

# Statistics from Data Science

## Statistics

1) Descriptive Statistics      2) Inferential Statistics

{ ① Measure of Central Tendency  
  ② Measure of Dispersion }

→ Z test  
t test

Anova test

CHISQUARI

Hypothesis testing.

Confidence Intervals

- 1) Gaussian Distribution.
- 2) Log normal Distribution.
- 3) Binomial Distribution
- 4) Bernoulli's Distribution
- 5) Pareto Distribution (Power law)
- 6) Standard Normal Distribution
- 7) Transformation and standardisation.
- 8) qq plot

Z table, t table,  
chi square table.

The two most important functions of descriptive stats are -  
1) communicate information

vivo V20 Pro (Bipolar) measures about data

Ques. what is Statistics?

Ans. Statistics is the science of collecting, organising and analyzing data,  
{ Better Decision making }

Data  $\Rightarrow$  Facts or pieces of information that can be measured. Data is information such as facts and numbers.

Example  $\Rightarrow$  The IQ of a class used to analyze something or make decision  
 $\{ 98, 97, 60, 55, 75, 65 \} \Rightarrow$  Data.

Ages of students of a class

$\{ 30, 25, 24, 23, 27, 28 \} \Rightarrow$  DATA

### Type of Statistics.

1) Descriptive Stats  $\Rightarrow$  It consists of organising and summarizing data.

2) Inferential Stats  $\Rightarrow$  It is a technique where in we used the data that we have measured to form conclusions.

### Example

v) Classroom of Maths Students (20) students.

Average Marks of 1st Semester

(84, 86, 78, 72, 75, 65, 80, 81, 92, 95, 96, 97 ...)

What kind of question may come in descriptive stats?

What kind of questions may

Come in inferential stats.

### Ex Descriptive Stats.

⇒ What is the average marks of the students in the class?

### Ex Inferential Stats.

⇒ Are the marks of the students of this classroom similar to the other Maths classroom in the college?

Similar to the other Maths classroom in the college?

## Population And Samples

### Exercises

Population ⇒ The entire group of subjects about which we want information.

Parameter ⇒ The quantity about the population we are interested in.

Sample ⇒ A part of the population from which we collect information.

Statistic ⇒ The quantity we are interested in as measured in the sample.

Population ( $N$ )

Sample ( $n$ )

Sampling Techniques  $\Rightarrow$

1) Simple Random Sampling  $\Rightarrow$  when performing.

Simple Random Sampling every member of the population ( $N$ ) has an equal chance of being selected from your Sample ( $n$ ).

2) Stratified Sampling  $\Rightarrow$  where the population ( $N$ ) is split into non-overlapping groups (strata)

Ex  $\Rightarrow$  Gender - [Male] There is no chance of overlapping.  
[Female]

Age Based

(0-10) (10-20) (20-40) (40-100)

Ex Profession  $\Rightarrow$  net PMP Python

But we can apply on Doctor, Engineer, Government job.

### 3) Systematic Sampling → Here from the

population ( $N$ ) we just pick up every  $n^{\text{th}}$  individual from this population.

loop of  $n$  and  $(N)$   $\rightarrow n^{\text{th}}$  individual.

loop of  $n$  and  $(N)$

Eg  $\Rightarrow$  Mall  $\rightarrow$  Survey (Covid)

$\hookrightarrow$  8<sup>th</sup> person  $\rightarrow$  survey.

(ii) and along with survey & population, below are

iii) Convenience Sampling  $\Rightarrow$  I am doing a survey

$\hookrightarrow$  survey,  $\leftarrow$  related to a specific topic.

$\hookrightarrow$  in this particular example like data & demo

$\hookrightarrow$  DATA SCIENCE Only those people

Variables  $\Rightarrow$  A variable is a property that can take on any value.

Eg = Height = 172 cm  
Weight = 65 kg

Two kinds of Variables; ex  $\Rightarrow$  Age, Height, weight

1) Quantitative Variables  $\rightarrow$  Measured Numerically  
 $\quad \quad \quad$  Add, subtract, multiple, divide

2) Qualitative / Categorical Variables

$\hookrightarrow$  Eg: Gender  $\begin{bmatrix} M \\ F \end{bmatrix}$  (Based to some characteristics we can divide some categorical variables)

Eg: Blood Group

## Quantitative Variables

### Discrete Variable

Eg → Whole Number.

① No. of Bank Accounts.

You can say I have 1, 2,  
or 3 bank account.

② Total no. of children  
in a family.  
⇒ 1, 2, 3, 5, 8 (Bargraph)

### Continuous Variable

Eg → Height - 178, 162

Weight = 100kg, 78.2kg

Account of = 1.25 inch,  
Rainfall = 1 inch

## \* Variable Measurement Scales

### 4 types of Measured Variables

- 1) Nominal. {categorical data} → Colour, Gender.
- 2) Ordinal.
- 3) Interval.
- 4) Ratio.

Ordinal ⇒ Order of the data matter, but values do not.

Eg	Students (Marks)
	100
	96
	43
	53
	89

Rank	→ Ordinal Data.
1	
2	
5	
4	
3	

Interval  $\rightarrow$  Order matters, Value also matters,  
natural zero is not present -

Eg. Temperature (F)

$$\begin{array}{c} 70-80^{\circ}\text{F} \\ + \\ 80-90^{\circ}\text{F} \\ + \\ 90-100^{\circ}\text{F} \end{array}$$

You have some range of value between them and the order also basically matters.

Ratio  $\rightarrow$

## \* Frequency Distribution

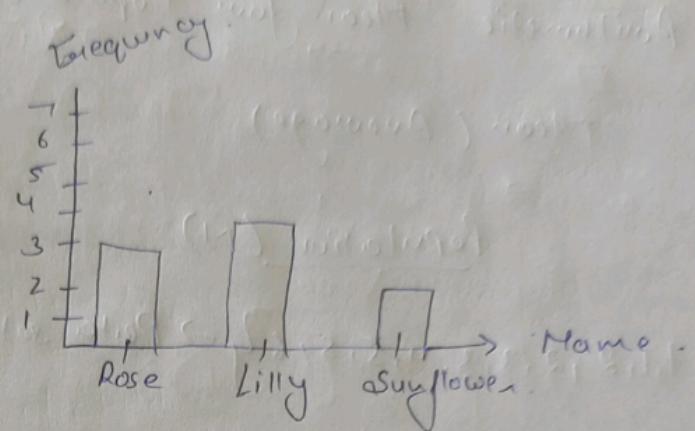
Sample dataset  $\rightarrow$  Rose, lilly, Sunflower.  
Rose, lilly, Sunflower,  
Rose, lilly, lilly.

<u>Flower</u>
Rose.
lilly.
Sunflower.

<u>Frequency</u>	<u>Cumulative Frequency</u>
3.	3
4.	7
2.	9
	Total 9

→ This is frequency.  
Total no. of above dataset.

\* Bar Graph (Discrete)



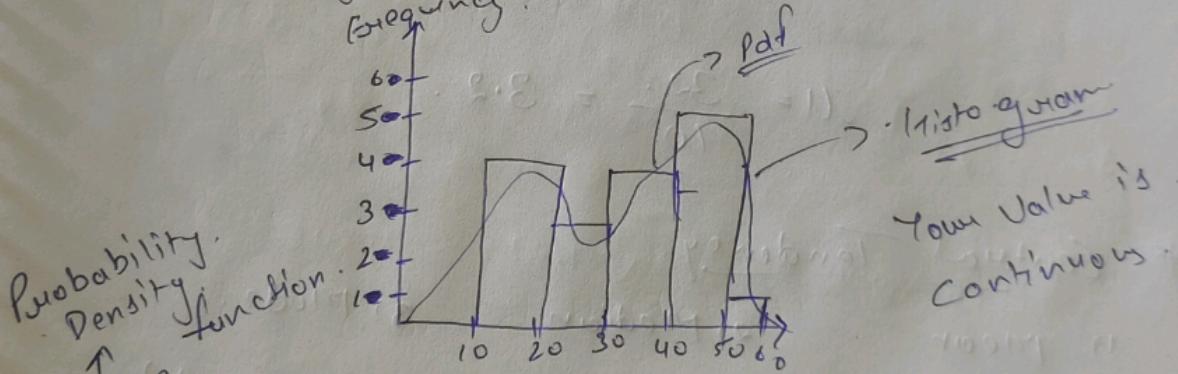
\* Histogram (Continuous)

Age = { 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 45, 50, 51 }.

Bins (Grouping)

Default = 10

Frequency



Pdf  $\Rightarrow$  Smoothing of histogram

## \* Arithmetic Mean from Population & Sample

Mean (Average)

<u>Population (N)</u>	<u>Sample (n)</u>
$n = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\} \rightarrow \text{Data set}$	$= \cancel{2}$

$$\text{mean } \bar{M} = \frac{\sum_{i=1}^N n_i}{N} = \frac{1+1+2+2+3+3+4+5+5+6}{10} = 3.2$$

$$\bar{M} = \frac{32}{10} = 3.2$$

## \* Central Tendency.

- 1) Mean
- 2) Median
- 3) Mode.

$\Rightarrow$  It refers to the measure used to determine the center of the distribution of data.

### Median.

$\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\} + 100$

$$\text{Mean} = \frac{32 + 100}{11} = \frac{132}{11} = 12. \quad \mu = 3.2 + \frac{100}{11}$$

$$\mu = 12$$

Median = {1, 1, 2, 2, 3, 3, 4, 5, 5, 16, 100}.

⇒ First Sort the numbers.

⇒ If the total count of the

numbers is odd number.

mode = 3  
median.

⇒ If the total count of numbers is even.



$$\text{Avg} = \frac{3+4}{2} = 3.5$$

{ Median works well with outliers }.

Mode = { 1, 2, 2, 3, 4, 5, 5, 6, 6, 6, 7, 8, 100, 200 }.

mode = { most frequent Element }.

Mode = 6 ⇒ Measure of central Tendency.

Type of flower. Petal length Petal width  
Rose. 10.2 10.5 10.5  
Lily. 10.2 10.5 10.5  
Sunflower. 10.2 10.5 10.5  
10% missing data

{ Missing Value } → most frequent number } Well worked on Categorical Value.

## Measure of Dispersion

Deviation

i) Variance . ii) Standard Deviation.

(Dispersion) = Spread.

i) Variance  $\Rightarrow$  Variance is the expectation of the squared deviation of a random variable from the population mean or sample mean. Population Variance.

Sample Variance

dataset

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (u_i - \mu)^2 = \frac{10.84}{6} = 1.81$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (u_i - \bar{u})^2$$

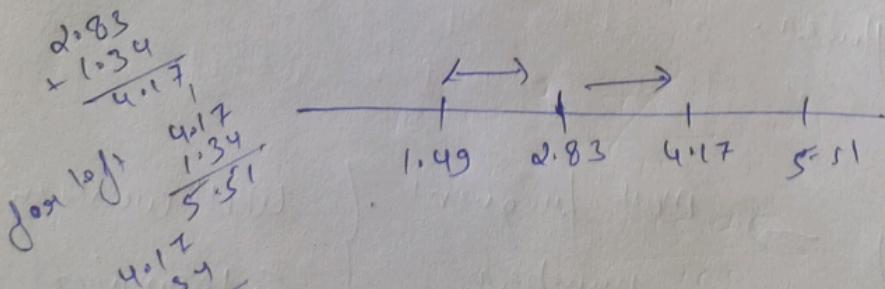
A

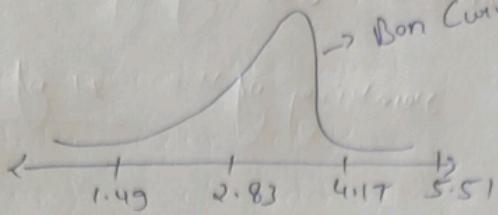
$u_i$	$\mu$	$u_i - \mu$	$(u_i - \mu)^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.68
3	2.83	-0.83	0.68
4	2.83	0.17	0.03
5	2.83	1.17	1.37
		2.17	4.71
			<u>10.84</u>
	$\mu = 2.83$		

$$\sigma = \sqrt{\text{Variance}} = \sqrt{1.81} = 1.345$$

One side Standard deviation  $\Rightarrow$  Elements are.

basically present between two points.





\* Percentiles And Quartiles {find outliers}.

Percentage: 1, 2, 3, 4, 5

% of the numbers that are odd?

? =  $\frac{\text{No. of numbers that are odd}}{\text{Total Numbers}}$

$$? = \frac{3}{5} = 0.6 = 60\%$$

Percentile  $\Rightarrow$  A percentile is a value below which a certain percentage of observations lie.

Example  $\Rightarrow$  Dataset  $\Rightarrow$  2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12.

$\Rightarrow$  What is the percentile ranking of 10?

Ans Percentile Rank of 10 =  $\frac{\# \text{ of values below } 10 \times 100}{\text{No. of values}}$

$$\begin{aligned} &= \frac{16}{20} \times 100 \\ &= \frac{80}{100} = 80\% \end{aligned}$$

Meaning that  $\Rightarrow$  80% of entire distribution is less than 10.

⇒ What is Percentile ranking of 11?

Ans

Percentile ranking of 11 =  $\frac{\# \text{ of values below } 11}{n} \times 100$

$$= \frac{17}{28} \times 100$$

$$= \frac{170}{2} = 85\%$$

Meaning = 85% of entire distribution is less than 11.

Q What value exists at percentile ranking of 25%?

Ans what is formula

$$\text{Value} = \frac{\text{Percentage}}{100} \times (n+1)$$

$$\text{Value} = \frac{25}{100} \times 21$$

$$= \frac{21}{4} = 5.25 \quad \begin{matrix} \text{Index} \\ \text{Position} \end{matrix}$$

So 5.25 index value = 5 is the value of 25%.

# What value exists at percentile ranking of 75%?

$$\text{Value} = \frac{75}{100} \times 21$$

$$= \frac{1575}{100}$$

So the value of 75% = 15.75