

Advance Statistics

① Sample vs Population

	<u>Sample</u> n	<u>Population</u> N
<u>numbers</u> →	Statistics (a subset of population)	Parameters (collection of all items of interest)

Why prefer sample

- easy to manage
- less time consuming
- less costly

A sample has to be random !

↓
(Each member is equally likely to be chosen)

② Covariance

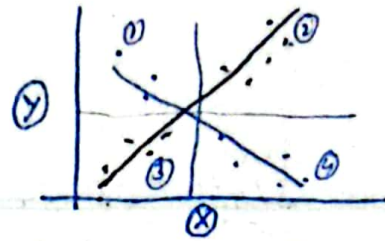
$$\text{Cov}(x, y) = E[(x - E(x)) \times (y - E(y))]$$

$$\text{Cov}_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

(It measures the direction of relationship) (+ive, -ive, no)
b/w two variables

① & ② → -ive Cov values

② & ③ → +ive Cov values



→ Covariance value doesn't tell us about steepness of the line or how close the data points are to the line.

(It just tells us about nature of relationship)

→ Cov value can be changed even if relationship doesn't (i.e. by changing scale value)

→ Used for PCA (Principle Component Analysis)

→ Just like correlation but its values can be varied more than (-1 to 1)

③ Pearson's Correlation (r)

- "Explain relationship b/w two numerical variables"

- Value → (-1 to 1)

- Doesn't depend on scale of data

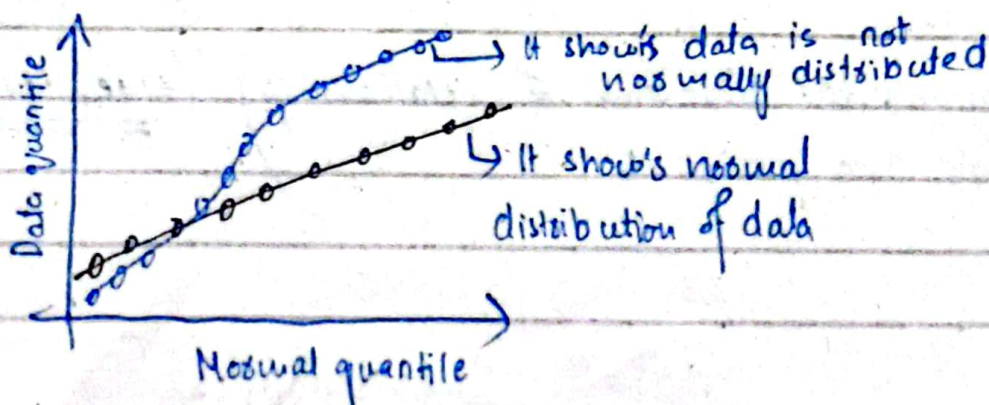
- Educated guesses \propto Amount of data points

- Educated guesses \propto $1/p$ -value

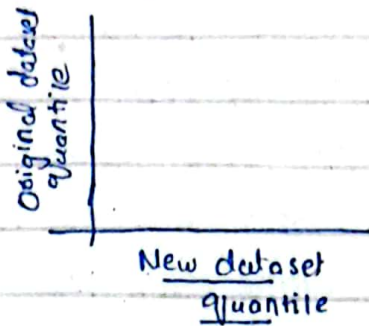
$$r = \frac{\text{Cov}(x, y)}{S_x \cdot S_y}$$

→ This denominator value ensures that the r value becomes independent of scale.

④ Quantile - Quantile Plots (Q-Q plots)



We can also compose two data sets:



⑤ Confidence Interval (p-value)

Value \rightarrow 0.05 or 5% (universally accepted)

understanding through bootstrapping

\rightarrow Take 15 books and weigh them

\rightarrow Calculate mean

\rightarrow Now take random sample of these 15 books which is again 15 in number (double or triple pick is OK!)

It is done to estimate actual mean of book weight of all available books in shop.

\rightarrow Bootstrapping # \rightarrow 10,000

Now here p-value of 0.05 or confidence interval of 95% means

"All the space values in which 95% of the mean's are present"

\Rightarrow Important in hypothesis testing ✓

(vi) Null hypothesis / Hypothesis testing

(vii) t-test / Chi-square test / Anova / \rightarrow when to use?

1. One categorical data set \rightarrow One sample proportion test
if it gives value less than p-value of 0.05 then
 \rightarrow reject null hypothesis and accept alternative hypothesis

Two categorical data set \rightarrow Chi-square test

One numerical data set \rightarrow T-test

Two numerical data set \rightarrow Correlation test (r)

One numerical & One or two categorical data set \rightarrow ANOVA

(further subsets
present in data)

If,

$P < \alpha$ H_0 is rejected

If,

$P \geq \alpha$ H_0 is accepted