

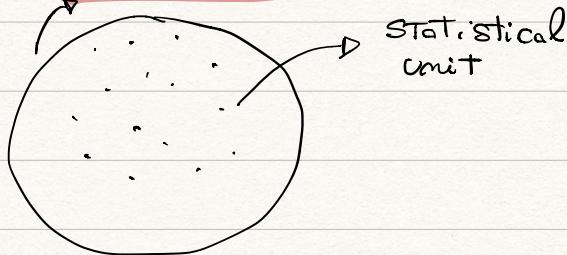
STATISTICS

IS A MATHEMATICAL SCIENCE PERTAINING TO THE COLLECTION, ANALYSIS, INTERPRETATION OR EXPLANATION, AND PRESENTATION OF DATA

DATASET is a sort of a collection of data and also the object of statistical analysis, representation, etc...

We fix a set of object that we want to describe

(**POPULATION**)



We can arrange the population in a list

U_1
 U_2
⋮
 U_i
⋮
 U_m

we can observe different attribute of each subject

Weight Height Age ... Sex

n_1 n_2 n_3 ... n_k

n_{i1} represent the attribute n_1 of the unit ;

m = number of units in the population

Set of unit + Set of attribute = DATASET

DATA SET X

| STATISTICAL UNITS | ATTRIBUTES | | | | |
|-------------------|------------|----------|----------|----------|----------|
| | X_1 | \dots | X_j | \dots | X_k |
| U_1 | x_{11} | \dots | x_{1j} | \dots | x_{1k} |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots |
| U_i | x_{i1} | \dots | x_{ij} | \dots | x_{ik} |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots |
| U_m | x_{m1} | \dots | x_{mj} | \dots | x_{mk} |

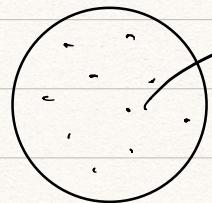
This can be considered as a matrix in a mathematical point of view

Can be seen as a table in an informatic point of view

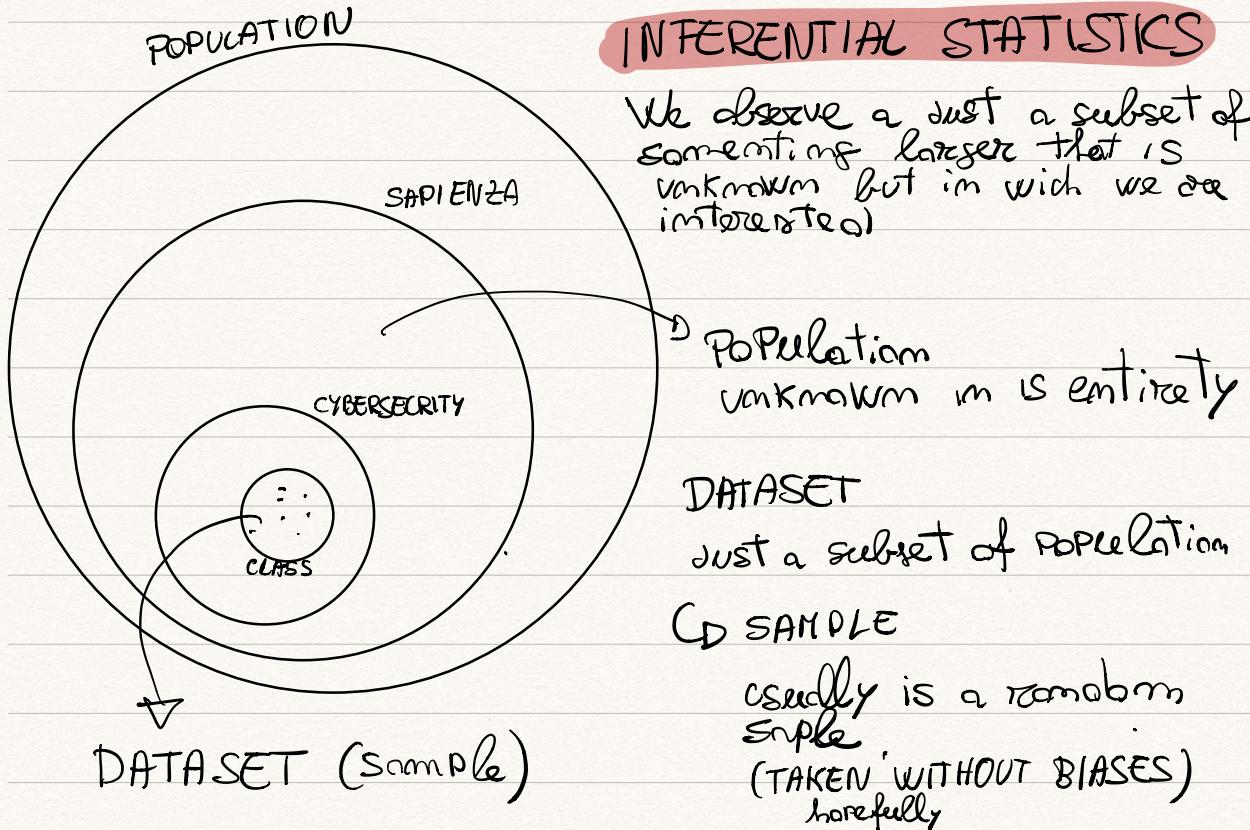
Population with some attribute so we have a dataset
this is the easiest kind of setup and is called

DESCRIPTIVE STATISTICS

Because it is just a description of my population.



D DATASET = POPULATION



This distinction can be applied to both inferential and descriptive statistics:

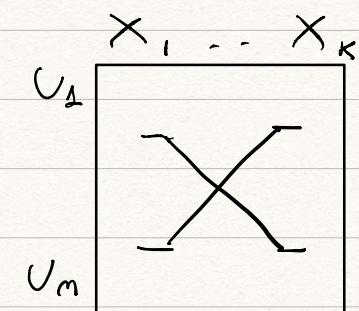
UNIVARIATE

The distinction is on the **structure of the dataset**

MULTIVARIATE

NORMALLY

x_{ij} com $j = 1 \dots m$ e $i = 1 \dots k$



When we consider
only 1 attribute

$\Rightarrow K=1$ we
talk about

UNIVARIATE
STATISTICS

Usually techniques developed
for univariate case can be
easily extended to multivariate
case using matrix mathematical
tools.

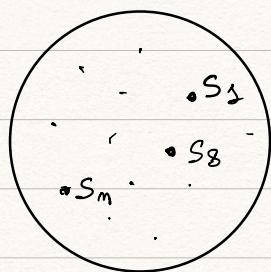
If we talk about
multiple attributes
 $K > 1$ we talk
about

MULTIVARIATE
STATISTICS

\hookrightarrow we can have
bivariate
trivariate
etc ...

EXAMPLE

POPULATION = class



1 attribute

$x_1 = \text{Avg grade}$

We have a univariate
DATASET

50
Students

↙
So we have
just a list
or a vector

DATASET

| | |
|-------------|------|
| Student 1 | 28.1 |
| | 27.9 |
| | : |
| | 29.5 |
| Student m | |

=>

| |
|-------|
| X |
| n_2 |
| : |
| n_m |

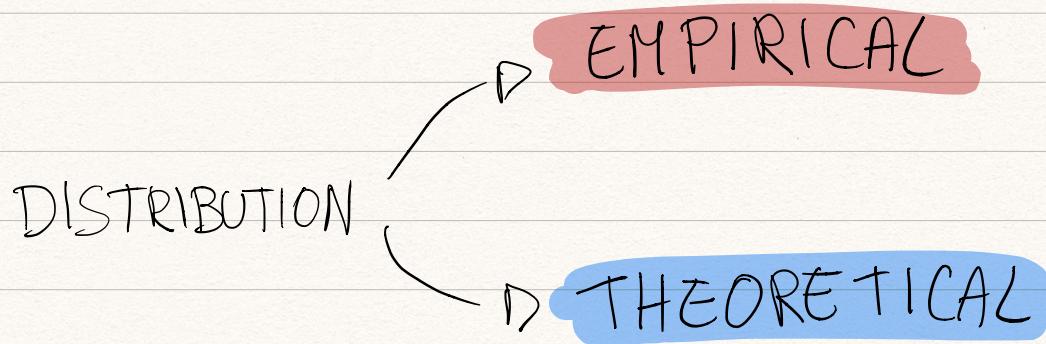
| X | Y |
|-------|-------|
| n_2 | y_1 |
| : | : |
| n_m | y_m |

in case of
MULTIVARIATE
case

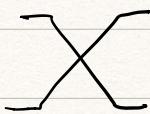
From a DATASET we can extract a DISTRIBUTION that is really crucial and central in statistics because from the dataset we have too much informations and we need to take conclusions.

The final goal of statistics is describe the DISTRIBUTION. This can be exact in the DESCRIPTIVE case and can be an assumption/deduction in the INFERENCEAL case.

Distribution can be done both in the UNIVARIATE case and in the MULTIVARIATE case



DATASET



EMPIRICAL
DISTRIBUTION

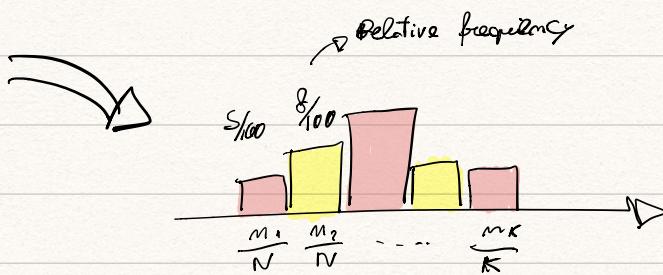
- Grouping
 - Counting
 - To get relative frequencies
- } operation to do
To create the empirical dist.

UNIVARIATE DATASET

| X | Weight |
|-------|--------|
| x_1 | 75 |
| x_2 | 82 |
| : | : |
| x_i | 65 |
| : | : |
| x_N | 92 |

Set of statistics units : students in statistic course

X : Weight

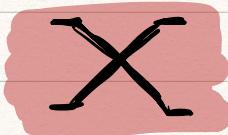


To have a better view of the dataset : create a distribution.

To create this there is an operation of grouping, group the stats units that have similar weight and counting how many there are of them.

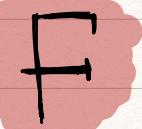
I'm essentially determining the relative frequencies

DATASET



| | |
|-------|-------|
| U_1 | y_1 |
| : | : |
| U_N | y_m |

EMPIRICAL DISTRIBUTION

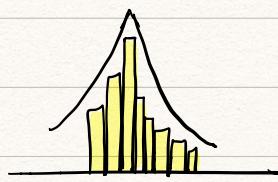


- Define intervals (closure)

- Compute the relative frequencies

| | |
|-------|---------|
| I_1 | m_1/N |
| I_2 | m_2/N |
| : | : |
| I_M | m_M/N |

Large is the DATASET less important are the association
to create the DISTRIBUTION



In statistics we are not interested in the association between units and the values.
This association is discarded when we compute the distribution.

We lose information To gain more knowledge
(we gain more Privacy)