



Correlation -

⇒ 1) Correlation

↳ Pearson

↳ Spearman

↳ Viz of Correlation

⇒ Relationship b/w 2 Variables -

1) Num vs Categorical → T test ($\leq 2 \text{ Grp}$)
→ Anova Test ($> 2 \text{ Grp}$)

2) Cat vs Cat → Chi Square

3) Num vs Num → Correlation

⇒ T Test → IQ score of S1 to IQ score of S2

Anova

→ Not Multiple Numeric field

→ But multiple groups of same numeric field.

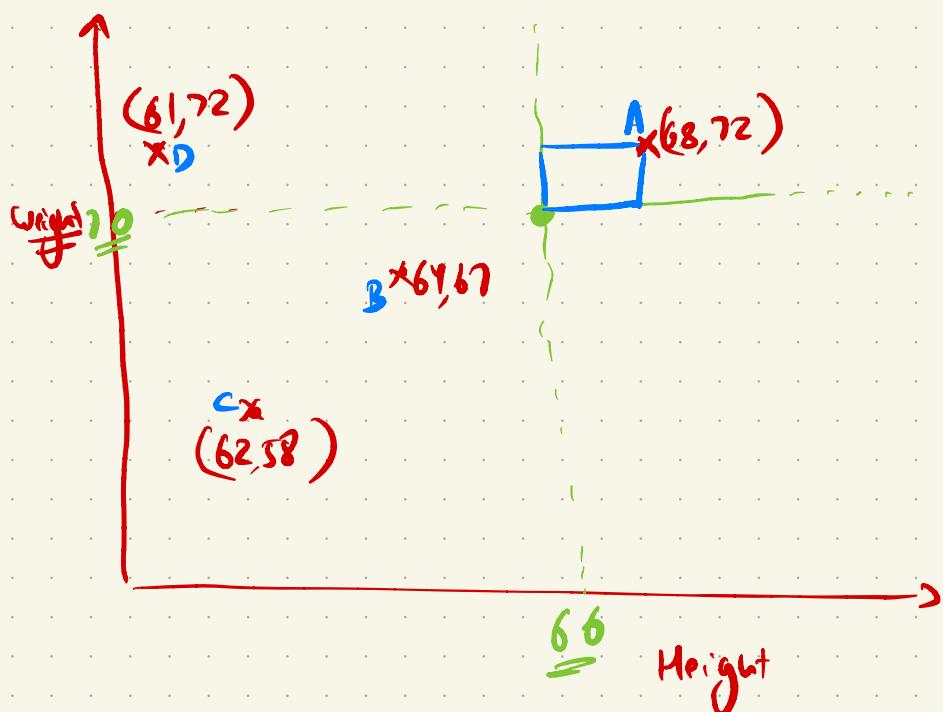
⇒ Num vs Num →

⇒ Height vs Weight →

Height (inches)	Weight (kg)
68	72
62	58
64	67
61	72
70	79
66	61
61	68
65	64
71	80
72	79
$\bar{h} = 66$	$\bar{w} = 70$

~~Height~~ & ~~weight~~ are related
 ↗ Num variable
 ↘ Num variable

⇒ Num vs Num relationship



$$A \Rightarrow (68 - 66)(72 - 70) = 2 \times 2 = 4$$

$$B \Rightarrow (64 - 66)(67 - 70) = -2 \times -3 = 6$$

$$C \Rightarrow (62 - 66)(58 - 70) = -4 \times -12 = 48$$

$$D \Rightarrow (61 - 66)(72 - 70) = -5 \times 2 = -10$$

Area → $(x - \mu_x)(y - \mu_y)$ → $\begin{cases} +ve \\ -ve \end{cases}$

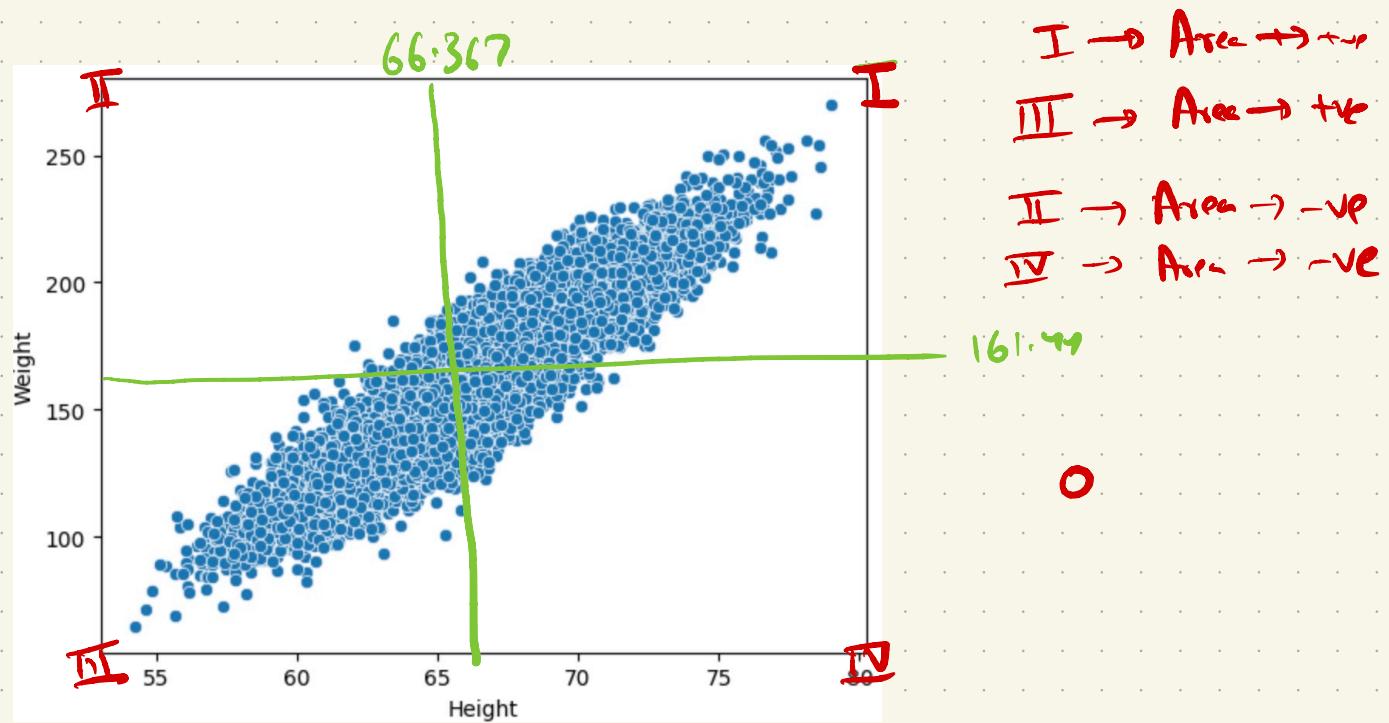
\Rightarrow +ve area \rightarrow +ve relationship

or Height increases when weight increases

-ve area \rightarrow -ve relationship

or Height decreases when weight increases

* \Rightarrow Area can be added \rightarrow slopes can't be added.



of I, III have \gg # of points
of II, IV have

$$\sum (x - \bar{x}_x)(y - \bar{y}_y)$$

Positive
Correlation

$$\frac{\sum_{i=1}^n (h_i - \bar{h})(w_i - \bar{w})}{n} \Rightarrow \text{tve number}$$

\Rightarrow tve if more pts in Ist & IIIrd Quadrant

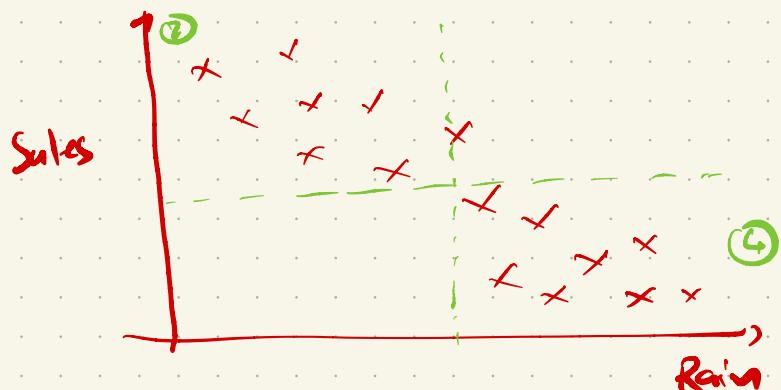
$$\text{Cov}(x, y) = E[(x - E(x))(y - E(y))]$$

$E(x)$ = avg of all x pts

$E(y)$ = avg of all y pts

$$E((x - E(x))(y - E(y))) = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{n}$$

\Rightarrow Ice Cream Sales vs Amount of Rain



Negative Correlation

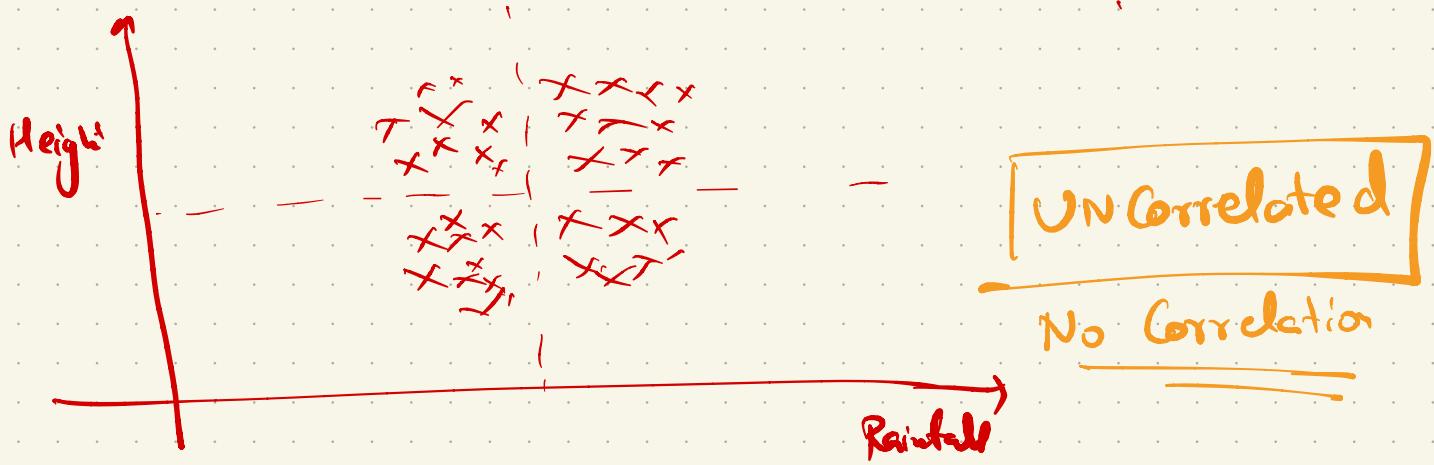
\Rightarrow Most data will lie in II & IV Quadrant \rightarrow

So Covariance \rightarrow ~~pos~~

-ve \checkmark

Rainfall & Ice cream Sales are -ve related.

\Rightarrow Height vs Rainfall



Covariance \rightarrow Sum of areas

positive area \approx negative area

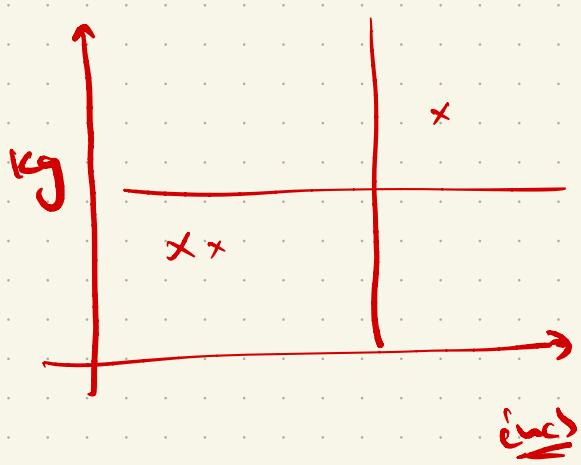
Covariance ≈ 0

\Rightarrow

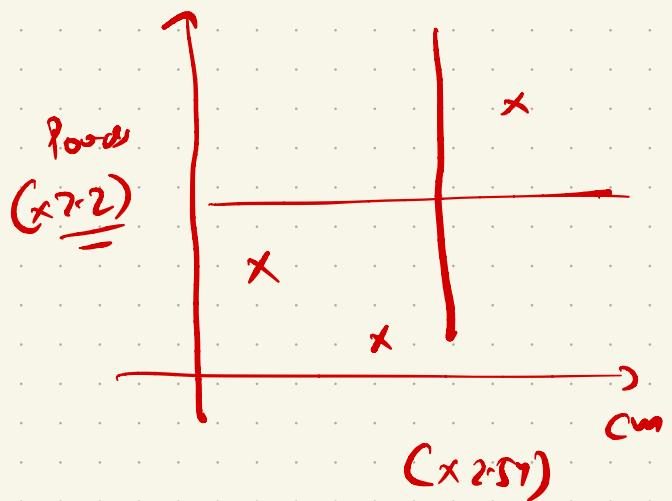
Pearson Correlation \rightarrow

Covariance

① Inch vs Kgs



② Cms vs Pounds



$$(h - \bar{h})(w - \bar{w})$$

$$h = 69 \text{ in}$$

$$\bar{h} = 66 \text{ in}$$

$$w = 72 \text{ kg}$$

$$\bar{w} = 70 \text{ kg}$$

$$\underline{(3)(2)}$$

$$h = 69 \times 2.54$$

$$\bar{h} = 66 \times 2.54$$

$$w = 72 \times 2.2$$

$$\bar{w} = 70 \times 2.2$$

$$\underline{(3 \times 2.54)(2 \times 2.2)}$$

Covariance \rightarrow dependent on the units \rightarrow

$1 \text{ kg} \approx 2.2 \text{ pound}$

$1 \text{ cm} \approx 2.54 \text{ cm}$

SPI
1

$$\frac{(h - \bar{h})}{\sigma_h} \frac{(\omega - \bar{\omega})}{\sigma_\omega}$$

$$\left(\frac{h - \bar{h}}{\sigma_h} \right) \left(\frac{\omega - \bar{\omega}}{\sigma_\omega} \right)$$

SPI
2

$$\frac{(h - \bar{h})}{\sigma_h} \quad \frac{(\omega - \bar{\omega})}{\sigma_\omega}$$

$$\text{Corr} = \frac{\text{Covariance}(x, y)}{\sigma_x \sigma_y} \rightarrow$$

$$f(\rho) = \frac{1}{n} \sum (h_i - \bar{h}) \frac{(\omega_i - \bar{\omega})}{\sigma_h \sigma_\omega}$$

↳ Correlation

Coefficient

Pearson's
Correlation

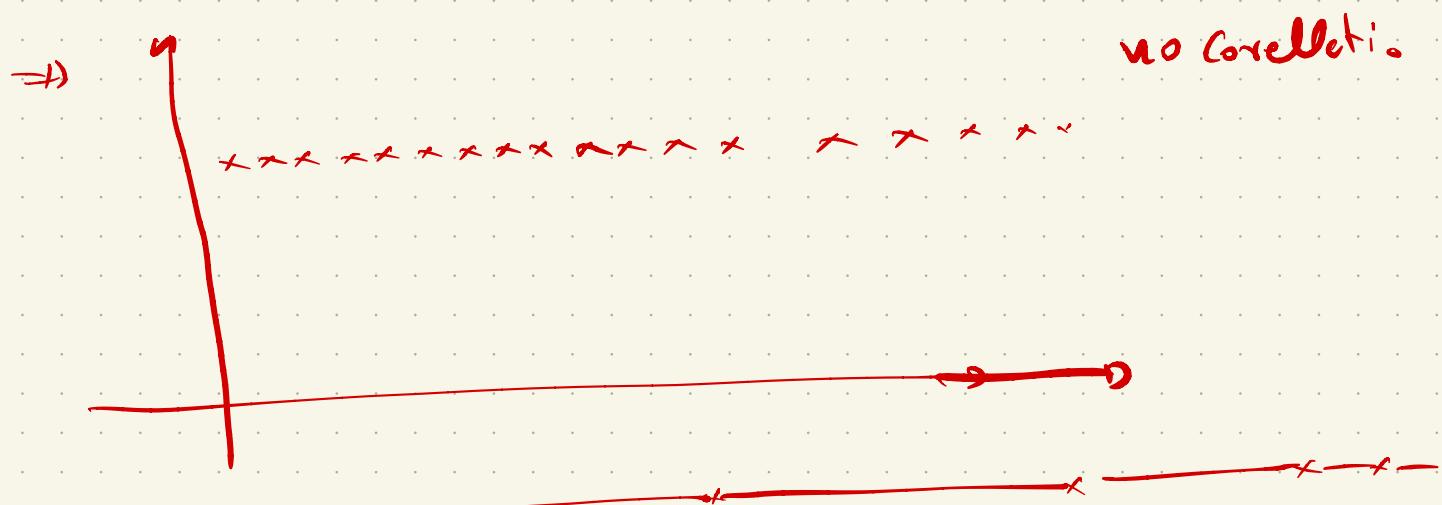
$$\rightarrow -1 \leq f \leq 1$$

$f \approx -1 \Rightarrow$ High -ve Correlt

$f \approx 1 \Rightarrow$ High +ve Correlt

$$\rho \approx 0 \Rightarrow \text{no Correlation}$$

$\text{Corr}(\rho)$ is a unit free measurement \rightarrow



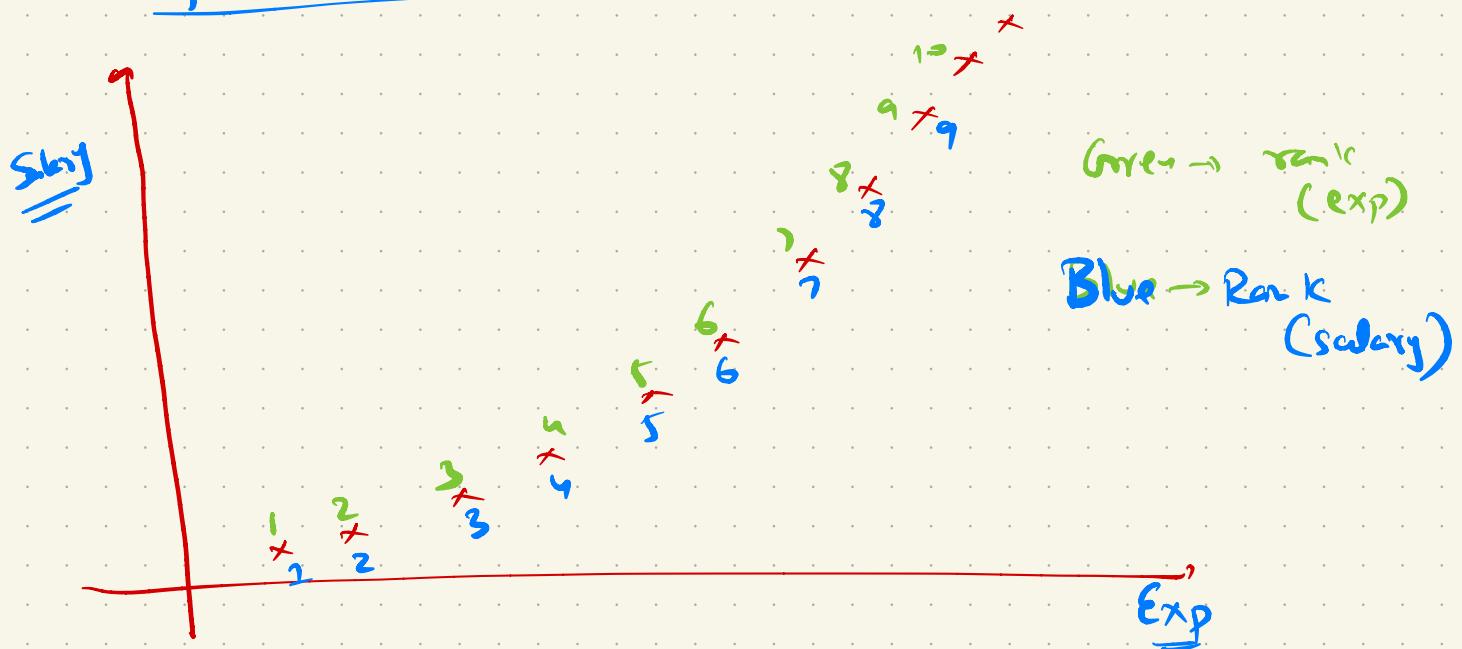
\Rightarrow Salary vs Experience \rightarrow



\Rightarrow Total area may be ≈ 0 or $-up$

\Rightarrow Pearson's Correlation \rightarrow works very well with linear dependencies

\Rightarrow Spearman Correlation \rightarrow



Rank (x)	Rank (r)
1	1
2	2
3	3
4	4
5	5
6	6

$S_{\text{rank } x \text{ rank } r}$

$S \approx 1 \rightarrow$

Spearman coefficient

\Rightarrow Spearman would be high

\Rightarrow Pearson could be low or even -ve.

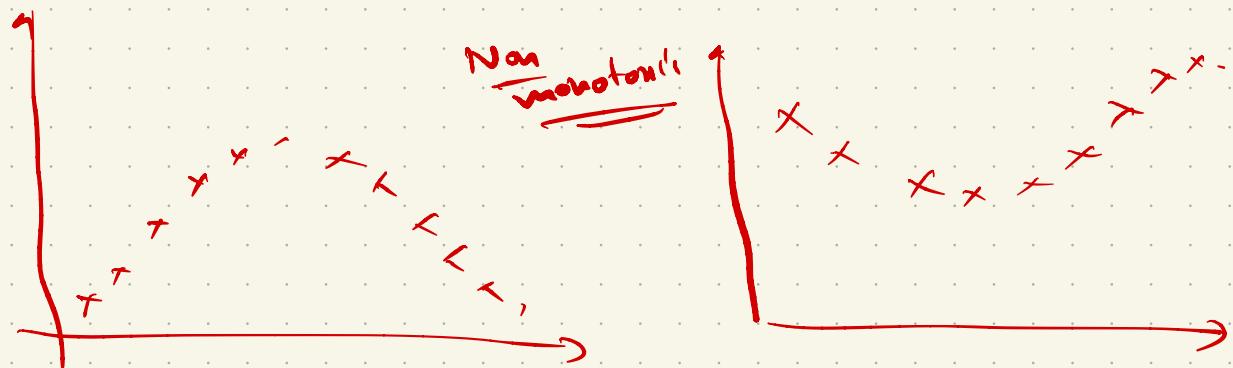
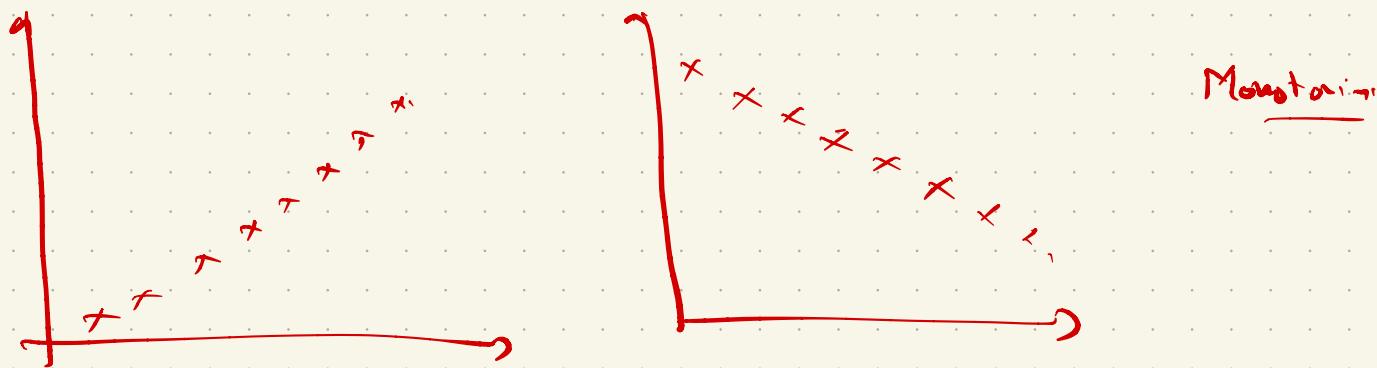


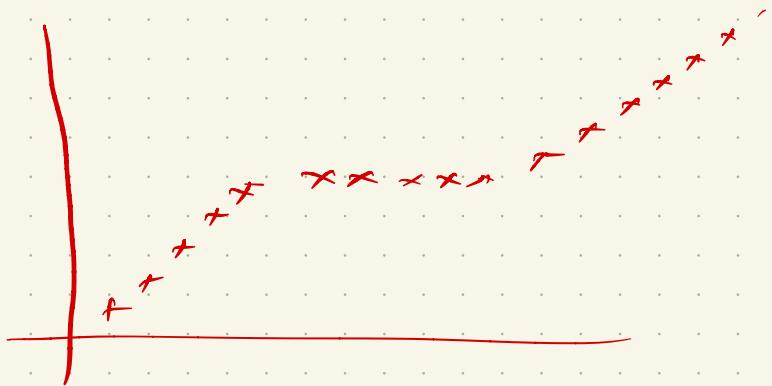
$$r \approx 0$$

Both Pearson & Spearman $r = 0$

Spearman Correlation \rightarrow Monotonic increasing/decreasing.

if x increases y should also increase \rightarrow





Ranks are not

Heights are not

Specimen → we have -ve corre.

Rank X	Rank Y
1	10
2	9
3	8
4	7
5	6

Hw → Calculate Specimen Correlation using basic python
and `•Corr()`.

⇒ a way to get rank → ??

\Rightarrow Hypothesis framework

$H_0 \rightarrow$ Fields are uncorrelated

$H_a \rightarrow$ one field is dependent on another.

Vizualize Correlation \rightarrow

\Rightarrow Parametric \rightarrow Some persons are assumed

Non parametric \rightarrow

\Rightarrow Pearson \rightarrow Parametric

Correlation Test \rightarrow Bivariate Normal distribution

\Rightarrow Spearman \rightarrow Non parametric

\rightarrow on the basis of rank

\Rightarrow Both Height & weight follow a normal dist'

\Rightarrow Bivariate \rightarrow X, Y are Random variables

for every real no. a, b

such
that

$$\textcircled{2} \quad \underline{\underline{ax + by}}$$

$Z \rightarrow$ follows a normal distribution

p-value \rightarrow tells that what are the chances that data is unrelated.

p-value $\approx 0 \rightarrow$ data is Highly Correlated

p-value $< 0.05 \rightarrow$ data is related

else
data is not related

Collab Link : <https://colab.research.google.com/drive/1CTwBoTB7E11KoYiDGlugDfq!HFkxdZHw?usp=sharing>