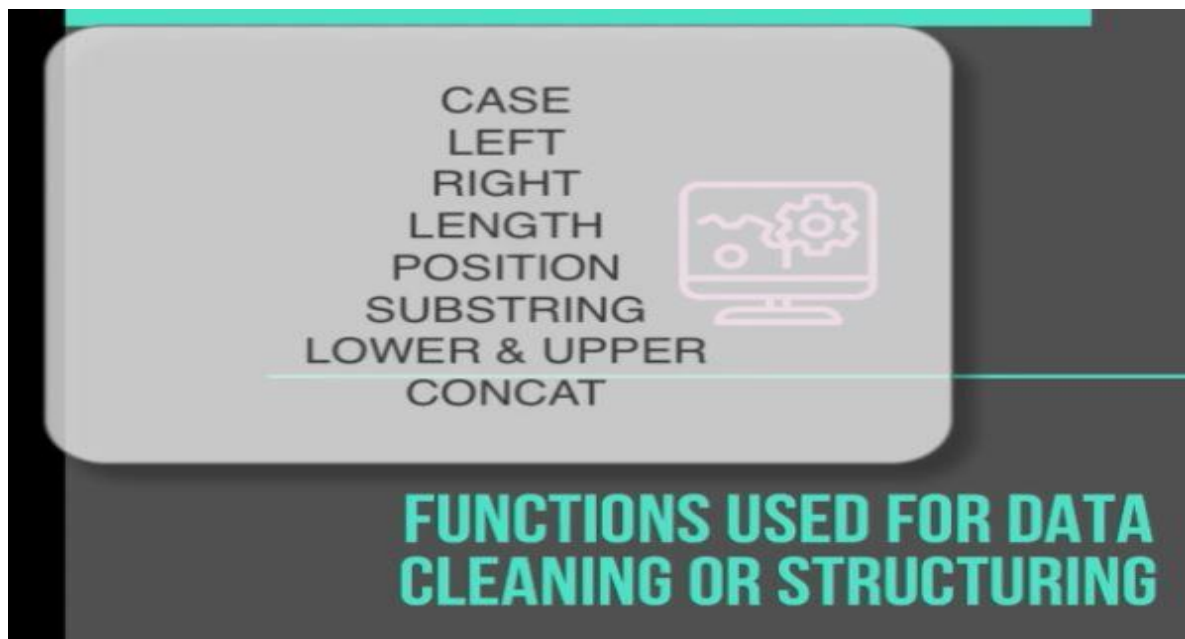


## DATA CLEANING AND STRUCTURING



### 1. CASE Statement

- The **CASE statement** is SQL's way of handling if/then logic.
- The **CASE statement** is followed by at least one pair of **WHEN and THEN** statements- SQL's equivalent of IF/THEN in Excel.
- Every **CASE statement** must end with the **END** statement.
- **ELSE** statement is optional, and captures values not specified in the **WHEN/THEN** statements.

#### General Form:

```
CASE WHEN (condition) THEN (result if condition is met)
      ELSE (result if previous condition not met)
      END AS new_attribute_name
```

#### Example 1:

```
Query Editor  Query History
1  /* Create a new column that categorizes customers based on the total
2     amount they have paid: BRONZE - $0-$70, BRONZE - $0-$70 AND
3     GOLD - greater than $150. Call this "customer_category" and sort
4     the results by the total amount paid. */
5
6  SELECT customer_id, SUM(amount),
7         CASE WHEN SUM(amount) BETWEEN 0 AND 70 THEN 'BRONZE'
8             WHEN SUM(amount) BETWEEN 71 AND 150 THEN 'SILVER'
9             ELSE 'GOLD'
10        END AS customer_category
11  FROM payment
12  GROUP BY 1
13  ORDER BY 2;
```

### Example 2:

Query Editor Query History

```
1  /* Count the inventory of each film and then categorizes each film as
2     'high' inventory if it's inventory is greater than or equal to 5
3     'low' inventory otherwise */
4
5  SELECT film_id, COUNT(inventory_id),
6         CASE WHEN COUNT(inventory_id) >= 5 THEN 'HIGH'
7              ELSE 'LOW'
8         END AS inventory_status
9  FROM inventory
10 GROUP BY 1;
```

Data Output Explain Messages Notifications

	film_id smallint	count bigint	inventory_status text
1	652	4	LOW
2	273	7	HIGH

### Example 3:

Query Editor Query History

```
1  /* Count the number of rentals for each film and categorize films
2     based on number of rentals: 'low' (less than 10 rentals). 'medium'
3     (between 10 to 20 rentals) and 'high' (greater than 20 rentals). */
4
5  SELECT f.film_id, title, COUNT(rental_id),
6         CASE WHEN COUNT(rental_id) < 10 THEN 'LOW'
7              WHEN COUNT(rental_id) BETWEEN 10 AND 20 THEN 'MEDIUM'
8              ELSE 'HIGH'
9         END AS demand_status
10 FROM film f
11 JOIN inventory i
12 ON f.film_id=i.film_id
13 JOIN rental r
14 ON i.inventory_id=r.inventory_id
15 GROUP BY 1;
```

Data Output Explain Messages Notifications

	film_id [PK] integer	title character varying (255)	count bigint	demand_status text
1	652	Pajama Jawbreaker	14	MEDIUM
2	273	Effect Gladiator	25	HIGH
3	51	Balloon Homeward	23	HIGH

## 2. LEFT, RIGHT & LENGTH Statement

- **LEFT** pulls a specified number of characters for each row in a column starting at the beginning (or from the left). **LEFT (string, number)**
- **RIGHT** pulls a specified number of characters for each row in a column starting at the end (or from the right). **RIGHT (string, number)**
- **LENGTH** provides the number of characters for each row of a specified column. **LENGTH (string)**

### Example 1:

Query Editor	Query History												
<pre>1  /* Extract each customer's local phone number including the area code 2     (the last 10 digits in the given phone number),e.g. 14033335568 3     Hint : Find the phone attribute in address table. */ 4 5  SELECT RIGHT(phone,10) AS local_phone_number 6  FROM address; 7 8  -- Result was changed from 14033335568 to 4033335568 in no.3</pre>													
Data Output	Explain Messages Notifications												
<table><thead><tr><th></th><th>local_phone_number text</th></tr></thead><tbody><tr><td>1</td><td></td></tr><tr><td>2</td><td></td></tr><tr><td>3</td><td>4033335568</td></tr><tr><td>4</td><td>6172235589</td></tr><tr><td>5</td><td>8303384290</td></tr></tbody></table>		local_phone_number text	1		2		3	4033335568	4	6172235589	5	8303384290	
	local_phone_number text												
1													
2													
3	4033335568												
4	6172235589												
5	8303384290												

### Example 2:

Query Editor

Query History

```
1  /* Extract each customer's country code from the given phone number.
2     Hint : The country code is either 1-digit or 2-digits If a phone
3     number does not include a country code, then write 'NULL'. */
4
5  SELECT phone, CASE WHEN LENGTH(phone) = 11 THEN LEFT(phone, 1)
6                     WHEN LENGTH(phone) = 12 THEN LEFT(phone, 2)
7                     ELSE 'NULL'
8                     END AS country_code
9  FROM address;
```

Data Output

Explain

Messages

Notifications

	phone character varying (20)	country_code text
1		NULL
2		NULL
3	14033335568	1
4	6172235589	NULL
5	28303384290	2
6	838635286649	83

### Example 3:

Query Editor		Query History
<pre>1  /* How many customers have first names starting with a vowel(a,e,i,o,u) 2  */ 3  SELECT SUM(firstname_vowel) 4  FROM (SELECT customer_id, 5           CASE WHEN LEFT(first_name, 1) IN ('A','E','I','O','U') THEN 1 6           ELSE 0 7           END AS firstname_vowel 8           FROM customer) sub1; 9 10 --Since the values under firstname_vowel are either 1 or 0, so summed up.</pre>		
Data Output		Explain Messages Notifications
	sum bigint	
1	83	

### 3. POSITION, SUBSTRING, UPPER & LOWER Statement

- **POSITION** provides the position of a string starting from the left. In other words, it provides an index where the character is for each row. Note that the index of the first position is 1 in SQL.

General form: POSITION ('string' IN attribute)

- Since the **POSITION** is case-sensitive, you may use either **LOWER** or **UPPER** to make all the characters in either **lowercase** or **uppercase**.

General form: LOWER (attribute)

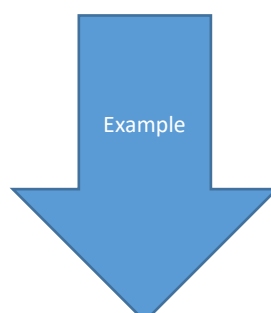
UPPER (attribute)

- **SUBSTRING** is used to extract a string containing a specific number of characters from a particular position of a given string.

General form: SUBSTRING (attribute, str\_pos, ext\_char)

Where,

- str\_pos – (optional) this is the position of the string where extracting starts. If omitted, the substring function will start at position 1.
- ext\_char – (optional) specifies the number of characters to be extracted from the starting position. If omitted, the substring function will extract until the last character of the string.



**Example 1:** The email attribute under the customer table is formatted as:

[firstname.lastname@sakilacustomer.org](mailto:firstname.lastname@sakilacustomer.org)



Create a new column showing the firstname.lastname only and call this new column email\_name.

```
1 SELECT LEFT(email, POSITION('@' IN email))
2 FROM customer;
```

Data Output	Explain	Messages	Notifications
 <b>left</b> text			
1	jared.ely@		
2	mary.smith@		

OR

```
1 SELECT LEFT(email, POSITION('@' IN email)-1)
2 FROM customer;
```


Data Output	Explain	Messages	Notifications
 <b>left</b> text			
1	jared.ely		
2	mary.smith		

**Example 2:** The email attribute under the customer table is formatted as:

[firstname.lastname@sakilacustomer.org](mailto:firstname.lastname@sakilacustomer.org)

Create a new column showing the last name only and call this new column email\_lastname.

```
SELECT SUBSTRING(email, POSITION('.') IN email)+1,
           POSITION('@' IN email) - POSITION('.') IN email)-1)
AS email_lastname
FROM customer;
```

Output	Explain	Messages	Notifications
<b>email_lastname</b> text			
ely			
smith			



**Example 3:** The email attribute under the customer table is formatted as:

[firstname.lastname@sakilacustomer.org](mailto:firstname.lastname@sakilacustomer.org)

Create a new column showing the first name only and call this new column email\_firstname.

```
1 SELECT LEFT(email, POSITION('.') IN email)-1)
2           AS email_firstname
3 FROM customer;
```

	email_firstname	
	text	
1	jared	
2	mary	

**Example 4:** The film's title attribute under the film table is formatted as:

["Chamber Italian"](#)

Grab the second word in the film's title (in lowercase) and call this new attribute film\_phrase.

```
1 SELECT RIGHT(lower(title), LENGTH(title)-POSITION(' ' IN title))
2           AS titile_phrase
3 FROM film;
```

	titile_phrase	
	text	
1	italian	
2	wonderful	

#### 4. CONCAT Statement

**CONCATE** combines values from several columns into one column.

General form: `CONCATE (column1, column2)`

Note: You may also insert a character or space in between the columns.

Instead of `CONCAT`, you may use `||` (pipes), e.g. `column1 || column2`

**Example 1:** Combine the first\_name and last\_name columns into one column full\_name. Make sure to put a space between the first name and last name.

```
1 SELECT CONCAT(first_name, ' ', last_name) AS full_name
2 FROM customer;
```

	full_name	
	text	
1	Jared Ely	
2	Mary Smith	

**Example 2:** Create a new version of email addresses (name it email\_address). It will be of the following format: [f.lastname@sakila.org](mailto:f.lastname@sakila.org)  
 'f' stands for the first letter of the first name. All characters should be in lowercase.

```
1 SELECT CONCAT(LOWER(LEFT(first_name,1)),',',LOWER(last_name),'@sakila.org')
2         AS email_address
3 FROM customer;
```

Data Output Explain Messages Notifications

	email_address text	
1	j.ely@sakila.org	
2	m.smith@sakila.org	

**Example 3:** Create a column showing the customer's district and country. Make sure to put a comma and a space in between.

For example: Alberta, Canada

```
1 SELECT district, country, CONCAT(district,',',' ',country)
2         AS dist_country
3 FROM address a
4 JOIN city ci
5 ON a.city_id = ci.city_id
6 JOIN country co
7 ON ci.country_id = co.country_id;
```

Data Output Explain Messages Notifications

	district character varying (20)	country character varying (50)	dist_country text	
1	Alberta	Canada	Alberta, Canada	
2	QLD	Australia	QLD, Australia	