

: DATASCIENCE - FOR - BEGINNERS :

Population

collection of all items of interest

→ N

→ The numbers we've obtained when using population is called parameters.

→ Hard to observe and contact.

Sample.

→ A subset of the population.

→ n

→ The numbers we've obtained using samples is called statistics.

Adv :-

- 1) Less time consuming
- 2) Less costly.

SAMPLE :-

Two defining characteristics

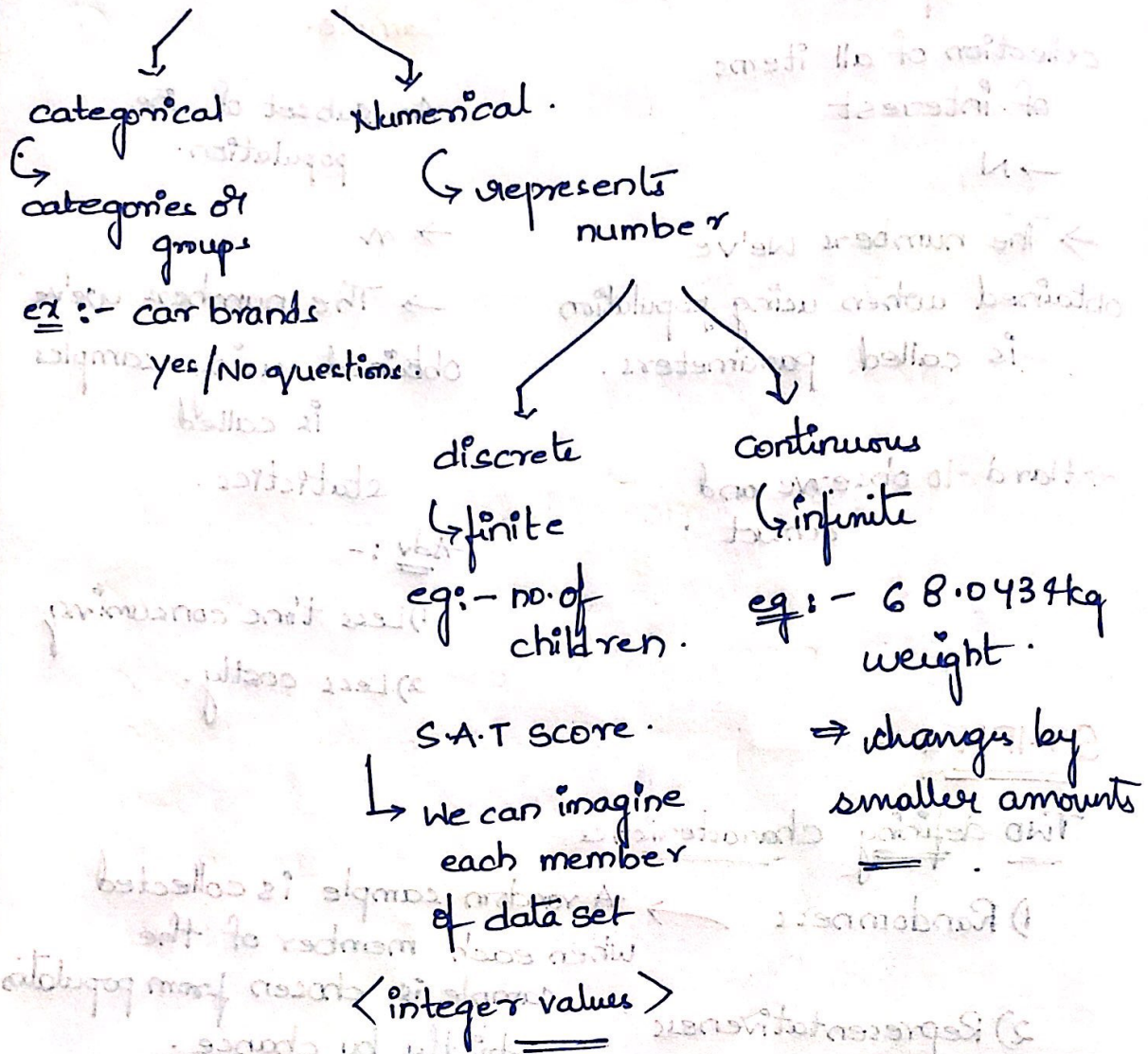
1) Randomness

→ A random sample is collected when each member of the sample is chosen from population strictly by chance.

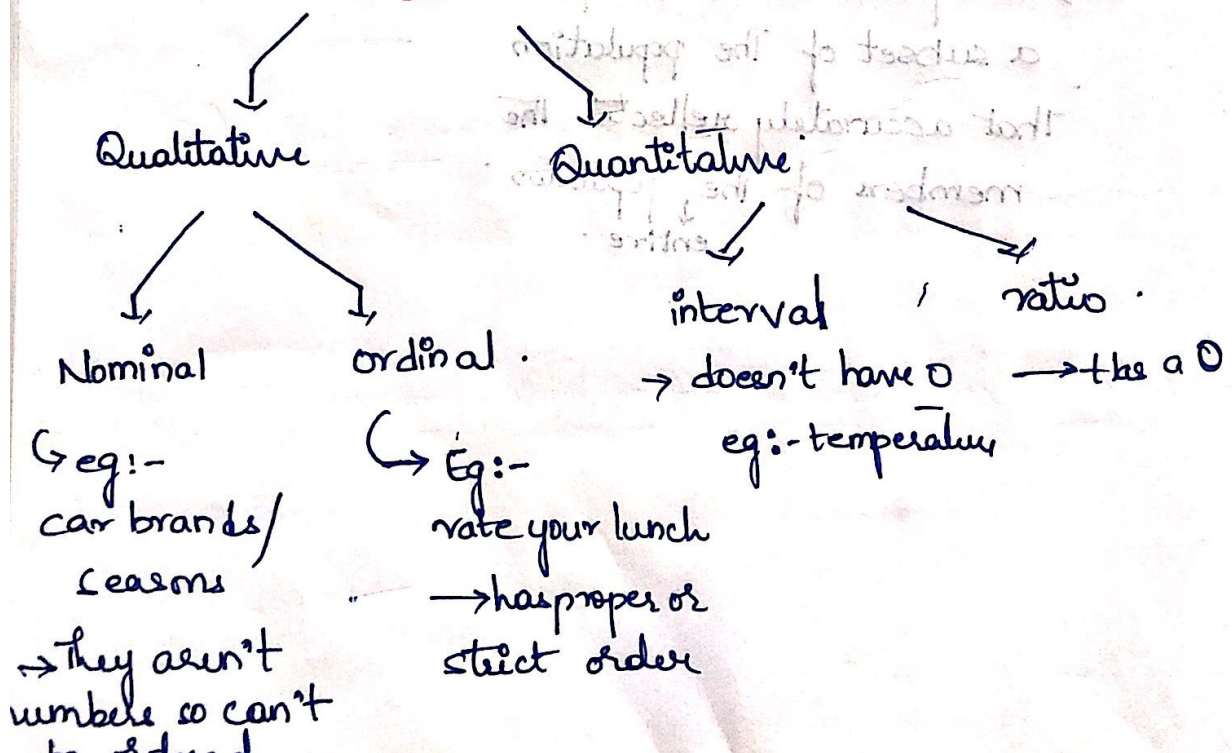
2) Representativeness

→ A representative sample is a subset of the population that accurately reflects the members of the population entire.

Types of Data:-



Measurement Levels:-



Visualizing techniques for categorical variables:-

1) Frequency distribution tables

2) Bar charts

3) Pie charts

4) Pareto Diagrams.

→ categories are shown in descending order of frequency.

→ show the line graph of cumulative frequency.

	Frequency	relative frequency	cummulative frequency.
Audi	124	37%	37%
Mercedes	113	33%	71%
BMW	98	29%	100%
	<u>335</u>	<u>99%</u>	

$$\text{Relative frequency} = \frac{124}{335} \times 100\% \\ = 37$$

$$= \frac{113}{335} \times 100 = 33\%$$

$$= \frac{98}{335} \times 100 = 29\%$$

$$\text{cummulative frequency} = \frac{\text{subgroup relative frequency}}{\text{Total of relative frequency}}$$

↓
Sum of relative frequency.

Visualizing techniques for numerical data

① Frequency Distribution Data with intervals

$$\text{Desired intervals} = \frac{\text{largest number} - \text{smallest number}}{\text{number of desired intervals}}$$

$$\text{eg:- } \frac{100-1}{5} = 19.8 \approx 20$$

1-21, 21-41, 41-61, 61-81, 81-101

A number is included in an interval if the number:

- 1) is greater than lower bound
- 2) is lower or equal to the upper bound.

interval frequency relative frequency

1-21 2 0.10

21-41 4 0.20

41-61 3 0.15

61-81 6 0.30

81-101 5 0.25

20

② Histogram charts :-

⇒ Histogram with unequal intervals.

Graphs to represent relationship between two variables :-

1. Cross tables

2. Scatter plots

→ used when we are representing two numerical

Measures of Central tendency:-

1/19/2022

Mean:-

Also known as simple average.

denoted by $\mu \rightarrow$ population

$\bar{x} \rightarrow$ sample.

- 1. Trimmed mean
- 2. weighted mean
- 3. Harmonic mean
- 4. Geometric mean

$$\text{Mean} = \frac{\sum_{i=1}^n x_i}{N}$$

} By adding up all the components and then dividing by the number of components.

$$(or) \frac{x_1 + x_2 + x_3 + x_4 + \dots + x_n}{N}$$

Disadvantages

\Rightarrow easily affected by outliers.

\Rightarrow Mean is not enough to make definite conclusions

Median :-

Median is middle number in ordered dataset.

\Rightarrow Order data in ascending order.

\Rightarrow

Median is the number at position $(n+1)/2$ in the ordered list

Mode :- \rightarrow it is ^{mostly} useful for categorical variable.

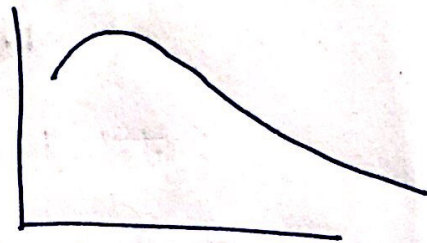
The mode is the value that occurs most often.

Skewness :-

skewness indicates whether the data is concentrated on one side.

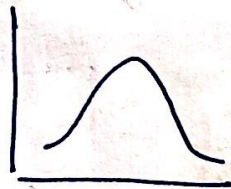
① Right skewness \rightarrow outliers are on the right side

\hookrightarrow mean $>$ median
 \hookrightarrow positive skew.



② Zero skew.

Mean = Median = Mode



③ Left skewness \rightarrow outliers are to the left.

\hookrightarrow mean $<$ median
 \hookrightarrow Negative skewness



\Rightarrow Skewness tells us where the data is situated.

\Rightarrow