

DAV-3

HYPOTHESIS TESTING

(Class starts
@ 9:08 PM)



power of test

		pred	
		0	1
Act	0	TN	FN
	1	FP	TP

Annotations:

- α (Type I error) is associated with the cell (Act 1, pred 0).
- β (Type II error) is associated with the cell (Act 0, pred 1).
- $(1 - \beta)$ (Power of Test) is associated with the cell (Act 1, pred 1).
- $(1 - \alpha)$ is associated with the cell (Act 0, pred 0).
- $FN \downarrow$ (False Negative) is associated with the cell (Act 0, pred 1).
- $TP \uparrow$ (True Positive) is associated with the cell (Act 1, pred 1).
- $(1 - \beta) \uparrow$ (Power of Test) is associated with the cell (Act 1, pred 1).

goal

(power of Test)

$(1 - \beta)$

Lecture 8: Correlation

#Agenda

- ① Parametric Vs Non parametric
- ② Covariance & Correlation
- ③ Pearson Correlation
- ④ Spearmann Correlation
- ⑤ Visualizing using Heatmap

Parametric vs Non-Parametric Test

Parametric Vs Non-parametric Test

Parametric

Check the underlying distribution of data before Test

T-test, ANOVA

Non-Parametric

No need to check the underlying distribution of data before test

KSTest, KWTest

Parametric HT

Assumptions:

- Parametric tests make specific assumptions about the population distribution from which the data is drawn.
 - Common assumptions include normality (data follows a normal distribution) and variance is constant across groups or conditions.
 - Parametric tests are typically used when the data reasonably follows the assumed distribution and other assumptions are met.
- ✓ Parametric tests tend to be more powerful (i.e., better at detecting true effects) than non-parametric tests when the assumptions are met.
- This is especially true when the sample size is large.

Non-Parametric HT

Assumptions:

- Non-parametric tests make fewer or no assumptions about the population distribution.
 - They are distribution-free or rely on fewer assumptions, such as independence of observations.
- ✓ Non-parametric tests are useful when the assumptions for parametric tests are violated.
- They are also suitable for data types that don't fit well with parametric assumptions, such as ordinal or skewed data.
 - Non-parametric tests are generally less powerful than parametric tests when data conforms to parametric assumptions.
- ✗ However, they can be more robust and appropriate when dealing with non-normally distributed data.

Parametric

Z test

paired T test

ANOVA

Lemene's Test

Non-parametric

Chi Sq Test

KSTest

KWTest

Types of Test

~~Sofar~~

✓ Numeric Vs Categorical → 2 (T-test)
→ >2 (ANOVA)

✓ Categorical Vs Categorical → Chi square test

✓ Numeric Vs Numeric → Correlation Test

- ① Pearson Correlation
- ② Spearman Correlation

Covariance And Correlation

Class starts at 9:05 pm

Agenda

- ① Covariance and Correlation
- ② Pearson Correlation
- ③ Spearman Correlation
- ④ Visualization using Heatmaps.

}

Example: Height Weight

Statistician gathers Height & Weight data

→ To find the relation b/w H & W.

(Army Camp)

Fit

Height	Weight

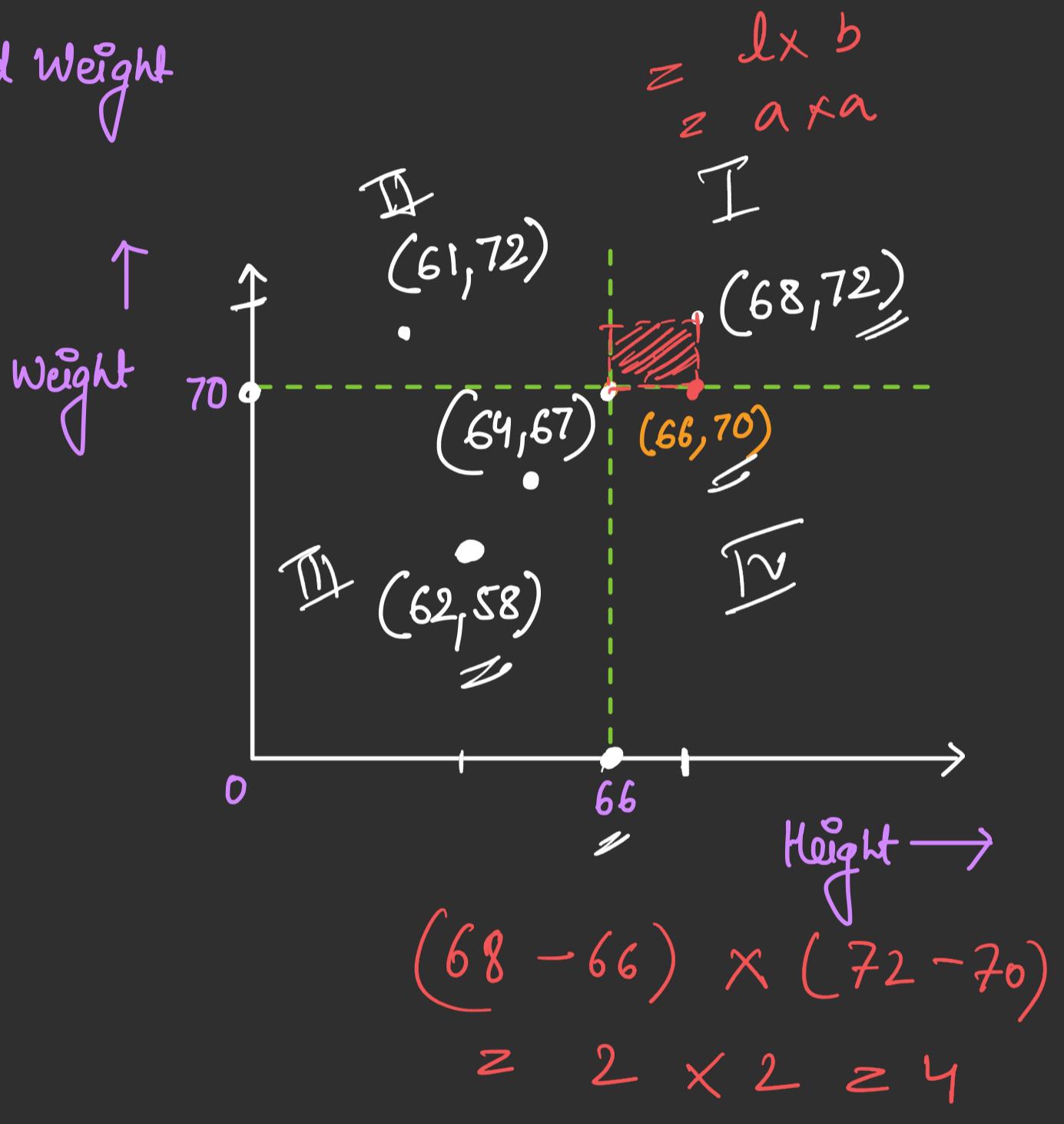
Height vs Weight → correlation

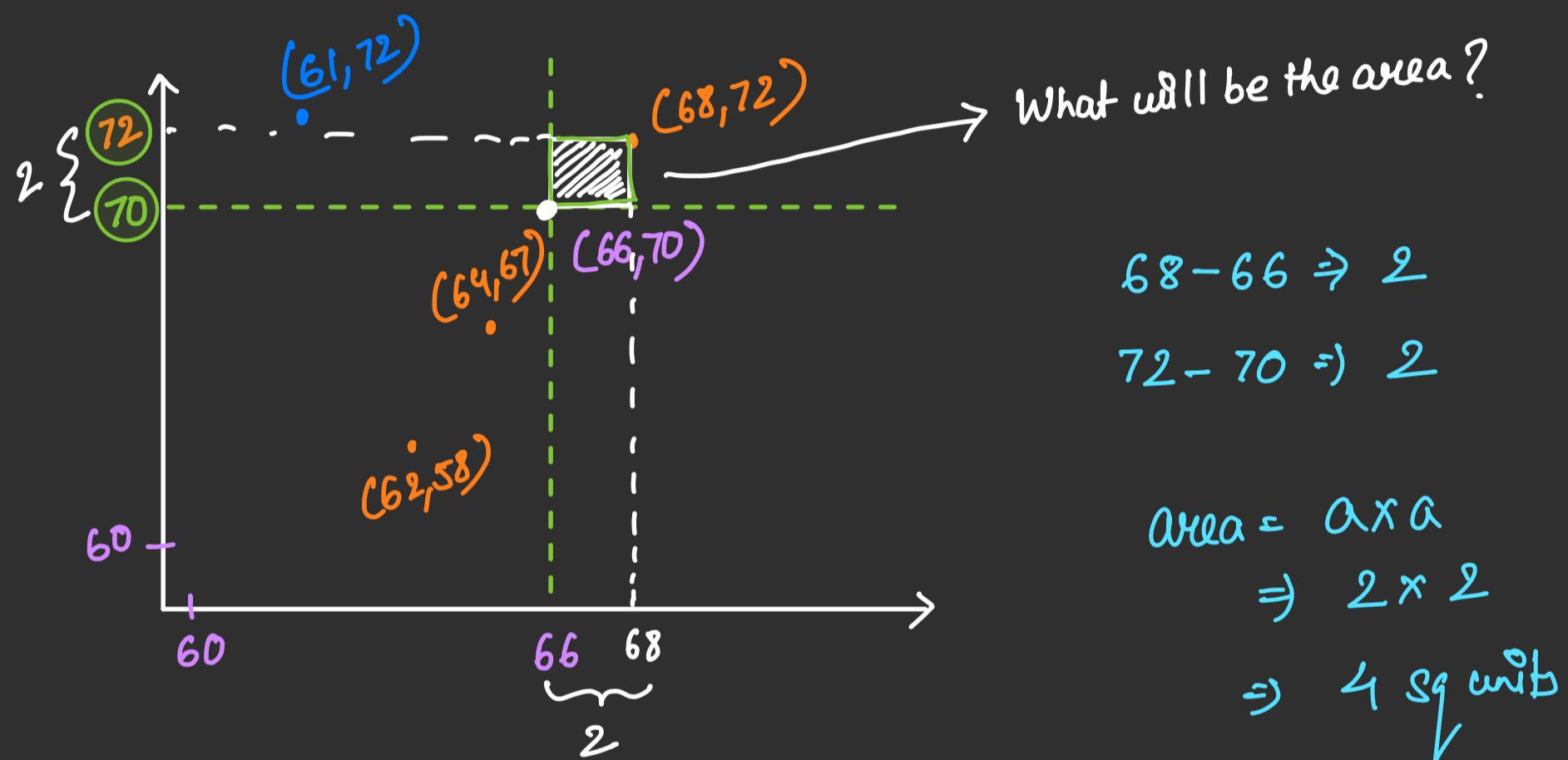
num num

Example 1: Height and Weight

Height (inches)	Weight (kg)
68	72
62	58
64	67
61	72
70	79
66	61
61	68
65	64
71	80
72	79

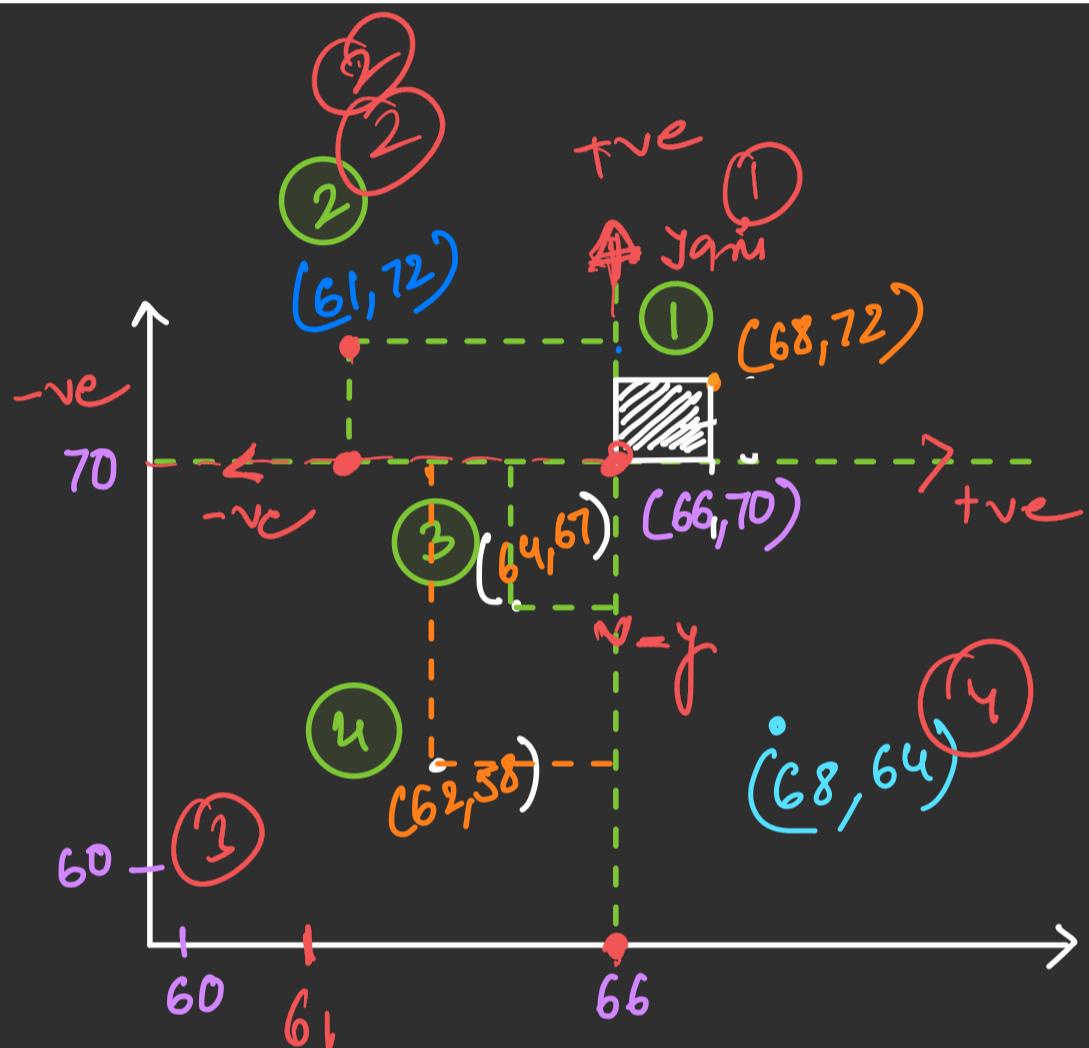
$\bar{h} = 66$ $\bar{w} = 70$





$$(68 - 66) \times (72 - 70) = 4$$

mean mean



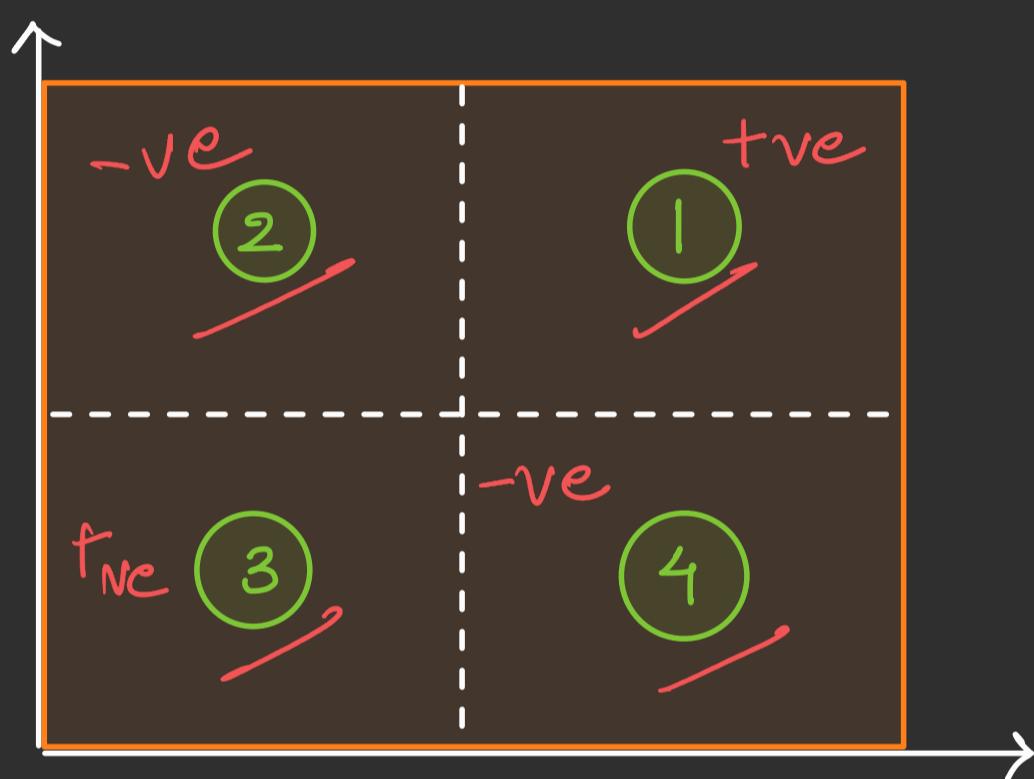
$$\textcircled{1} (68 - 66) \times (72 - 70) = 4$$

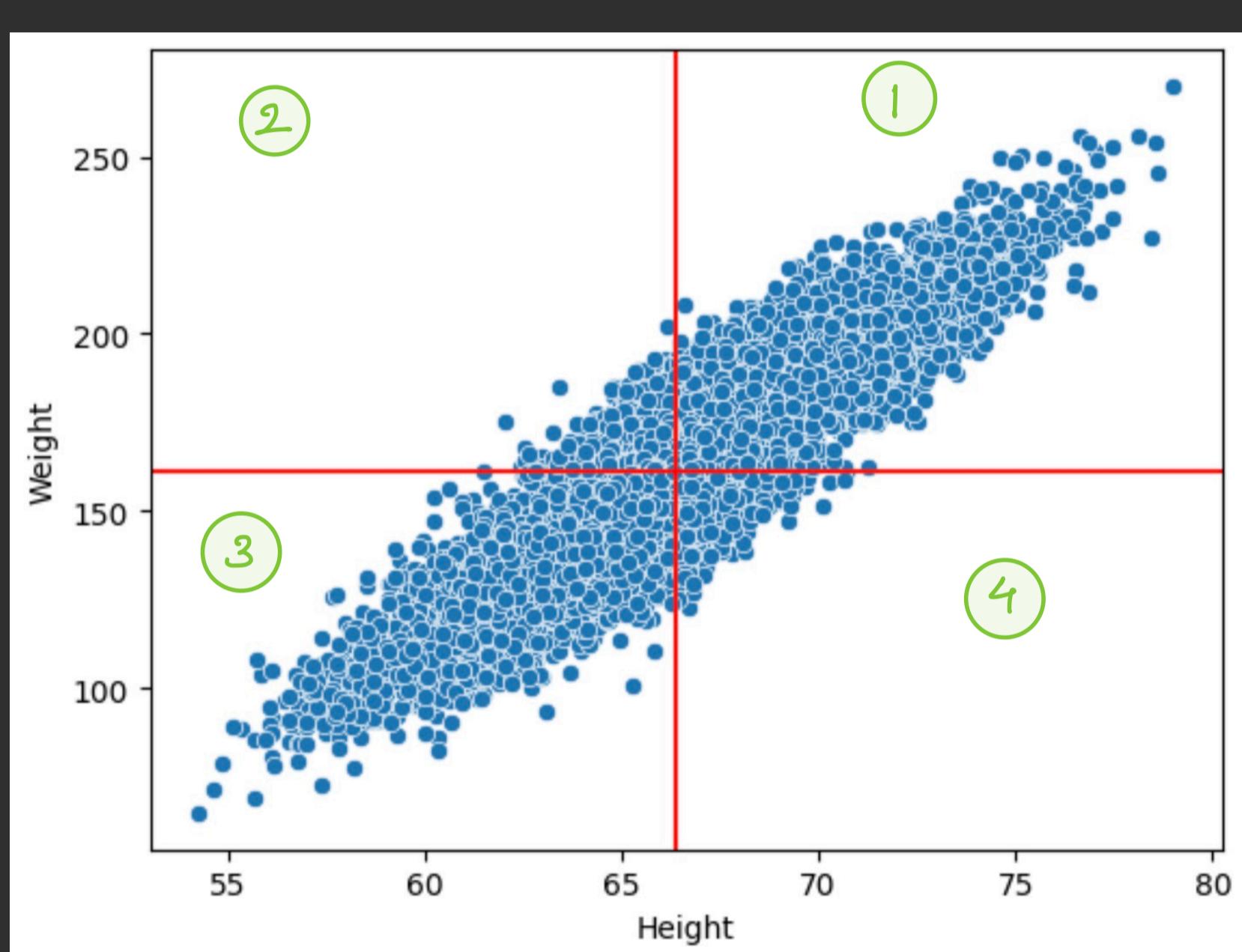
$$\cancel{\textcircled{2}} (61 - 66) \times (72 - 70) \\ \cancel{- 10}$$

$$\textcircled{3} (64 - 66) \times (67 - 70) \\ \cancel{6}$$

$$\textcircled{4} (62 - 66) \times (58 - 70) \\ \cancel{48}$$

$$(68 - 66) \times (64 - 70) \\ +ve(2) r - 6 = -12$$





Average area for the points falling in area $(1, 3)$
will be positive

Also known as "positively correlated"

Average area for the points falling in area $(2, 4)$
will be negative

Also known as "negatively correlated"

$$\textcircled{1} \quad (68 - 66) \times (72 - 70) = 2 \times 2 = 4$$

$$\textcircled{2} \quad (62 - 66) \times (58 - 70) = (-4) \times (-12) = 48$$

$$\textcircled{3} \quad (64 - 66) \times (67 - 70) = (-2) \times (-3) = 6$$

$$\textcircled{4} \quad (61 - 66) \times (72 - 70) = (-5) \times (2) = -10$$

$$\begin{matrix} [4, 48, 6, -10] \\ \xrightarrow{\text{Avg } \underline{12}} \\ \cancel{\text{Avg }} [\text{Covariance}] \end{matrix}$$

Formula for Covariance

~~~~~ in ~~~~~

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

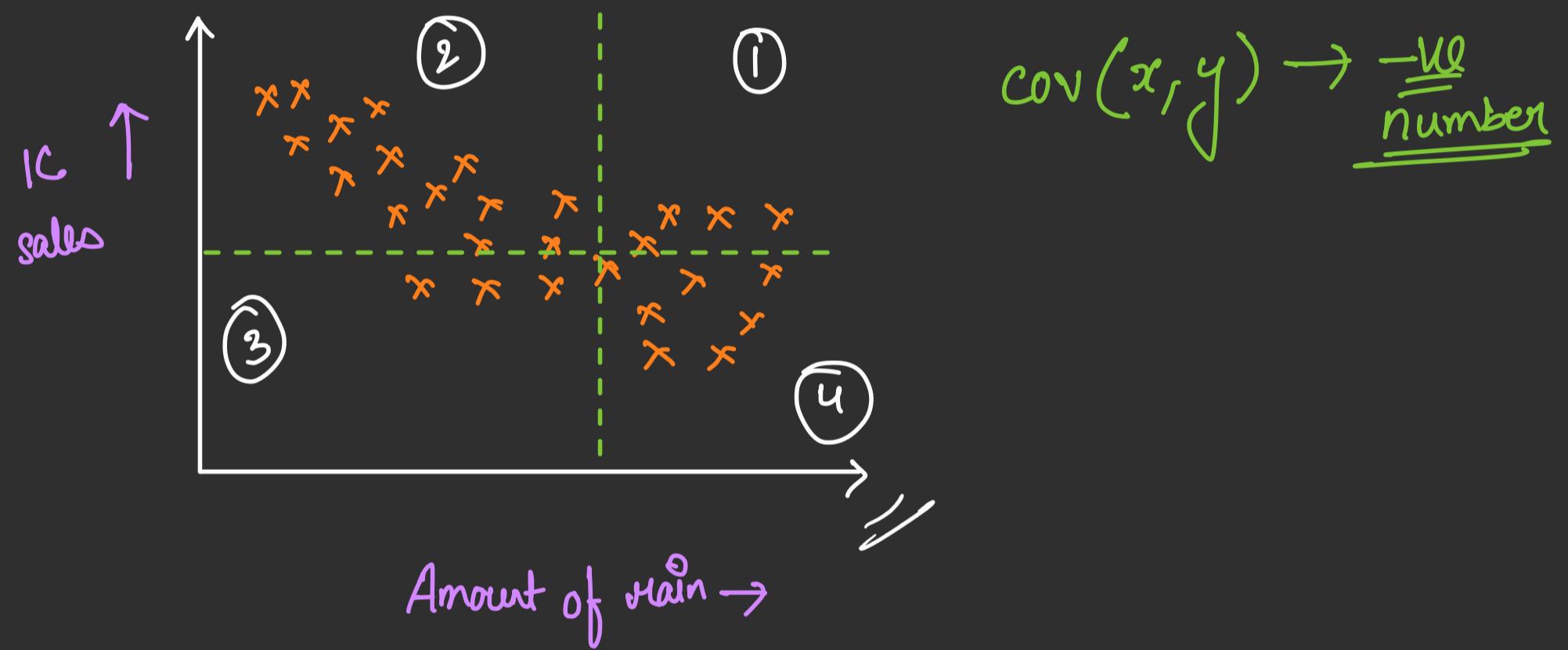
Annotations:

- $x$  dtpt  $n \rightarrow$  Total no of dt pt
- mean
- mean
- $y$  dtpt

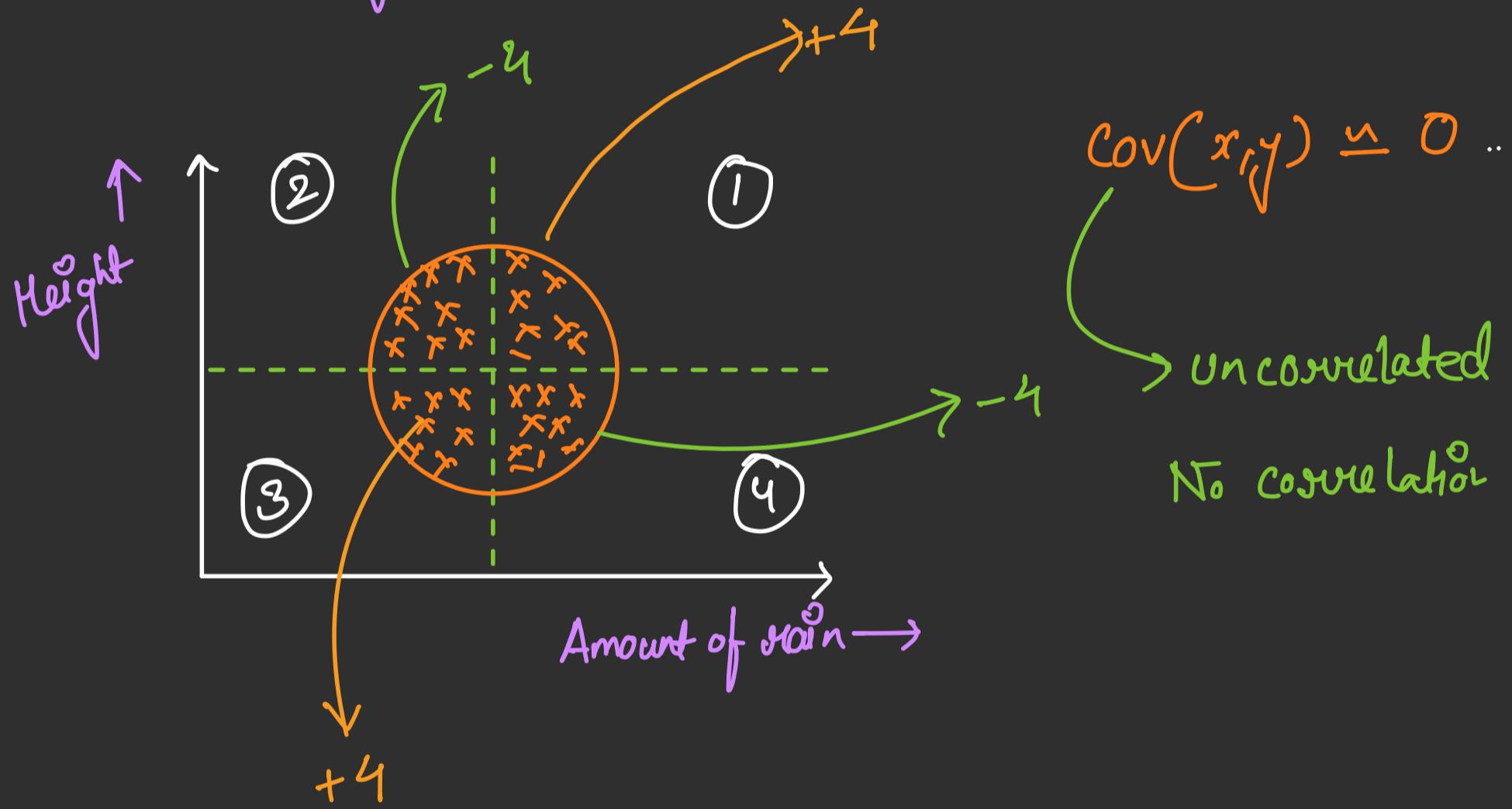
Covariance is a mathematically entity to understand correlation

Co variance  $\rightarrow x \uparrow y \uparrow$

# Example 2: Ice cream sales Vs Amount of Rain



# Example 3: Height Vs Rain

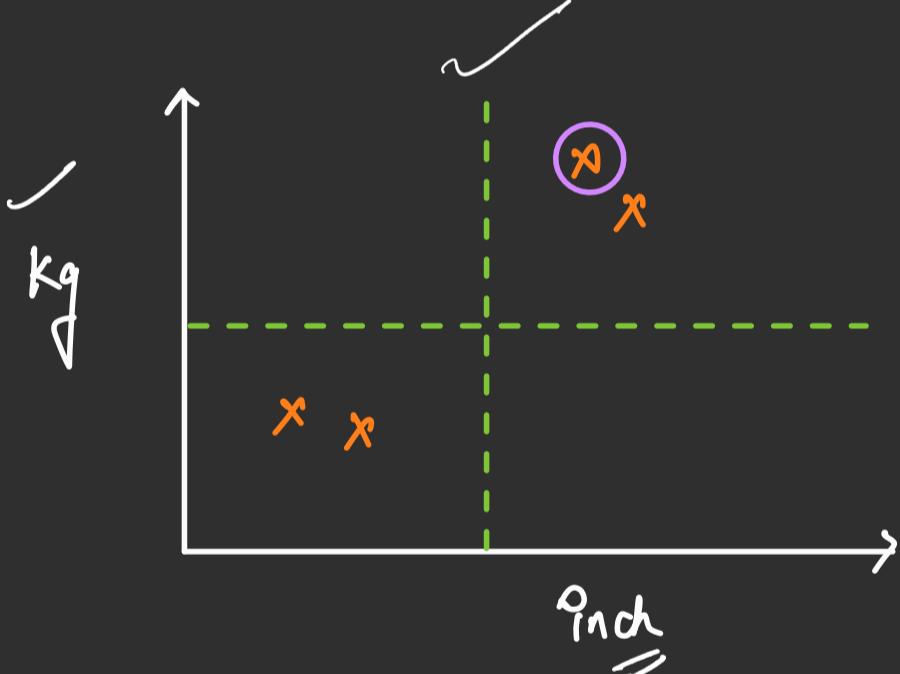


# Pearson Correlation

# Pearson Correlation

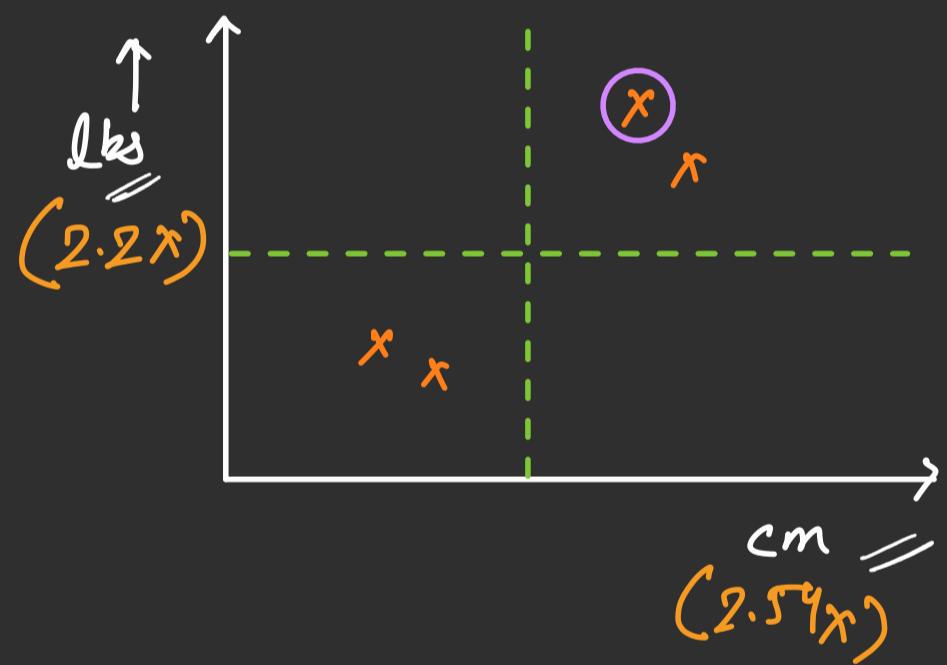
lets take the example of Height & Weight

① Inch & kg ( $S_1$ )



$$(1 \text{ inch} = 2.54 \text{ cm})$$

② Cm & pound ( $S_2$ )



$$(1 \text{ kg} = 2.2 \text{ pounds})$$

$$\text{cov}_{S_1}(x, y) = (2.54)(2.2) \times \text{cov}_{S_2}(x, y)$$

$$\text{cov}_{S_1}(x, y) = \underline{\underline{2.39}}$$

$$\text{cov}_{S_2}(x, y) = \underline{\underline{48.39}}$$

$$X_{\text{cm}} = \frac{2.54}{2.2} X_{\text{inch}}$$

$$Y_{\text{lb}} = \frac{2.2}{2.54} Y_{\text{kg}}$$

$$\begin{aligned}\text{cov}_{S_2} &= (2.54) \times (2.2) \\ &= 5.58 \times \text{cov}_{S_1}\end{aligned}$$

S<sub>1</sub>

H  
(inch)  
—

W  
kg  
—

—  
—  
—

men      ne

inch      kg

S<sub>2</sub>

H  
(cm)  
—

W  
(pounds)  
—

—  
—  
—

cm      pounds

## Pearson Correlation

$$PC = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

where  $\sigma_x$  &  $\sigma_y$  are std of  $x$  &  $y$

Two cols

$$PC = \frac{1}{n} \sum \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \times \left( \frac{y_i - \bar{y}}{\sigma_y} \right)$$

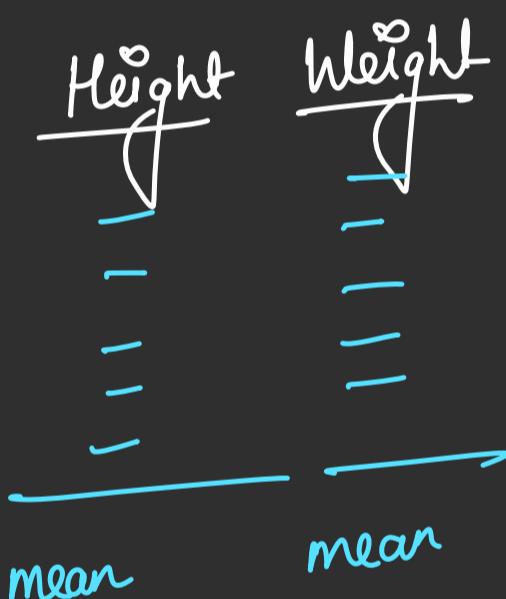
"Independent of scale"

Z score

$$[-1 \leq \rho \leq 1]$$

np.cov

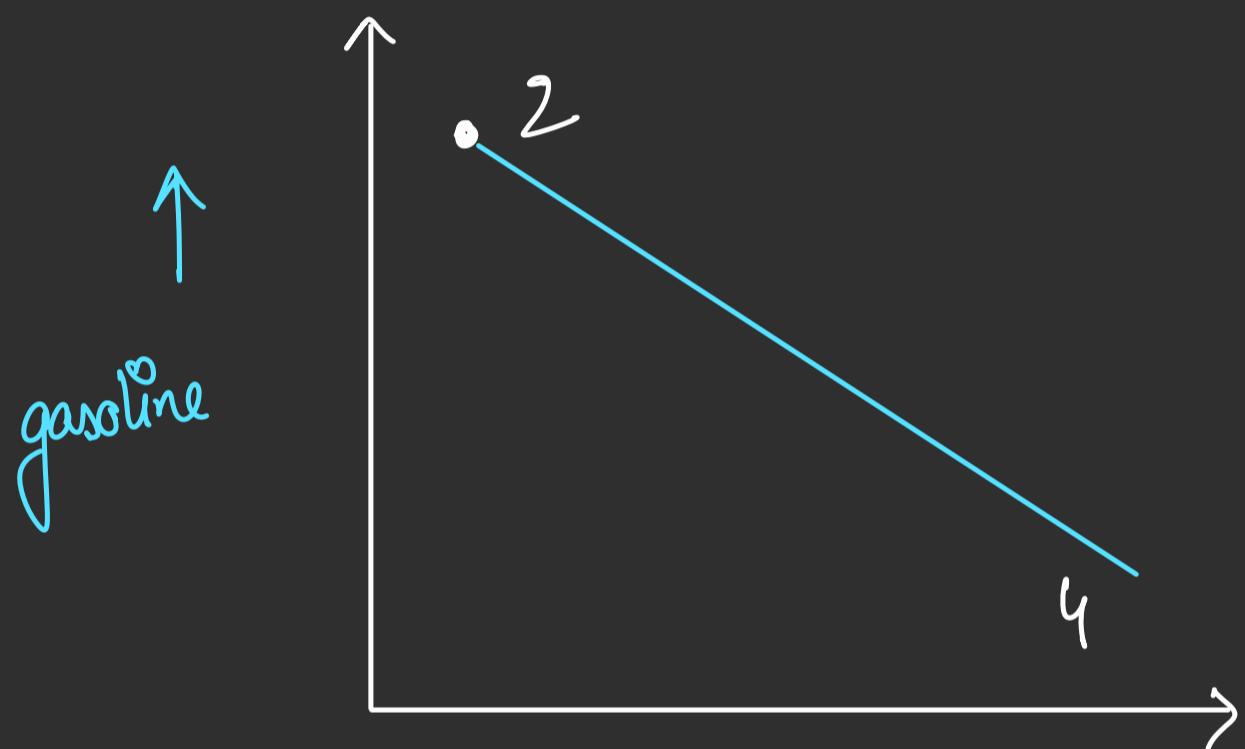
(  
Two Col)



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

$$\text{np.cov} = \begin{bmatrix} \text{var}(x) & \text{cov}(y,x) \\ \text{cov}(x,y) & \text{var}(y) \end{bmatrix}$$

$$\checkmark \text{cov}(x,y) = \text{cov}(y,x)$$

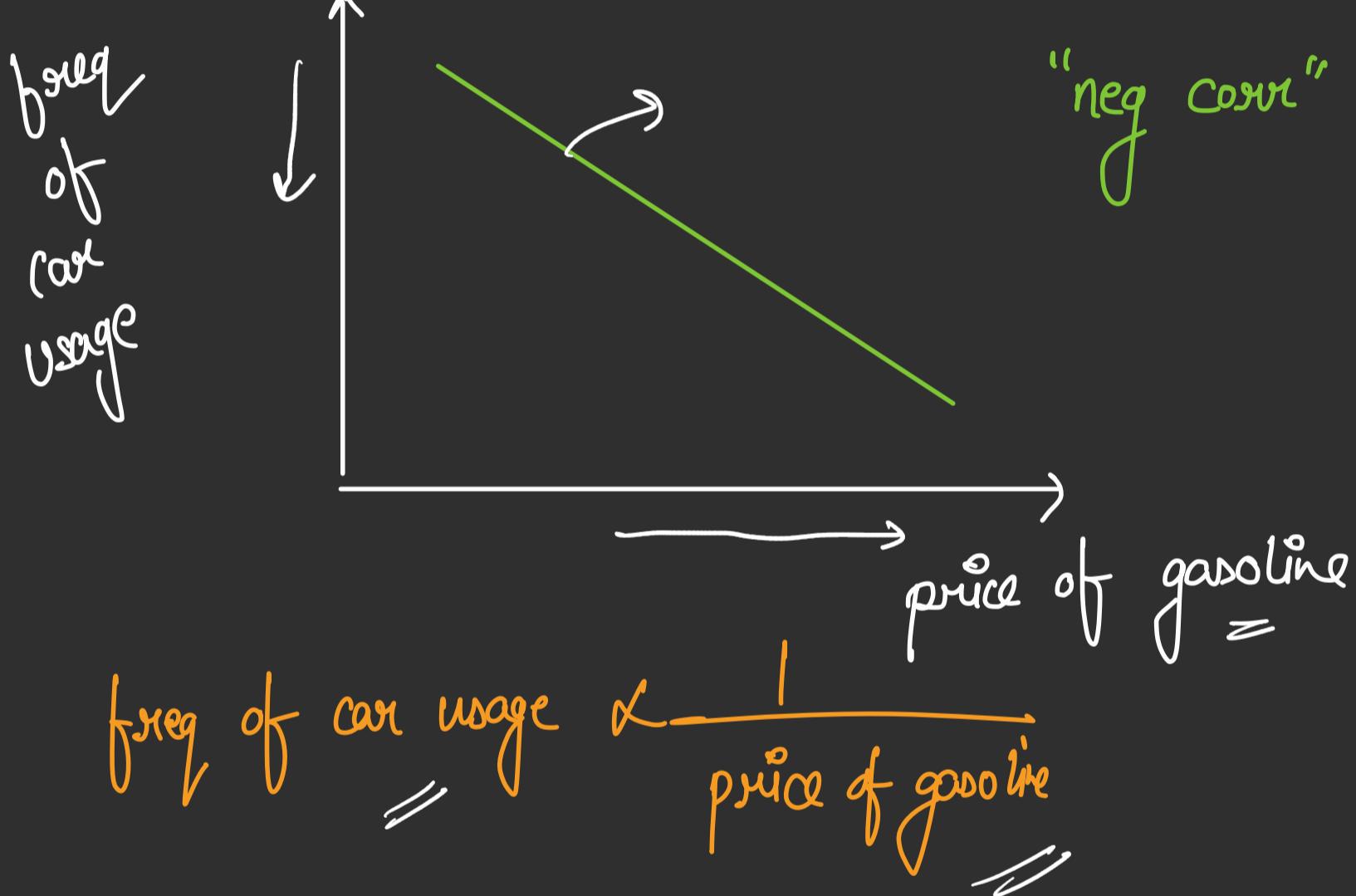


$$\text{cov}(x, y) = -\frac{\text{var}}{m}$$

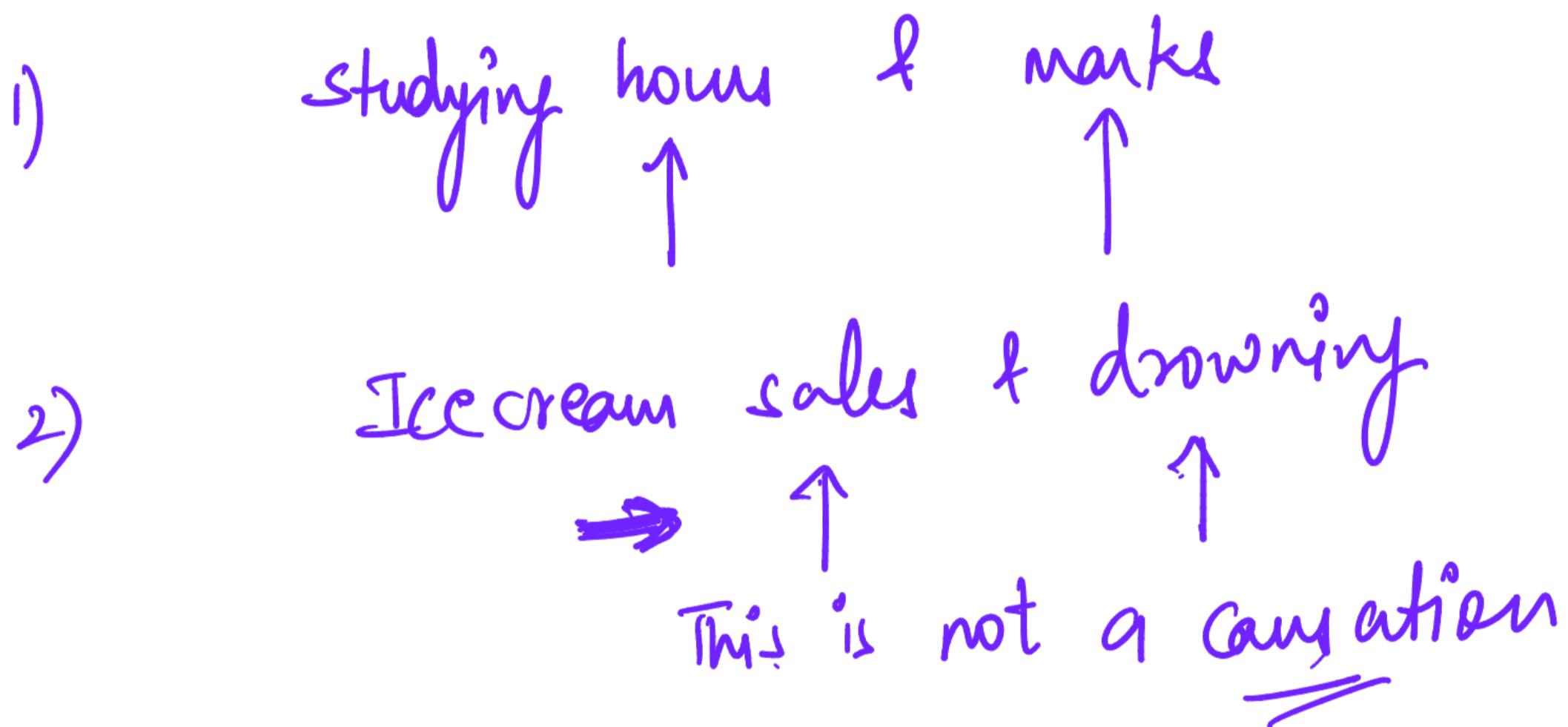
"negatively correlated"

freq of car  
usage ↓

## Quiz

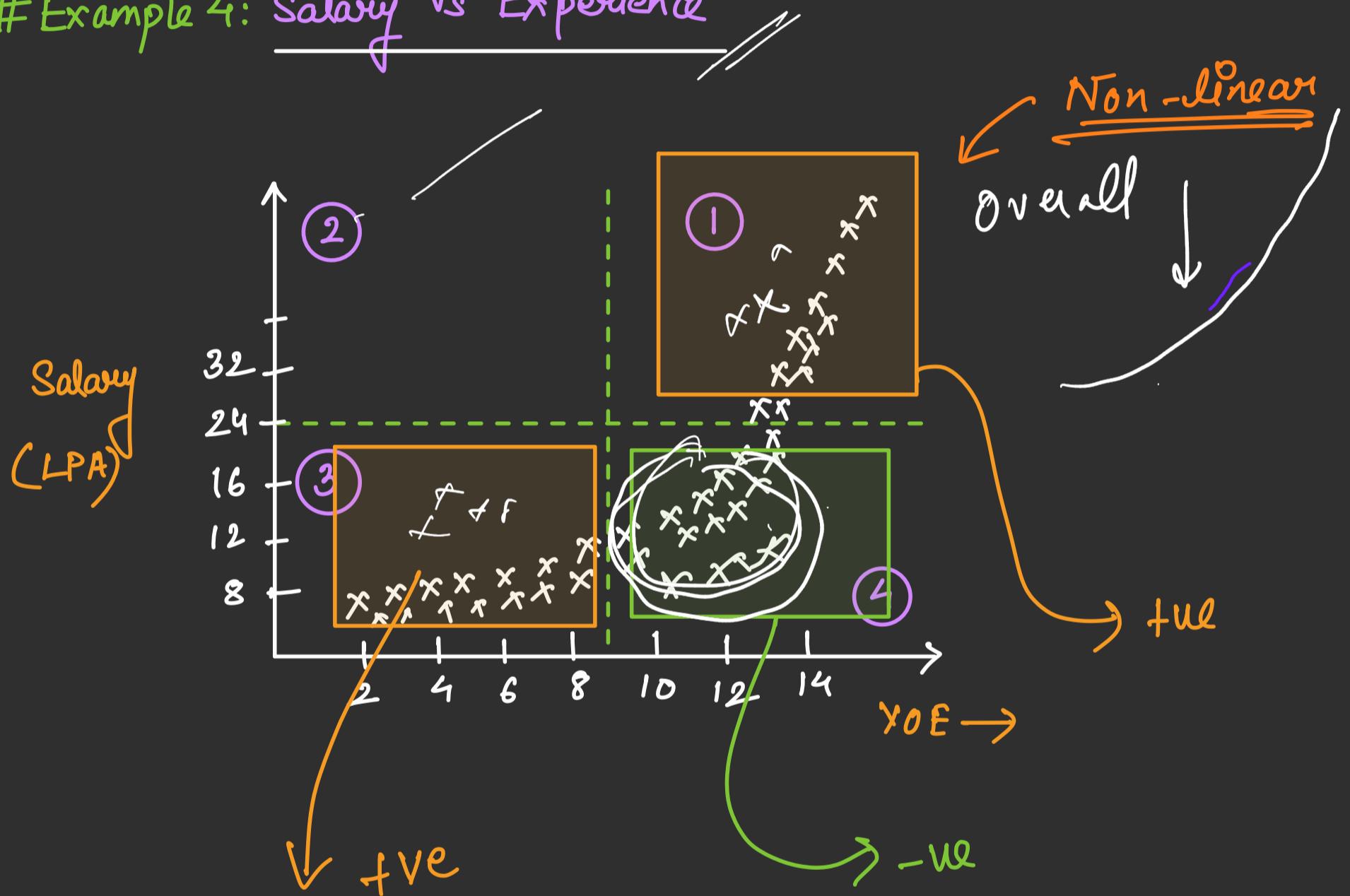


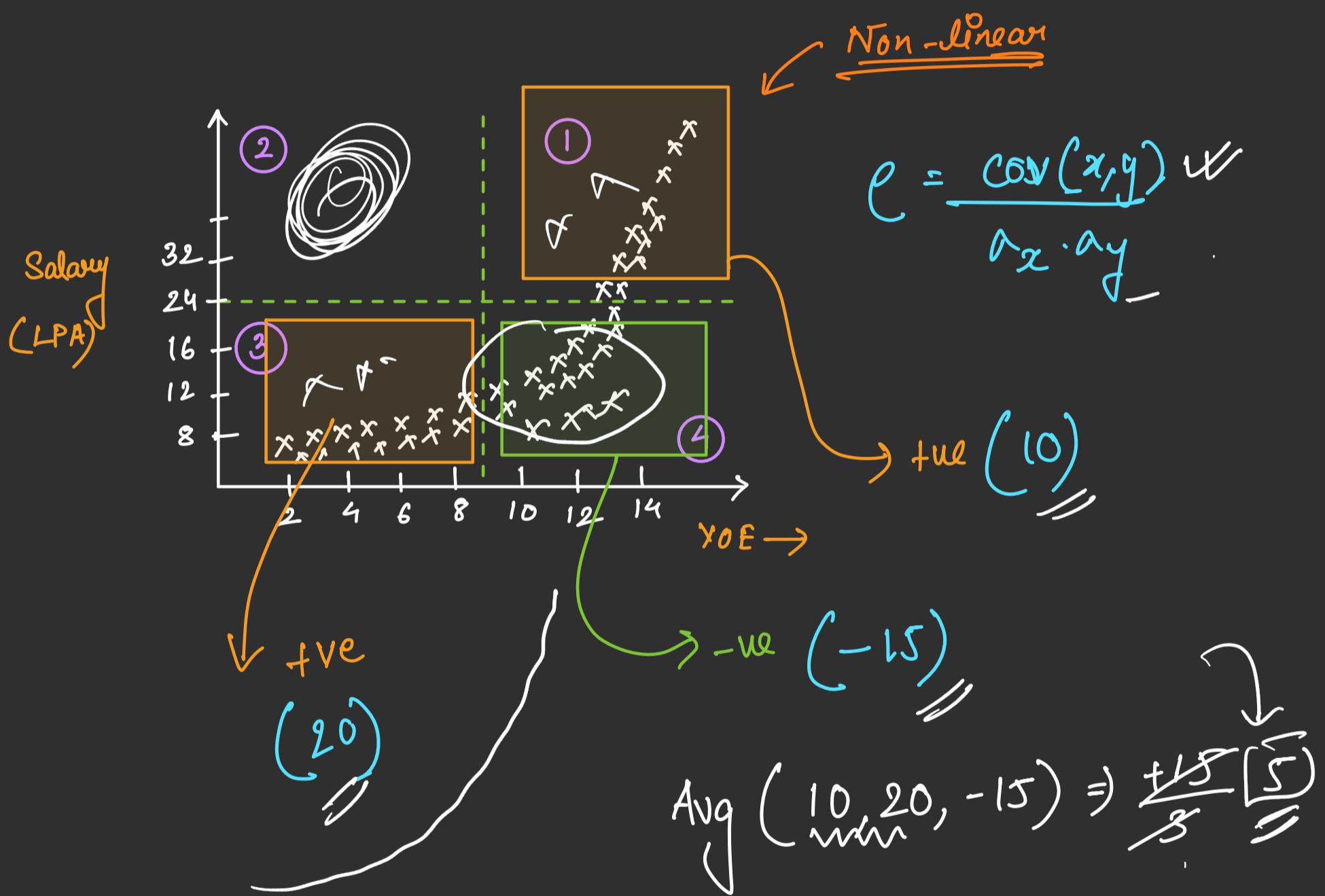
$X$  is the reason for change in  $Y$

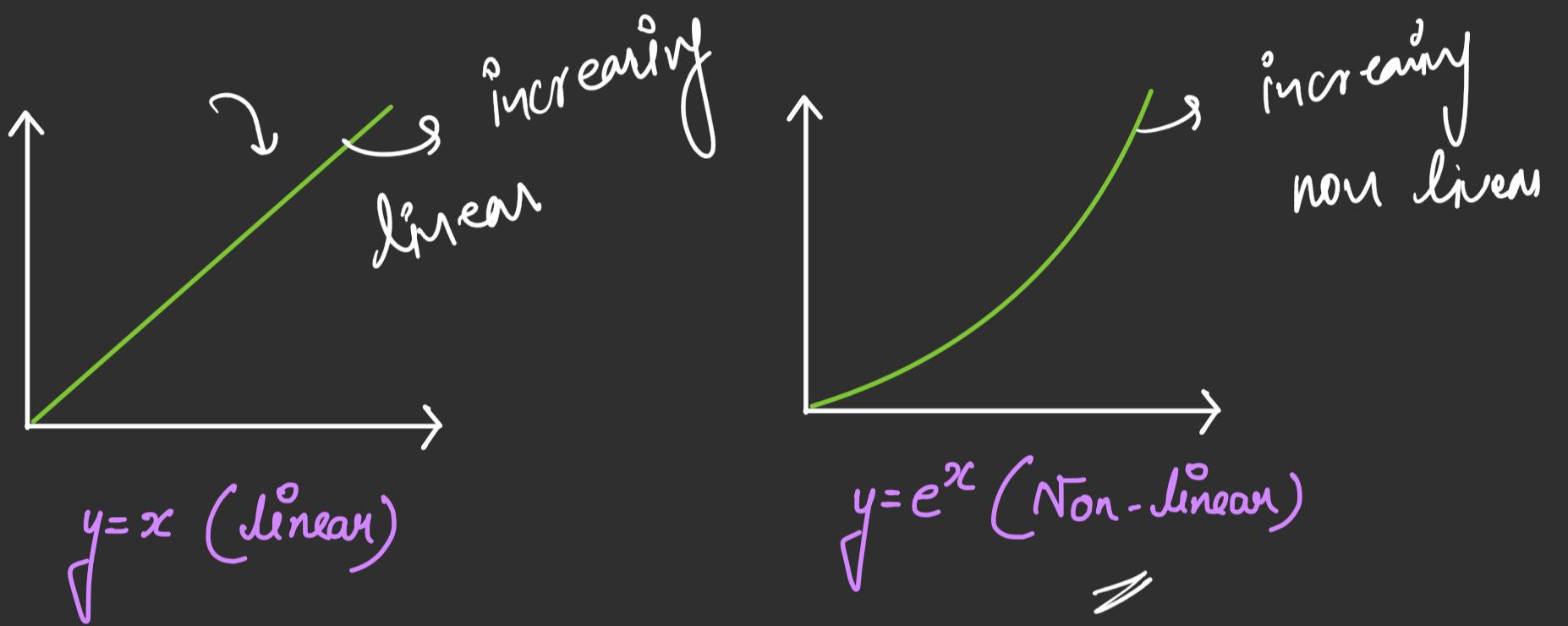


# Spearman Correlation

## # Example 4: Salary Vs Experience

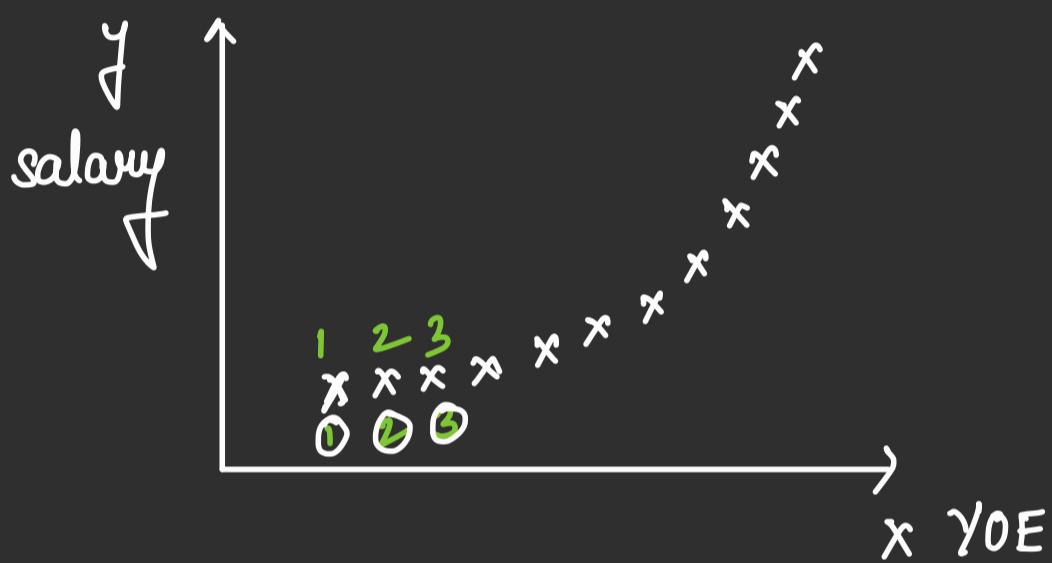






Pearson Correlation  $\rightarrow$  Cannot Handle Non-linear data

## # Spearman Correlation



| $x$ | $y$ |
|-----|-----|
| 32  | 20  |
| 42  | 36  |
| 100 | 48  |
| 46  | 31  |

| Rank X | Rank Y |
|--------|--------|
| 1      | 1      |
| 2      | 2      |
| 3      | 3      |
| 4      | 4      |
| :      | :      |
| 10     | 10     |

linear

## # Example

| Students | Maths | Science |
|----------|-------|---------|
| A        | 35    | 24      |
| B        | 20    | 35      |
| C        | 49    | 39      |
| D        | 44    | 48      |
| E        | 30    | 45      |

difference in ranks

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

↑

→ Spearman Correlation

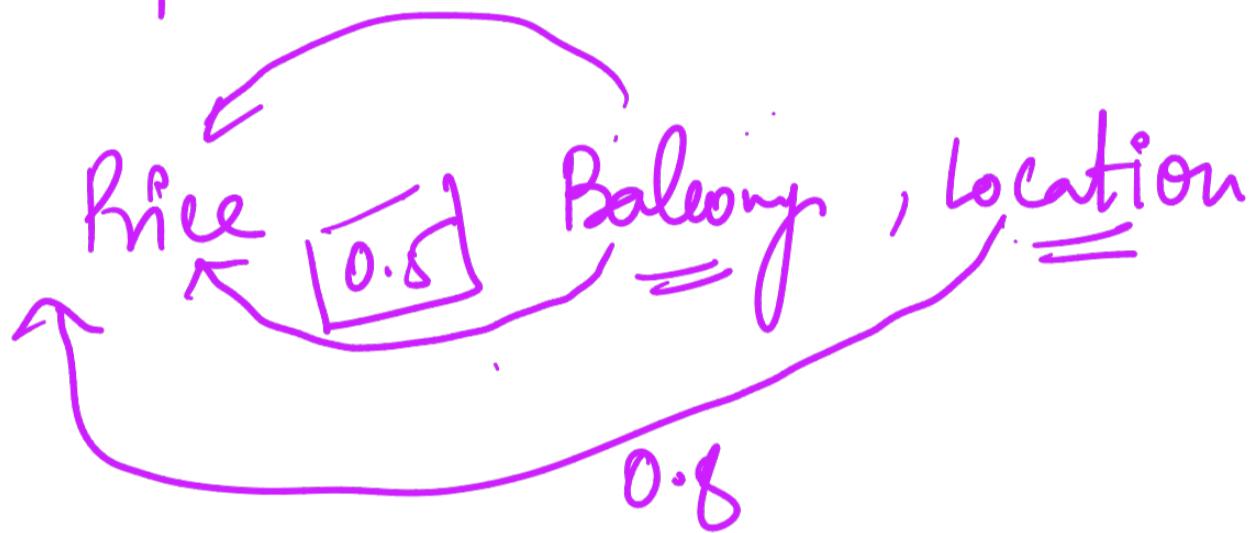
| Students | Maths | Ranks M | Science | Ranks S | Difference d | $ d ^2$ |
|----------|-------|---------|---------|---------|--------------|---------|
| A        | 35    | 3       | 24      | 5       | -2           | 4       |
| B        | 20    | 5       | 35      | 4       | 1            | 1       |
| C        | 49    | 1       | 39      | 3       | -2           | 4       |
| D        | 44    | 2       | 48      | 1       | 1            | 1       |
| E        | 30    | 4       | 45      | 2       | 2            | 4       |

$$SC = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} \Rightarrow 1 - \frac{6 \times (14)}{5(25-1)} \Rightarrow \underline{\underline{0.32}}$$

14

non linear data  $\xrightarrow{\text{rank}}$  linear  
 $\hookrightarrow$  Pearson correlation

$$\text{Spearman}(X, Y) = \text{Pearson}(\text{rank}(X), \text{rank}(Y))$$



$$\text{norm. cof}(z) = P(Z \leq z)$$

(left tail)

$$\text{Right tail} = 1 - \text{norm. cof}(z)$$

$$\rightarrow \text{norm. cof}(\text{abs}(z)) \rightarrow p \text{ value}$$

$$p = 2 \times P(Z \geq |z|) \hookrightarrow \text{Two tailed test}$$

↑ A

↑ B

(0.98) c

### pearson corr

- 1) Covariance
- 2)  $C = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y}$

linear data

### spear corr

$$1) SC = 1 - \frac{6 \sum d^2}{n(n^2-1)}$$

(Non-linear data  
as well as linear data)

Nominal Data  
~~~~~ ~~~~

35
24
36
22



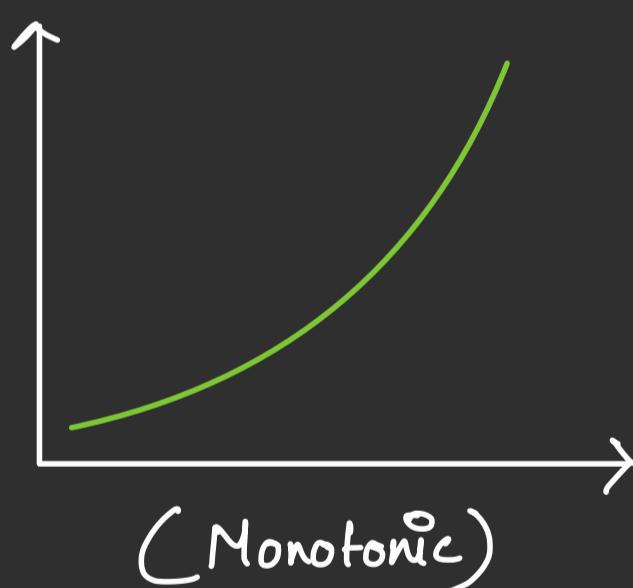
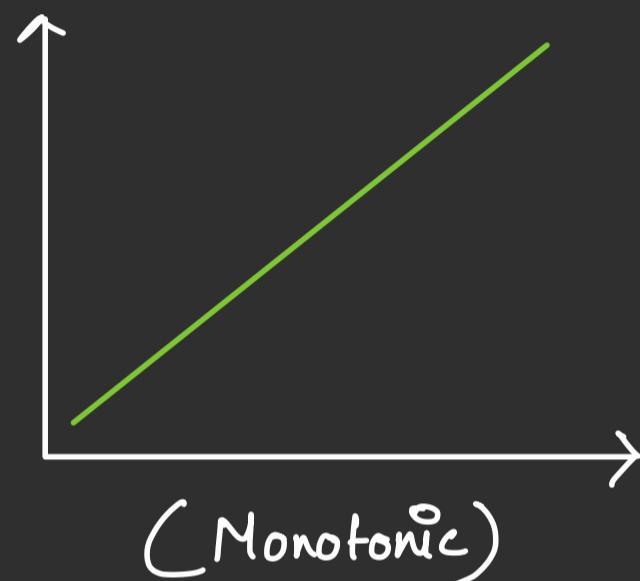
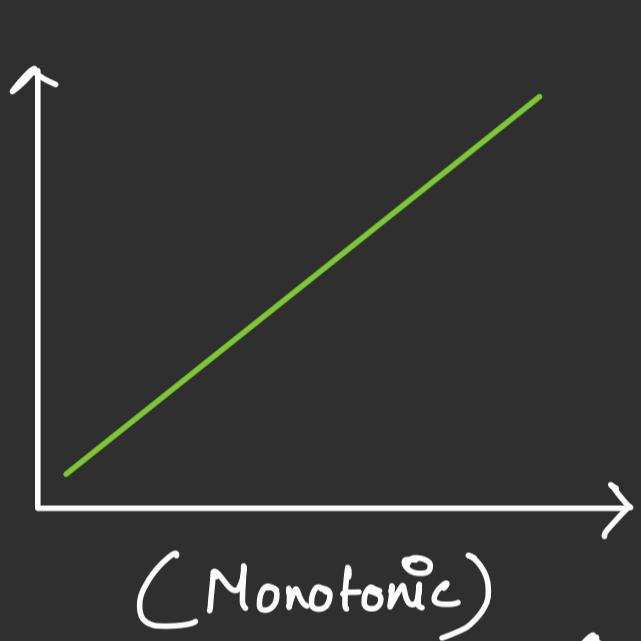
Ordinal Data
~~~~~ ~~~~

2  
3  
1  
4

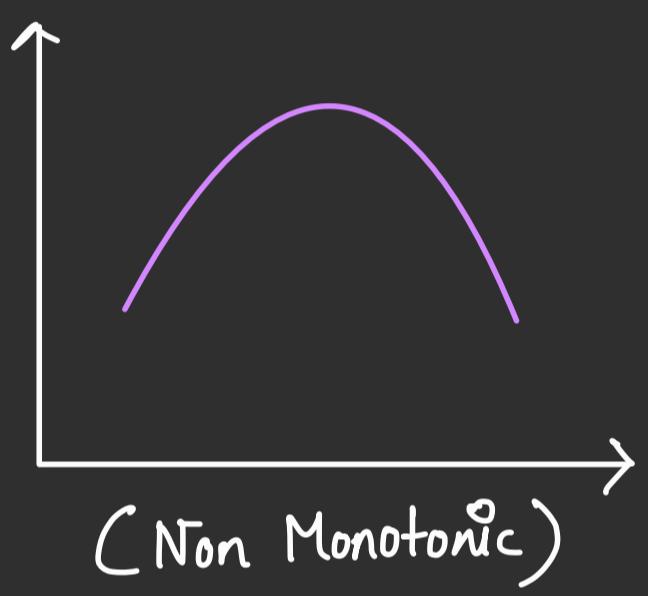
(Ordered data)

Rank

## # Monotonic Vs Non-Monotonic



(Any graph that is  
strictly increasing or  
Monotonic decreasing)



(Non Monotonic)

