Statistics for Data Science & Business Analysis

Population                              Sample
↓¡
collectⁿ of Items of          ↳ Subset of populatⁿ
interest                          → denoted by "n"
→ "N"                            ↳ statistics
→ Parameter
    (The number we'have    → less time consuming
    obtained when using a  → less costly
    populatⁿ are called p)   ↳
→ populatⁿ is hard to        ┌ statistical data  work
                             └ with incompleteness of
define and hard to                        data
observe in real life                    ⁽ⁱⁱ⁾
                                     work with
(EX: NYU                              sample data
      Campbes)                          & make
                                       data-driven
→ Hard to observe    → easy to observe   decisⁿ]

→ Hard to contact    → easy to contact

                     → statistics usually based on
                     Sample data - to accurate
                     statistical insights
                         2 characteristics.
                             i) Randomness
                             ii) Representativeness

1

Randomness:-
   It is collected when each member of the sample
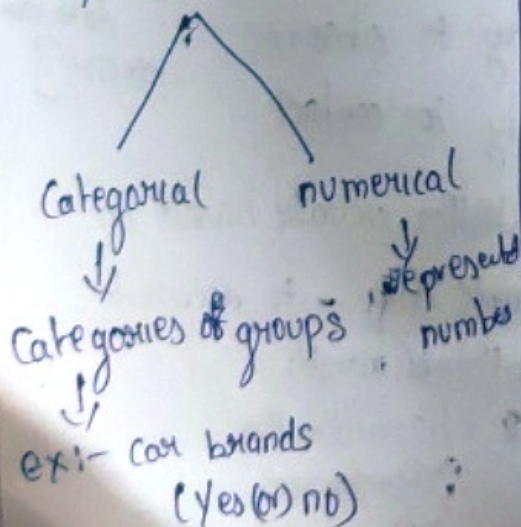is chosen from the population strictly by chance.

Representativeness:-
      It is a subset of the population that accurately
reflects the members of the entire population

[Sample is not random But it is representativeness]
                                          ↓
                                    Ex- Cateen

We can classify data in two main ways:-

      Types of  Data                 measurement level



Categorical         numerical
    ↓                   ↓
Categories of groups  represent
    ↓                  numbe
ex:- car brands
      (Yes or no)

Numerical
   /
discrete
   ↓
finite man
ex:- Grad
     univer
→ no of be

levels of

     Q

   /
nomin
   ↓
Categori
like BM
mercec
not n
Cannot b

2

Numerical data

discrete          continuous

↓                    ↓
finite manner        infinite & impossible to
                     count

Ex:- Grade of        Ex:- height, area, time. etc.
universities

→no. of bottles                    (72·12345)

                                              Ex:- Temperature { no true zero
                                                              Celsius & F

levels of measurements

            Qualitative    Quantitative ─┐─ interval ┐ both vepresented
                                          │           } by numbers
nominal    ordinal                        └─ Ratio    -but one major
                                                       difference
  ↓           ↓                                      Ratio has true zeo
Categories    strict order                          & 2 apples & 8 ap
                                                       areb of 6 & 2 u
like BMW      Ex tasty, delicious                           3
              -ve to +ve (goto)          Kelvin
mercedes                                   ↓
                                         true zero
not numbers
Cannot be ordered

Types of Data

↙    ↘

categorical    Numerical

↓
representation ⟶ frequency distribut$^n$ tables,
⟶ bar charts
pie charts (% calculate chayall)
pareto diagram ⟶ 

↓
It is a special type of bar chart where categories are shown in descending order of frequency

↓

* a curve on the Same graph showing the cumulative frequency

↓
It is the sum of the relative frequencies

It is the no. of occurences of each item.

* pareto principle :
:
⟶ 80 - 20 Rule
⟶ 80% of effect Come from 20% of causes

## Numerical

$$\frac{\text{largest number} - \text{smallest number}}{\text{No. of desired intervals}}$$

frequency distribution
table
calculate

Construct frequency distribution
table

* A Number is included in an interval if that number
  i) is Greater than the lower bound
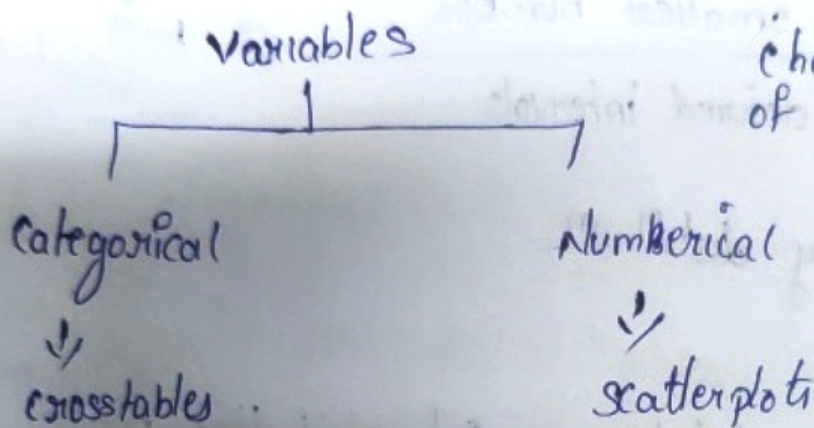  ii) is lower our equal to the upper bound

The relative frequency is the frequency of a given interval as
Part of total

$$\text{relative frequency} = \frac{\text{Frequency}}{\text{Total frequency}}$$

$\left[\frac{8}{L} = \frac{8}{16} = 0.10\right]$  $\frac{\frac{20}{110}}{\frac{600}{110}}$  $\frac{1}{6}$  i) $10(0.1$  $\frac{6}{40}$

* The most common graph is used create numerical data
  is Histogram
      ↓
  create with unequal intervals

# Crosstables & Scatter plots

Variables

Categorical                    Numberical

crosstables                    scatterplots

The side by side bar chart is a variation of the bar chart

→ are used when we are representing 2 numerical variables

→ represents bts & bts of observ -atn

→ outliers are data pts that go - against the logic of the whole dataset