

For each attribute we can have different level of measurement

A **VARIABLE** is a sort of **OPERATIONALIZED** version of the attribute

An **ATTRIBUTE** is just an intuitive general concept  
The **VARIABLE** also defines the set of all possible values that it can take

The **OPERATIONALIZATION** is the process of defining the values which a variable can take

### LEVEL OF MEASUREMENT

4 LEVELS.

- NOMINAL
- ORDINAL

Qualitative variables  
(Tall, short)

- INTERVAL
- RATIO

Quantitative variables  
(173 cm)

## QUALITATIVE ( CATEGORICAL )

### • NOMINAL

When we can't establish a relation of order between variables we can say that the scale of measurement is nominal.

### • ORDINAL

derive from order.

is when we are observing variable that can be ordered.

When we can establish a relation of order between values we say that the scale of measurement is ordinal.

$$V_1 < V_2 \leq V_3$$

ex: short < tall

The order is defined as a particular relation that has defined the property of reflexivity ( $a \leq a$ )

antisymmetry (if  $a \leq b$  and  $b \leq a$  then  $a = b$ )

transitivity (if  $a \leq b$  and  $b \leq c$  then  $a \leq c$ )

## QUANTITATIVE (NUMERICAL)

not all numerical variable are quantitative variable because not all number derive from a measurement

for example ID, credit card number are qualitative variables.

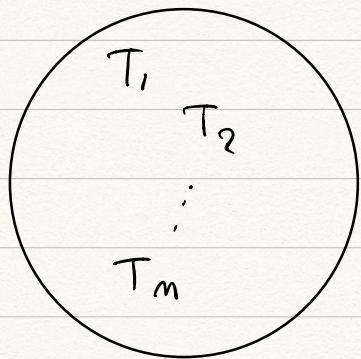
### • INTERVAL

The ratio does not make sense, but is important  
the differences between value

### • RATIO

The ratio between possible variables can make sense

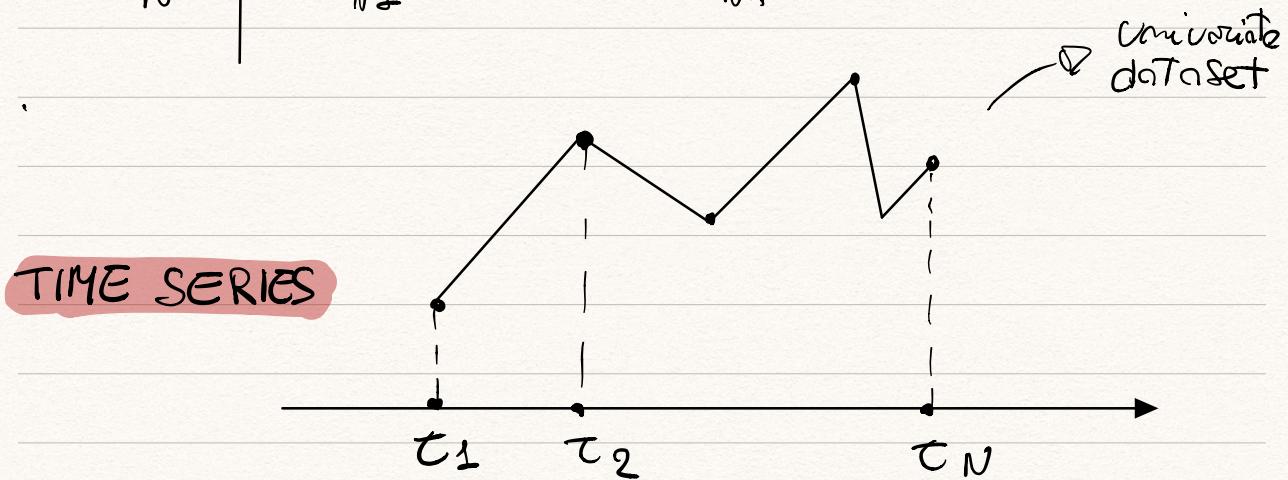
# TIME SERIES ANALYSIS



Instead of a set of items we have a set of **INSTANTS OF TIME**

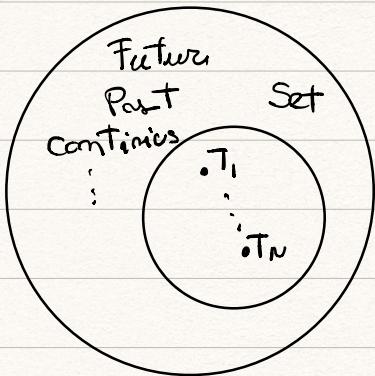
For each instant of Time we have a set of observation (DATA POINT) which would be a the observations of some variable  $(x_1, \dots, x_k)$  at an instant

time	$x_1$	...	$x_k$
$t_1$	$x_{11}$	...	$x_{1k}$
:	:		:
$t_N$	$x_{N1}$	...	$x_{Nk}$



A Time Series can be constituted by ~~least~~ one variable or multiple values.

We can order instances of time. So the first thing that differentiate the TIME SERIES SETUP from other setup is that the statistical units have defined a VARIATION OF ORDER



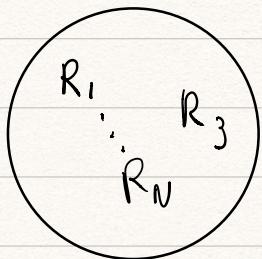
We can consider our set as a Sample

( $\Rightarrow$  INFERENTIAL nature intrinsically)

It is important of a financial instrument

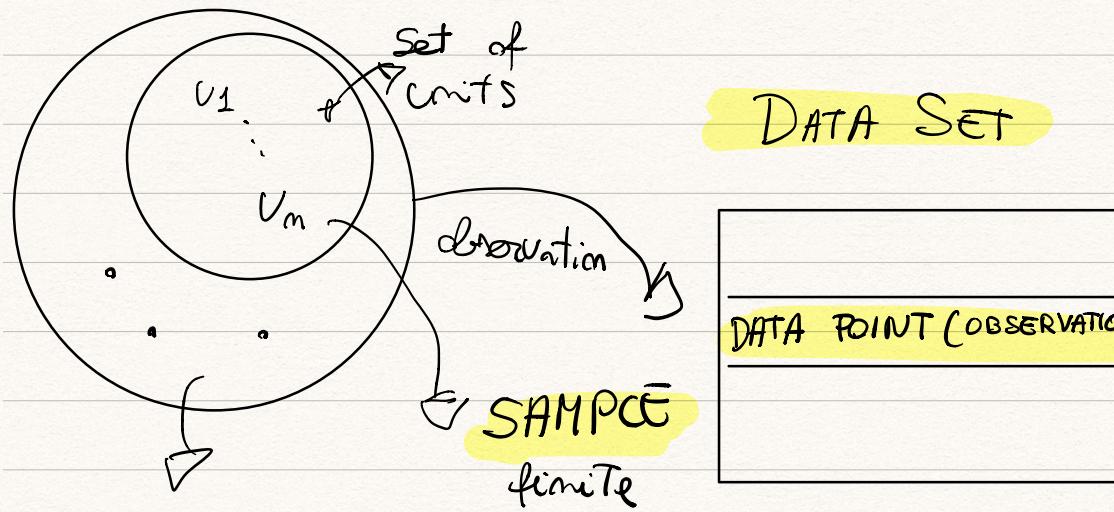
# SPATIAL DATA ANALYSIS

In Spatial Data analysis we have Region of space as statistical units



The regions of space haven't a material order, but we can do assumption to introduce some order.

We could order for longitude, latitude etc...



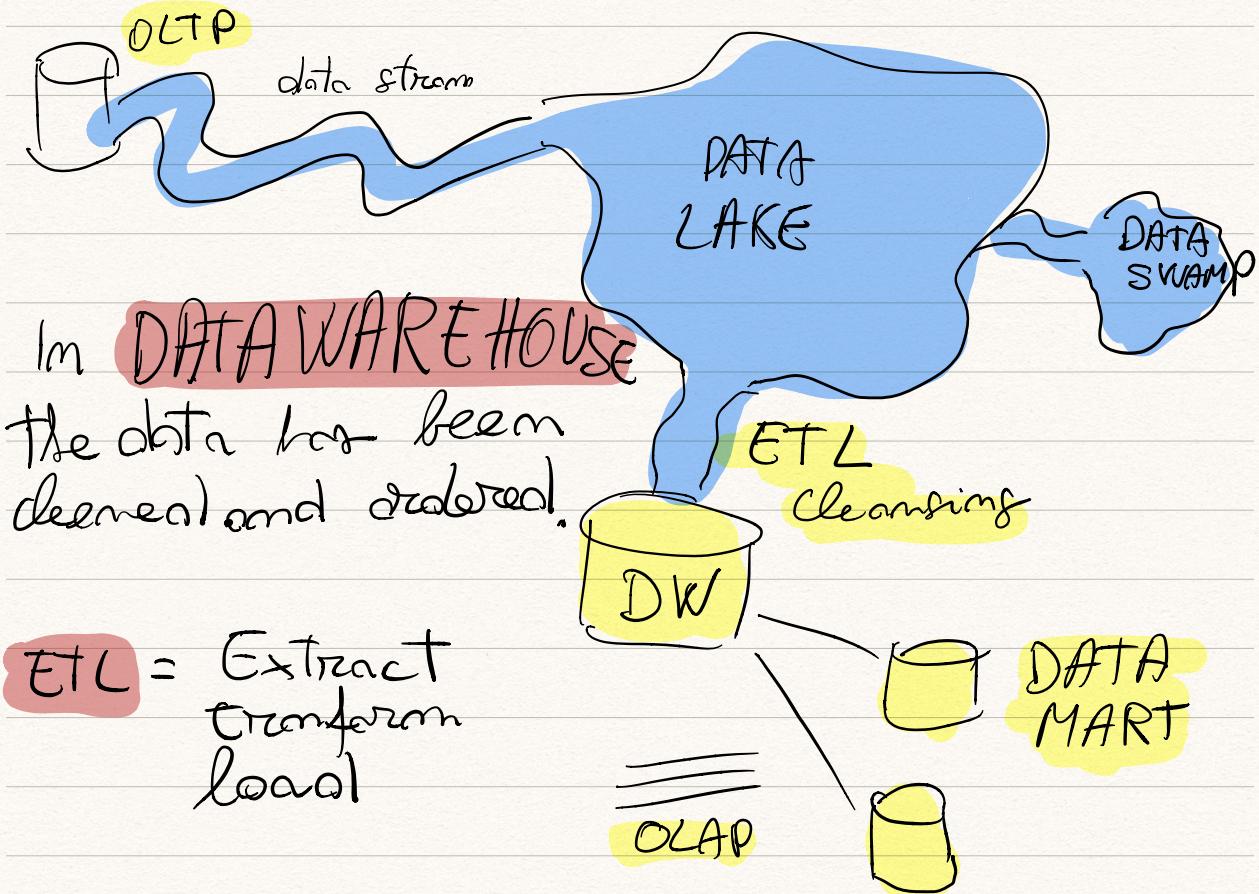
**POPULATION**  
finite or infinite

In Real world we consider a **STREAM** of **DATA** of set of units of observation.



even if i have already collected the units ; always think as a stream establishing an order .

The data stream are stored in  
DATA LAKE or Data WAREHOUSE



ETL = Extract  
transform  
load

DATA MART are special data warehouse  
specialized

OLAP = Online Analysis Processing  
are operations of reporting, analysis and  
mining done to create DATA MART

DATA SWAMP = Some data more confused  
and less useful.

**OLTP** = Online Transactional Processing  
Operative type of system where can collect data

## Process DATA STREAM

		1) Store the data
Ex 1	30	2) Data Processing
Ex 2	28	
Ex 3	20	$\begin{array}{ c } \hline 39 \\ 28 \\ \vdots \\ 25 \\ \hline \end{array}$ $\frac{30 + 28 + 25}{\text{Count}}$

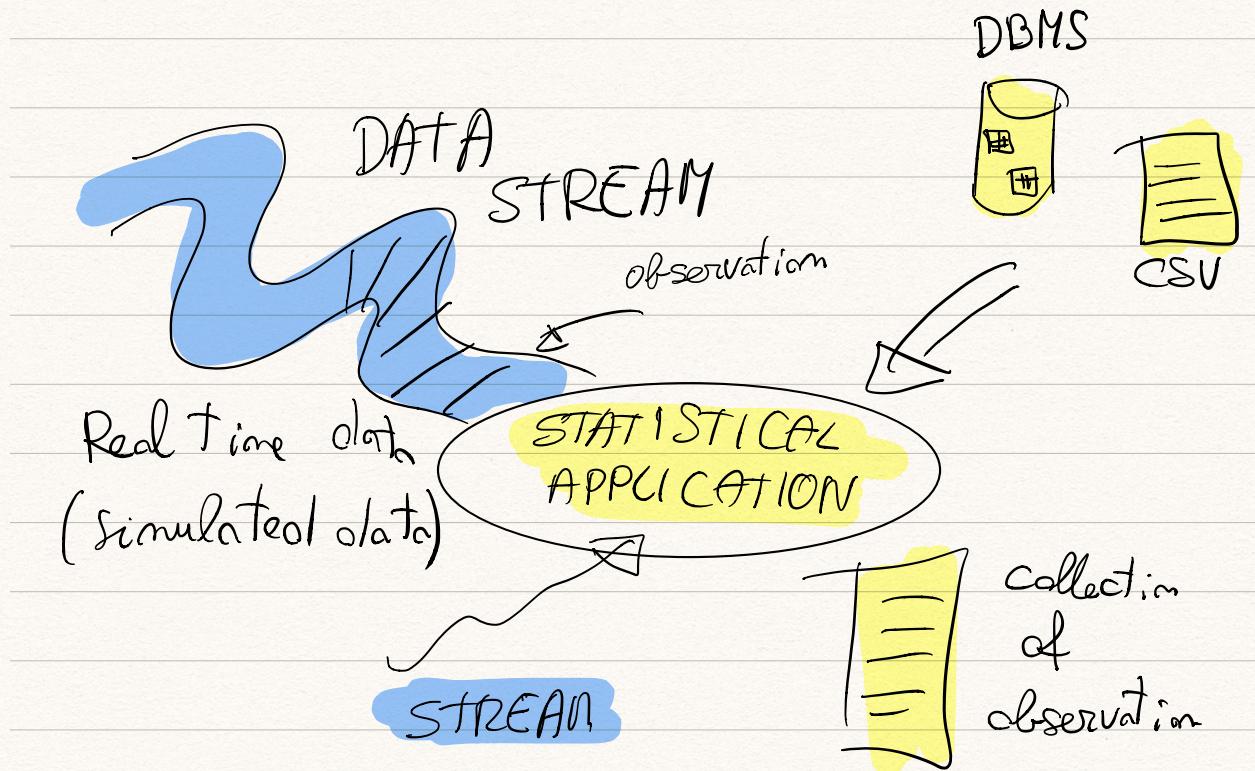
OR

## STREAM PROCESSING

- 1) RECEIVE DATA
- 2) PROCESS AND UPDATE THE RESULT

$$\text{Avg} = \frac{\bar{x} + x_m \cdot m}{m+1} \quad \bar{x} = \text{last data received}$$

If we have an infinite stream of data we can't store all data and then process it, we need a streaming algorithm to do streaming processing



An application can work on data stream or within an external collection of observation

- We can store all dataset in memory → **BATCH PROCESSING**  
(offline)
- We reserve a limit amount of data received from a stream → **STREAM PROCESSING**  
(online)

When we read only once the data we say that the algorithm is **ONE PASS**, otherwise the algorithm is **MULTI PASS**

↓  
single pass

This distinction is applicable to both Batch processing and stream processing

## FOUNDATIONAL CONCEPTS

- **RANDOM OBJECT**  
Generator of pseudo random numbers

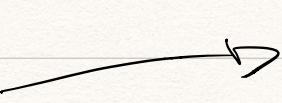
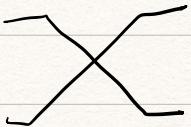
**COLLECTION of object**  
↳ collection of types

- **TIMER**  
Permit to generate event spaces by a deterministic interval of time

- Array
- List
- Dictionary
- Queue
- Stack
- Sorted list

:

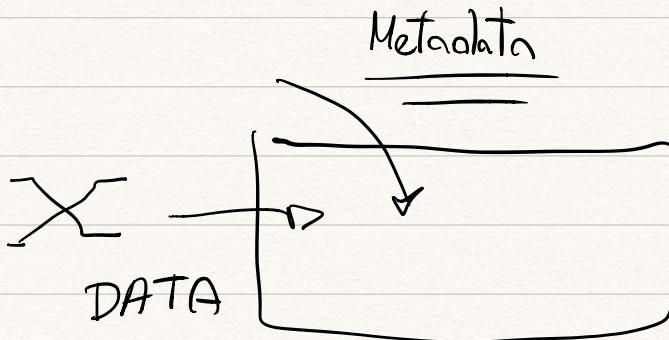
Summarizing



The whole dataset into a unique "representative" value

for example in the case of an arithmetic mean we store the sum of the units as representative value.

**METADATA** is about the meaning of the DATA and they are also important.



Metadata will change the interpretation of the data, so it is important to store it.

## PSEUDOCODE

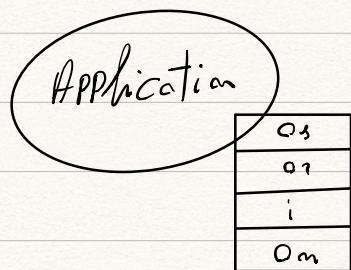
IS a plain language description of how the algorithm is going to work

We need to define a **CLASS** that represent the statistical units and within the class we can define the **PROPERTIES** that define the variable

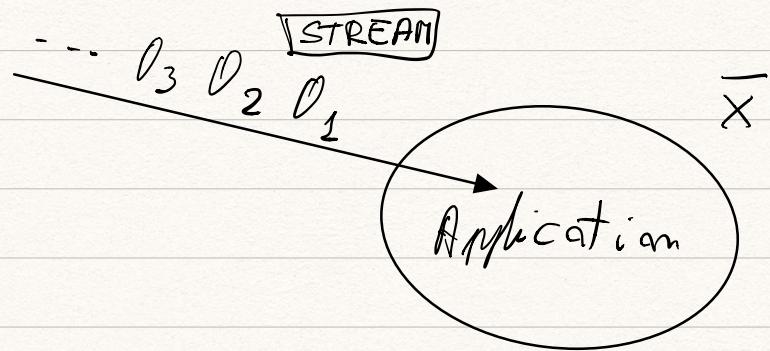
↳ the instance of the class will be the **DATA POINT**

The **DATA** is just the content of the object  
The **METADATA** is represented by the **CLASS** itself.

In the **BATCH** processing the application will hold all the set of object



In STREAM Processing we feed The application with a continuous stream of object



$$\bar{x}_i = \bar{x}_{i-1} + \text{update Term}$$

Online algorithm (On The fly)