



# Day\_1 Statistics

Statistics :-

Q- What is statistics?

↳ It is the science of collecting, organizing and analyzing data. { Better Decision Making }

• Data: facts or pieces of information that can be measured.

ex: Ages of student of a class  
{ 30, 25, 27, 22, 28, 27 }

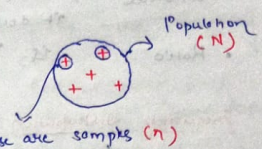
Type:-

- # Descriptive stats: It consist of organizing and summarizing data.
- # Inferential stats: Technique we use the data that we have measured to form conclusions.

ex: Q- What is the average age of student of a class : Descriptive  
Q- Are the ages of the student of this classroom similar to the Maths classroom in the College? Inferential

\* Population And Samples:

for ex. UP & Goa Population we have to predict Exit poll



Population (N)  
These are samples (n)

# Sampling Techniques

- Simple Random Sampling  
↳ Every member has equal chance of being get selected.
- Stratified Sampling  
↳ where the population is split into non-overlapping groups  
for ex: { Male survey }  
          { female }  
Age group: (10-20) (20-30) (30-40)
- Systematic Sampling  
(N)  $\rightarrow$  (nth) individual  
eg. Mall  $\rightarrow$  survey (Covid)  
      ↳ Every 8th person

- Convenience Sampling } only those who interested and have knowledge of it.  
     ↳ survey for ex: Data Science

\* Variables: A variable is a property that can take on any value  
 ex. Height = {102, 142, 160, 150, 160}  
 Weight = {60, 55, 75, 82}

Two kinds of Variables:-

- Quantitative Variable: → Measured Numerically, {Add, sub, multiply}
- Qualitative / Categorical Variable: → Gender {Male, Female}  
     ex: Blood group

Four Types of Measured Variables:

- Nominal data {categorical data} → classes
  - Ordinal → Order of the data matters, value does not
  - Interval → order matters, value also matters but it does not include zero value
  - Ratio → ex: Temperature {Fahrenheit}
- for ex: rank of students.

# Frequency Distribution:

Sample data: {Rose, Lily, Sunflower, Rose, Lily, Rose, Lily}

ex

Flower	Frequency	Cumulative freq
Lily	3	3
Rose	3	6

# \* Bar graph Vs Histogram  
discrete                      continuous

used for

$$M = \frac{1+1+2+2+3}{5} = \frac{9}{5}$$

# Arithmetic Mean for Population and Sample  
 Average

$$X = \{1, 1, 2, 2, 3, 3, 4, 5\}$$

$$M = \sum_{i=1}^N \frac{X_i}{N}$$

focus on notation used

$$\bar{X} = \sum_{i=1}^n \frac{X_i}{n}$$



## # Central tendency

- ① Mean    ② Median    ③ Mode

Refers to the measure used to determine the centre of the distribution of data.

ex { 1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100 }

$$\text{Mean} = \frac{33 + 100}{10 + 1} \Rightarrow \mu = 3.2$$

These are called  
outliers

$$\Rightarrow \frac{132}{11} \Rightarrow \mu = 12$$

- Median: first step sort the data.

{ Median works <sup>good</sup> with outliers well }

As we add outlier but it will only show slight change rather than mean

- Mode: we can use this with categorical variable

we can replace our data which is missing to more recurring data.

# Measure of Dispersion means Spread { how well your data is spread }

- Variance:

Def:

Population Variance

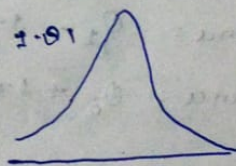
$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

Sample Variance

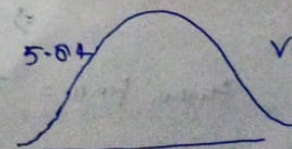
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

ex:

$x$	$\mu$	$x - \mu$	$(x - \mu)^2$
1	2.03	-1.03	3.34
2	2.03	-0.03	0.0009
2	2.03	+0.03	0.0009
3	2.03	0.17	0.03
4	2.03	1.17	1.37
5	2.03	2.17	4.71
$\mu = 2.03$			10.04



Variance  
More



Variance

Spread More



# Percentiles And Quartiles { find Outliers }

$\gamma$  is a value below which a certain percentage of observation lie.

ex: Dataset: 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12  
 $n=20$

what is the percentile of 10?

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

$\Rightarrow \frac{16}{20} \times 100 = 80\%$

Q. what value exists at percentile ranking of 25%?

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

$\Rightarrow \frac{25}{100} \times (21) = 5.25$  This is index value not value

# Five Number Summary { for Removing Outliers }

- Minimum
- First Quartile ( $Q_1$ )
- Median
- Third Quartile ( $Q_3$ )
- Maximum

• Removing the Outliers

{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 9, 27 }

[ Lower fence  $\longleftrightarrow$  Higher Fence ]

Lower fence =  $Q_1 - 1.5(IQR)$

$IQR \Rightarrow$  Interquartile Range

Upper fence =  $Q_3 + 1.5(IQR)$

$Q_3 - Q_1$   
 (75%) (25%)

Lower fence =  $3 - 1.5(4)$   
 $\Rightarrow 3 - 6 = -3$

Higher fence =  $7 + 1.5(4)$   
 $7 + 6 = 13$

} Remove other data



Remaining data:

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 9, ~~10~~

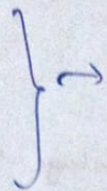
Minimum = 1

$Q_1 = 3$

Median = 5

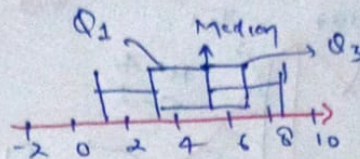
$Q_3 = 7$

Max = 9



5 Number Summary

Box Plot



# Application of Box Plot

It gives a visualisation to see where an outlier is present

# Gaussian / Normal Distribution

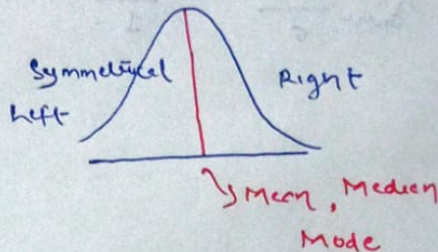
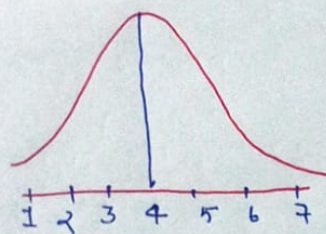
Empirical formula

68-95-99.7% Rule

e.g.: Height

- weight
- IRIS dataset

$\mu = 4$   $\sigma = 1$



→ To find how far from  $\mu$

$$Z_{\text{score}} = \frac{x_1 - \mu}{\sigma}$$

$$\Rightarrow \frac{4.75 - 4}{1} = 0.75 \text{ sd}$$

how many sd left or right

# Practical Application

Dataset

(years)	(Rs)	(kg)
Age	Salary	Weight
24	40k	70
25	50k	80
26	60k	55
27	70k	65

To bring all data in single form as we have these different units

Standardization

⇒ we can apply zscore and convert into Standard normal distribution

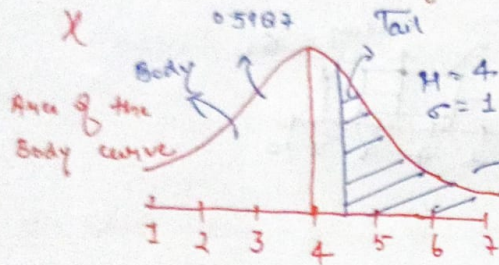
Normalization:

$$\{ \mu=0, \sigma=1 \}$$

↳ convert data in range (0 to 1)

MinMax Scaler  $\rightarrow$  (0 to 1)

Star Stats Interview Question



Ques: what percentage of scores falls above 4.25?

$$(1 - \text{left Area})$$

$$\Rightarrow 1 - 0.5987$$

$$\Rightarrow 0.4013 \Rightarrow 40\%$$

$$Z = \frac{X - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

