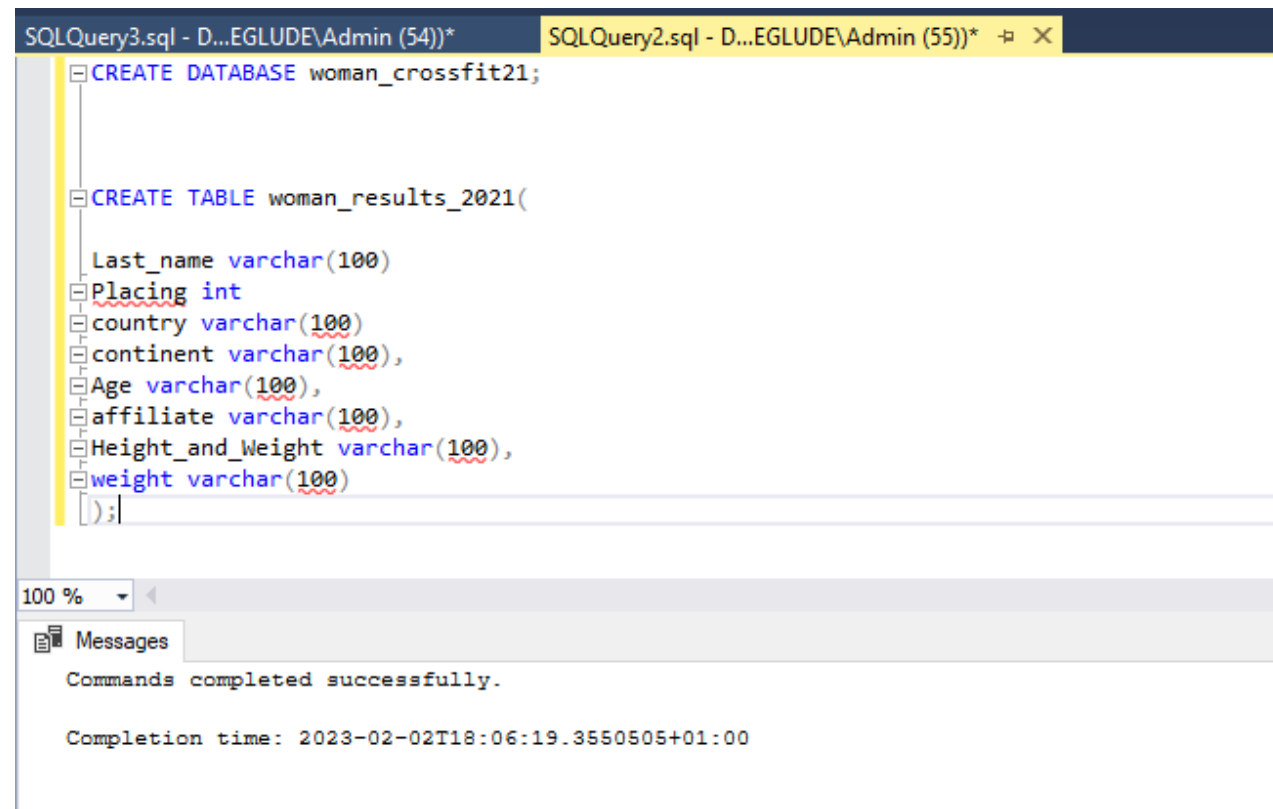# Cleaning Data in SQL:
# Women Ranking in Crossfit Games 2021

In order to present this cleaning up job, I am going to use a dataset from the Crossfit games 2021, where we can find the women's ranking details.

- First, let's create the Database and Table

# • **Import the data from .CSV to SQL Server**

## 1) Running a python script

# • **Import the data from .CSV to SQL Server**

2) Running a SSIS package

# • Verify the records loaded & start to clean the data

# Missing Values



- We have some missing records, for example in the "affiliate" column
- But we I count the null records, it gave me 0. So, it is not considering as a NULL.

- To avoid some troubles in ''JOIN'' statements for example, we can fill the blank (in this case) or null records with some useful description. In this case, I use UPDATE to set new record values:



```
QLQuery5.sql - D...EGLUDE\Admin (53))*  ⊣ ✕  SQLQue
⊟UPDATE woman_results_2021

  SET affiliate = 'No data'
  WHERE affiliate = ' ';

⊟UPDATE woman_results_2021
  SET Height_and_Weight = 'No data'
  WHERE Height_and_Weight = ' ';
```

- The outputs:



```
SQLQuery5.sql - D...EGLUDE\Admin (53))*        SQLQuery4.sql - D...EGLUDE\Admin (54))*  ⊣ ✕  SQLQuery2.sql - D...EGLUDE\Adm
⊟SELECT
  SUM(CASE WHEN affiliate= 'No data' THEN 1 ELSE 0 END) AS Affiliate_Replaced
     , COUNT(affiliate) AS Affiliate_Non_Null_Values,
  SUM(CASE WHEN Height_and_Weight= 'No data' THEN 1 ELSE 0 END) AS Height_and_Weight_Replaced
     , COUNT(Height_and_Weight) AS Height_and_Weight_Non_Null_Values
     FROM woman_results_2021;
```

100 %

⊞ Results  📄 Messages

| | Affiliate_Replaced | Affiliate_Non_Null_Values | Height_and_Weight_Replaced | Height_and_Weight_Non_Null_Values |
|---|---|---|---|---|
| 1 | 15074 | 108688 | 58281 | 108688 |

# Now that we deal with missing data, we can review the content

- First, I going to work in the "Age" column because I am only interesting in the number of years. I'm going to use RIGHT statement to keep only the number and re-type the column into int.

- The output:

# In "Height_and_Weight" column we have to work a bit deeper.

- First, let's create the new columns called 'height' and 'weight'. Previously, I renamed the original column 'weight' as 'lbs_lifted'.

- Then, we want to split the 'Height_and_Weight' column in two. I worked out the splitting using *CHARINDEX* to map the "|" character and then used *SUBSTRING* to keep the data from the left side into height column and the data from the right into weight column.

- Finally, drop the old column.

- The Output:

**1**

```
SQLQuery7.sql - D...EGLUDE\Admin (54))*        SQLQuery6.sql - D...EGLUDE\Admin (69))*    ⊹ ✕    SQLQuery5.sql - D...EGLUDE\Admin (
UPDATE woman_results_2021
  set height=CASE WHEN CHARINDEX('|', Height_and_Weight) > 0 THEN SUBSTRING(Height_and_Weight, 1,
CHARINDEX('|', Height_and_Weight) - 1)
              ELSE 'No data' END;

UPDATE woman_results_2021
  SET weight =  CASE WHEN CHARINDEX('|', Height_and_Weight) > 0 THEN SUBSTRING(Height_and_Weight,
CHARINDEX('|', Height_and_Weight) + 1, LEN(Height_and_Weight))
              ELSE 'No data' END;
```
100 %

Messages

```
SELECT height, weight FROM woman_results_2021
```
00 %

**2**

Results | Messages

| | height | weight |
|---|---|---|
| 1 | 163 cm | 58 kg |
| 2 | 64 in | 140 lb |
| 3 | 162 cm | 158 lb |
| 4 | 62 in | 132 lb |
| 5 | 63 in | 145 lb |
| 6 | 165 cm | 64 kg |
| 7 | 66 in | 168 lb |
| 8 | No data | No data |
| 9 | 66 in | 150 lb |
| 10 | 67 in | 140 lb |
| 11 | 174 cm | 68 kg |
| 12 | 170 cm | 75 kg |
| 13 | 170 cm | 150 lb |
| 14 | 169 cm | 150 lb |

**3**

```
ALTER TABLE woman_results_2021
  DROP COLUMN Height_and_Weight;
```

# "Lbs_lifted" column

- Now, let's move on into the "lbs_lifted" column and just maintain the actual weight lifted. To do that, I used the function *SUBSTRING*.

- The syntax is:
  SUBSTRING(*string*, *starting character position*, *# of characters*)

- When I was trying to change the datatype into *int,* I noticed that there were some spaces in column. To remove them, use the *TRIM (or RTRIM, LTRIM)* function.



```
UPDATE woman_results_2021

SET lbs_lifted = SUBSTRING(lbs_lifted,16,3);

SELECT lbs_lifted  FROM woman_results_2021;
```

| | lbs_lifted |
|---|---|
| 1 | 230 |
| 2 | 218 |
| 3 | 232 |
| 4 | 217 |
| 5 | 217 |
| 6 | 224 |
| 7 | 236 |
| 8 | 216 |
| 9 | 211 |
| 10 | 202 |
| 11 | 203 |
| 12 | 221 |
| 13 | 210 |
| 14 | 206 |

```
UPDATE woman_results_2021
SET lbs_lifted = REPLACE(LTRIM(RTRIM(lbs_lifted)), ' ', '')
```

# Capturing Insights

- Now, we can querying to discover some insights as your requests.

- For this project, I will connect the SQL Server Database with Power BI in order to extract the data and gain insights. After analyzing the data, I will create visualizations to present in the final report.

```sql
SELECT country,
    Age,
    count(Age) AS Participants_per_Age_and_Country
    FROM woman_results_2021
GROUP BY country, Age
ORDER BY country
```

Results   Messages

| country | Age | Participants_per_Age_and_Country |
|---------|-----|----------------------------------|
| Albania | 39 | 1 |
| Albania | 36 | 2 |
| Albania | 33 | 1 |
| Albania | 28 | 2 |
| Albania | 46 | 1 |
| Albania | 29 | 2 |
| Albania | 23 | 1 |
| Algeria | 38 | 1 |
| Algeria | 24 | 1 |
| Algeria | 32 | 1 |
| Algeria | 20 | 1 |
| Algeria | 45 | 1 |
| Algeria | 25 | 1 |
| Andorra | 42 | 2 |
| Andorra | 33 | 1 |