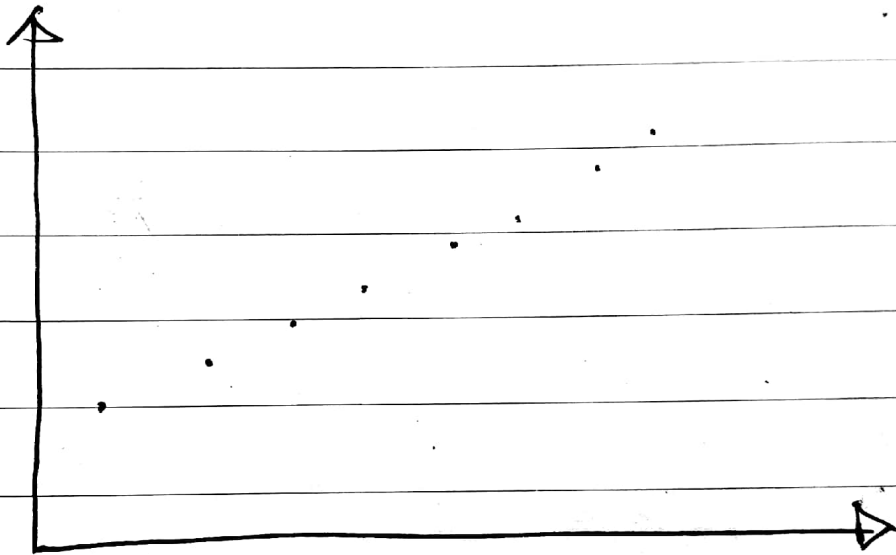


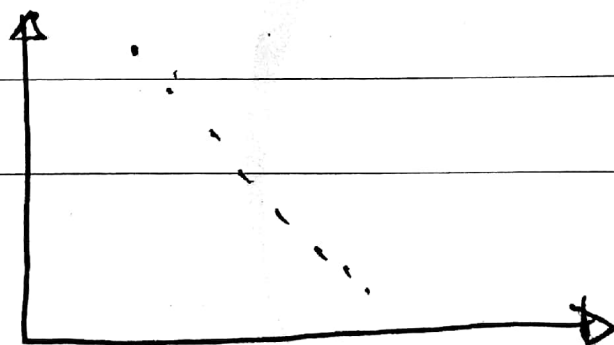
## Notes

Correlation: is usually defined as a measure of the linear relationship between two quantitative variables (e.g: height and weight) → is a statistical measure that indicates the extent to which two or more variables fluctuate together.

\* ~~Posi~~ A positive correlation indicates the extent to which those variables increase or decrease in parallel.



\* A negative correlation indicates the extent to which one variable increases as the other decreases.



→ A correlation coefficient is a statistical <sup>Notes</sup> measure of the degree to which changes to the value of one variable predict change to the value of another.

Pearson's Product moment coefficient  $r$ :

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n \sigma_A \sigma_B}$$
$$= \frac{\sum_{i=1}^n (a_i b_i) - n \bar{A} \bar{B}}{n \sigma_A \sigma_B}$$

Here,

$n$  = no. of tuples

$\bar{A}$  and  $\bar{B}$  are the respective means

$\sigma_A$  and  $\sigma_B$  are the respective S.D of A and B

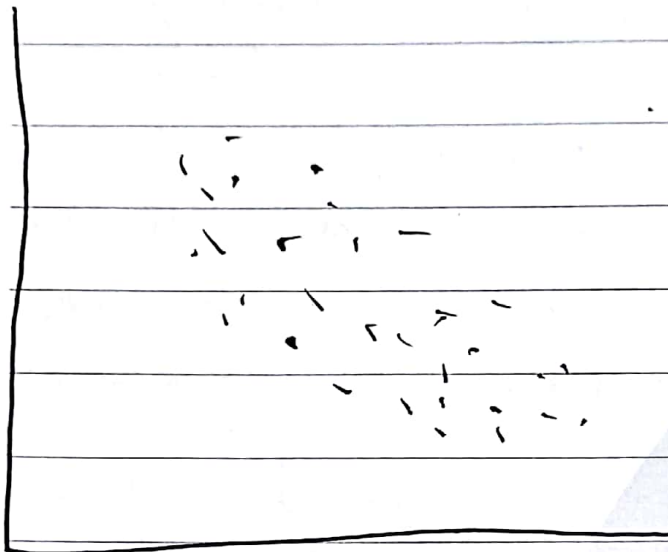
→ If  $r_{A,B} > 0$ , A and B are positively correlated (A's values increase as B's). The higher, the stronger correlation.

→  $r_{A,B} = 0$ ; independent;

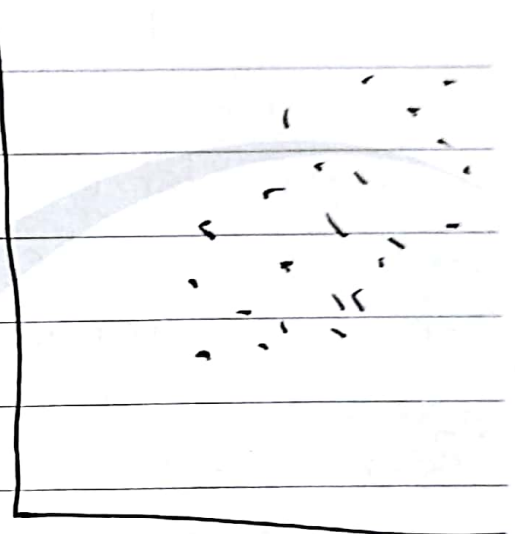
→ There is no correlation

Notes  $r_{AB} < 0$  : negatively correlated  
~~One attribute~~

- A correlation of 1 indicates a perfect positive correlation.
- A correlation of -1 indicates a perfect negative correlation.



weak negatively correlated data



Weak positively correlated data

age	23	23	27	27	39	41
% fat	9.5	<del>26.5</del> 26.5	<del>17.8</del> 17.8	<del>31.7</del> 17.8	31.4	25.9
47	49	50	52	54	54	56
<del>27.4</del> 27.4	27.2	31.2	34.6	42.5	28.8	33.4

## Regression:

~~Notes~~

There is single response variable  $Y$ , also called the dependent variable, which depends on the value of a set of input, also called independent variables  $x_1, x_2, \dots, x_n$ .

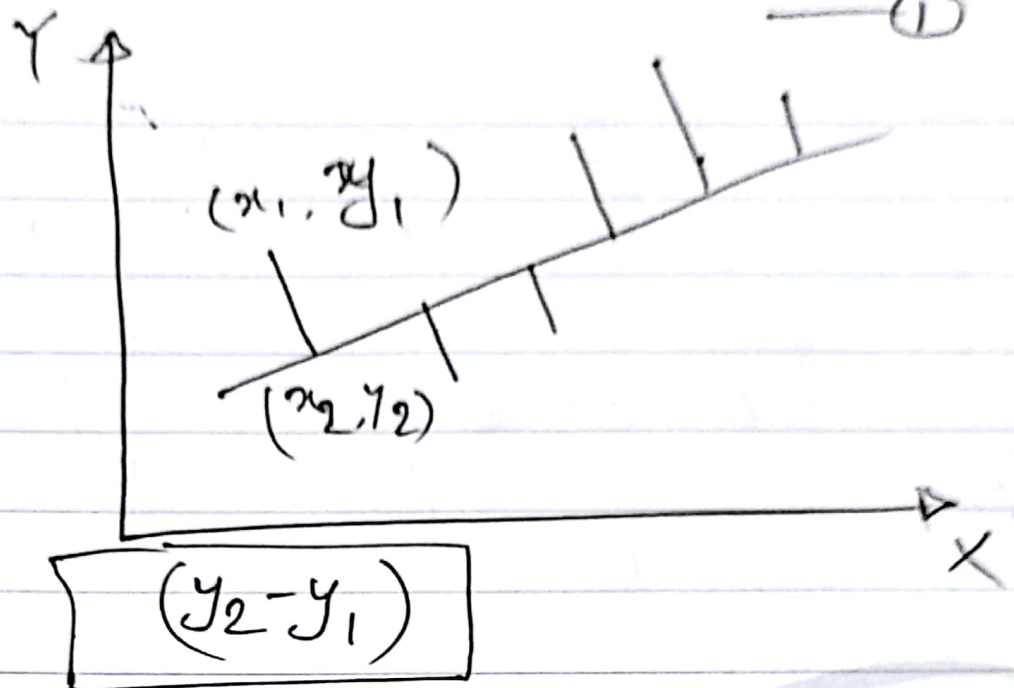
→ A technique for determining the statistical relationship between two or more variables where a change in a dependent variable is associated with and depends on, a change in one or more independent variables.

\* Regression analysis is a powerful statistical method that allows you to examine the relationship between two or more variables of interest.

The simplest type of relationship between the dependent variable  $Y$  and the input variables  $x_1, \dots, x_n$  is a linear relationship for some constraints  $\beta_0, \beta_1, \dots, \beta_n$ .

Notes

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n \quad \text{--- (I)}$$



$\epsilon$  = random error

For error,

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \epsilon \quad \text{--- (II)}$$

$\hookrightarrow$  is called a linear regression equation.  $\beta_0, \beta_1, \dots, \beta_n$  are called regression coefficients

Q: Difference between Correlation and Regression. (Next class)

$\hookrightarrow$  Self study



Notes

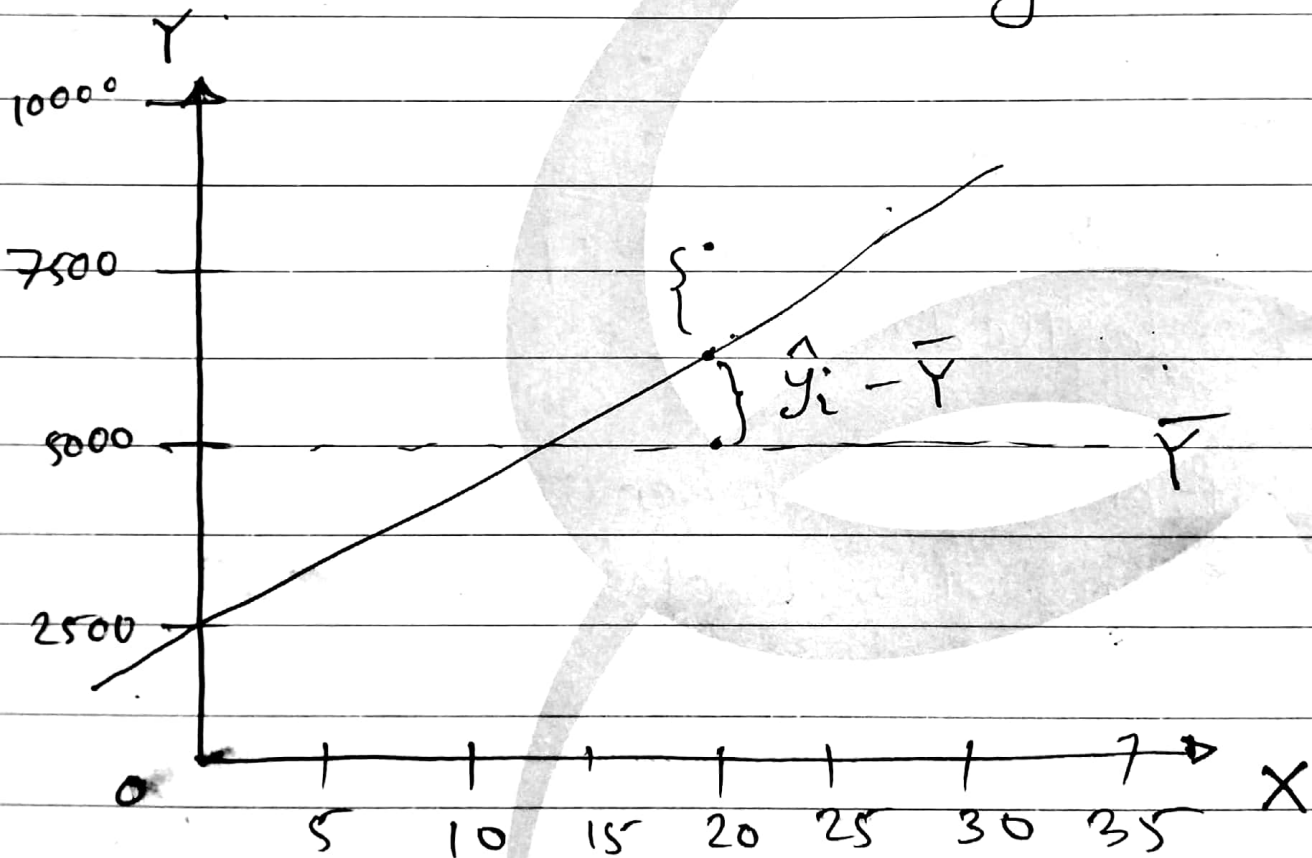
Simple linear Regression: A regression equation containing a single independent variable (that is, one in which  $n=1$ ) is called a simple linear regression equation.

$$Y = \beta_0 + \beta_1 x + \epsilon \rightarrow \text{Population regression function}$$

$\beta_0 \rightarrow$  unknown intercept

$\beta_1 \rightarrow$  Slope parameter

multiple<sup>linear</sup> Regression: A regression equation containing many independent variables is called a multiple linear regression equation.



$$\text{Explained} = \hat{Y}_i - \bar{Y}$$

Notes

$$\text{Unexplained} = Y_i - \hat{Y}_i$$

$$\begin{aligned} SSR &= \text{Sum of squares due to regression} \\ &= \sum (\hat{Y}_i - \bar{Y})^2 \end{aligned}$$

$$\begin{aligned} SSE &= \text{Sum of squares due to error/residual} \\ &= \sum (Y_i - \hat{Y}_i)^2 \end{aligned}$$

$$\begin{aligned} SST &= \text{Sum of squared total} \\ &= SSR + SSE \\ &= \sum (Y_i - \bar{Y})^2 \end{aligned}$$

$$R^2 = \frac{SSR}{SST}$$

→ Proportion of total variation that is explained

Linear Regression Function:

$$y = a + bx$$

Slope of Regression line:

$$b = r \cdot \frac{s_y}{s_x}$$

$$\begin{aligned} r &= \text{Pearson's Correlation Coefficient} \\ &= \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}} \end{aligned}$$

Notes

Y - intercept of regression line:

$$a = \bar{y} - b\bar{x}$$

$x$	$y$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})y_j$	$(x - \bar{x})^2$	$(y - \bar{y})^2$
17	94					
13	73					
12	59					
15	80					
16	93					
14	85					
16	66					
16	79					
18	77					
19	91					

$\bar{x} = \bar{y} =$  1) Find the regression eqn and draw the regression line

$$s_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n-1}}$$

$$s_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

The co-efficient of determination & the sample correlation coefficient:

To measure the amount of variation in the set of response variables



<sup>Notes</sup>  
 $Y_1, Y_2, \dots, Y_n$  of the amount of variation in a set of values  $Y_1, \dots, Y_n$  is given by,

$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$SSR = \sum (\hat{Y}_i - \bar{Y})^2$$

$$SSE = \sum (Y_i - \hat{Y}_i)^2$$

$Y_i = \text{observed}$

$\hat{Y} = \text{predicted}$

$$R^2 = \frac{SSR}{SST}$$

$R^2$  represents the proportion of the variation in the response variables that is explained by the different types of input values. It's called the coefficient of determination.

$$0 \leq R^2 \leq 1$$

If  $R^2 \rightarrow 1$ , it indicates that most of the variation of the response data is explained by the different input values.

$R^2 \rightarrow 0$ , it indicates that little of the variation is explained by the different input values.

Now the sample correlation co-efficient:

Notes

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

As  $r$  is a unitless coefficient, it provides a measure of the degree to which high value of  $x$  are paired with high values of  $y$  and low values of  $x$  with low values of  $y$ .

$$-1 \leq r \leq 1$$

$r \rightarrow 1$  strongly positively correlated

$r \rightarrow -1$  large  $x$  values are strongly associated with small  $y$  values and ~~small~~ small  $x$  values ~~are~~ with large  $y$  values.

Thus  $(r)^2 = \text{similar to } R^2$

$$\therefore |r| = \sqrt{R^2}$$

Q: The grades of a class of 9 students on a midterm report ( $x$ ) and on the final exam ( $y$ ) are as follows:

www.sslwireless.com



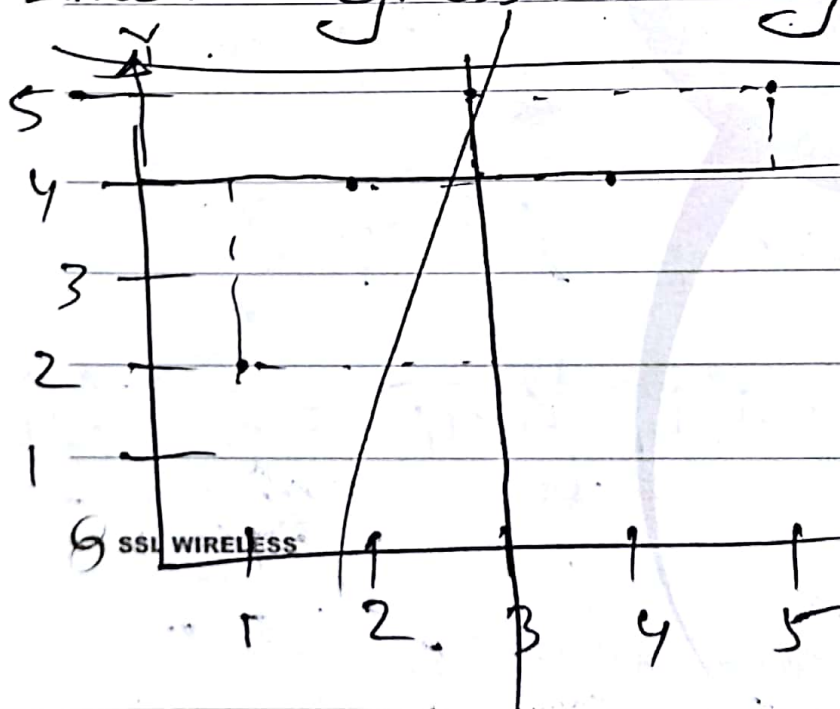
Notes

$x$	$y$
85	75
82	95
71	78
72	84
81	47
94	8.5
96	99
99	99
67	68

a) Estimate the line

b) Estimate the final exam grade of a student who received a grade of 85 on the midterm report.

Linear regression using least square method



$x$	$y$
1	2
2	4
3	5
4	4
5	5
mean	3

# Least Square Method:

Notes

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

$$\frac{\partial(SSE)}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i)$$

$$\frac{\partial(SSE)}{\partial b_1} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) x_i$$

$$\textcircled{1} \sum_{i=1}^n (y_i - b_0 - b_1 x_i) = 0$$

$$\Rightarrow \textcircled{2} \sum_{i=1}^n y_i - n b_0 - b_1 \sum_{i=1}^n x_i = 0$$

$$\Rightarrow n b_0 + b_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \quad \checkmark$$

Again,

$$\sum_{i=1}^n (y_i - b_0 - b_1 x_i) x_i = 0 \quad \text{--- (i)}$$

$$\Rightarrow \sum_{i=1}^n x_i y_i - b_0 \sum_{i=1}^n x_i - b_1 \sum_{i=1}^n x_i^2 = 0$$

$$\Rightarrow b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \quad \text{--- (ii)}$$



Notes

$$nb_0 + b_1 \sum_{i=1}^n x_i$$

Q

$$nb_0 = \sum_{i=1}^n y_i - b_1 \sum_{i=1}^n x_i$$

$$\Rightarrow nb_0 = \frac{\sum_{i=1}^n y_i - b_1 \sum_{i=1}^n x_i}{n}$$

$$= \bar{y} - b_1 \bar{x}$$

$$b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - b_0 \sum_{i=1}^n x_i$$

$$\Rightarrow b_1 = \frac{\sum_{i=1}^n x_i y_i - \bar{x} \sum_{i=1}^n y_i}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}$$

$$(i) \times \sum_{i=1}^n x_i \quad (i) \times n \quad - \quad (i) \times \sum_{i=1}^n x_i$$

$$nb_0 \sum_{i=1}^n x_i + nb_1 \sum_{i=1}^n x_i^2$$

$$- nb_0 \sum_{i=1}^n x_i + b_1 \left( \sum_{i=1}^n x_i \right)^2 = n \sum_{i=1}^n x_i y_i$$

$$- \sum_{i=1}^n x_i \sum_{i=1}^n y_i$$

Notes

$$\Rightarrow b_1 \left\{ n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2 \right\}$$

$$= n \sum x_i y_i - \sum x_i \sum y_i$$

$$\Rightarrow b_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

(1)  $\Rightarrow$

$$b_0 = \frac{\sum y_i - b_1 \sum x_i}{n}$$

$$= \bar{y} - b_1 \bar{x}$$