

→ Pearson Correlation coefficient

1. Co-Variance $\rightarrow \text{Cov}(x, y)$

$$= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) \cdot (y_i - \mu_y)$$

2. Pearson Correlation $= \rho_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y}$

3. Spearman^{rank} Correlation

1. Covariance

Size Price

1200sqm \$100k

1800sqm \$200k

1500sqm \$150k

Is there a relationship between size and Price?

Quantifying the relationship

$$\underset{\substack{\downarrow \\ x}}{\text{Cov}(\text{size}, \text{Price})} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) \cdot (y_i - \mu_y)$$

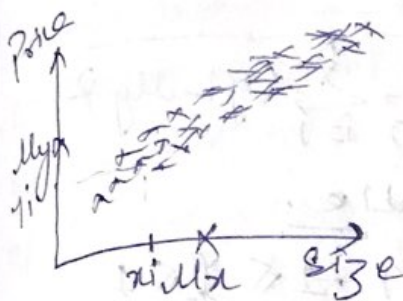
$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) \cdot (x_i - \mu_x)$$

if x increases and y increase as x increases and y decreases

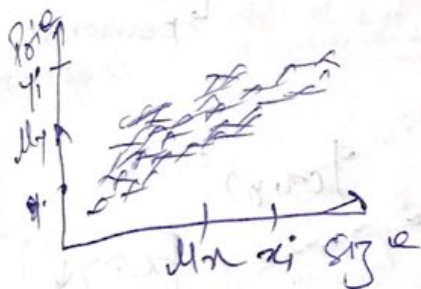
$$x \uparrow \text{ and } y \uparrow = \boxed{} +ve$$

$$x \uparrow \text{ and } y \downarrow = \boxed{} -ve$$



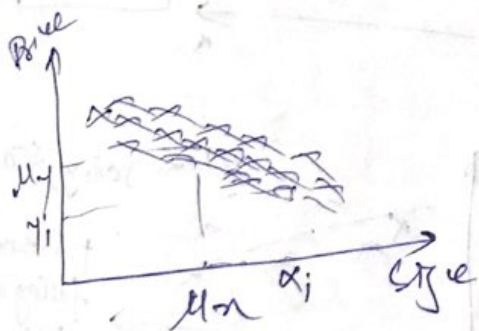
$$\frac{1}{n} (x_i - \mu_x)(y_i - \mu_y)$$

$$\frac{1}{n} (-ve) \times (-ve) = +ve$$



$$\frac{1}{n} (x_i - \mu_x)(y_i - \mu_y)$$

$$\frac{1}{n} (+ve) \times (+ve) = +ve$$



$$\frac{1}{n} (x_i - \mu_x)(y_i - \mu_y)$$

$$\frac{1}{n} (+ve) \times (-ve) = -ve$$

So, Covariance tells us about the direction of relationship between two random variable. But doesn't tell ~~us~~ about the degree to which it is +vely related or -vely related.

For that we have Correlation.

→ Pearson correlation coefficient - $r_{xy} = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$

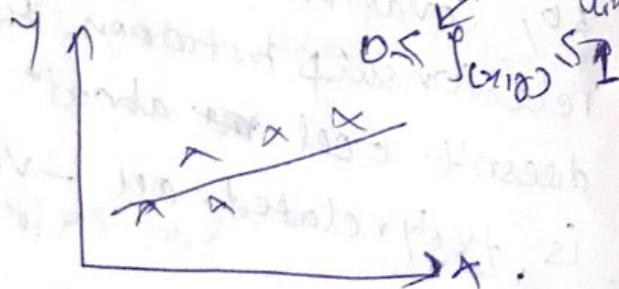
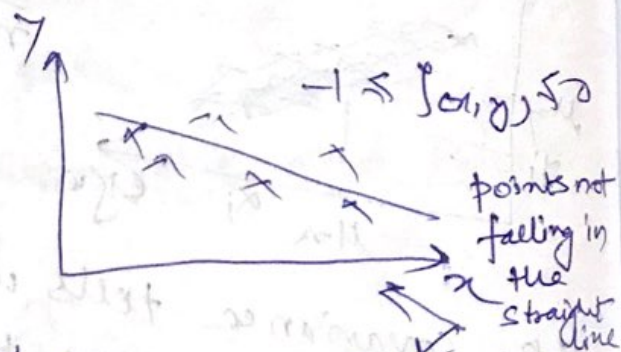
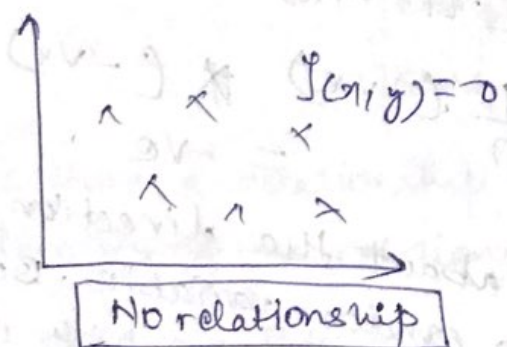
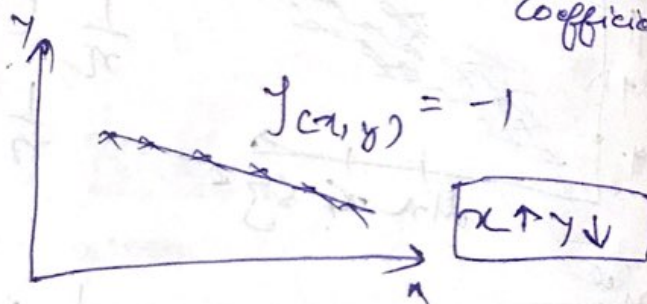
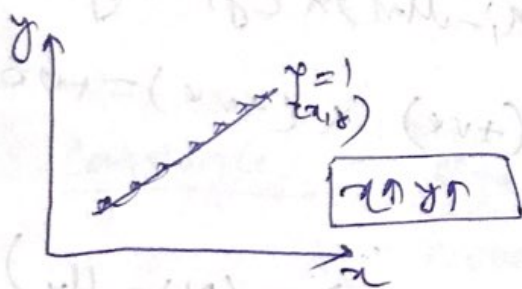
σ_x → σ_y
 stan. dev. → stan. dev.

$$= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \mu_y)^2}}$$

Tells the strength of two variables.

the value ranges between $-1 \leq r_{xy} \leq 1$

↳ Pearson's Coefficient



→ Spearman rank correlation coefficient

$$\rho_{(x,y)} = \frac{\text{Cov}(r_{gx}, r_{gy})}{\sigma_{gx} \sigma_{gy}}$$

↳ we find the correlation between rank of x and rank of y .

In this example, the raw data in the table below is used to calculate the correlation between the IQ of a person with the no. of hours spent in front of TV per week.

<u>IQ x_i</u>	<u>Hours of TV/week y_i</u>
106	7
86	2
86	2
101	50
99	28
103	29
97	20
113	12
112	6
110	7

1. sort the data by the 1st column (x_i). create a new column x_i and assign it the ranked values, 1, 2, 3, ..., n .
2. create a 4th column and sort the y_i value and place it accordingly against x_i value.
3. create a 5th col d_i to hold the differences between two rank columns (x_i and y_i)
4. create one final column d_i^2 to hold the value of column d_i^2 .

<u>ID x_i</u>	<u>Hours TV/week y_i</u>	<u>rank x_i</u>	<u>rank y_i</u>	<u>d_i</u>	<u>d_i^2</u>
86	0	1	1	0	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	29	6	9	-3	9
106	7	7	3	+4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

$$r = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} \quad \sum d_i^2 = 194$$

$$= 1 - \frac{6 \times 194}{10(10^2-1)} = -0.17575$$

- Pearson focuses mainly on linear relationship linear aspect of data.
- where as Spearman rank correlation focuses on the outliers also and gives the correlation value on non-linear relationships also.
- In python heatmap it uses this technique only.