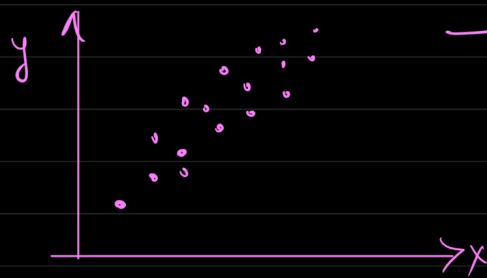


* Covariance and Correlation

$\left\{ \begin{array}{l} Y \uparrow X \uparrow \\ Y \uparrow X \downarrow \\ Y \downarrow X \uparrow \\ Y \downarrow X \downarrow \end{array} \right.$

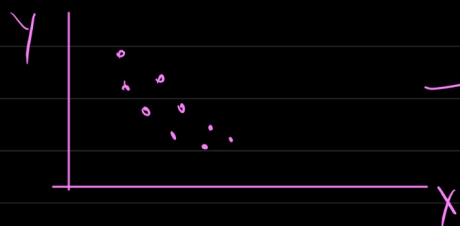
(X)	(Y)
transaction amount	transaction count
—	—
—	—
—	—
—	—

Understanding the relationship



→ relationship b/w x, y is direct

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

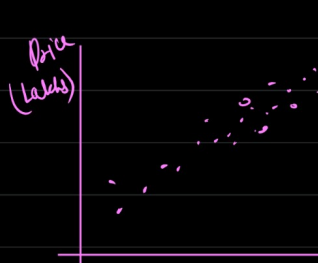


→ indirect relationship

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

* Example (direct)

Predict price of house based on area of house



Area of house (sqft)

Price

$A \uparrow P \uparrow$

$A \downarrow P \downarrow$

1100

80

1200

85

—

—

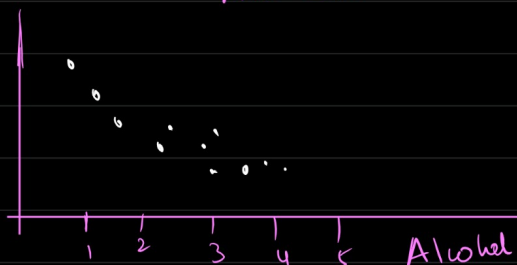
—

—

Area is sqft

* Example (indirect)

No. of years to live



$A \uparrow YTL \downarrow$

$A \downarrow YTL \uparrow$

Alcohol consumption (daily)

① Covariance

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

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

Variance was spread of data \Rightarrow relationship of a feature with itself

\Downarrow
Cov mean, you are trying to understand the relationship of a feature with respect to other features.

\rightarrow Covariance \rightarrow relationship b/w two variable.

\rightarrow

$x \uparrow y \uparrow$

$x \downarrow y \downarrow$

(+ve)
Cov

or

$x \uparrow y \downarrow$

$x \downarrow y \uparrow$

(-ve)
Cov

X	Y
2	3
3	5
6	6
1	8
$\bar{x} = 3$	$\bar{y} = 5.5$

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

$$(2-3)(3-5.5) + (3-3)(5-5.5) + (6-3)(6-5.5) + (1-3)(8-5.5)$$

$$\Rightarrow \frac{(-1)(-2.5) + 0 + 3 \times 0.5 + (-2) \times 2.5}{3}$$

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

\rightarrow The two features x and y are negatively related.

X	Y
2	3
4	5
6	7
$\overline{x} = 4$	$\overline{y} = 5$

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

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

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

→ x & y are having a positive relation.

Advantage

→ Relationship b/w X, Y

+ve or -ve

* Disadvantage

① X Y

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



∞ to ∞

A B

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

→ No comparison of Strength of relationship in Covariance

→ No any standardized scale to interpret the strength.

② Covariance has dimension

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

x → height → ft

y → wt → kg.

Cov. ⇒ ft · kg

(x)

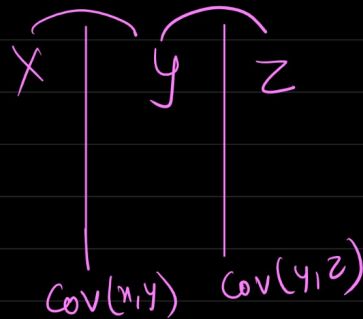
(y)

(z)

transaction Amount (Rs)	height (ft)	weight (kg)
-------------------------------	----------------	----------------

$$\begin{aligned} \text{Cov}(\text{tr Amt}, \text{ht}) &\Rightarrow \text{Rs} \cdot \text{ft} \Rightarrow 450 \text{ Rs} \cdot \text{ft} \\ \text{Cov}(\text{height}, \text{wt}) &\Rightarrow \text{ft} \cdot \text{kg} \Rightarrow 600 \text{ ft} \cdot \text{kg} \end{aligned} \quad \left. \vphantom{\begin{aligned} \text{Cov}(\text{tr Amt}, \text{ht}) \\ \text{Cov}(\text{height}, \text{wt}) \end{aligned}} \right\}$$

→ We can not compare two different dimension



Not comparable → different dimensions.

* Soln

→ -1 to 1
→ dimensionless Quantity

② Pearson Correlation Coefficient [-1 to 1]

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



← the more -ve
Correlated the
features are.

→ the more +vely correlated the
features are

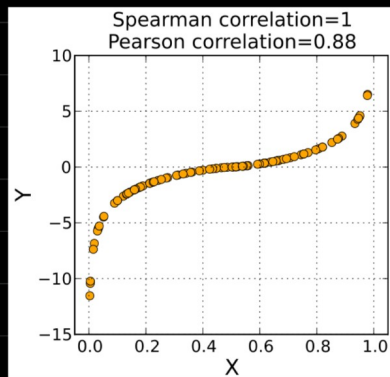
$$\rho_{x,y} = 0.4$$

1

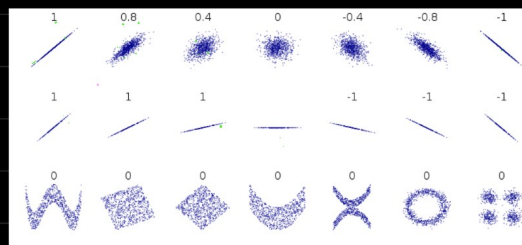
$$\rho_{A,B} = 0.8$$

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

→ Pearson correlation coefficient always measures the linear relationship



Pearson correlation
↓
linear relationship
 $x \uparrow y \uparrow$



What to do for Non linear relationship

Spearman Rank Correlation

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma_{R(x)} * \sigma_{R(y)}}$$

$R(x)$ — Rank of x
 $R(y)$ — Rank of 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

5th 4th 3rd 2nd 1st
↑ ↑ ↑ ↑ ↑
1, 2, 5, 7, 8

$x \rightarrow 5, 7, 8, 1, 2 \rightarrow$ Sort the value $\rightarrow 1, 2, 5, 7, 8$

Sort the no \rightarrow Highest no will be rank 1.

Dataset \rightarrow 1000 feature

\checkmark X_1 X_2 X_3 X_4 X_5 X_6 \dots X_{1000} $\overset{\text{opt } (y)}{\uparrow}$
Price

\checkmark
 $\left\{ \begin{array}{l} X_1 - y \\ X_2 - y \\ X_3 - y \\ X_4 - y \end{array} \right.$

$\text{Corr} \approx 0$