

Introduction to Stats:-

① Basics to Advance

Descriptive stats
Inferential Stats

What is Statistics?

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

Data: facts (or) pieces of info that can be measured.

Descriptive Stats:

It consists of organizing and summarizing data

Inferential Stats: Eg: mean, median, mode etc..

Technique where in we use the data that we have measured to form conclusion.

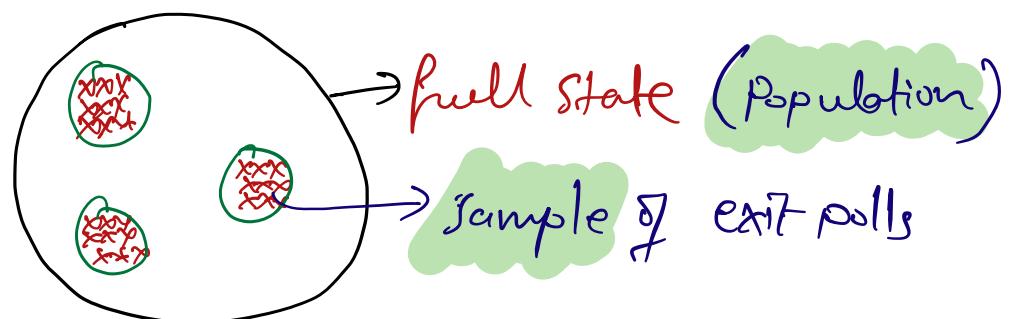
e.g. Are the marks of the students of this classroom similar to the marks of the Maths classroom in the college--

eg(1):

Class room - Sample
College - population

eg(2):

election



Population $\rightarrow N$

Sample $\rightarrow n$

Note:

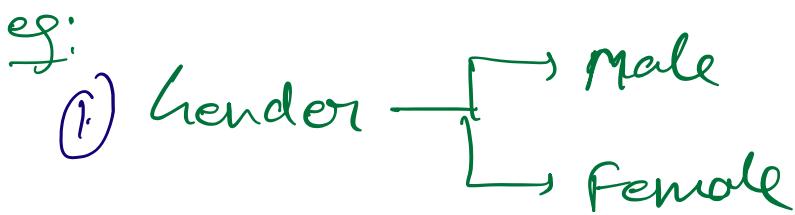
we selected the samples randomly, to make it more systematic we use various sampling techniques

① Simple Random Sampling:

Every member of the population (N) has an equal chance of being selected as a sample (n)

② Stratified Sampling:

where the population (N) is split into non-overlapping groups



- ② Age Group (0-10) (10-20) (20-40) (40-100)
- ③ Profession Doctors, Eng, Actors

③ Systematic Sampling:

from population (n) → we keep every (n^{th}) individual

eg: Mall → survey (Covid)

↳
8th person → Survey

④ Convenience Sampling

Sample is taken from a group of people easy to contact (or) reach.

It is also called Grab Sampling (or) availability Sampling.

eg: Data Science.

↳ we consider a certain domain

20

Variable Measurement Scales

- (1) Nominal data → categorical data["]
- (2) Ordinal → order of data matters, value does not
- (3) Interval → order matters, value also matters,
natural zero is not present
- (4) Ratio

Frequency Distribution

e.g.: Sample data - Rose, lily, sunflower, Rose, lily, sunflower, Rose, lily, lily

<u>Flower</u>	<u>Frequency</u>
Rose	3
Lily	4
Sunflower	2

From this we
can plot diff
graphs.

Cumulative Frequency

$$\text{Rose} - \textcircled{3} \Rightarrow \text{lily} - \textcircled{7} \Rightarrow \text{sunflower} - \textcircled{9}$$

Note:

- if Variable is discrete we draw "BarChart"
- if Variable is continuous we draw "Histogram"

Agedai

pdf is smoothing
of histogram

- ① Measure of Central Tendency
- ② Measure of dispersion
- ③ Gaussian Dist
- ④ Z-score
- ⑤ Standard Normal Dist

① Arithmetic Mean for population & Sample :

Mean (Avg)

population (N)

Sample (n)

e.g. $X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$

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

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$

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

Central Measure Tendency:

① Mean

② median

③ Mode

① \Rightarrow Refers to the measure used to determine the centre of the distribution of data.

② \Rightarrow (1) Sort the numbers
(2) the central element } Good for outliers

③ Most frequent element } Bad for outliers

Measure of dispersion: \rightarrow Spread

① Variance

- to understand how two dispersions are diff

Population

Variance

Sample Variance

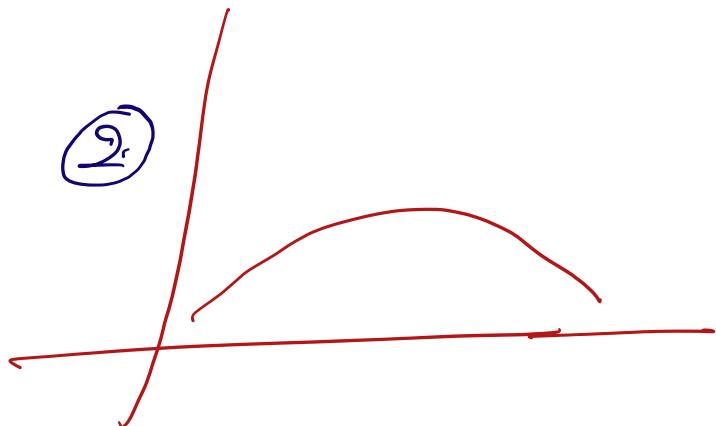
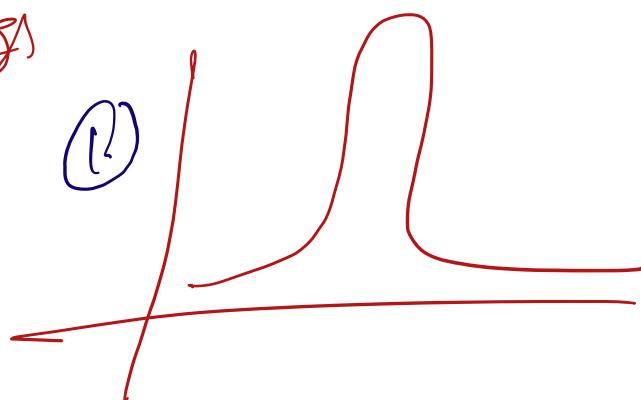
$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

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

Q1

x	μ	$x-\mu$	$(x-\mu)^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
2	2.83	+0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71
<hr/>			
$\mu = 2.83$			<hr/> <u>10.84</u>

Q2



Variance is more in ②, as the distribution is more, the space b/w the points from avg should be more...

Standard deviation:

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

Percentile:

percentage: (eg) 1, 2, 3, 4, 5

x. If nms that are odd?

$$3/6 \times 100 = 50\%$$

percentiles: A percentile is a value below which a certain percentage of observation lie

five Number Summary: [To remove outlier]

- ① Minimum
- ② First Quartile (Q_1)
- ③ median
- ④ Third Quartile (Q_3)
- ⑤ Maximum

$\frac{\text{no. of values less than that number}}{\text{Total samples}} \times 100$

IQR = Inter Quartile Range

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

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

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

ex1

$\{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 10\}$

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

(total elements + 1)

$Q_3 - Q_1$

$$IQR = 7 - 3 = 4$$

$\Rightarrow Q_1 = 3$

by $Q_3 = 7$

25%

This formula is derived from percentile.

$$\begin{aligned} \text{lower fence} &= 3 - 1.5(4) \\ &= 3 - 6 \\ &= -3 \end{aligned} \quad \begin{aligned} \text{higher} &= 7 + 6 \\ &= 13 \end{aligned}$$

$$\text{range} = (-3 \leftrightarrow 13)$$

$$\Rightarrow \min = 1$$

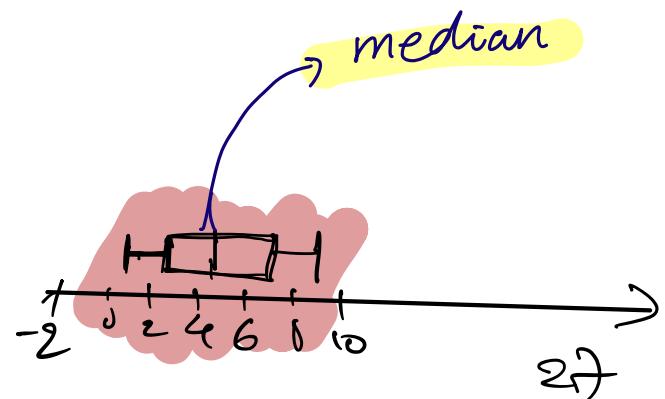
$$Q_1 = 3$$

$$\text{median} = 5$$

$$Q_3 = 7$$

$$\max = 9$$

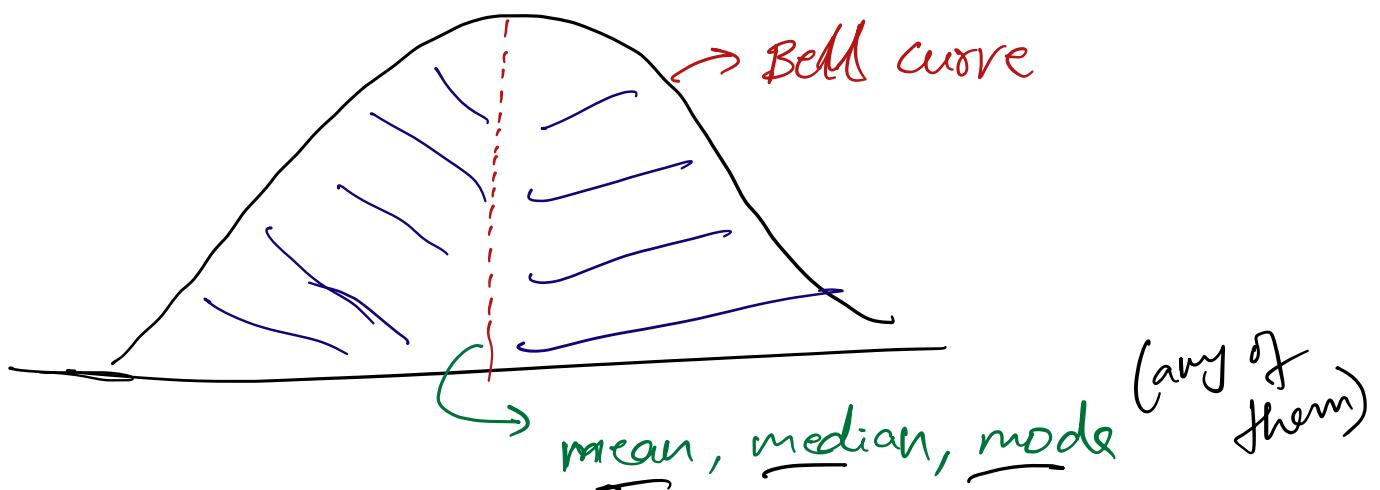
boxplot



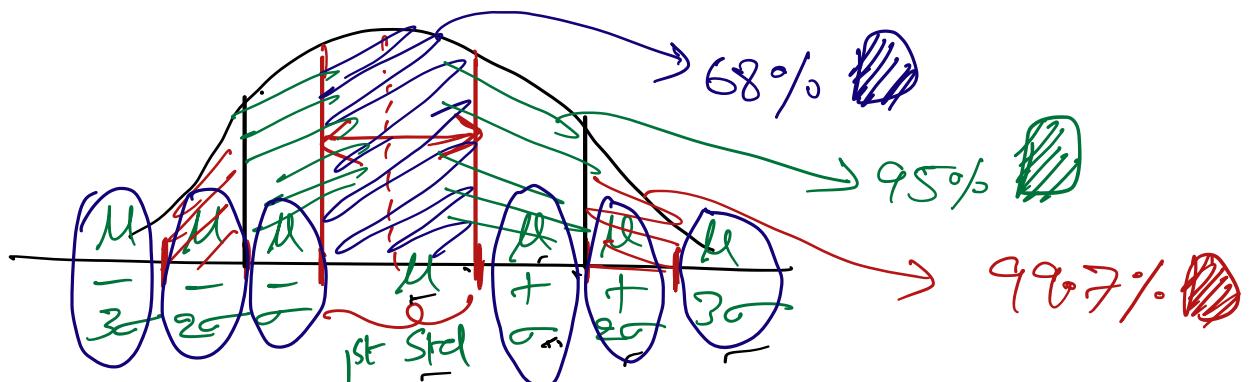
Distributions:

- ① Normal Distr / Gaussian Distr
- ② Standard Normal Distr
- ③ Z score
- ④ Log normal Distr }
⑤ Bernoulli Distr
⑥ Binomial Distr }

Gaussian / Normal Distribution:



→ Normal distribution (or) Gaussian Distribution



Empirical formula: $68 - 95 - 99.7\%$ Rule

%) of distribution

e.g:

ht, wt, Iris dataset etc.....

mean $\mu = 4$ $\sigma = 1$ SD



68,

4.5 → where

① it will fall
in terms of
 SD

+0.5 SD

② $4.75 \Rightarrow +0.75 SD$

To find this we use

Z-Score tells us how much

SD is the point far from the
mean

~~Z-Score~~ 2)



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

- Let, this is +ve number then it is on right
 - if it is -ve then it will be on left
- ⇒ After applying Z-Score then the values will change---

$$Z(1) = \frac{1-4}{1} = -3, \quad Z(2) = \frac{2-4}{1} = -2, \quad Z(3) = \frac{3-4}{1} = -1$$

$\Rightarrow \{ -3, -2, -1, 0, 1, 2, 3 \}$

they look like Standardized values...

They are called as Standard Normal Dist

$$\Rightarrow [\mu = 0, \sigma = 1]$$

Eg:

<u>Age</u>	<u>Salary</u>	<u>Weight</u>
24	40K	70
25	80K	80
26	60K	55
27	70K	45

(yrs) (Rs) (kg)

If we want to make $\mu = 0$ and $\sigma = 1$

unit

We call this as **Standardization**

Z-score formula is applied internally

Normalization:

↳ If we want to change all the values in b/w range of 0 and 1 then we use normalization

↳ we use **MinMax Scaler** for this process.

↳ (0 to 1)

MinMax Scaler:

- Transform features by scaling each feature to a given range.
- This estimator scales and translates each feature individually such that it is in the given range on the training set.

eg $0 - 1$

①

$$X_{\text{std}} = \frac{(X - X_{\text{min}}(\text{axis}=0))}{(X_{\text{max}}(\text{axis}=0) - X_{\text{min}}(\text{axis}=0))}$$

②

$$X_{\text{scaled}} = X_{\text{std}} * (\text{max} - \text{min}) + \text{min}$$

where min, max = feature-range.

Note:

This transformation is always used as an alternative to $\text{mean}=0$,
 $\text{Variance}=1$

Ex: ODI Series

2021

Series Avg 2021 = 250

S.D = 10

V.K Avg Score = 260

2020

Series Avg 2020 = 260

S.D = 12

V.K Avg Score = 245

Compared to both the series in which year V.K score was better?

Sol:

To find out this value, we find the z -score

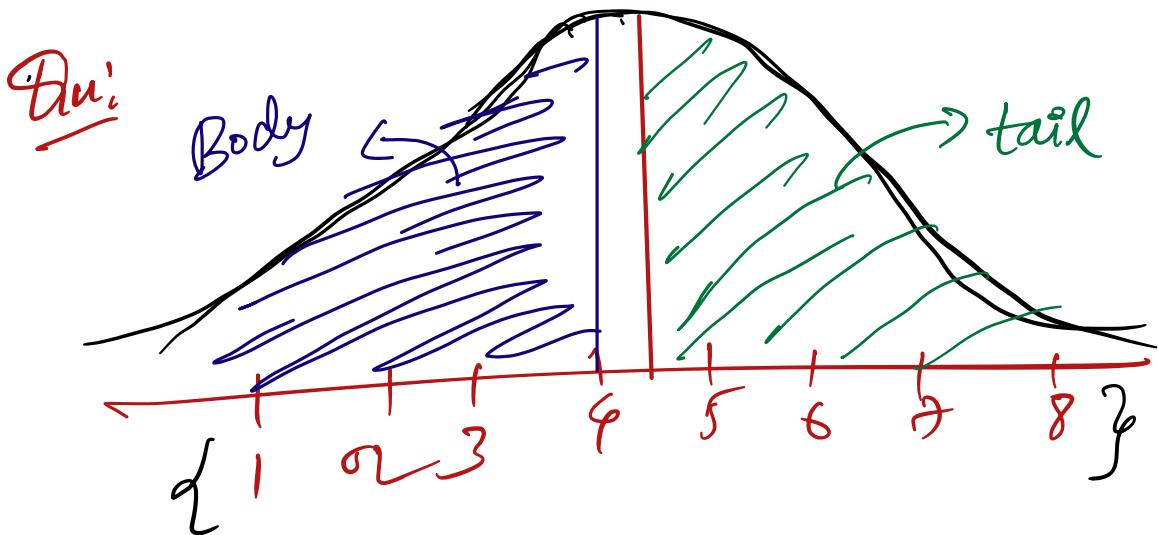
$$Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

$$\Rightarrow \boxed{2021} = \frac{260 - 250}{10} = 1$$

$$\Rightarrow \frac{245 - 260}{12} = -1.25$$

2020

↓ (Best)



what percentage of scores fall above
4.25?

Note:

Z-Score finds the area of the
body curve

$$z = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

help!

In online we have values for all kinds of z-score and then you can find the area required.

0.4013 from table

⇒ 40%

Agenda:

- ① Probability
- ② P & C
- ③ Confidence Interval
- ④ P value
- ⑤ Hypothesis testing

Probability:

- What is probability
- Addition Rule

- ↳ Mutual Exclusive
- ↳ Non-Mutual Exclusive

→ Multiplication Rule

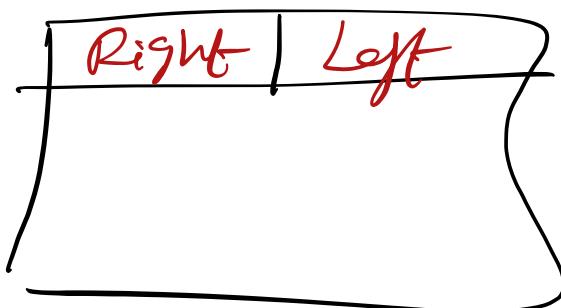
- ↳ Independent \Rightarrow (Conditional Prob)
- ↳ dependent
 - \hookrightarrow (Naive Bayes)

Per & Com:

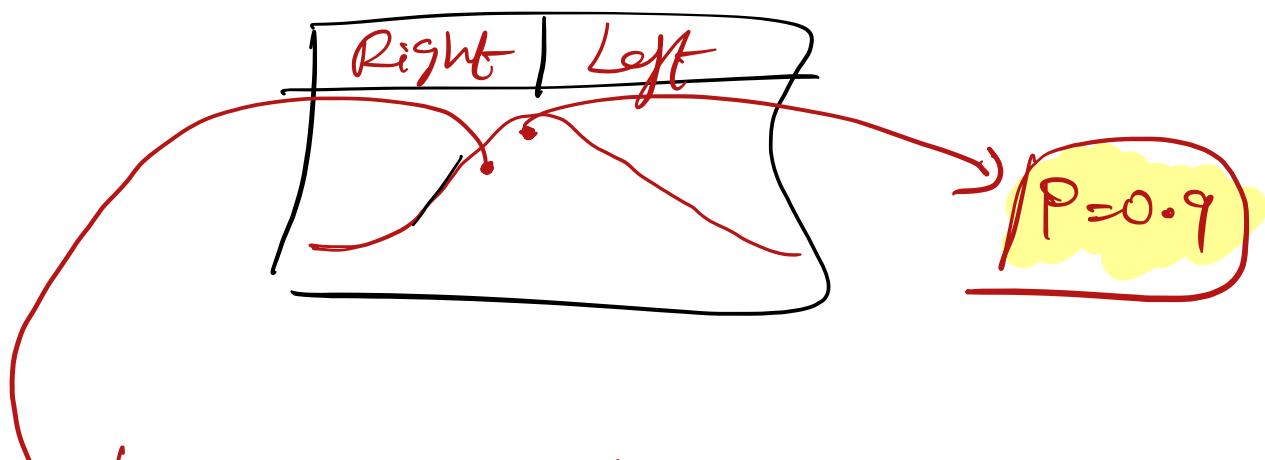
n_p , n_C , Circular arrangement, repetitions,
etc. ---

P value:

Let's consider a laptop mouse pad.



- ① we will use right click (or) left click, when we want to select something.
- ② Else most of the time we just play our finger on the bottom area.
- ③ So if we plot the graph for the area we touch mostly it will look smtg like,



lets, consider the value of

$P=0.8$

⇒ Every 100 times we touch the pad, there can be chance of 80 times we touch this specific region.

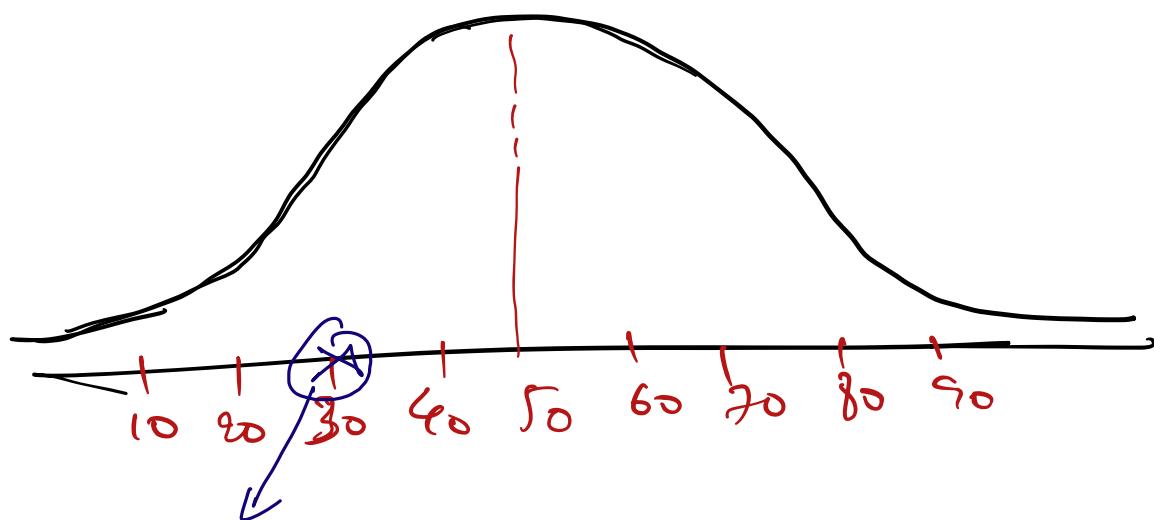
Hypothesis Testing:

- ① Null Hypothesis: Coin is fair
- ② Alternate hypothesis: Coin is unfair
- ③ Experiment
- ④ Reject or Accept the null hypothesis

ex: lets take a coin example, and we tossed it 100 times

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

If we are drawing a SD graph.



If we say we got 30 times head
Is that fair coin (or) an unfair coin?

Note:

- ⇒ our value should be as close as possible to the mean value $\boxed{50}$
- ⇒ for this, to determine, tell how far the points accepted from the mean, it is calculated by using.

Significance Value

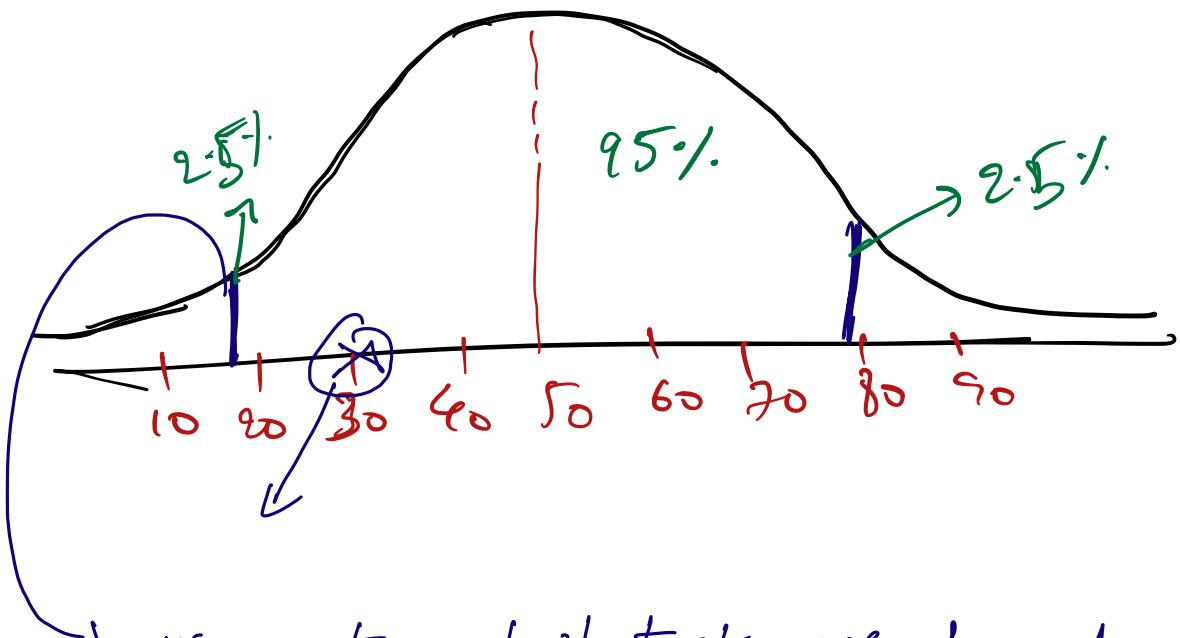
Rep by $\boxed{\alpha}$

bt,
 $d = 0.05$ $\Rightarrow 5\%$ {Domain Expert}

$$100 - \delta_2 = 95\%$$

Confidence interval

on graph



→ Using two tail test we found out
these are our 2.5% from front
and back

Note:

→ If the experiment fall in the 95%
confidence interval we call it a fair

confidence interval we call it a fair coin

→ Else it is a unfair coin

e.g. we got 10 heads out of 100 tails

It is outside of CI

⇒ it is unfair

- (x) Null hypothesis $\Rightarrow (H_0)$
(x) Alternate hypothesis $\Rightarrow (H_1)$

Reality Check:

Null hypothesis is True (x)
Null hypothesis is False

Decision:

H_0 is true (x) H_0 is false

Outcome 1 :-

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

↳ **Yes** ↳ **No**

O/P 2:

We **reject** the null hypothesis,
when in reality it is **true**

↳ **No**

Type 1 Error

Op 3:

We Accept the null hypothesis,
when in reality it is false

No

Type 2 Error

Op 4:

We Accept the null hypothesis,
when in reality it is true

Yes

Good

	P	N
T	TP	TN
F	FP	Fn

Type 2

Type 1

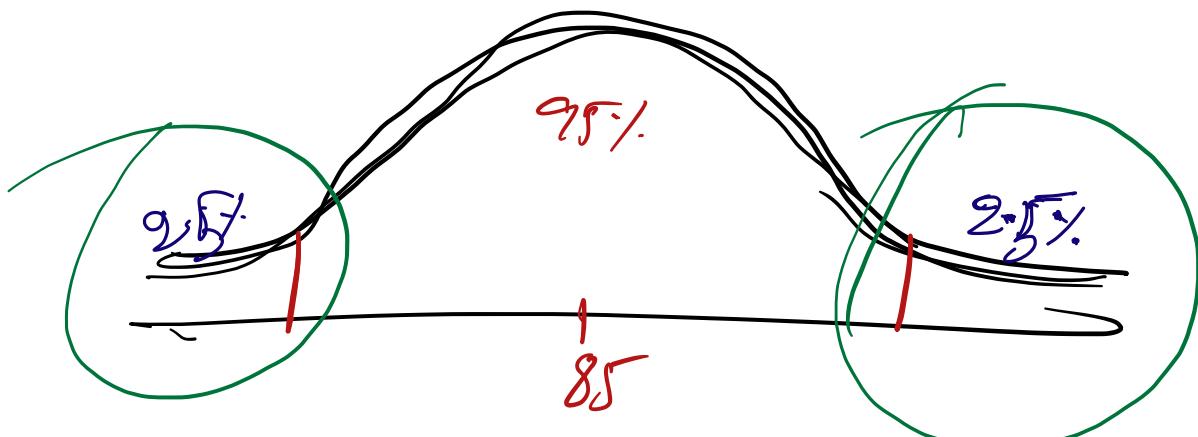
1 tail and 2 tailed test:

e.g.: Colleges in Karnataka have an 85% placement ratio. A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88% with a S.D (Rey.) . Does this college have a diff placement rate compared to other colleges?

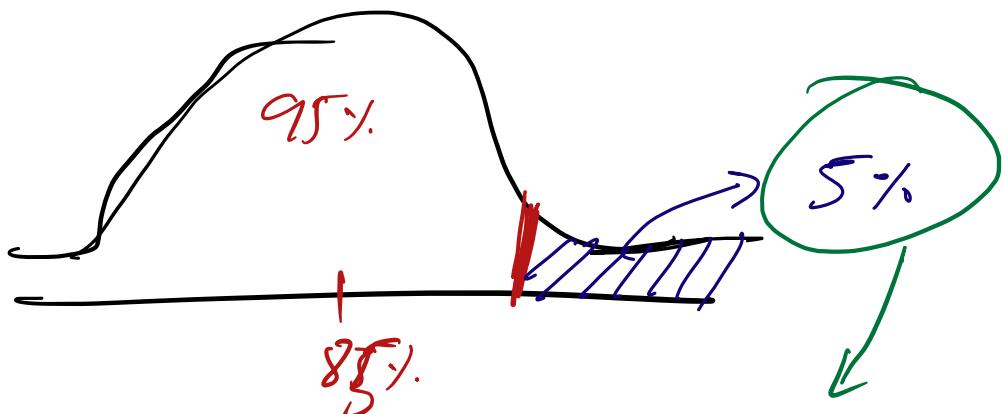
Let, $\alpha = 0.05$ (Sig. Val)

Then $C.I \Rightarrow 95\%$

two tailed



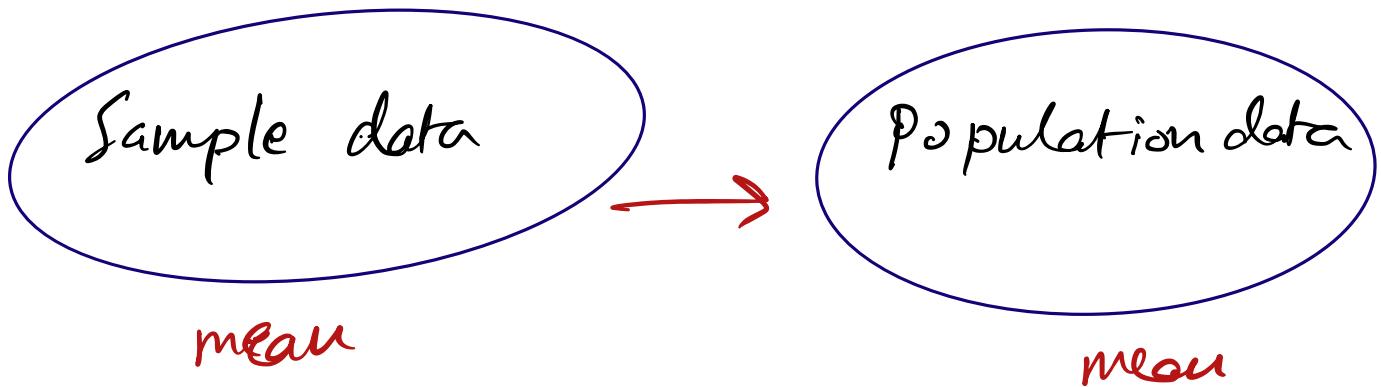
⑧ Does the college have a placement rate greater than 85%?



Confidence Interval:

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

Inferential Stats:

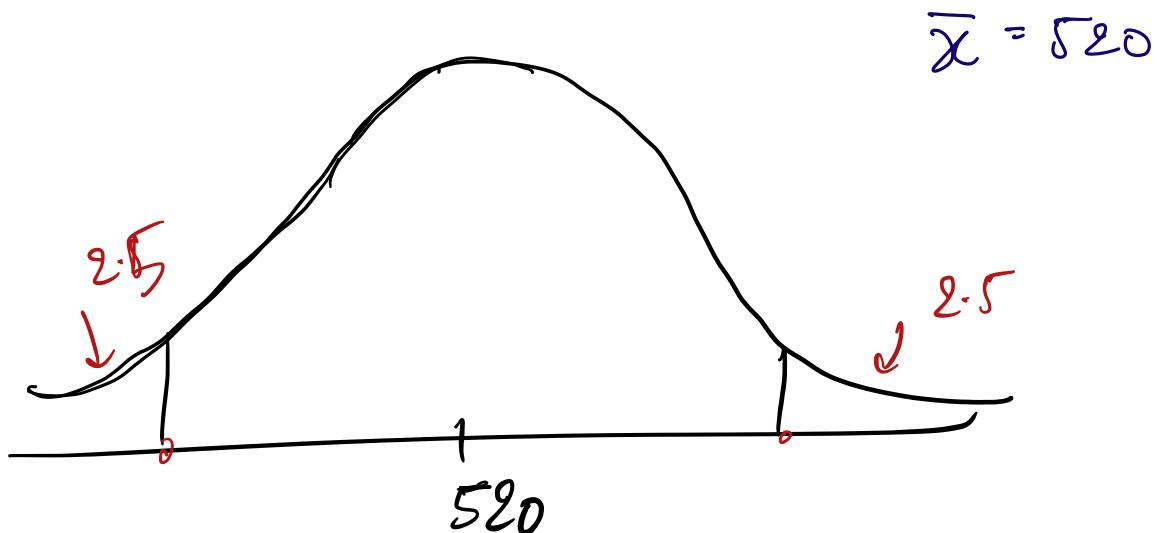


→ with the help of sample data, we will estimate population data

$$CI = \text{Point Estimate} \pm \text{Margin of Error}$$

ex: On the Quant test of CAT Exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520 score. Construct a 95% CI about the mean?

Sol: $\sigma = 100$, $n = 25$, $\alpha = 0.05$,



→ Point Estimate \pm Margin of error

$n \geq 30$
⇒

generally

$$\bar{x} \pm z_{\alpha/2}$$

$$\frac{\sigma}{\sqrt{n}}$$

Standard
Error

Upper bound = $\bar{x} + z_{\frac{0.05}{2}} \frac{100}{\sqrt{25}}$

Lower bound = $\bar{x} - z_{\frac{0.05}{2}} \frac{100}{\sqrt{25}}$

Note!

$$z_{\frac{0.05}{2}} = \sqrt{z_{0.025}}$$

you can find
this from
online table

of z-score.

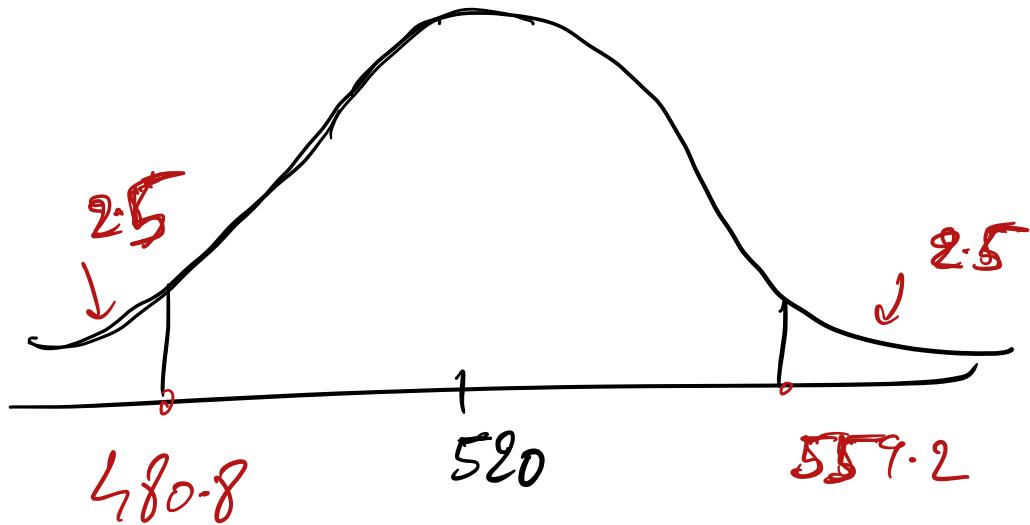
$$\alpha = 0.025$$

$$1 - \alpha = 0.975$$

↳ from table it lies at 1.96

$$\boxed{\text{Upper}} = 520 + 1.96(20) = \boxed{559.2}$$

$$\boxed{\text{Lower}} = 520 - 1.96(20) = \boxed{480.8}$$



(Q.) On the Quant test of CAT exam, a sample of 25 test takers has a mean of 520 with a SD of 80. Construct 95% CI about the mean.

Ans: here, population std is not given

↳ hence, we use t-test

= point estimate \pm margin of error

$$= \bar{x} \pm t_{\frac{\alpha}{2}} \left(\frac{s}{\sqrt{n}} \right) \quad \text{standard error}$$

Upper bound = $\bar{x} + t_{\frac{0.05}{2}} \left(\frac{s}{\sqrt{n}} \right)$

Degree of freedom $\Rightarrow n-1 \Rightarrow 25-1 = 24$

$$= 520 + 2.064 \left(\frac{80}{5} \right)$$

$$= 553.024$$

Lower Bound = $\bar{x} - t_{\frac{0.05}{2}} \left(\frac{s}{\sqrt{n}} \right)$

$$= 520 - 2.064 \left(\frac{80}{5} \right)$$

$$= 486.97$$

$$\left[486.97 \longleftrightarrow 553.024 \right]$$

One Sample Z-test:

① In the population, the average IQ is 100 with a SD of 15. Researchers wants to test a new medication to see if there is a positive (or) negative effect on intelligence, (or) no effect at all. A sample of 30 participants who have taken the medication has a mean of 120. Did the medication effect the intelligence?

St. ② Define Null hypothesis

$$H_0 \Rightarrow \mu = 100$$

$$\alpha = 0.05$$
$$S_{\bar{X}} = 9.57$$

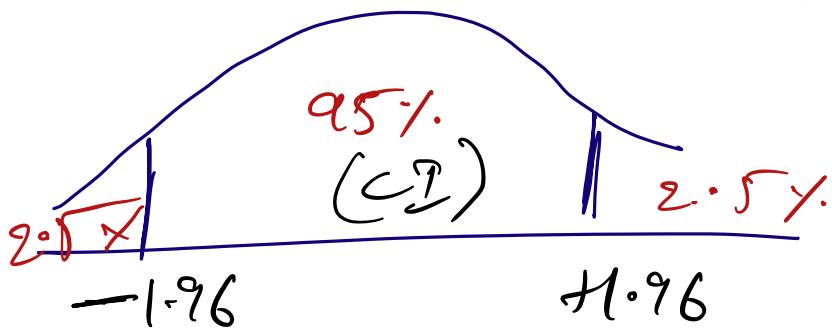
③ Alternate hypothesis

$$H_1 \rightarrow \boxed{\mu \neq 100}$$

④ State Alpha $\alpha = 0.05$

⑤ State Decision Rule

(2 tail test)



5. Calculate Z test:

$$Z = \frac{\bar{x} - u}{\left(\frac{\sigma}{\sqrt{n}} \right)}$$

note:

- ① for sample value we take standard error
- ② for population value we take only (~~error~~) not standard error.

$$\rightarrow \frac{140 - 100}{\sqrt{30}} = 14.60$$

6. State our Decision

$$\rightarrow 14.60 > 1.96$$

If z is $<$ than -1.96 (or) > 1.96 ,
reject the null hypothesis.

② did the medication improve the intelligence (or) decrease.

↳ [Improved]

One Sample T-test:

$\left[\begin{array}{l} z\text{-test} \Rightarrow \text{population std} \\ t\text{-test} \Rightarrow \text{unknown " " } \end{array} \right]$

③ population the average IQ = 100

$$n=30 \quad \bar{x}=140 \quad s=20$$

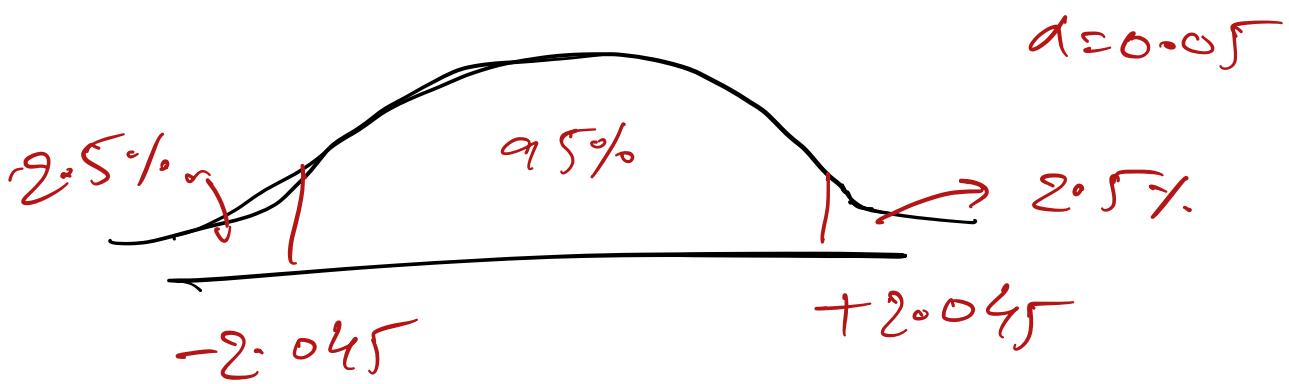
Did the medication affect the int?

$$\beta d=0.05$$

- Sol.
- ① $H_0 \Rightarrow \mu = 100$
 - ② $H_1 \Rightarrow \mu \neq 100$
 - ③ Calculate the degree of freedom

$$n-1 = 30-1 = 29$$

- ④ State Design rule



- ⑤ T test

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

$$\bar{x} = 140, \mu = 100, s = 20, n = 30$$

$$= 10.96$$

$t = 10.96 > 2.045$

⇒ Reject null hypothesis

$$P \leq \text{sig value}$$

[increased the intelligence]

Agenda:

- ① CHI Square
- ② Covariance
- ③ Pearson Correlation Coefficient
- ④ Spearman Rank Correlation
- ⑤ F test (ANOVA)

Chi Square:

- Chi Square test claims about population proportions
- It is a non-parametric test that is performed on categorical variables can be either 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

Using $\alpha = 0.05$, would you conclude the population distribution of ages has changed in the last 10 yrs?

Sol: note: when we have data in terms of proportion then we need to

implement non-parametric test.

< 18	$18 - 35$	> 35	
121	288	91	(Obs)
100	150	250	(Exp)

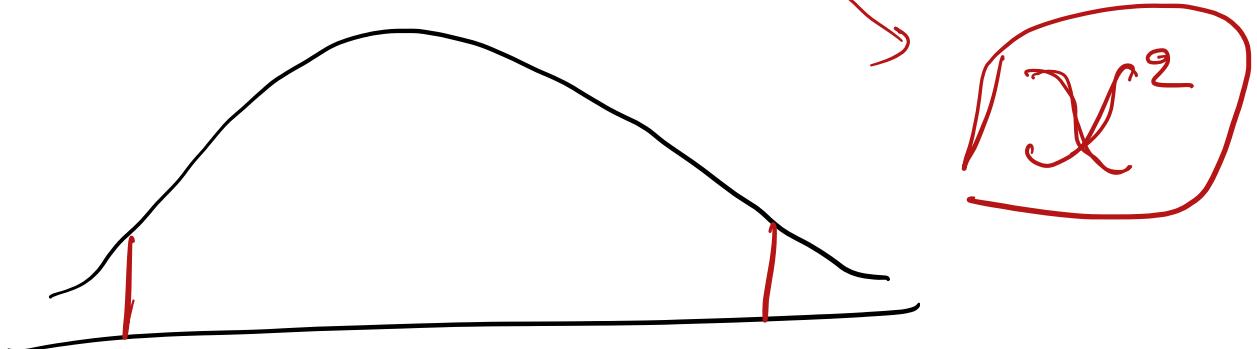
① H_0 = The data meets the distribution
of 2000 census

H_1 = The data does not meet the

② $\alpha = 0.05$ (95% CI)

③ Degree of freedom $\Rightarrow n-1$
 $\Rightarrow 3-1 = 2$

④ Check in the Chi Square table to
find Decision Boundary



If χ^2 is greater than 5.99 then reject the H_0

⑤ Calculate the test

Statistic

$$\chi^2 = \sum \frac{(f - f_e)^2}{f_e}$$

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

$$= 232.94$$

$$\chi^2 = 232.94 > 5.99$$

Reject the Null hypothesis

note:

If P-Value < Significance Value

then we Reject the null hypothesis

Covariance:

Relation b/w values

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

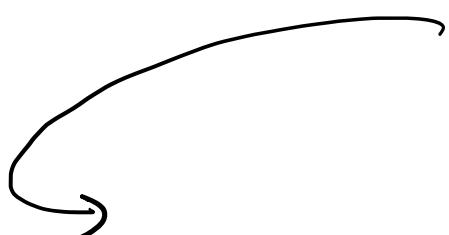
= +ve \Leftrightarrow -ve \Leftrightarrow 0

Note:

This value can be outg in terms of magnitude.

To put a constraint for this we use

pearson Correlation



it restricts all the values b/w [-1 to +1]

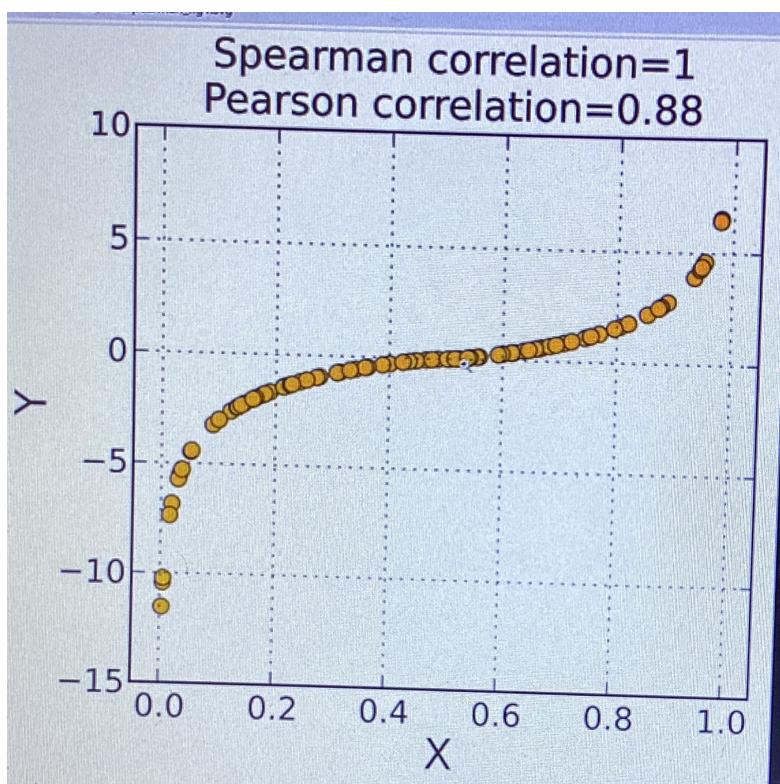
The more it is towards +1 \Rightarrow **+ve Correlation**
The more .. -1 \Rightarrow **-ve Correlation**

$$\rho(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y}$$

$\{-1 \text{ to } +1\}$
range

→ it captures linear properties well.

Spearman's rank correlation:



→ along with linear, non-linear

Properties can also be captured well.

$$S_p(x,y) = \frac{\text{Cov}(R(x), R(y))}{\sqrt{R_x} \sqrt{R_y}}$$

ex)

x (HT)	y (WT)	$R(x)$	$R(y)$
120	75	2	2
160	62	3	3
150	60	4	4
145	55	5	5
180	85	1	1