

IV Covariance & Correlation :-

Suppose we have 2 features on 2 Random Variables (X, Y) :-

X	Y
2	3
4	5
6	7
8	9

So, can we find a relationship b/w X and Y. (ie like if $X \uparrow$ then $Y \uparrow$ etc.)

Plotting: $\begin{matrix} X \uparrow, Y \uparrow \\ X \downarrow, Y \downarrow \end{matrix} \equiv \begin{matrix} x & x & x & x & x \\ & x & & & \\ & & x & & \end{matrix}$

Plotting: $\begin{matrix} X \downarrow, Y \uparrow \\ X \uparrow, Y \downarrow \end{matrix} \equiv \begin{matrix} x & & & & \\ & x & & & \\ & & x & & \\ & & & x & \\ & & & & x \end{matrix}$

5.1 Covariance (X, Y) or $\text{Cov}(X, Y)$:-

now, covariance w/ random variables X and Y is given by formulae :-

$$\text{Cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$$

Note \rightarrow what the difference b/w Variance and Covariance

Variance (X^2) = $\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$ & Covariance (X, Y) = $\frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$
(signifies spread)

So; Variance is nothing but Covariance of random variable X with itself i.e $\text{Cov}(X, X)$

i.e

$$\text{Cov}(X, X) = \frac{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}{n-1} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} = \text{Var}(X^2)$$

1) For ($X \uparrow, Y \uparrow / X \downarrow, Y \downarrow$) Relationship: Covariance = +ve

X	Y
2	3
4	5
6	7
$\bar{X}=4$	$\bar{Y}=5$

$$\left. \begin{matrix} \text{Cov}(X, Y) = \frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{3-1} \\ = \frac{4 + 0 + 4}{2} = \boxed{4} \text{ (+ve)} \end{matrix} \right\}$$

\therefore X and Y are having a +ve covariance

Qy Example where Covariance can be used?

Sol \rightarrow Suppose we have a dataset where Price and size of different houses is given, and In such case we have to predict the price of the house

So, Covariance can be used to find the relationship b/w the size and prices of the house.

Advantages and Disadvantages of Covariance:-

\rightarrow Advantage:- Using covariance we can find a relation b/w X and Y. It can be +ve or -ve

\rightarrow Disadvantage:- Covariance doesnot have a specific limit value

Ex:- Suppose we have X, Y, Z

if $\text{Cov}(X, Y) = +100$ & $\text{Cov}(X, Z) = +500$

it does not mean X is heavily correlated to Z as compared to Y coz there is no specific limit to Cov value to compare high or low

So,

in order to overcome this advantage of covariance - we specifically use another correlation technique i.e. "Pearson Correlation Coefficient"

5.24 Pearson Correlation Coefficient (P) [range: -1 to +1]

As its value is restricted b/w -1 and +1 \therefore , now we can differentiate which variables are highly correlated and which are not.

$$\text{Correlation}(X, Y) (\rho_{X,Y}) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = [-1 \text{ to } +1]$$

- more the value is towards +1, the more +vely correlated X & Y are.
- more the value of $\rho_{X,Y}$ is towards -1, the more -vely correlated X & Y are

Qy Where can this correlation concept is applied in data science project?

Soln Suppose a dataset with 1000 features, now, for an ML project it is not feasible to take all of the features

Ex: dataset =

Size of house	No of Rooms	Location	No. of ppl staying	Price
Independent features				Dependent feature

no. using Feature Selection :-

if $\text{Cov}(\text{No. of ppl, Price}) \approx \text{or near } 0$ then it can be removed from our ML model.

5.24 Spearman Rank Correlation :- (γ_s)

Sometime it is considered to be better than Pearson Correlation. here instead of considering the value of X & Y we consider the rank of both X and Y i.e

$$\gamma_s = \frac{\text{Cov}[R(X), R(Y)]}{\sigma_{R(X)} \sigma_{R(Y)}}$$

<u>ex:</u>	if	x	y	$R(x)$	$R(y)$
		1	2	5	5
		3	4	4	4
		5	6	3	3
		7	8	2	1
		0	7	6	2
		8	1	1	6