

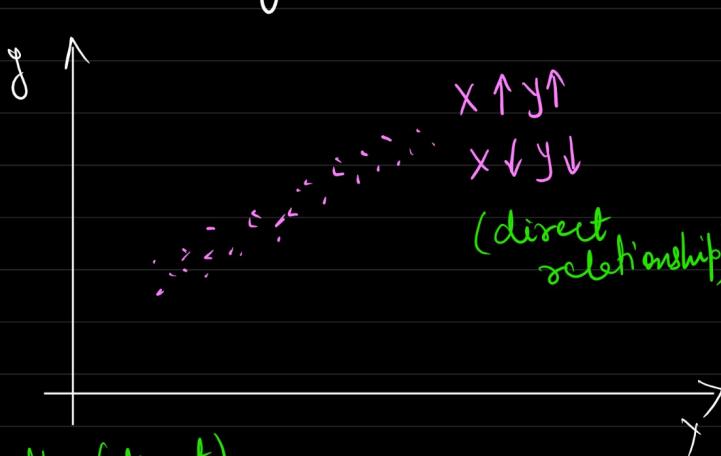
Covariance and Correlation

$$\left\{ \begin{array}{l} y \uparrow x \uparrow \\ y \uparrow x \downarrow \\ y \downarrow x \uparrow \\ y \downarrow x \downarrow \end{array} \right\}$$

⇒ relationship b/w two features.

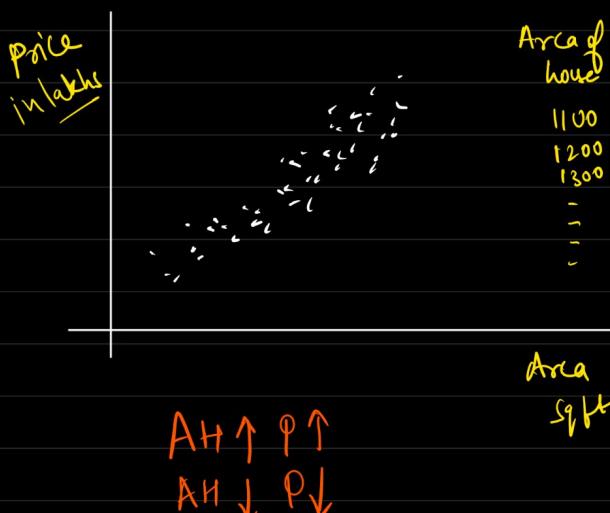
(x)	(y)
total transaction count	5
1 lakh	4
50K.	-
-	-
-	-
-	-

* Understanding the relationship

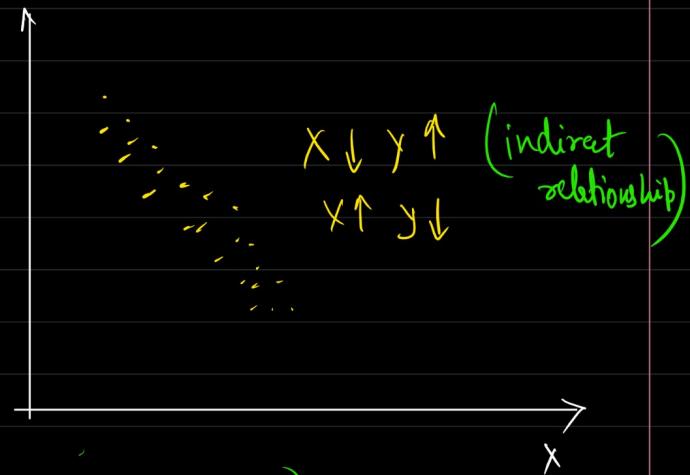


Example (direct)

* Predict price of house based on area of house.

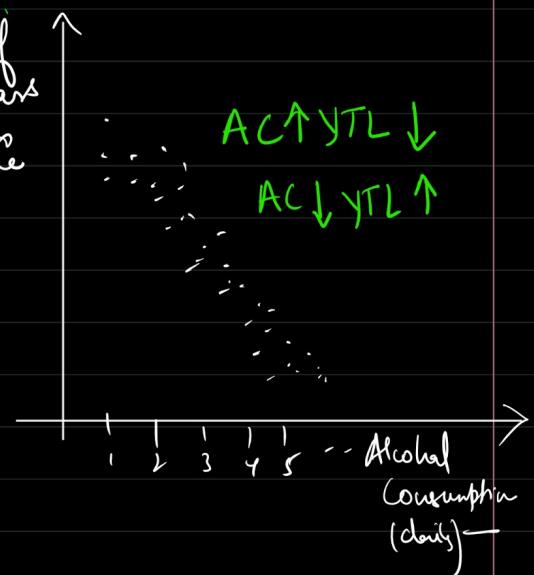


Area of house	Price (lakh)
1100	80
1200	85
1300	90
-	-
-	-
-	-



Example (Indirect)

AC ↑ YTL ↓
AC ↓ YTL ↑



* To quantify / measure the relationship

① Covariance

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

where x &
 y are
two features

Cov - Variance
 ↓
 Existing together
 ↓
 Spread of data

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

↓
 Spread of data
 ↓
 relationship of a
 feature with itself.

* Covariance means, you are trying to understand the relationship of a feature with respect to other feature.

* Covariance - relationship b/w two variables.

→ $X \uparrow Y \uparrow$
 $X \downarrow Y \downarrow$

or

$X \uparrow Y \downarrow$
 $X \downarrow Y \uparrow$

(+ve covariance)

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

(-ve covariance)

X	Y	$x - \bar{x}$	$y - \bar{y}$	
-2	3	(-2 - 3)	(3 - 5.5)	$\Rightarrow (-2-3)(3-5.5) + (3-3)(5-5.5) + (6-3)(6-5.5) + (-3)(8-5.5)$
-3	5	(-3 - 3)	(5 - 5.5)	<hr/>
-6	6	(-6 - 3)	(6 - 5.5)	$4-1$
-1	8	(-1 - 3)	(8 - 5.5)	$\Rightarrow (-1)(-2.5) + 0 + 3 \times 0.5 + (-2) \times 2.5$
		$\bar{x} = 3$	$\bar{y} = 5.5$	<hr/>

$= \frac{2.5 + 0 + 1.5 - 5}{3} = -\frac{1}{3} = 0.33$

$-0.33 \rightarrow$ interpretation \rightarrow The two features $x \& y$ are negatively selected.

$$\begin{array}{cc}
 \checkmark & \checkmark \\
 2 & 3 \\
 4 & 5 \\
 6 & 7 \\
 \hline
 \bar{x} = 4 & \bar{y} = 5
 \end{array}
 \quad \text{Cov}(x,y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$= \frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{3-1}$$

$$= \frac{4+0+4}{2} = 4 \text{ (+ve)}$$

$$\text{Cov}(x,y) = 4$$

interpretation $\rightarrow x \& y$ are having a positive relationship.

* Advantage

\rightarrow Relationship b/w x, y
+ve, -ve.

* Disadvantage of Covariance

$$\begin{array}{c|c}
 \checkmark & \checkmark \\
 \hline
 \text{Cov} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} & x | 1 \\
 & - | 2 \\
 & - | 3 \\
 & - | 4 \\
 & - | 5 \\
 & - | 6 \\
 & - | 7 \\
 & - | 8 \\
 & - | 9 \\
 & - | 10
 \end{array}$$

① $x \quad y$

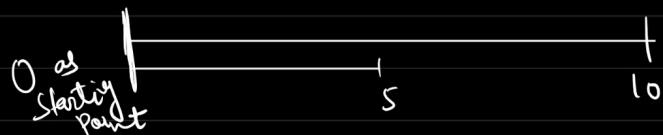
A B

$$\text{Cov}(x,y) = 50$$

$$\text{Cov}(A,B) = 100$$

Can I say the relationship b/w x and y is twice of A and B ?? \Rightarrow No

$\text{Cov}(x,y)$
 \Downarrow
range $\rightarrow -\infty$ to ∞ \rightarrow You cannot compare the strength of relationship.



→ No comparison of strength of relationship in covariance.

→ No any standardised scale to interpret the strength.

② Covariance has dimension

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

n is no. of rows.

$$\text{Cov}(\text{transaction}, \text{ht}) = \text{Rs. ft} \Rightarrow 450 \text{ Rs. ft.}$$

$$\text{Cov}(\text{ht}, \text{wt}) = \text{ft} \cdot \text{kg} \Rightarrow 600 \text{ ft} \cdot \text{kg.}$$

x_1 height	x_2 Weight	x_3 transaction amount
(ft)	(kg)	(Rs)
-	-	-
-	-	-
-	-	-
-	-	-

* we can not compare two different dimension.



$$\frac{\text{Cov}(x, y)}{\text{Rs ft}}$$

$\frac{\text{Cov}(y, z)}{\text{ft kg}}$
→ Not comparable → different dimension.

Soln

→ -1 to 1
→ dimensionless

② Pearson Correlation Coefficient [-1 to 1]

$$f(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \underline{[-1 \text{ to } 1]}$$

→ dimensionless quantity

↑
No relationship
—1 .0 1.
 ← →
 the more negatively correlated the features are.
 (indirect)
 The more positively correlated features are.
 (direct)

$$\rightarrow \int_{x,y} = 0.4 , \quad \int_{A,B} = 0.8$$

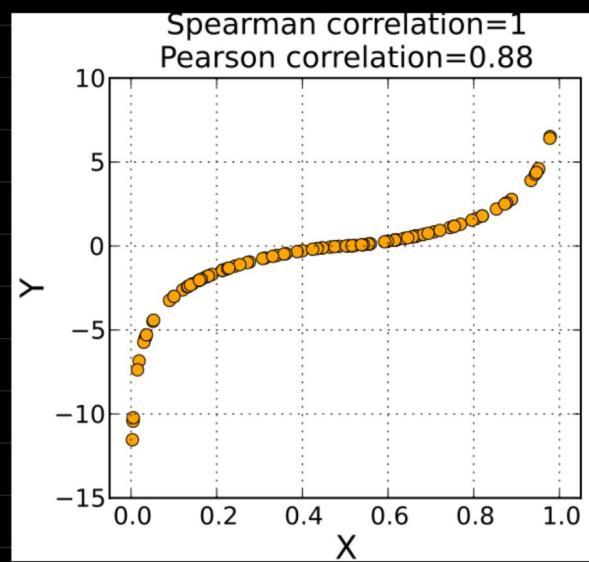
→ feature A, B is highly correlated as compared to x, y.

$$\rightarrow \int_{x,y} = -0.2 \quad \int_{A,B} = -0.5$$

$x \uparrow y \downarrow$ $A \uparrow B \downarrow$
 $y \downarrow x \uparrow$ $A \downarrow B \uparrow$

A, B is more negative correlated as compared to x, y

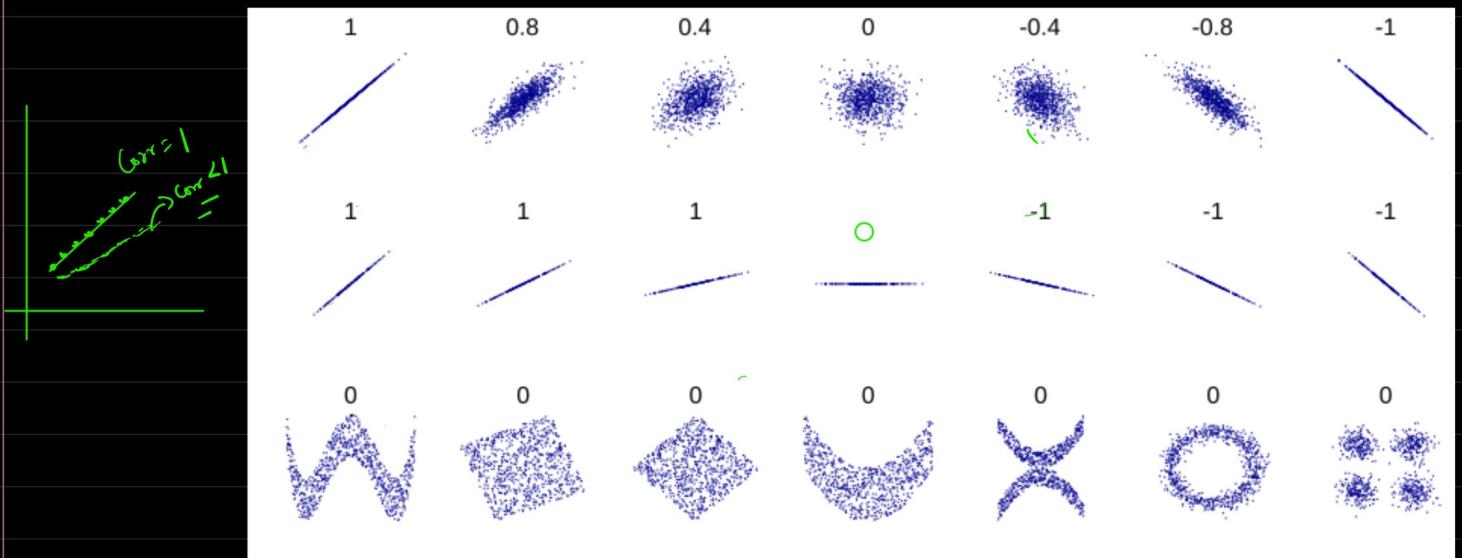
* Pearson correlation coefficient always measures the linear relationship



Person correlation

$$\left\{ \begin{array}{l} x \uparrow y \uparrow \\ x \downarrow y \downarrow \\ x \uparrow y \downarrow \\ x \downarrow y \uparrow \end{array} \right.$$

$\Rightarrow x \uparrow y \uparrow$



- * Correlation is not slope
- * Non linear data will have Pearson correlation 0.
- * What to do for understanding Non-linear relationship

Spearman Rank Correlation

$$\rho_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

$R(x)$ - Rank of x
 $R(y)$ - Rank y .

x	y	$R(x)$	$R(y)$
-5	6	3	1
-7	4	2	2
-8	3	1	3
1	1	5	5
2	2	4	4

$x \rightarrow -5, -7, -8, 1, 2 \Rightarrow$ Sort the value $1, 2, 5, 7, 8$
 Highest no will be rank 1

* Use Case

$x_1 \quad x_2 \quad x_3 \quad \dots \quad x_{1000}$ $\stackrel{y}{\text{Price}}$

$x_1 - y$

$\text{Corr} \approx 0$

→ Correlation is
Used for feature
Selection in
ML modelling

$x_2 - y$

$x_3 - y$

$x_4 - y$