

*# MODULE 1 : DATA EXPLORATION*



# DATA-MANIPULATION-1

Wrangling data to extract meaningful insights



Club Informatique & Télécom  
Data Cell

# Data Manipulation - 1

- What is data wrangling ?
- Pandas data structures
- Viewing / inspecting data
- Selecting data
- Cleaning data
- Today's lab : Olympic Data

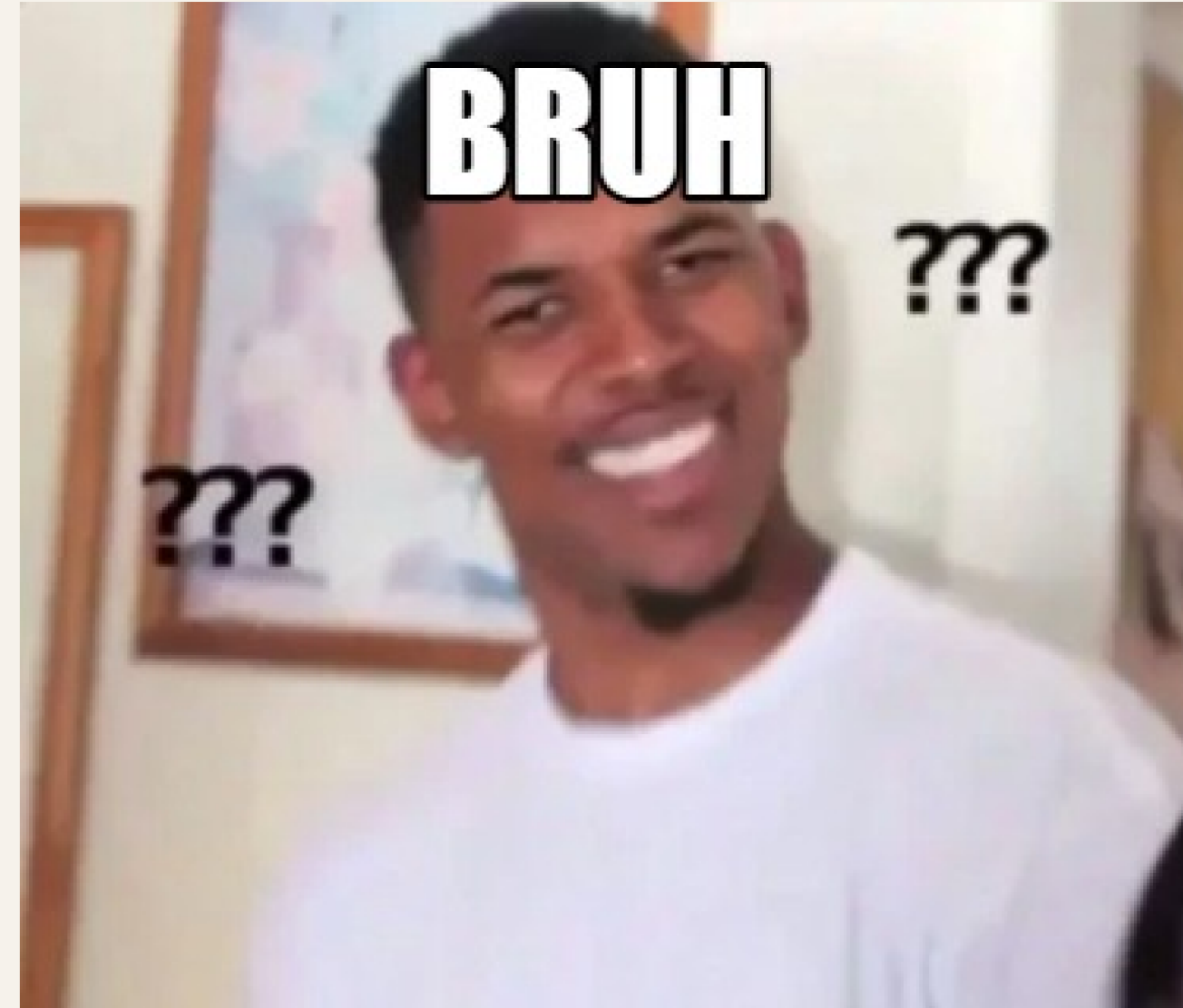


# Data Manipulation

Data wrangling is the process of transforming and structuring data from one raw form into a desired format with the intent of improving data quality and making it more consumable and useful for analytics or machine learning. It's also sometimes called data munging.

[alteryx.com](https://www.alteryx.com)





# Data Wrangling



# Pandas Library



# Pandas Data Structures

pd.Series

	apples
0	3
1	2
2	0
3	1

+

pd.Series

	oranges
0	0
1	3
2	7
3	2

=

pd.DataFrame

	apples	oranges
0	3	0
1	2	3
2	0	7
3	1	2



# Viewing / inspecting data

<code>df.head(n)</code>	First n rows of the DataFrame
<code>df.tail(n)</code>	Last n rows of the DataFrame
<code>df.shape</code>	Number of rows and columns
<code>df.info()</code>	Index, Datatype and Memory information
<code>df.describe()</code>	Summary statistics for numerical columns
<code>s.value_counts(dropna=False)</code>	View unique values and counts
<code>df.apply(pd.Series.value_counts)</code>	Unique values and counts for all columns





# Selecting data

<code>df[col]</code>	Returns column with label col as Series
<code>df[[col1, col2]]</code>	Returns columns as a new DataFrame
<code>s.iloc[0]</code>	Selection by position
<code>s.loc['index_one']</code>	Selection by index
<code>df.iloc[0,:]</code>	First row
<code>df.iloc[0,0]</code>	First element of first column



# Data Cleaning

<code>pd.isnull()</code>	Checks for null Values, Returns Boolean Array
<code>df.dropna()</code> / <code>df.dropna(axis=1)</code>	Drop all rows / columns that contain null values
<code>df.fillna(x)</code>	Replace all null values with x
<code>s.fillna(s.mean())</code>	Replace all null values with the mean
<code>s.astype(float)</code>	Convert the datatype of the series to float
<code>s.replace([2,3],['two', 'three'])</code>	Replace all 2 with 'two' and 3 with 'three'
<code>df.rename(columns={'old': 'new'})</code>	Selective renaming
<code>df.set_index('column_one')</code>	Change the index



# Time to practice !

