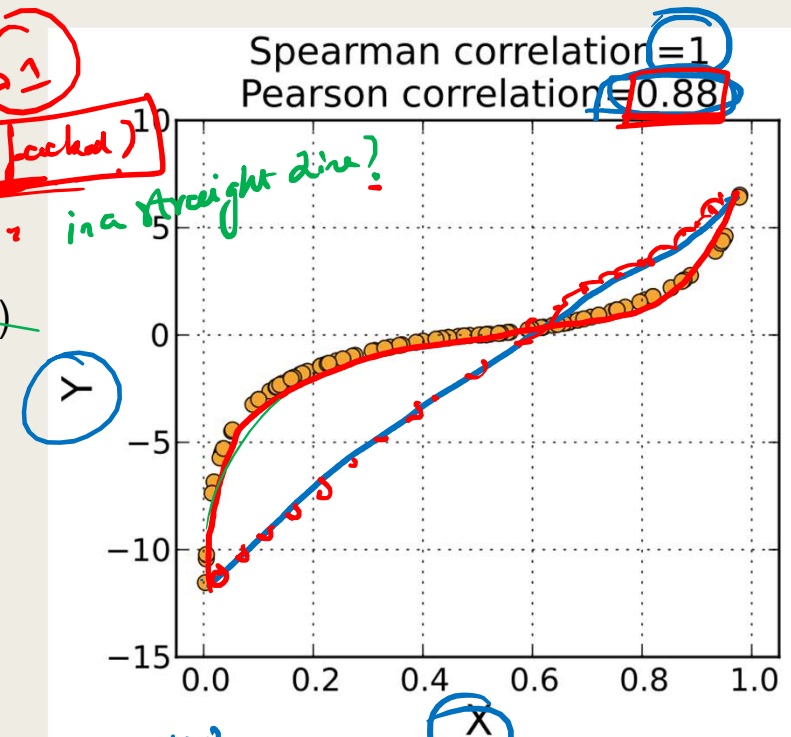


Spearman Rank Correlation

- Used for Non-Linear Variables
- Spearman Corr Coeff = Pearson Corr coeff (rank variables)
- In Python: `DataFrame.corr(method='spearman')`

Non-linear
correlation coeff.

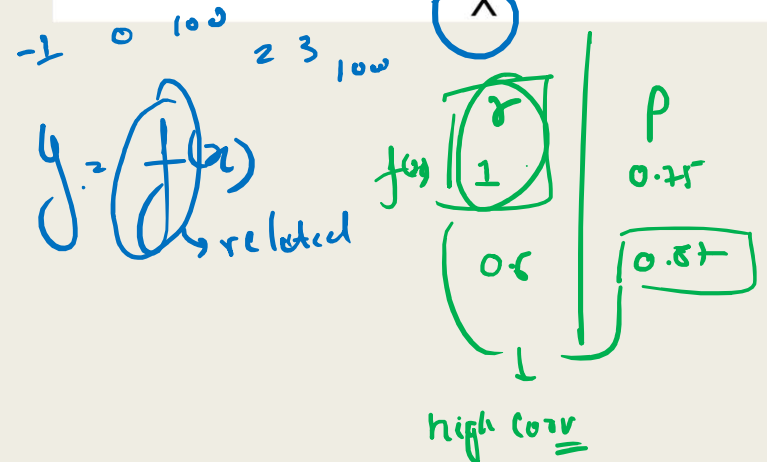
Are they closely linked?
Are they in a straight line?



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

Denoted by rho.

Correlation → linear
→ non-linear



Steps for Spearman Correlation Coefficient

1. Create a new column for rank(x) and assign the rank of each variable.
2. Assign the rank of 2nd variable in a new column rank(y).
3. Calculate the difference in rank of both the variables = d.
4. Calculate the d-squared.
5. Add up d-squared score.
6. Put in the formula provided:

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

of observations.

diff. in Rank.
g values

Question: The scores for 10 students in English and Maths are as follows:

	Marks									
English	56	75	45	71	62	64	58	80	76	61
Maths	66	70	40	60	65	56	59	77	67	63

Compute the Spearman rank correlation.

Solution:

Step 1, 2, 3 and 4:

$H \rightarrow L$

$L \rightarrow H$

$\geq 0.75 \rightarrow n$

< 0.75 with

English (mark)	Maths (mark)	Rank (English)	Rank (maths)	d	d ²
56	66	9 2	4 7	5 -5	25
75	70	3 8	2 9	1 -1	1
45	40	10 1	10 1	0 0	0
71	60	4 7	7 4	-3 3	9
62	65	6 5	5 6	1 -1	1
64	56	5 6	9 2	-4 4	16
58	59	8 3	8 3	0 0	0
80	77	1 10	1 10	0 0	0
76	67	2 9	3 8	-1 1	1
61	63	7 4	6 5	-1 1	1

Solution Contd.

Step 5:

$$\sum d_i^2 = 25 + 1 + 9 + 1 + 16 + 1 + 1 = 54$$

Step 6:

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

$$\rho = 1 - \frac{6 \times 54}{10(10^2 - 1)}$$

$$\rho = 1 - \frac{324}{990}$$

$$\rho = 1 - 0.33$$

$$\rho = 0.67$$

Weak +ve non linear corr
 \uparrow Eng \uparrow math \downarrow

$\begin{matrix} \uparrow & \uparrow \\ x & y \end{matrix}$
 $\begin{matrix} x & y \\ \downarrow & \downarrow \end{matrix}$
 $\begin{matrix} x \uparrow & y \downarrow \\ \underline{\quad} & \underline{\quad} \end{matrix}$
 $\begin{matrix} x \downarrow & y \uparrow \\ \underline{\quad} & \underline{\quad} \end{matrix}$
 $r > \rho \rightarrow \text{linear}$
 $r \leq \rho \rightarrow \text{Non-linear}$

Hence, the Spearman Rank Coefficient is 0.67.

Points to Ponder (Correlation):

Correlation does not ~~Causation~~

Example: A few years ago a survey of employees found a strong positive correlation between "Studying an external course" and Sick Leave Days.

Does this mean?

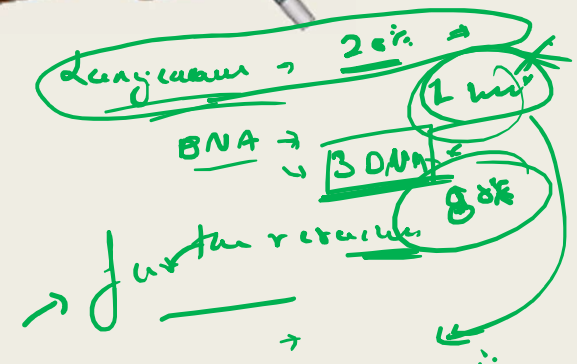
- Studying makes them sick?
- Sick people study a lot?
- Or did they lie about being sick so they can study more?

Without further research we can't be sure why. :-D

Example: Poor suburbs are more likely to have high pollution.

Why?

- Do poor people make pollution?
- Are polluted suburbs the only place poor people can afford?
- Is it a common link, such as factories with low paying jobs and lots of pollution?





Recap – Descriptive Statistics

- Statistics? Its Importance. Population vs Sample.
- Types of variable – (Quantitative, Categorical), - (Ordinal, Nominal), (discrete, continuous).
- Types of charts – Pie, Donut, Line, Scatterplot, Histogram, Bar chart, Box-plot
- Descriptive Stats:

Measure of central tendency – Mean , Median & Mode

Measures of Dispersion/spread – Standard Deviation, Variance, range & IQR

Measures of symmetry – Skewness and Kurtosis

01:01 PM

↳ sampling

techniques

- 5 Number Summary – Box Plot (Box and Whiskers)
- Effect of transformation on central tendency and spread.
- Outliers? How to detect? Modified Boxplot. (IQR Method)
- Covariance, Correlation. Pearson's Correlation Coefficient. Nature & Strength of Correlation. How to calculate Pearson's Correlation Coefficient and Spearman's Rank Correlation Coefficient.