complete,
5 language_1d,
6 rental_duration,
7 rental_trate,
8 length,
9 replacement_cost,
11 last_update,
12 special_teatures,
13 fulltoxt,
14 description,
15 read_ouration,
16 rental_trate,
18 length,
19 replacement_cost,
21 special_teatures,
22 special_teatures,
23 description,
24 replacement_cost,
25 reating,
26 last_update,
27 special_teatures,
28 language_1d,
29 rental_crate,
20 last_update,
20 rental_trate,
21 rental_duration,
22 rental_furation,
23 language_1d,
24 replacement_cost,
25 reating,
26 last_update
27 special_teatures,
28 follows
29 special_teatures,
30 follows
20 special_teatures,
31 language_1d,
32 rental_crate,
33 language_1d,
34 replacement_cost,
35 reating,
36 last_update
37 special_teatures,
38 follows
38 language_1d,
39 rental_duration / metal_crate
39 last_update
30 last_update
30 last_update
30 last_update
30 last_update
31 language_1d,
32 language_1d,
33 language_1d,
34 replacement_cost
35 reating,
36 last_update
36 language_1d,
37 rental_trate
37 special_teatures
38 language_1d,
39 rental_duration / metal_crate
38 language_1d,
39 rental_trate
39 language_1d,
30 rental_trate
30 language_1d,
31 rental_duration_1d,
32 language_1d,
33 language_1d,
34 rental_duration_1d,
35 rental_trate
30 language_1d,
30 rental_trate
30 l
```

Output: No duplicates

Script

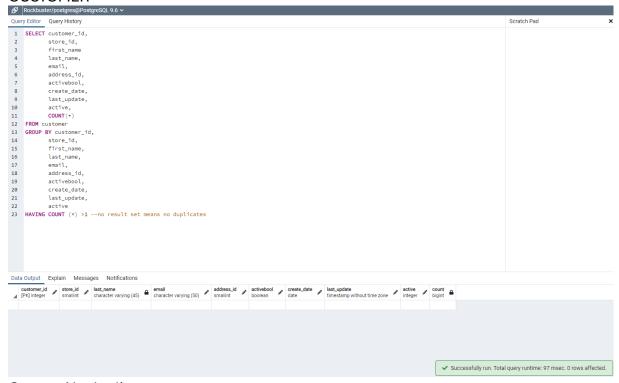
```
SELECT film_id,
        title,
        description,
        release_year,
        language_id,
        rental_duration,
        rental_rate,
        length,
        replacement_cost,
        rating,
        last_update,
        special_features,
        fulltext,
        COUNT(*)
FROM film
GROUP BY film_id,
        title,
        description,
        release_year,
        language_id,
        rental_duration,
        rental_rate,
```

length,

```
replacement_cost,
rating,
last_update,
special_features,
fulltext
```

HAVING COUNT (*) >1 --no result set means no duplicates

CUSTOMER



Output: No duplicates

Script

```
SELECT customer_id,
        store_id,
        first_name
        last_name,
        email,
        address_id,
        activebool,
        create_date,
        last_update,
        active,
        COUNT(*)
FROM customer
GROUP BY customer_id,
        store_id,
        first_name,
        last_name,
        email,
        address id,
```

```
activebool,
create_date,
last_update,
active

HAVING COUNT (*) >1 --no result set means no duplicates
```

CLEANING DUPLICATES:

After identifying whether or not duplicates exist, I would then create a view to flag the duplicates and be able to delete them like so:

```
CREATE VIEW AS duplicate_customers
SELECT
customer_id,
store_id,
first_name,
last_name,
email,
address_id,
activebool,
create_date,
last_update,
active,
```

ROW_NUMBER() OVER (-- ROW_NUMBER creates a column that assigns a number to each row on the view per the established partition, meaning that if there are more than one records of the same, the second record will be labeled as '2' and so forth. Therefore, a record that has no duplicates is only labeled as '1'

PARTITION BY store_id, first_name, last_name, email --PARTITION BY would determine how to group the records in the view and single out the duplicates. In this case, the partition would organize the records by store_id, the customer name, and email, thus only focusing on main identifiers of the customer.

ORDER BY last_update DESC – ORDER BY in this case would order the records in descending order by their last update, meaning that the newest records are at the top and are therefore labeled as '1' per the ROW_NUMBER command. I.e.: the older record will be the duplicate.

```
) AS row_num FROM customer;
```

After this you would delete the records like so:

```
DELETE FROM duplicate_customers --this is deleting the records from the view, not the base table

WHERE customer_id IN (
    SELECT customer_id
    FROM duplicate_customers
```

WHERE row_num > 1 --this part of the script is the one that scans for the row_num that is larger than 1, i.e. the duplicates);

Finally, you run the WHERE x>1 script again to check for duplicates and make sure you got rid of them all.

SELECT *
FROM duplicate_customers
WHERE row_num > 1;

1.2 NON-UNIFORM DATA

FILM

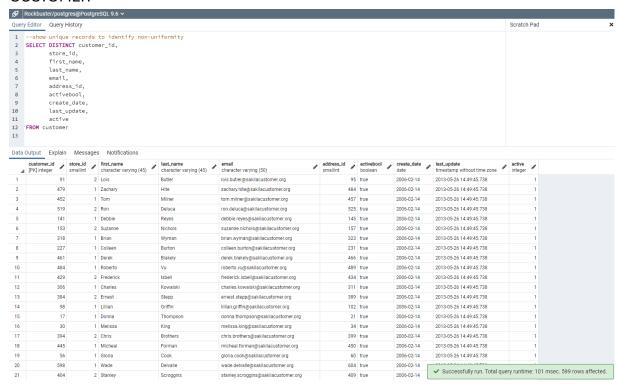
	Editor Q	uery H		Scratch Pad					
			records to identify						
	ELECT D								
3		itle.	- /						
4			ption,						
5	r	eleas	e_year,						
5	l	angua	ge_id,						
7			_duration,						
3			_rate,						
9 a		ength	ement cost,						
1		ating	- /						
2		_	pdate.						
3	S	pecia	l_features,						
4	f	ullte	xt						
5 F	ROM fil	.m							
ata O	utput Ex	xplain	Messages Notifications						
	film_id		title	description	release_year _	language_id	. rental_duration	_ rental_rate	
	[PK] intege	er 🎤	character varying (255)	text	integer	smallint	smallint "	numeric (4,	,2)
1		er 1	character varying (255) Academy Dinosaur	text A Epic Drama of a Feminist And a Mad Scientist who must Battle a Teacher in The Canadian Rockies	integer 2006	smallint	smallint *	numeric (4,	,2) 0.99
1 2		er 1 2	character varying (255) Academy Dinosaur Ace Goldfinger	text A Epic Drama of a Feminist And a Mad Scientist who must Battle a Teacher in The Canadian Rockies A Stounding Epistie of a Database Administrator And a Explorer who must Find a Car in Ancient China	integer 2006 2006	smallint	smallint '	numeric (4, 6	,2) 0.99 4.99
1 2 3		1 2 3	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes	text A Epic Drama of a Feminist And a Mad Scientist who must Battle a Teacher in The Canadian Rockies A Astounding Epistle of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baboon Factory	2006 2006 2006	smallint	smallint f	numeric (4, 6 3 7	0.99 4.99 2.99
1 2 3 4		1 2 3 4	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes Affair Prejudice	text A Epic Drama of a Feminist And a Mad Scientist who must Battie a Teacher in The Canadian Rockies A Astounding Episte of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Bakoon Factory A Fanciful Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shark Tank	2006 2006 2006 2006 2006	smallint	smallint f	numeric (4, 6 3 7 5	,2) 0.99 4.99 2.99 2.99
1 2 3 4 5		1 2 3 4 5	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes Affair Prejudice African Egg	text A Epic Drama of a Feminist And a Mad Scientist who must Battie a Teacher in The Canadian Rockies A Astounding Epistel of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Babon Factory A Fanciful Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shark Tank A Fast-Paced Documentary of a Pastry Chef And a Dentist who must Pursue a Forensic Psychologist in The Gulf of Mexico	2006 2006 2006 2006 2006 2006	smallint "	smallint ' 1 1 1 1 1 1 1	numeric (4; 6 3 7 5 6	0.99 4.99 2.99 2.99 2.99
1 2 3 4 5		1 2 3 4 5 6	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes Affair Prejudice African Egg Agent Truman	text A Epic Drama of a Feminist And a Mad Scientist who must Battle a Teacher in The Canadian Rockies A Astounding Epistle of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baboon Factory A Fanchful Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Post York And a Demistrative hormust Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Past-York And a Demistrative hormust Pursue a Forensic Psychologist in The Gulf of Mexico A Intrepid Panorama of a Robot And a Boy who must Escape a Sumo Wiestler in Ancient China	2006 2006 2006 2006 2006 2006 2006	smallint "	smallint ' smallint ' 1 1 1 1 1	numeric (4; 6 3 7 5 6 3	0.99 4.99 2.99 2.99 2.99 2.99
1 2 3 4 5 6 7		1 2 3 4 5 6 7	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes Affair Prejudice African Egg Agent Truman Airplane Sierra	text	2006 2006 2006 2006 2006 2006 2006 2006	smallint **	smallint " smalli	numeric (4; 6 3 3 7 5 6 6 3 6 6	0.99 4.99 2.99 2.99 2.99 2.99 4.99
1 2 3 4 5 6 7 8		1 2 3 4 5 6 7	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes Affair Prejudice African Egg Agent Truman	text A Epic Drama of a Feminist And a Mad Scientist who must Battle a Teacher in The Canadian Rockies A Astounding Epistle of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baboon Factory A Fanchful Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Post York And a Demistrative hormust Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Past-York And a Demistrative hormust Pursue a Forensic Psychologist in The Gulf of Mexico A Intrepid Panorama of a Robot And a Boy who must Escape a Sumo Wiestler in Ancient China	2006 2006 2006 2006 2006 2006 2006	smallint **	smallint ' smallint ' 1 1 1 1 1	numeric (4, 6 3 3 7 5 6 6 3 3 6 6 6 6	0.99 4.99 2.99 2.99 2.99 2.99
1 2 3 4 5 6 7 8		1 2 3 4 5 6 7 8	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes Affair Prejudice African Egg Agent Truman Airplane Sierra	text	2006 2006 2006 2006 2006 2006 2006 2006	smallint "	smallint " smalli	numeric (4; 6 3 3 7 5 6 6 3 6 6	0.99 4.99 2.99 2.99 2.99 2.99 4.99
1 2 3 4 5 6 7 8		1 2 3 4 5 6 7 8 9	character varying (255) Academy Dinosaur Ace Goldfinger Adaptation Holes Affair Prejudice Affrican Egg Agent Truman Airplane Sierra Airport Pollock	text A Epic Drama of a Ferninist And a Mad Scientist who must Battie a Teacher in The Canadian Rockies A Astounding Episte of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baloon Factory A Fanchul Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shark Tank A Fast-Paced Documentary of a Pastry Chef And a Dentist who must Pursue a Forensic Psychologist in The Gulf of Mexico A Interplid Panorama of a Robot And a Boy who must Eccape a Sumo Westler in Ancient China A Touching Saga of a Hunter And a Buttler who must Discover a Butter in A Jet Boat A Epic Tale of a Moose And a Girl who must Confront a Monkey in Ancient India	101eger 2006 2006 2006 2006 2006 2006 2006 2006	smallint **	smallint '	numeric (4, 6 3 3 7 5 6 6 3 3 6 6 6 6	0.99 4.99 2.99 2.99 2.99 2.99 4.99 4.99
1 2 3 4 5 6 7 8 9		1 2 3 4 5 6 7 8 9 10	character varying (255) Academy Dinosaur Ace Goldringer Adaptation Holes Affair Prejudice African Egg Agent Truman Airplane Sierra Airport Poliock Alabama Devil	text A Epic Drama of a Feminist And a Mad Scientist who must Battie a Teacher in The Canadian Rockies A Astounding Episte of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Babon Factory A Fanciful Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shark Tank A Fast-Paeced Documentary of a Pastry Chef And a Dentist who must Pursue a Forensic Psychologist in The Gulf of Mexico A Intrepid Panorama of a Robot And a Boy who must Escape a Subset in Ancient China A Touching Saga of a Hunter And a Butter who must Discover a Butter in A Jet Bost A Epic Tale of a Moose And a Girl who must Confront a Monkey in Ancient India A Thoughtful Panorama of a Database Administrator And a Mad Scientist who must Outgun a Mad Scientist in A Jet Bost	100 integer 2006 2006 2006 2006 2006 2006 2006 200	smallint "	smallint " smalli	numeric (4, 6 3 3 7 5 6 6 3 3 6 6 6 6 3 3	,2) 0.99 4.99 2.99 2.99 2.99 4.99 4.99 2.99
1 2 3 4 5 6 7 8 9		1 2 3 4 5 6 7 8 9 10 11	character varying (255) Academy Dinosaur Academy Dinosaur Academy Dinosaur Adaptation Holes Affair Prejudice Affair Prejudice Affair Prejudice Affair Prejudice Affair Prejudice Alipone Sierra Alipone Sierra Alipone Deli	text A Epic Drama of a Feminist And a Mad Scientist who must Battle a Teacher in The Canadian Rockies A Astounding Epistle of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baboon Factory A Fanciful Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Pastry Chef And a Dentist who must Puse a Forensic Psychologist in The Gulf of Mexico A Intrepic Panorama of a Rock And a Boy who must Escape a Sumo Wirestler in Ancient China A Touching Saga of a Hunter And a Butter who must Discover a Butler in A Jet Boat A Epic Tale of a Moose And a Gif who must Confront a Monkey in Ancient India A Thoughtful Panorama of a Database Administrator And a Mad Scientist who must Cutfun a Mad Scientist in A Jet Boat A Action-Packed Tale of a Man And a Lumberjack who must Reach a Feminist in Ancient China	Integer 2006	smallint	smallint '	numeric (4, 6 3 3 7 5 6 6 3 3 6 6 6 3 3 6	,2)
1 2 3 4 4 5 5 6 7 7 8 8 9 10 11		1 2 3 4 5 6 7 8 9 10 11 12	character varying (255) Academy Dinosaur Aca Goldringer Adaptation Holes Affair Pequidice Affair Pequidice Affair Regulate Affair Pequidice Affair Regulate Airport Politock Alabama Devil Aladidin Calendar Alamo Videotape	text	2006 2006 2006 2006 2006 2006 2006 2006	smallint	smallint '	numeric (4, 6 3 3 7 5 6 6 3 3 6 6 6 6 6	,2)
1 2 2 3 4 4 5 5 6 6 7 7 8 8 9 9 110 111 112 113		1 2 3 4 5 6 7 8 9 10 11 12 13	character varying (255) Academy Dinosaur Academy Dinosaur Aca Coldringer Adaptation Holes Affair Prejudice Affaira Reg Agent Truman Airpaine Sierra Airport Pollock Alabama Devil Aladdin Calendar Alamo Videotape Alaska Phantom	A Epic Drama of a Feminist And a Mad Scientist who must Battie a Teacher in The Canadian Rockies A Astounding Episte of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Bakon Factory A Fanchiu Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shark Tank A Fast-Paced Documentary of a Pastry Chef And a Dentist who must Pursue a Forensic Psychologist in The Gulf of Mexico A Integrid Panorama of a Robot And a Boy who must Escape a Sumo Westler in Ancient China A Touching Saga of a Hunter And a Buttler who must Discover a Butter in A Jet Boat A Epic Tale of a Moose And a Girl who must Confront a Monkey in Ancient India A Thoughtful Panorama of a Database Administrator And a Mad Scientist who must Outgun a Mad Scientist in A Jet Boat A Rotion-Packed Tale of a Man And a Lumberjack who must Reach a Femiliar in Ancient China A Boring Epistle of a Butter And a Cat who must Flight a Pastry Chef in A MySoll, Convention A Fanciful Saga of a Hunter And a Pastry Chef who must Vanquish a Boy in Australia	Integer	smallint	smallint	numeric (4, 6 3 3 7 5 5 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6	,2)
1 2 3 3 4 4 5 5 6 6 7 8 8 9 10 11 12 13		1 2 3 4 5 6 7 8 9 10 11 12 13 14	character varying (255) Academy Dinosau Academy Dinosau Academy Dinosau Academy Dinosau Adaptation Holes Affair Prejudice Affair Dinosau Affair Prejudice Affair Prejud	text A Epic Drama of a Feminist And a Mad Scientist who must Battle a Teacher in The Canadian Rockies A Astounding Episte of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baloon Factory A Fanchit Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Pastry Orde And a Demitt who must Pursue a Foreisn-Psychologist in The Gulf of Mexico A Intrepid Panorama of a Robot And a Boy who must Escape a Sumo Wiestler in Ancient China A Touching Saga of a Hunter And a Butler who must Discover a Butler in A yet Boat A Epic Tale of a Moose And a Gif who must Confront a Monkey in Ancient India A Thoughtful Panorama of a Database Administrator And a Mad Scientist who must Outgun a Mad Scientist in A Jet Boat A Action-Packed Tale of a Man And a Lumberjack who must Resch a Feminist in Ancient China A Boring Epistie of a Butler And a Cat who must Fight a Pastry Chef in A MySQL Convention Fanciful Saga of a Hunter And a Pastry Chef who must Stattle a Feminist in The Canadian Rockies	Integer	smallint	smallint (numeric (4, 6 6 3 3 6 6 6 6 6 6 4 4	0.99 4.99 2.99 2.99 2.99 4.99 2.99 4.99 4
1 2 2 3 3 4 4 5 5 6 6 7 7 8 8 9 9 0 0 1 1 2 2 3 3 4 4 4 5 5 6 6 6 1 1 2 2 3 3 4 4 4 5 5 7 8 7 8 7 8 7 8 7 8 7 8 7 8 7 8 7 8		1 2 2 3 4 4 5 5 6 6 7 7 8 9 10 11 12 13 13 14 15	character varying (255) Academy Dinosaur Academy Dinosaur Academy Dinosaur Adaptation Holes Affair Prejudice Affair Prejudice Affair Prejudice Affair Prejudice Affair Prejudice Affair Prejudice Alabema Devil Alabema Devil Alabema Devil Alabema Devil Alamon Varientar Alamon Vari	text A Epic Drama of a Feminist And a Mad Scientist who must Battie a Teacher in The Canadian Rockies A Astounding Epistle of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baboon Factory A Fanciful Documentary of a Frisbee And a Lumberjack who must Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Pastry Chef And a Dentist who must Pursue a Forensic Psychologist in The Gulf of Mexico A Intregiol Panorama of a Robot And a Boy who must Escape a Sumo Wrestler in A Jet Boat A Touching Saga of a Hunter And a Butler who must Discover a Butler in A Jet Boat A Epic Tale of a Moses And a Girl who must Confront a Monkey in Ancient India A Thoughtuff Panorama of a Pastabase Administrator And a Mad Scientist who must Outgun a Mad Scientist in A Jet Boat A Action-Packed Tale of a Man And a Lumberjack who must Reach a Feminist in Ancient China A Boring Epistler of a Butler And a Cat who must Fight a Pastry Chef in A MySQL Convention A Fanciful Saga of a Hunter And a Pastry Chef who must Vanquish a Boy in Australia A Action-Packed Drama of a Dentist And a Croccolle who must Battle a Feminist in The Canadian Rockies A Emotional Drama of a A Shark And a Database Administrator who must Starquish a Pioneer in Soviet Georgia	Integer	smallint	smallint (numeric (4, 6 6 3 3 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6	,2)
1 2 3 4 5 6 7		1 2 3 3 4 5 6 6 7 7 8 9 10 11 12 13 13 14 15 16	character varying (255) Academy Dinosaur Academy Dinosaur Aca Goldringer Adaptation Holes Affrica Fegulice Affrican Egg Agent Truman Airplane Sierra Airport Politock Alabama Devil Aladdin Calendar Alamo Videotape Alaska Phantom Ali Forever Alice Fartistaia Allen Center	text A Epic Drama of a Feminist And a Mad Scientist who must Battie a Teacher in The Canadian Rockies A Astounding Episte of a Database Administrator And a Explorer who must Find a Car in Ancient China A Astounding Reflection of a Lumberjack And a Car who must Sink a Lumberjack in A Baloon Factory A Fanciful Documentary of a Prisbee And a Lumberjack who must Chase a Monkey in A Shank Tank A Fast-Paced Documentary of a Pastry Chef And a Dentist who must Pursue a Forensic Psychologist in The Gulf of Mexico A Intregicif Panorama of a Robot And a Boy who must Escape a Sumo Westler in Ancient China A Touching Saga of a Hunter And a Butler who must Discover a Butler in A Jet Bost A Epic Tale of a Mosee And a Girl who must Confront a Monkey in Ancient India A Thoughtful Panorama of a Database Administrator And a Mad Scientist who must Outgun a Mad Scientist in A Jet Bost A Action-Packed Tale of a Man And a Lumberjack who must Resch a Feminist in Ancient China A Boring Epistle of a Butler And a Cat who must Fight a Pastry Chef in A MySQL Convention A Fancific Saga of a Hunter And a Pastry Chef who must Sattle a Feminist in The Canadian Rockies A Action-Packed Drama of a Dentist And a Croccolle who must Battle a Feminist in The Canadian Rockies A Emotional Drama of a Cat And a Mad Scientist who must Battle a Feminist in The Canadian Rockies A Brilliant Drama of a Cat And a Mad Scientist who must Battle a Feminist in A MySQL Convention	Integer	smallint	smallint 4	numeric (4, 6 6 3 3 6 6 6 6 6 6 4 4 6 6 5 5	0.99 4.99 2.99 2.99 2.99 4.99 4.99 4.99 0.99 0.99 0.99

Script

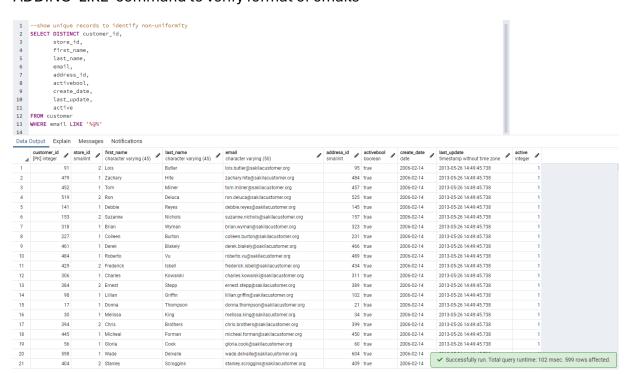
rating, last_update, special_features, fulltext

FROM film

CUSTOMER



ADDING 'LIKE' command to verify format of emails



Output after LIKE statement: Uniform values (same number of records as without LIKE statement)

SCRIPT

FROM customer

WHERE email LIKE '%@%' --adding like command to verify that all unique values have a correct email format

CLEANING NON-UNIFORM DATA

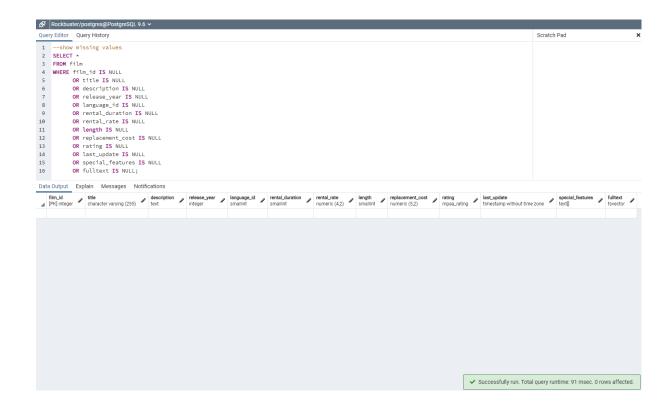
Like in the example in the reading, after identifying all possible values you would have to UPDATE them to only one:

However, you would first need to make sure that the data is inconsistent to begin with using LIKE statements like the one that was used in the script for the customers table. In this table I wanted to make sure that all email addresses were consistent and I tested this by comparing the number of records to the number of records if I applied the LIKE statement searching for the '@' sign.

However, this can only be done with certain columns where the data has categories or structure. For example, if it is numerical data, you can make sure by figuring out if all the data falls within the acceptable range by using "BETWEEN" or "IN" statements. For unstructured data like names and addresses it is more complex.

1.3 NULL VALUES

FILM



Script

--show missing values

SELECT *
FROM film
WHERE film_id IS NULL
OR title IS NULL
OR description IS NULL
OR release_year IS NULL
OR language_id IS NULL
OR rental_duration IS NULL
OR rental_rate IS NULL
OR length IS NULL
OR replacement_cost IS NULL
OR rating IS NULL
OR last_update IS NULL
OR special_features IS NULL
OR fulltext IS NULL;

CUSTOMER



Script

--show missing values

SELECT*

FROM customer

WHERE customer_id IS NULL

OR store_id IS NULL

OR first_name IS NULL

OR last name IS NULL

OR email IS NULL

OR address id IS NULL

OR activebool IS NULL

OR create_date IS NULL

OR last_update IS NULL

OR active IS NULL;

CLEANING MISSING DATA

Cleaning missing data would be the hardest because we would have to corroborate the data at the source, depending on what type of data is missing. If it is in the customer's table for example, and a name is missing but we have everything else, we can look up in older records with its customer id, or with the email or address. If it's something like the language id in the film table, we can corroborate by consulting the film directly.

Depending on the type of information that is missing the input is either manual or in batches.

If it is in batches, it would be like it was illustrated in the reading:

--imputing missing values with the AVG value

```
UPDATE tablename
SET = AVG(col1)
WHERE col1 IS NULL
```

If it were manual and one by one (as with a name, for example), you would have to do like so

--imputing missing value manually

```
UPDATE customer

SET first_name = 'Emma'

SET last_name = 'Ratri'

WHERE customer id = 23;
```

2. SUMMARIZING DATA

2.1 FILM NUMERICAL VALUES

```
SELECT MIN(release_year) AS min_release_year,
           MAX(release_year) AS max_release_year,
 2
 3
           AVG(release_year) AS avg_realease_year,
 4
           COUNT(release_year) AS count_release_year_values,
 5
           MIN(rental_duration) AS min_rental_duration,
 6
           MAX(rental_duration) AS max_rental_duration,
 7
           AVG(rental_duration) AS avg_rental_duration,
           COUNT(rental_duration) AS count_rental_duration_values,
           MIN(rental_rate) AS min_rent,
9
10
           MAX(rental_rate) AS max_rent,
11
           AVG(rental_rate) AS avg_rent,
12
           COUNT(rental_rate) AS count_rent_values,
           MIN(length) AS min_length,
13
14
           MAX(length) AS max_length,
           AVG(length) AS avg_length,
15
16
           COUNT(length) AS count_length_values,
           MIN(replacement_cost) AS min_replacement_cost,
17
           MAX(replacement_cost) AS max_replacement_cost,
18
           AVG(replacement_cost) AS avg_replacenemt_cost,
19
           COUNT(replacement_cost) AS count_replacement_cost_values,
20
21
           COUNT(*) AS count_rows
22
    FROM film;
```

min_release_ year	max_rele ase_year	avg_realea se_year	count_relea se_year_valu es	min_renta l_duration	max_rental_du ration	avg_rent al_durati on
2006	2006	2006.0	1000	3	7	4.985
count_rental_ duration_valu es	min_rent	max_rent	avg_rent	count_ren t_values	min_length	max_len gth
1000	0.99	4.99	2.98	1000	46	185
avg_length	count_le ngth_valu es	min_repla cement_c ost	max_replace ment_cost	avg_repla cenemt_c ost	count_replace ment_cost_val ues	count_ro ws
115.272	1000	9.99	29.99	19.984	1000	1000

NON-NUMERICAL VALUES

```
1 SELECT MODE() WITHIN GROUP (ORDER BY title)
           AS mode_title,
2
3
           MODE() WITHIN GROUP (ORDER BY description)
           AS mode_description,
4
           MODE() WITHIN GROUP (ORDER BY rating)
5
6
           AS mode_rating,
           MODE() WITHIN GROUP (ORDER BY last_update)
7
           AS mode_last_update,
8
           MODE() WITHIN GROUP (ORDER BY special_features)
9
           AS mode_special_features,
10
           MODE() WITHIN GROUP (ORDER BY fulltext)
11
           AS mode_fulltext
12
13 FROM film;
```

mode_title	mode_des cription	mode_rating	mode_last _update	mode_special_ features	mode_fulltext
	A Action-				
	Packed				
	Character				
	Study of a				
	Astronaut				
	And a				
	Explorer				'baloon':19
	who must				'confront':14
	Reach a		2013-05-		'documentari':5
	Monkey in		26	{Trailers,Comm	'feminist':8,11,16
Academy	A MySQL		14:50:58.9	entaries,"Behin	'mile':2 'must':13 'spi':1
Dinosaur	Convention	PG-13	51	d the Scenes"}	'thrill':4

2.2 CUSTOMER

NUMERICAL

```
1 SELECT MIN(active) AS min_active,
2 MAX(active) AS max_active,
3 AVG(active) AS avg_active,
4 COUNT(active) AS count_active,
5 COUNT(*) AS count_rows
6 FROM customer;
```

min_activ	max_acti	ve	avg_active	count_active	count_rows
)	1	0.97495826377295492487	599	599

NOTE: There is one non-active customer, meaning a record is probably flagged as deleted.

NON NUMERICAL

```
SELECT MODE() WITHIN GROUP (ORDER BY first_name)
1
2
          AS mode_firstname,
3
           MODE() WITHIN GROUP (ORDER BY last_name)
4
          AS mode_lastname,
5
           MODE() WITHIN GROUP (ORDER BY email)
           AS mode_email,
6
7
           MODE() WITHIN GROUP (ORDER BY activebool)
           AS mode_activebool,
8
9
           MODE() WITHIN GROUP (ORDER BY create_date)
10
          AS mode_create_date,
          MODE() WITHIN GROUP (ORDER BY last_update)
11
          AS mode_last_update
12
13 FROM customer;
```

mode_first name	mode_last name	mode email	mode_activ ebool	mode_create date	mode_last_u pdate
		aaron.selby@sakilacus			2013-05-26
Jamie	Abney	tomer.org	TRUE	14/02/2006	14:49:45.738

3. 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.

I believe SQL is a faster and more convenient way of quickly summarizing the data. However, SQL requires that the analyst knows a priori what they want to know—They need to know what to ask. You cannot see the dataset as it comes, so there is no way to make empirical observations before you start programming other than what you may observe at random by running SELECT * FROM table.

This lack of immediate interaction with the data can create a certain bias for whatever the results of this first summarization are.

In brief, SQL is faster and more convenient, but you need to spend more time getting acquainted with the dataset and cleaning it to make sure you are not getting bad results. There is less of an opportunity to randomly catch a mistake without having to troubleshoot the whole process.

Excel, on the other hand, is very interactive and one of its greatest advantages is that you can already visualize some of the summarization with the app's graphs and pivot tables.

However, it takes considerably more time to arrive at the same results. You still spend a lot of time cleaning data as well.

In this sense, both programs are useful for certain purposes, but I think that once you overcome the learning curve, SQL is more convenient.