

## Introduction to statistics

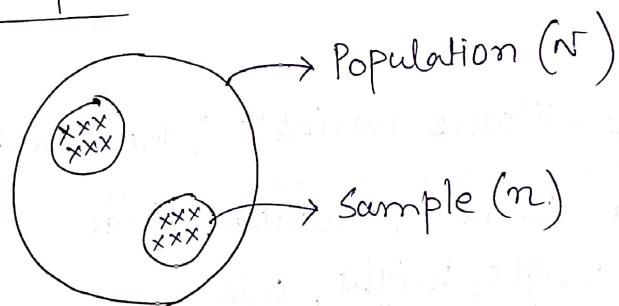
defn :⇒ Statistics is the science of collecting, organising, and analyzing the data. This will imply better decision making.

Types :

> Descriptive stats : Consist of organizing and summarizing of the data. ex: Mean, mode, median, S.D, Var etc.

> Inferential Stats : It's a technique where in we use the data that we have measured to form conclusions.  
ex: Are the age of students of this classroom is similar to the age of math dept students.

Population vs Sample :



Sampling Techniques :

i) Simple Random Sampling : Just go and randomly pick a sample. Every member of the population has equal chance of being selected for your sample ( $n$ ).

ii) Stratified Sampling : Population ( $N$ ) is split into non overlapping groups.

Ex: Gender

male } male, female different.  
female } Not overlapped.  
Ex2: Profession → .net } Sometimes .net developer may  
Python } know Python then overlapping arises. Sometimes it's  
stratified, sometimes it's not.

i) Systematic Sampling: Based on  $n^{\text{th}}$  individual.

Ex: for COVID survey we are checking for every 7<sup>th</sup> person. (No reason for 7<sup>th</sup> person, it's our choice)

iv) Convenience Sampling:

A domain expert will take part in specific surveys.

Ex: A ML engineer will survey ML & literacy in a population.

Note: Technique of sampling to be used depends on the use case.

Variables: It can take values.

Ex: weights = { 78, 75, 50, 52 .. }

Types:

i) Quantitative variable: Measured numerically.

We can add, subtract, multiply, divide.

Ex: Age, weight, height, salary.

ii) Qualitative/categorical variables:

Based on some characteristics we can have categorical variable.

Ex: IQ: (0 - 10) = Less IQ } No addition,  
(10 - 50) = (medium IQ) } subtraction,  
(50 - 100) = (Good IQ) } mul, div  
possible

## Quantitative

discrete variable

e.g: total no of children in family:

e.g: 2, 3, 4, 5, 6, 7

continuous variable

Any value it can have

e.g: height: 170 cm, 180.5 cm

weight: 65 kg, 63.2 kg.

Note:

variable	kind of variable
Gender	categorical (Qualitative)
Marital status	categorical
River length	continuous
Population of state	Discrete
Song length	continuous
Blood pressure	continuous

## Variable Measurement Scales :

i) Nominal Data: categorical data (Qualitative)

e.g: Gender, Type of flower

ii) Ordinal: Here order of the data matters but not the value.

ex:	marks	Rank	This is ordinal data. Rank has order.
	100	1	
	96	2	
	90	3	

iii) Interval: Order and value both matters

Natural zero not present.

e.g: Temperature :  $(70 - 80)^\circ\text{F}$   $(81 - 90)^\circ\text{F}$   $(91 - 100)^\circ\text{F}$   
but  $0^\circ\text{F}$  is not here.

Ratio Data: It is measured along a numerical scale that has equal distances b/w adjacent values and a true zero. (or meaningful zero)

Ex: weight (kg)  $\Rightarrow$  50, 70, 90

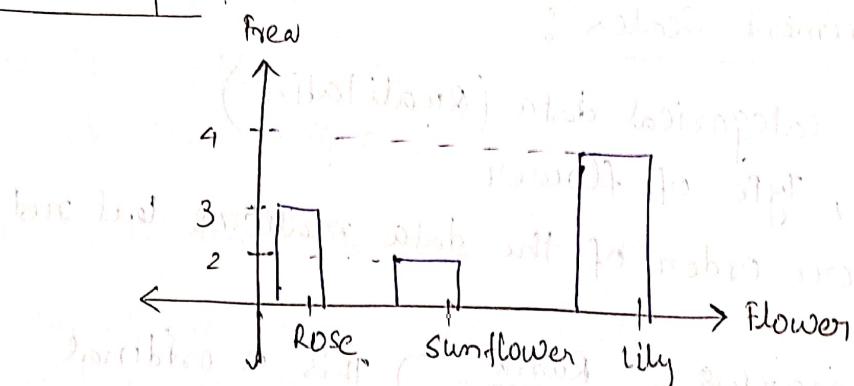
No of Staff  $\Rightarrow$  10, 20, 30

### Frequency Distribution:

Dataset: Rose, lily, sunflower, rose, rose, rose, lily, sunflower, lily, lily

Flower	Frequency	Cumulative Frequency
Rose	3	3
Sunflower	2	5
Lily	2	9

### Bar Graph:



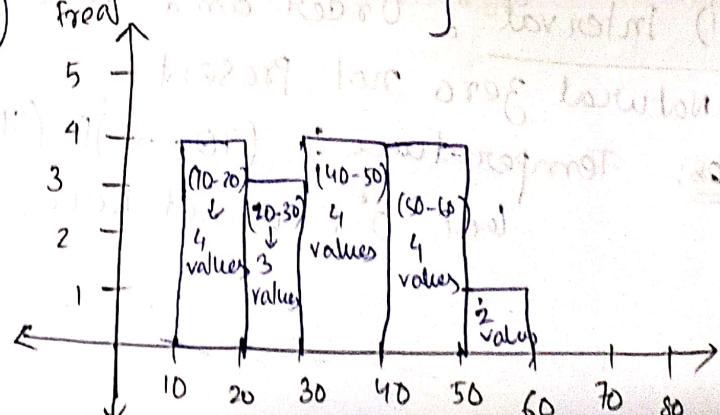
### Histogram: Data should be continuous.

Ex: Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37,

Bins = 10 (Grouping)

Default Size

> Pdf's can be built from here.



## Bar vs Histogram

→ Bar graph is discrete. Histogram is continuous.

### Measure of Central Tendency:

#### Mean :

$$x = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$

$$\mu = \sum_{i=1}^N \frac{x_i}{N} \quad (\text{Population } (N))$$

$$\mu = \frac{32}{10} = 3.2$$

Now for Sample ( $n$ ) mean is denoted by  $\bar{x}$ .

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n} = 3.2$$

Now if in the data set ( $x$ ) a point '100' is added.

$$x' = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$$

$\mu' = 12$  huge difference.

Previously  $\mu = 3.2$  now it's 12. So, median is better.

So, '100' is an outlier.

#### Median :

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

always sort the data in case of median.

$$\text{Mode} = 3$$

Now if 112 also added.

$$x' = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 112\}$$

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

\*\* Adding outliers also the median change is not high.

But in case of avg the difference is much higher.

So, in case of outlier median works better.

Mode: It's defined as most frequent element.

$$x = \{1, 2, 2, 3, 4, 5, 6, 6, 6, 7, 8, 100, 200\}$$

Mode ( $x$ ) = 6 (most frequent)

Significance: Suppose in a dataset there are many missing values. we can use Mode to fill those missing values. It is very useful for Categorical data.

Measure of dispersion:

'Dispersion' means spread.

Ex:  $x_1 = \{1, 1, 2, 2, 4\}$   $\mu_1 = 2$

$$x_2 = \{2, 2, 2, 2, 2\} \quad \mu_2 = 2$$

$\mu_1 = \mu_2$ . But  $x_1$  and  $x_2$  are differently distributed.

> Variance

Population variance:

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

Sample variance:

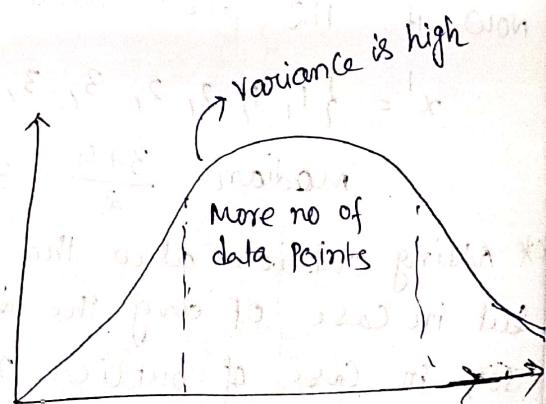
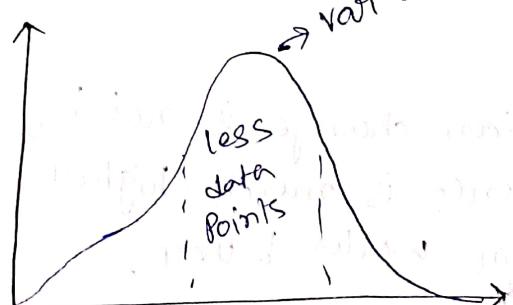
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

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

$$\mu = 2.83$$

$$\sigma^2 = \frac{(1-2.83)^2 + (2-2.83)^2 + \dots + (5-2.83)^2}{6}$$

$$\sigma^2 = 1.81$$

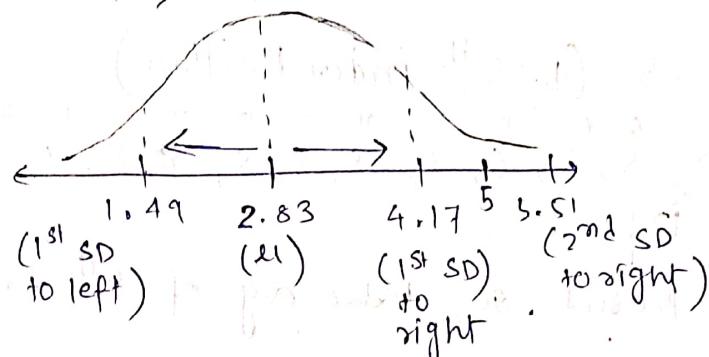


Standard deviation (S.D)( $\sigma$ ):

$$SD(\sigma) = \sqrt{\text{var}(x)}$$

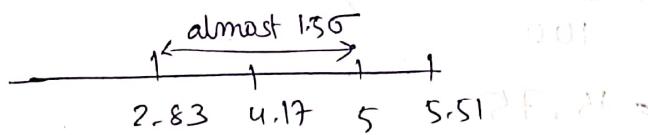
Previous data set  $S.D = \sqrt{10.81} = 1.34$

$$\mu = 2.83, SD = 1.34$$



$$\begin{array}{r} 2.83 \\ + 1.34 \\ \hline 4.17 \\ + 1.34 \\ \hline 5.51 \end{array}$$
$$\begin{array}{r} 2.83 \\ - 1.34 \\ \hline 1.49 \end{array}$$

- , in b/w 2.83 and 4.17 the values of the dataset falls in 1<sup>st</sup> SD to the right of mean
- , in b/w 1.49 to 2.83 the values of the dataset falls in 1<sup>st</sup> SD to the left of the mean
- > 5 falls almost 1.50 to the right of the mean



### Percentile and Quartiles (Outlier)

Percentile: A Percentile is a value below which a certain Percentage of observations lie.

Ex:  $x = \{2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11\}$

Percentile table of 10?  
 $\text{Percentile}(x) = \frac{\# \text{ values below } x}{n} * 100 \%$ .

$$\text{Percentile}(10) = \frac{16}{20} * 100 \% = 80 \text{ Percentile}$$

so, 80% of the dataset less than 10.

$$\text{Percentile}(11) = \frac{17}{20} * 100 \% = 85 \text{ Percentile}$$

What value exists at percentile ranking of 25%?

$$\text{value} = \frac{\text{Percentile}}{100} * (n+1)$$

$$\Rightarrow \text{value} = \frac{25}{100} * (20+1)$$

$$\Rightarrow \text{value} = 5.25 \text{ (5.25}^{\text{th}} \text{ index Position)}$$

$$i=1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16$$
$$x = \{2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12\}$$

5.25 index not present. So, take avg of 5<sup>th</sup> and 6<sup>th</sup> index value

$$\frac{x[5] + x[6]}{2} = \frac{5+5}{2} = \frac{10}{2} = 5 \text{ (ans)}$$

value of 75 percentile?

$$\text{value} = \frac{75}{100} * (21)$$
$$= 15.75$$

$$\frac{x[15] + x[16]}{2} = \frac{9+9}{2} = 9$$

Five no. Summary: Five no. summary includes

- ① minimum ② first quartile ( $Q_1$ ) ③ median
- ④ third quartile ( $Q_3$ ) ⑤ maximum

$$\text{Ex! } \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27\}$$

$$\boxed{IQR = Q_3 - Q_1}$$

$$Q_1(25\%) = \frac{25}{100} * (19+1) = 5^{\text{th}} \text{ index} = 3$$

$$Q_3(75\%) = \frac{75}{100} * (19+1) = 15^{\text{th}} \text{ index} = 7$$

$$IQR = Q_3 - Q_1 = 7 - 3 = 4$$



$$\text{Lower fence (LF)} = Q_1 - 1.5(IQR)$$

$$\text{High fence (HF)} = Q_3 + 1.5(IQR)$$

$$LF = 3 - 1.5(4) = -3$$

$$HF = 7 + 1.5(4) = 13$$

so, anything less than -3 and greater than 13 is outlier. Remove them.

$$\text{An outlier} = \begin{cases} x < LF & \text{or} \\ x > HF & \end{cases}$$

27 is removed from the data set. Remaining

$$\min = 1$$

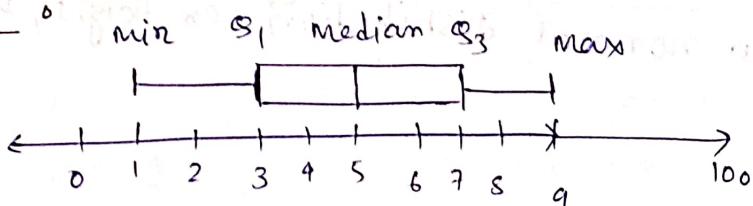
$$Q_1 = 3$$

$$\text{median} = 5$$

$$Q_3 = 7$$

$$\max = 9$$

Box-Plot:



Box-plot is a visualization technique to see if there is any outlier.

## Population variance vs Sample variance?

$$\Rightarrow \sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} \quad \text{this is Population variance}$$

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad \text{this is Sample variance}$$

Here  $(n-1)$  is called as degree of freedom.

In Sample data there is lack of information about majority of the data. So,  $(n-1)$  is used instead of  $N$ . which acts as an automatic ~~stop~~ shield for the analyst and decision maker.

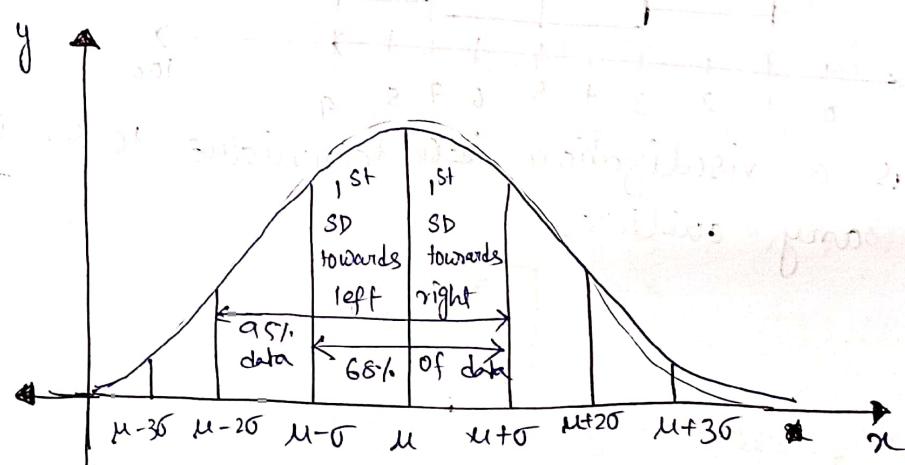
$$\text{value}(s^2) > \text{value}(\sigma^2)$$

so, Sample variance is ~~more~~ higher than population variance.

## Distributions :

### Gaussian Distribution :

It is also normal distribution. ex: Height, weight



Bell curve

$(\mu - \sigma)$  to  $(\mu + \sigma)$  = 68% of data

$(\mu - 2\sigma)$  to  $(\mu + 2\sigma)$  = 95% of data

$(\mu - 3\sigma)$  to  $(\mu + 3\sigma)$  = 99.7% of data

Suppose for a dataset

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

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

$$z(4.75) = \frac{4.75 - 4}{1} = 0.75 \Rightarrow 0.75 \text{ SD towards right of mean}(\mu).$$

$$z(3.75) = \frac{3.75 - 4}{1} = -0.25 \Rightarrow -0.25 \text{ SD towards left of mean}(\mu).$$

Ex! Dataset  $x = \{1, 2, 3, 4, 5, 6, 7\}$   $\mu = 4, \sigma = 1$

$$\Rightarrow z(1) = \frac{1-4}{1} = -3 \quad z(4) = \frac{4-4}{1} = 0$$

$$z(2) = \frac{2-4}{1} = -2 \quad z(5) = \frac{5-4}{1} = 1$$

$$z(3) = \frac{3-4}{1} = -1 \quad z(6) = \frac{6-4}{1} = 2$$

$$z(7) = \frac{7-4}{1} = 3$$

$$x = \{1, 2, 3, 4, 5, 6, 7\}$$

↓

Z-Score

$$z(x) = \{-3, -2, -1, 0, 1, 2, 3\}$$

$z(x) \Rightarrow \{\mu = 0\}$  This is a property of standard normal distribution.

$\sigma = 1$

$z(x) = \text{standardized normal distribution.}$

standardization: Performing Z-score to all the datapoints of the dataset making  $\{\mu = 0, \sigma = 1\}$  is called standardization.

Basically the process to get standard normal distribution is called standardization.

(Q) ① 2020 ODI series :-

$$\text{Avg Score} = 260$$

$$S.D = 12$$

$$\text{Pant Score} = 245$$

2021 ODI Series

$$\text{Avg Score} = 250$$

$$S.D = 10$$

$$\text{Pant} = 240$$

Which year Pant played better?

1) ~~2020 (68)~~

$$Z_{2020}(245) = \frac{245 - 260}{12} = -1.25$$

$$Z_{2021}(240) = \frac{240 - 250}{10} = -1$$

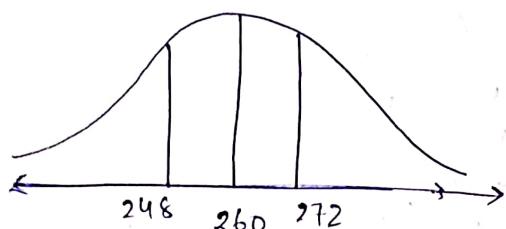


fig: 2020

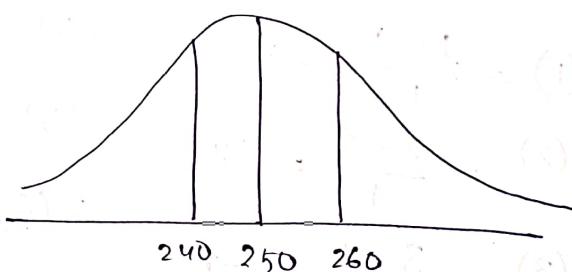
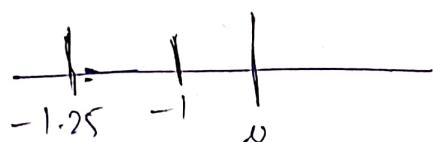


fig: 2021



-1 is close to avg than -1.25. In 2021 Pant played better compared to 2020

2) In maths you get 83 and in Physics = 68

$\mu_{\text{math}} = 83$ ,  $SD_{\text{math}} = 10$  which subject you did

$\mu_{\text{Physics}} = 62$ ,  $SD_{\text{phy}} = 6$ . better?

$$2(\text{Math}) = \frac{80 - 83}{10} = -0.3$$

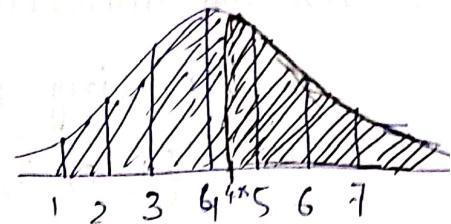
$$2(\text{Phy}) = \frac{68 - 62}{6} = 1$$

In Physics the marks is ~~more~~ 1 SD more than avg. You did physics well.

Q)  $\mu = 4$ ,  $\sigma = 1$ , what % of scores falls above 4.25?

$$\text{Total area} = 1$$

$$z(4.25) = \frac{4.25 - 4}{1} = 0.25$$



$\boxed{\quad}$   $\Rightarrow$  Required area

$$\begin{aligned}\boxed{\quad} &= 1 - \boxed{\quad} \quad \text{from } z\text{-table} \\ &= 1 - z(0.25) \\ &= 1 - 0.5987 \\ &= 0.4013\end{aligned}$$

= 40.13% of the data are over 4.25 (ans)

Q) Avg S.Q = 100, SD = 15. What % of the population would you expect to have an S.Q lower than 85?

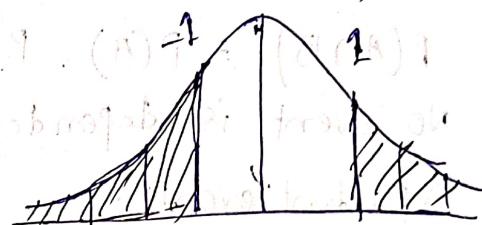
$$\Rightarrow z(85) = \frac{85 - 100}{15} = -1$$

~~It's left part.~~ Directly check  
left part of 85 = right part 115.  
bcz of symmetry.

$$\begin{aligned}\text{Area} &= 1 - z(1) \\ &= 1 - 0.84134\end{aligned}$$

$$= 0.15866$$

= 15.86% of the data (ans)



## Probability

It's a measure of likelihood of an event.

Eg: Rolling of dice.

$$S(\Omega) = \{1, 2, 3, 4, 5, 6\}$$

$$P(\text{even}) = \frac{3}{6} = \frac{1}{2} = 50\%$$

Mutual exclusive event:

Two events can't happen at the same time.

Eg: Tossing a coin. We can get either head or tail. Not both at the same time.

Non mutual exclusive event:

multiple events can happen at the same time.

Eg: Deck of Cards.

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Independent events:

$$P(A \cap B) = P(A) \cdot P(B)$$

No event is dependent on other event.

Dependent events:

$$P(A/B) = \frac{P(A \cap B)}{P(B)}$$

$$\Rightarrow P(A \cap B) = P(B) \cdot P(A/B)$$



P. value :

Here hypothesis testing, confidence interval, significance value etc all are included.

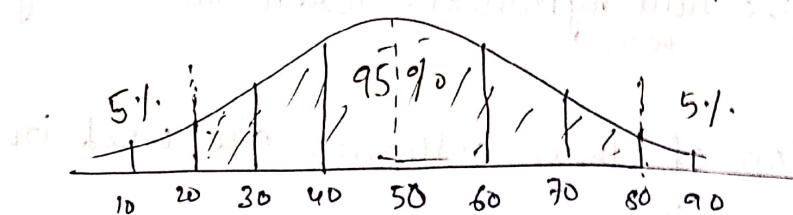
problem: Test whether the given coin is a fair coin or not.

$$\text{we know } P(H) = 0.5 = P(T)$$

Hypothesis testing:

- ① Null hypothesis : Coin is fair.
- ② Alternate hypothesis : Coin is unfair.
- ③ Experiment (Z-test, T-test or anything)
- ④ Reject or Accept the null hypothesis.

Procedure: Tossing of 100 times suppose we get 30 heads.



No of heads :  $\mu = 50$  } suppose  
 $\sigma = 10$  }

$$\text{Significance value } (\alpha) = 0.05$$

so,  $\alpha = 0.05 = 5\%$ . ( $\alpha$  is set by domain expert)

$$\text{Confidence interval} = 1 - \alpha$$

$$\begin{aligned} &= 100 - 5 \\ &= 95\% \end{aligned}$$

Now suppose 20 to 80 covers 95% of the area.  
Then if no of heads is in b/w 20 to 80 then it's a fair coin. otherwise not.

Decision :

Null hypothesis is true or Null hypothesis is false.

Outcome -1 :

We reject the null hypothesis, when in reality it is false.

Outcome -2 :

We reject the null hypothesis, when in reality it is true  $\Rightarrow$  Type 1 error.

Outcome -3 :

We accept the null hypothesis, when in reality it is false.  $\Rightarrow$  Type 2 error.

Outcome 4 :

We accept the null hypothesis, when in reality it is true.

Ex: Application of these outcomes are used in confusion matrix

		Predicted $\rightarrow$		Actual $\rightarrow$
		1	0	
T	TP	TN	FP	1 0
	FP	FN	TP	
F	0	FN	TN	TP

O 1

FT

TF

10

1 0

1 TP FP

0 FN TN

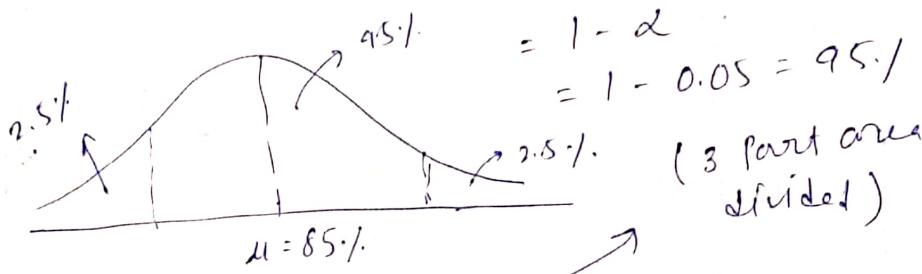
(10)  
(10)



## 1 tail and 2 tail test :

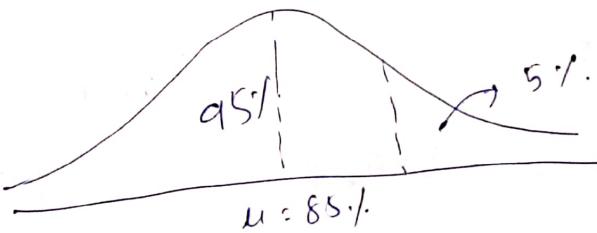
Suppose a college has 85%. placement record. Now you are given a new college. It is found that at the new college a sample of 150 students had a placement rate of 88%. with a standard deviation of 4%. Does this new college has a different placement rate? significant value ( $\alpha$ ) = 0.05.

confidence interval (CI)



2 tail test

Now if the problem statement was : does this new college has better placement than that old college .



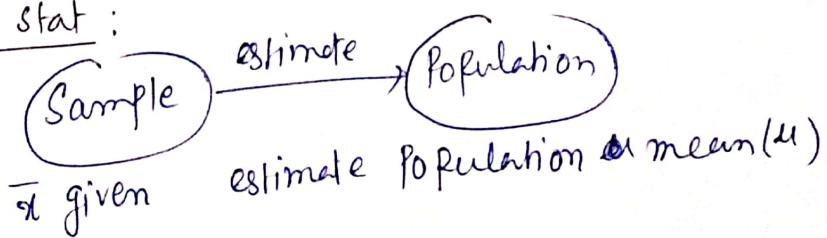
1 tail test .

Here better is told so,  $> 85\%$ . P.H. S of  $\mu$ . 1 tail test because area is divided in two parts

## Point estimate :

The value of any statistic that ~~estimates~~ estimate the value of a parameter.

## Inferential stat :



## Confidence Interval:

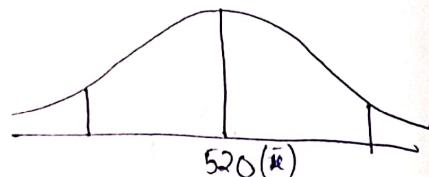
Point estimate  $\pm$  Margin of error.

- ⑧ \* In CAT exam. std = 100. A Sample of 25 test takers has a mean of 520 score. Construct a 95% CI about the mean. ( $\alpha = 0.05$ )

$$\Rightarrow \sigma = 100, n = 25, \alpha = 0.05, \bar{x} = 520$$

CI = Point estimate  $\pm$  margin of error

$$CI = \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$



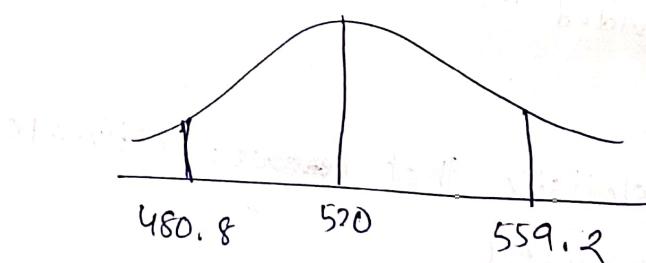
$$\text{Standard error} = \frac{\sigma}{\sqrt{n}}$$

$$CI = \begin{cases} \text{upperbound} = \bar{x} + Z_{0.05} \frac{\sigma}{\sqrt{n}} & \frac{100}{\sqrt{25}} \\ \text{lowerbound} = \bar{x} - Z_{0.05} \frac{\sigma}{\sqrt{n}} & \frac{100}{\sqrt{25}} \end{cases}$$

$$CI = \bar{x} + Z_{0.025} \frac{\sigma}{\sqrt{n}} \quad (Z_{0.025} = 1.96 \text{ from z-table})$$

$$CI_{\text{upper}} = 520 + 1.96 * 20 = 559.2$$

$$CI_{\text{lower}} = 520 - 1.96 * 20 = 480.8$$



⑨ In CAT exam a sample of 25 test takers has a mean of 520 with a std = 80. Construct 95% CI about the mean.

∴ Hence population std not given. Sample std is given. So, Perform T-test. ( $\alpha = 1 - 95\% = 0.05$ )

$$\alpha = 0.05, \bar{x} = 520, \sigma = 80$$

$$CI = \bar{x} \pm t_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right)$$

$$\text{Standard error} = \frac{\sigma}{\sqrt{n}}$$

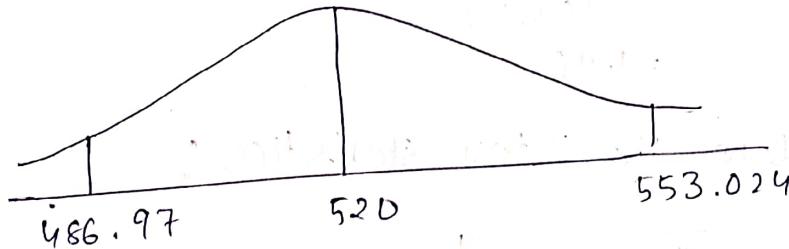
$$\boxed{\text{Degree of freedom} = n - 1 = 25 - 1 = 24}$$

\*  $t_{\alpha/2}$  is dependent on degree of freedom.

$$\text{upper bound} = \bar{x} + \frac{t_{0.05}}{2} \left( \frac{\sigma}{\sqrt{n}} \right) = 520 + 2.064 \left( \frac{80}{\sqrt{25}} \right) = 553.024$$

$$\text{lower bound} = \bar{x} - \frac{t_{0.05}}{2} \left( \frac{\sigma}{\sqrt{n}} \right) = 486.97$$

$(t_{0.025} = 2.064; \text{ from t table with degree of freedom } 24)$



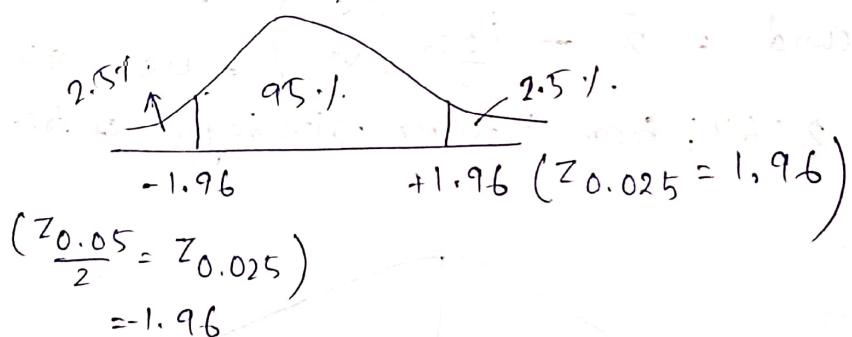
⑧ In the population avg IQ = 100. with std = 15. Researchers to test new medicine to see if there is positive or negative effect on intelligence or not effect. A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect the intelligence? ( $\alpha = 0.05$ )

→ ① Null hypothesis ( $H_0$ ) :  $\mu = 100$

② Alternate hypothesis ( $H_1$ ) :  $\mu \neq 100$

$$③ \alpha = 0.05$$

④ State decision rule : here it is said 'did the medication affect the intelligence?' It means it may increase or decrease. So, perform 2 tail test



⑤ Calculate the Z test statistics :

$$\boxed{Z = \frac{\bar{x} - \mu}{(\sigma/\sqrt{n})}} = \frac{140 - 100}{15/\sqrt{30}} = 14.60$$

⑥ State our decision :

Value has to be in b/w  $-1.96$  to  $+1.96$  but  $14.60 > 1.96$ .

So, Null hypothesis ( $H_0$ ) is rejected. That means avg IQ  $\neq 100$ .

$14.60 > 1.96 \Rightarrow$  improved intelligence.

Note:

i) Z-test  $\Rightarrow$  Population std.

ii) T-test  $\Rightarrow$  unknown population std.

③ Population avg  $\bar{x} = 100$ ,  $n = 30$ ,  $\bar{x} = 140$ ,  $s = 20$

$\alpha = 0.05$ . Did the medication affect?

①  $H_0: \mu = 100$  (Population std not given. So,

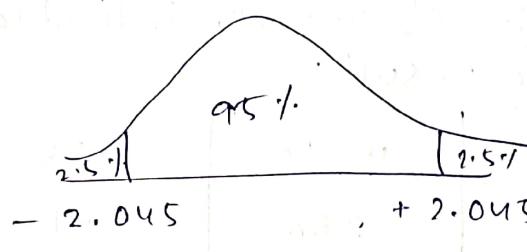
②  $H_1: \mu \neq 100$  (T-test)

③ Degree of freedom =  $n-1 = 30-1 = 29$

④ Decision Rule:

$$t_{\frac{\alpha}{2}} = 2.045$$

(Degree of freedom  
 $= 29$ )



⑤ T-test:

$$T = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

$$\Rightarrow T = \frac{140 - 100}{20/\sqrt{30}} = 10.96$$

⑥ State our decision:

$$10.96 > 2.045$$

So,  $\mu \neq 100$ ,  $H_0$  is rejected.

IQ has increased.



Scanned with OKEN Scanner

## Chi-Square Test :

- > Chi-Square test claims about Population Proportion
- > It is a non parametric test that is performed on Categorical (nominal or ordinal) data.

Q In the 2000 Indian Census, the age of the individual in a small town were found to be the following:

less than 18	18 - 35	> 35
20%	30%	50%

In 2010, age of  $n = 500$  individuals were sampled. Below are the results:

< 18	18 - 35	> 35
121	288	91

$\alpha = 0.05$ . Would you conclude the population distribution of ages has changed in the last 10 years?

From Census of 2000 Expected Population

$$< 18 : 500 * 0.2 = 100$$

$$18 - 35 : 500 * 0.3 = 150$$

$$> 35 : 500 * 0.5 = 250$$

Observed: ~~100~~ 121 288 91

Observed: 121 288 91

Expected: 100 150 250

Step-1:  $H_0$  = The data meets the distribution of 2000  
 $H_1$  = The data does not meet the distribution.

Step-2:  $\alpha = 0.05$

Step-3: Degree of freedom =  $n - 1 = 3 - 1 = 2$

Step-4 : Find chi square ( $\chi^2$ ) value at  $\alpha$  and degree of freedom.

Hence  $\alpha = 0.05$ , degree of freedom = 2

$$\chi^2_{\alpha=0.05} = 5.99$$

$$df = 2$$

so, if  $\chi^2$  is greater than 5.99 then reject  $H_0$

Step-5 : Calculate the  $\chi^2$  value for this problem.

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e} \quad \left( \begin{array}{l} f_o = \text{observed} \\ f_e = \text{expected} \end{array} \right)$$

$$= \frac{(121 - 100)^2}{100} + \frac{(288 - 150)^2}{150} + \frac{(91 - 250)^2}{250}$$

$$= 232.94$$

$$\chi^2 = 232.94 > 5.99 \quad (\text{Reject the null hypothesis } H_0)$$

so,  $H_1$  is true

### Covariance %

weight	Height
50	160
60	170
70	180
75	181

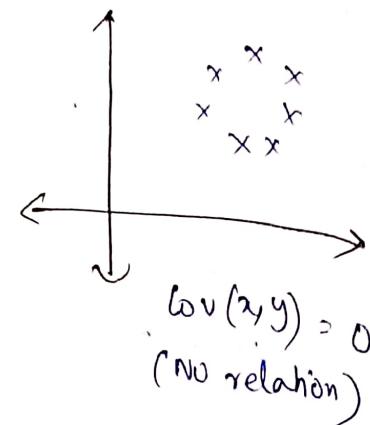
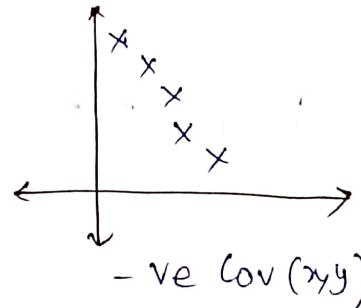
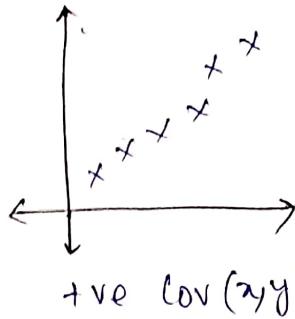
$$\begin{matrix} x \uparrow & y \uparrow \\ x \downarrow & y \downarrow \end{matrix}$$

Study Time	Play Time
2	6
3	4
5	2

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

if  $\text{Cov}(x, y) = +ve \Rightarrow \begin{matrix} x \uparrow & y \uparrow \\ x \downarrow & y \downarrow \end{matrix} \Rightarrow x \propto y$

$\text{Cov}(x, y) = -ve \Rightarrow \begin{matrix} x \downarrow & y \uparrow \\ x \uparrow & y \downarrow \end{matrix} \Rightarrow x \propto \frac{1}{y}$



Disadvantage of  $\text{Cov}(X, Y)$ :

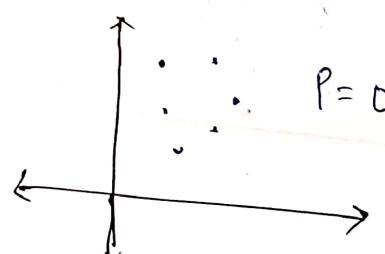
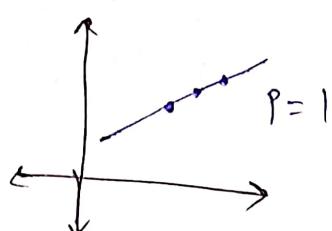
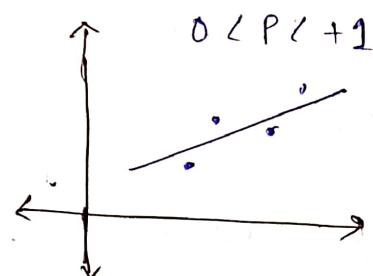
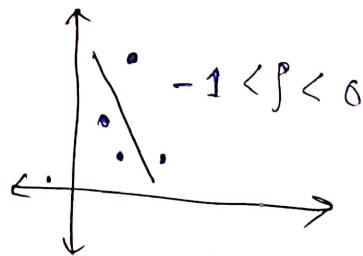
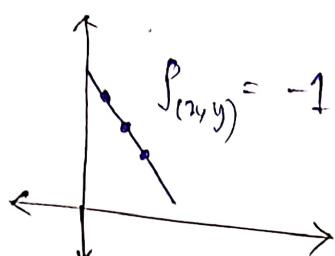
We can tell that it's +ve  $\text{Cov}(X, Y)$  or -ve  $\text{Cov}(X, Y)$ .  
But the value can't be told.

Soln ⇒

Pearson Correlation Co-efficient:

It ranges b/w -1 to +1. The more towards +1 more positive correlation. The more towards -1 more negative co-related.

$$\rho_{(X,Y)} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

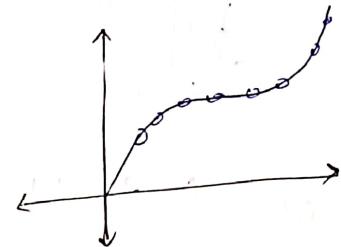


Pearson Co-relation Co-efficient is of linear type.  
For non linear Co-efficient we can use Spearman Co-relation Co-efficient.

$$\text{Spearman Correlation} : \frac{\text{Cov}(R(x), R(y))}{R_{\sigma_x} * R_{\sigma_y}}$$

$r(i)$  = Rank of i

	Height	weight	Rank(x)	Rank(y)
	170	75	2	2
	160	62	3	3
	150	60	4	4
	145	55	5	5
	180	85	1	1



$r(x) = 1$  for height = 180 (highest)

Rank 2 = 170 height.

Spearman Correlation captures the non-linear properties.

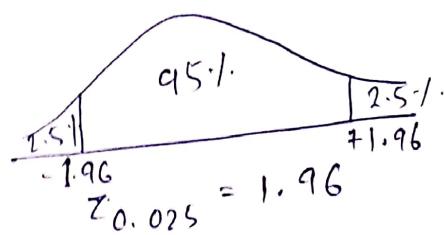
- (1) The avg weight of all residents in Bangalore city is 168 pounds with a std = 3.9. We take a sample of 36 persons and mean = 169.5 pounds. CI = 95%.
- $\therefore \mu = 168, \sigma = 3.9, \bar{x} = 169.5, n = 36, \alpha = 0.05$

(1)  $H_0$ :  $\mu = 168$

$H_1$ :  $\mu \neq 168$

(2)  $\alpha = 0.05$

(3)



(4)  $\alpha = 0.025$

$$Z_{\text{test}} = \frac{\bar{x} - \mu}{(\sigma/\sqrt{n})}$$

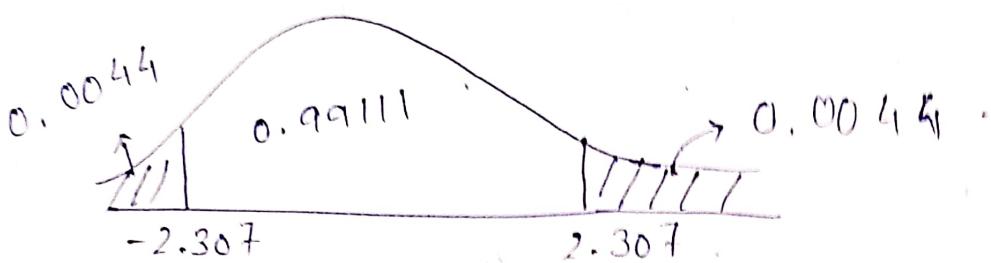
$$Z_{\text{test}} = \frac{169.5 - 168}{3.9/\sqrt{36}} = 2.307$$

(5)  $Z_{\text{test}} > 1.96$ , so,  $H_0$  is rejected.

P-value :

Previous Problem :  $Z_{0.025} = 1.96$

$$Z_{\text{test}} = 2.307$$



$$Z_{2.307} = 0.99111$$

$$\Rightarrow \text{Area} = 1 - Z_{2.307}$$

$$= 1 - 0.99111$$

$$= 0.00889$$

$$\text{Area}(\boxed{\text{II}}) = \frac{0.00889}{2} = 0.0044$$

$$P_{\text{value}} = 0.0044 + 0.0044$$

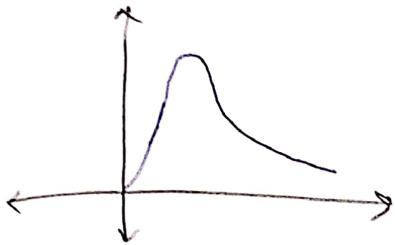
$$P_{\text{value}} = 0.0088$$

If  $P_{\text{value}} < \alpha$  = Reject  $H_0$

$P_{\text{value}} > \alpha$  = Accept  $H_0$

Here  $0.0088 < 0.05 \Rightarrow \text{Reject } H_0$

## Log normal distribution



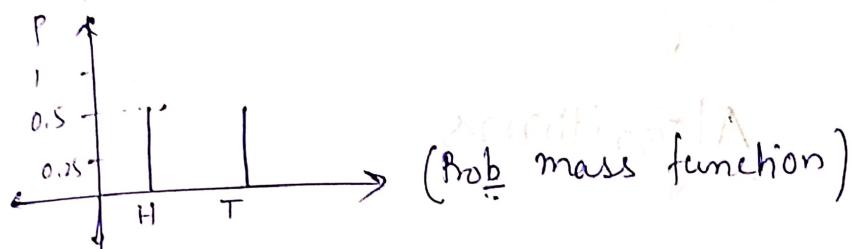
e.g.: wealth distribution

$y$  is random variable and it belongs to ~~not~~ log normal distribution if  $\log(y)$  follows normal distribution.

## Bernoulli Distribution :

It has two outcomes: ( $P$  and  $Q$ ;  $Q = 1 - P$ )

$$P(H) = 0.5 = P(T)$$

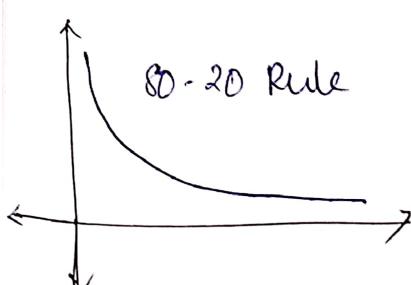


It is for single trial.

## Binomial Distribution :

It is for multiple trials. Each trial follows Bernoulli Distribution.

## Power-law distribution :



e.g.: 80% of the wealth is distributed 20% of people.