

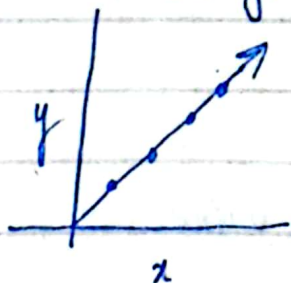
★ Scatterplots

★ b/w two numerical data to see relationship status

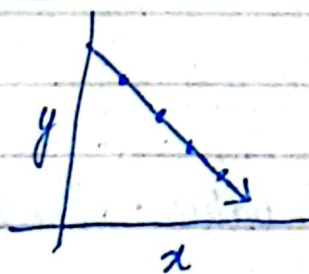
- Linearity (tive, -ive) or non-linear scenario
- Strength of relationship (how close points/data are present)
- Outliers
- Clusters (Distinguished regions of closely related datapoints)

★ Correlation Coefficient (Explain the relationship between two numerical variables)

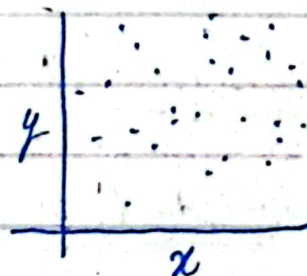
Range = $+1 \rightarrow -1$



$r = +1$



$r = -1$



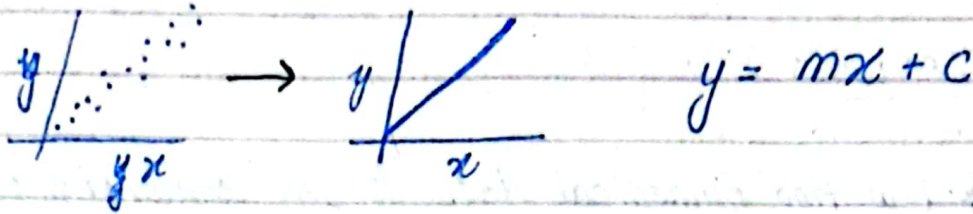
$r = 0$

- Calculation of "r"

($n \rightarrow$ data points)

$$r = \frac{1}{n-1} \sum \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

★ Introduction to trend lines (line graph)



- Now from given data, we can predict future data by just manipulation the x any y values.
- This whole concept is known as linear regression

- The line of best fit (!)

- Estimation with linear models ($y = mx + c$)

↳ It can be done through Eq as well as by simply extrapolating the line.

★ Least Squares regression Eq.

residuals = Actual value - Predicted value

+ive value → the predicted value is less (Actual high)

-ive value → The predicted value is high (Actual below)

(The aim of regression is to minimise the sum of residuals)

→ "A residual is a measure of how well a line fits an individual data point"

- Calculating the Eq. of a regression line

For given data set we have (\bar{x}, \bar{y}) and s_x, s_y
Linear regression Eq. is

Here,

$$y = mx + c$$

$$m = r \frac{s_y}{s_x}$$

And,

$$\bar{y} = m\bar{x} + c$$

$$c = \bar{y} - m\bar{x}$$

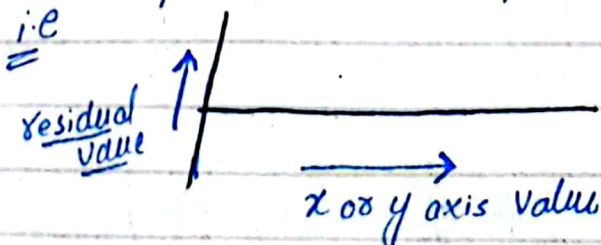
Thus,

$$\hat{y} = mx + c \quad \text{--- (1)}$$

* Assessing the fit in least-squares regression.

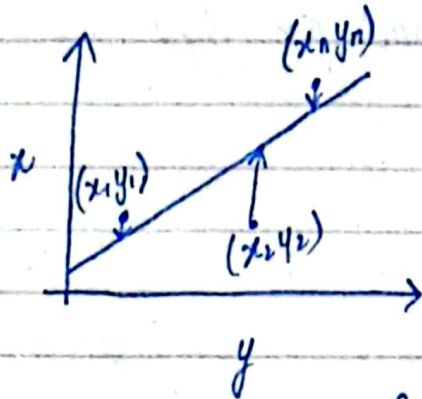
- Residual plots

→ Calculate the residual value of each data point and then plot it for x and y axis values.



If the points on residual plots are evenly scattered then relation is authentic (linear) if not then some non-linear relation should be suggested.

r (Correlation coefficient) vs r^2 (coefficient of determination)



$$SE_{LINE} = y_1 - (mx_1 + b)^2 + y_2 - (mx_2 + b)^2 \dots y_n - (mx_n + b)^2$$

Q: What %age of total variation in y is described by the variation in x ?
to answer this:

$$\text{Total variation in } y = SE \text{ from } \bar{y} = (y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + (y_n - \bar{y})^2$$

Q: What %age of variation is not described by the variation in x or by the regression line

$$\frac{SE_{LINE}}{SE_{\bar{y}}}$$

$$r^2 \Rightarrow 1 - \frac{SE_{LINE}}{SE_{\bar{y}}} \Rightarrow \left(\begin{array}{l} \text{What \%age of total variation} \\ \text{is described by the line} \\ \text{or by the variation in } x \end{array} \right)$$

↓
coefficient
of determination

if $SE_{LINE} \rightarrow \text{small} \rightarrow \text{line is a good fit}$ (r^2 value close to 1)

\Rightarrow It will mean a lot of total variation in y is described by the variation in x (which is good in relation terms)

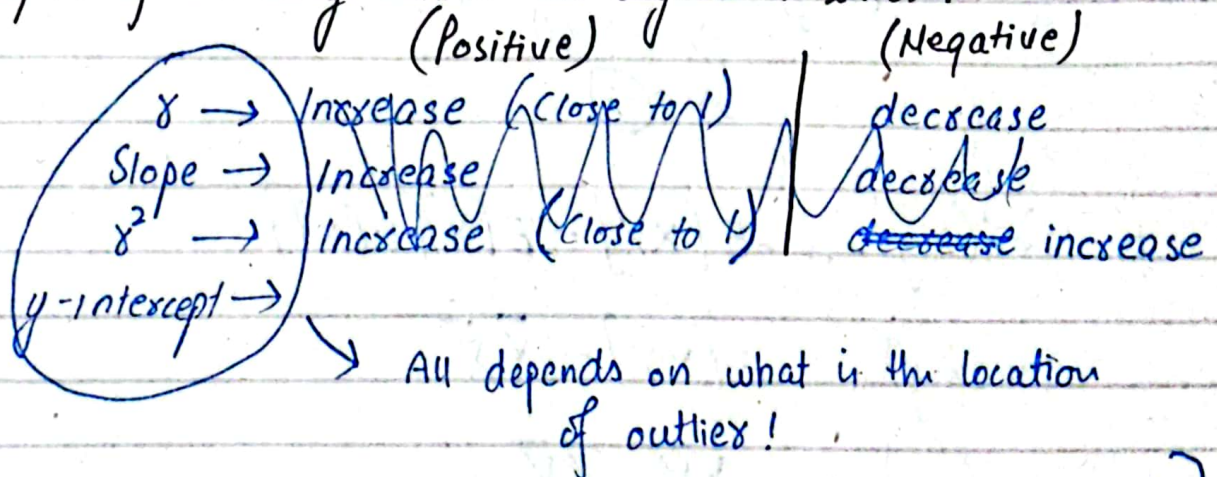
if $SE_{LINE} \rightarrow \text{large}$ ('opposite to above')

- Std of residuals / Root mean square error (RMSE)

$$RMSE = \sqrt{\frac{(\text{residual}_1)^2 + (\text{residual}_2)^2 + (\text{residual}_3)^2 + \dots + (\text{residual}_n)^2}{n-1}}$$

The lower this value is, the better fits the model !

- Impact of removing outliers on regression lines !



★ SE on regression line \Rightarrow (Sum is done because otherwise +ive and -ive residuals with counter effect each other)

"The whole idea of perfect regression line is to minimize the sum of squared error"

The line with least error !

From SE_{LINE} we get,

$$m = \frac{\overline{XY} - \bar{X}\bar{Y}}{(\bar{X})^2 - \bar{X}^2} \quad \left(m = \frac{\overline{XY} - \bar{X}\bar{Y}}{\bar{X}^2 - (\bar{X})^2} \right)$$

$$b = \bar{Y} - m\bar{X}$$

Covariance and regression line

$$\text{Cov}(X, Y) = E[(X - E(X)) \cdot (Y - E(Y))]$$

$$\text{i.e.} \quad \begin{array}{cc} X \rightarrow 1 & Y \rightarrow 2 \\ E(X) \rightarrow 3 & E(Y) \rightarrow 1 \end{array}$$

$$\begin{aligned} \text{Cov}(X, Y) &= (1-3)(2-1) \\ &\Rightarrow (-2)(1) \\ &\Rightarrow (-2) \end{aligned}$$

$$E[E(X)] \Rightarrow E(X)$$

$$\Rightarrow \text{Cov}(X, Y) = E[XY] - E[Y]E[X]$$

$$\text{Cov}(X, Y) = E[XY] - \bar{Y}\bar{X}$$

$$\text{Cov}(X, Y) = \overline{XY} - \bar{X}\bar{Y}$$

$$m = \frac{\text{Cov}(X, Y)}{\overline{X^2} - (\bar{X})^2} = \frac{\text{Cov}(X, Y)}{\overline{X^2} - \bar{X} \cdot \bar{X}} \Rightarrow \boxed{\frac{\text{Cov}(X, Y)}{\text{Var}(X)}}$$