# ALL BRANCHES
## ME,CE,EC,EE,CS

Probability and Statistics

Lecture No- 06

By- Vinay sir

# Revision

(Successes are Time Independent)

Discrete $\longrightarrow$ Binomial distribution $\longrightarrow$ Bernoulli Trials

$\quad\quad\quad\quad\quad\quad\longrightarrow$ Poisson distribution.

(Success are Time dependent).

When 'n' is large, Binomial distribution approaches to Poisson distribution.

$$P(x=\lambda) = n_{C_\lambda} \cdot p^\lambda \cdot q^{n-\lambda} \quad ; \quad\quad P(x=\lambda) = e^{-\lambda} \cdot \frac{\lambda^\lambda}{\lambda!}$$

$\quad\quad\quad\quad\quad\longrightarrow$ Binomial $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ Poisson

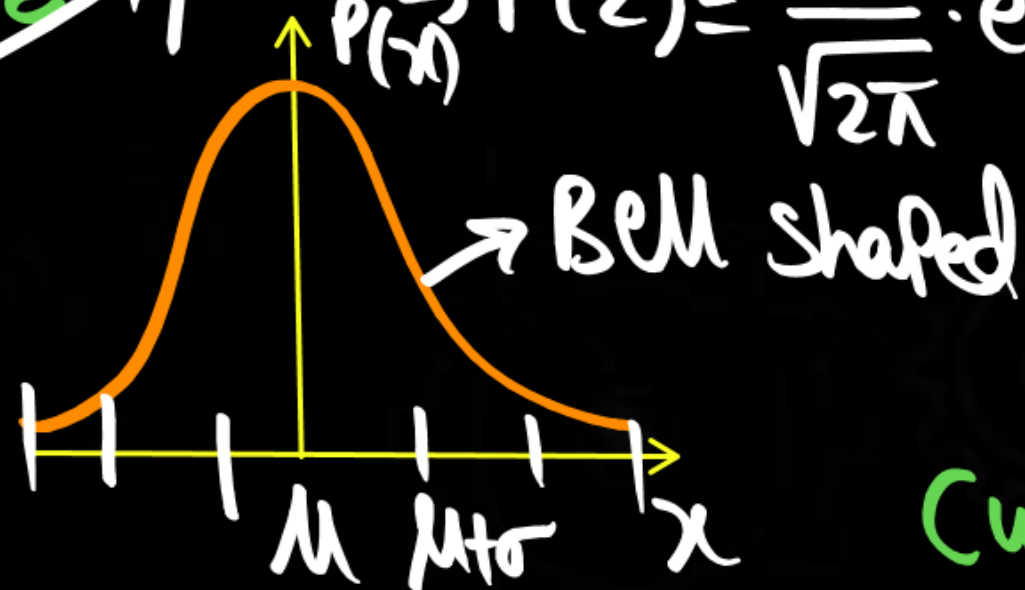**Binomial:** Mean $= nP$; Variance $= nPq$; $\sigma = \sqrt{nPq}$.

**Poisson:** Mean $=$ Variance $= \lambda$; $\sigma = \sqrt{\lambda}$.

**Normal/Gaussian distribution:** $P(x) = \dfrac{1}{\sigma \cdot \sqrt{2\pi}} \cdot \exp^{\dfrac{-(x-\mu)^2}{2\sigma^2}}$

$z = \dfrac{x-\mu}{\sigma}^{\nearrow 0}_{\searrow 1}$

$\Rightarrow P(z) = \dfrac{1}{\sqrt{2\pi}} \cdot e^{-z^2/2}$

$P(-1 \leq z \leq 1) = 0.68$

$P(-2 \leq z \leq 2) = 0.95$

$P(-3 \leq z \leq 3) = 0.997$



$\rightarrow$ Bell shaped

Cumulative function is 's' shaped.

# Exponential Distribution

The Probability density function of a Continuous Random Variable that is exponentially distributed is given as

$$P(x) = \begin{cases} \lambda \cdot e^{-\lambda x} & ; \ x \geq 0 \\ 0 & ; \ x < 0 \end{cases}$$

Where Mean of the Variable $= \dfrac{1}{\lambda}$.

$$\text{Mean} = \int_{-\infty}^{\infty} x \cdot P(x) \, dx$$

$$= \int_{-\infty}^{0} x \cdot (0) \, dx + \int_{0}^{\infty} x \cdot \left(\lambda \cdot e^{-\lambda x}\right) dx$$

$$= \int_{0}^{\infty} (\lambda x) \cdot \left(e^{-\lambda x}\right) dx$$

let $\lambda x = t$

$\Rightarrow \lambda \cdot dx = dt \Rightarrow dx = \dfrac{dt}{\lambda}$

L.L: $\lambda x = t \Rightarrow \lambda(0) = t \Rightarrow t = 0$

U.L: $\lambda x = t \Rightarrow \lambda(\infty) = t \Rightarrow t \to \infty$

$$\Rightarrow \text{Mean} = \int_{0}^{\infty} t \cdot e^{-t} \cdot \frac{dt}{\lambda}$$

$$= \frac{1}{\lambda} \cdot \int_{0}^{\infty} e^{-t} \cdot t^{2-1} \cdot dt$$

$$= \frac{1}{\lambda} \cdot \sqrt{2} = \frac{1}{\lambda} \cdot 1! = \frac{1}{\lambda}$$

$$\boxed{\therefore \text{Mean} = \frac{1}{\lambda}}$$

→ Variance: $\sigma^2$.

$\sigma^2 = E(x^2) - (E(x))^2$.

$E(x^2) = \int\limits_{-\infty}^{\infty} x^2 \cdot P(x)\, dx$

$= \int\limits_{0}^{\infty} x^2 \cdot \lambda \cdot \bar{e}^{\lambda x}\, dx$

let $\lambda x = t \Rightarrow x = \dfrac{t}{\lambda} \Rightarrow dx = \dfrac{dt}{\lambda}$.

$\dfrac{L.L.}{U.L.}$ $\lambda x = t \Rightarrow \lambda(0) = t \Rightarrow t = 0$

$\lambda x = t \Rightarrow \lambda(\infty) = t \Rightarrow t \to \infty$

$\therefore E(x^2) = \int\limits_{0}^{\infty} \dfrac{t^2}{\lambda^2} \cdot \lambda \cdot \bar{e}^{t} \cdot \dfrac{dt}{\lambda}$

$= \dfrac{1}{\lambda^2} \cdot \int\limits_{0}^{\infty} \bar{e}^{t} \cdot t^2\, dt$

$= \dfrac{1}{\lambda^2} \cdot \int\limits_{0}^{\infty} \bar{e}^{\lambda} \cdot t^{3-1}\, dt = \dfrac{1}{\lambda^2} \cdot \overline{|3}$

$= \dfrac{2!}{\lambda^2} = \dfrac{2}{\lambda^2}$.

$\therefore E(x^2) = \dfrac{2}{\lambda^2}$

$\therefore \sigma^2 = E(x^2) - (E(x))^2 = \dfrac{2}{\lambda^2} - \dfrac{1}{\lambda^2}$

$\boxed{\Rightarrow \sigma^2 = \dfrac{1}{\lambda^2}}$

∴ For a exponentially distributed Variable;

$$\text{Mean} = \text{Standard deviation} = \frac{1}{\lambda}$$

**Statistical Attributes:** The Parameters that governs the distributi~~on~~

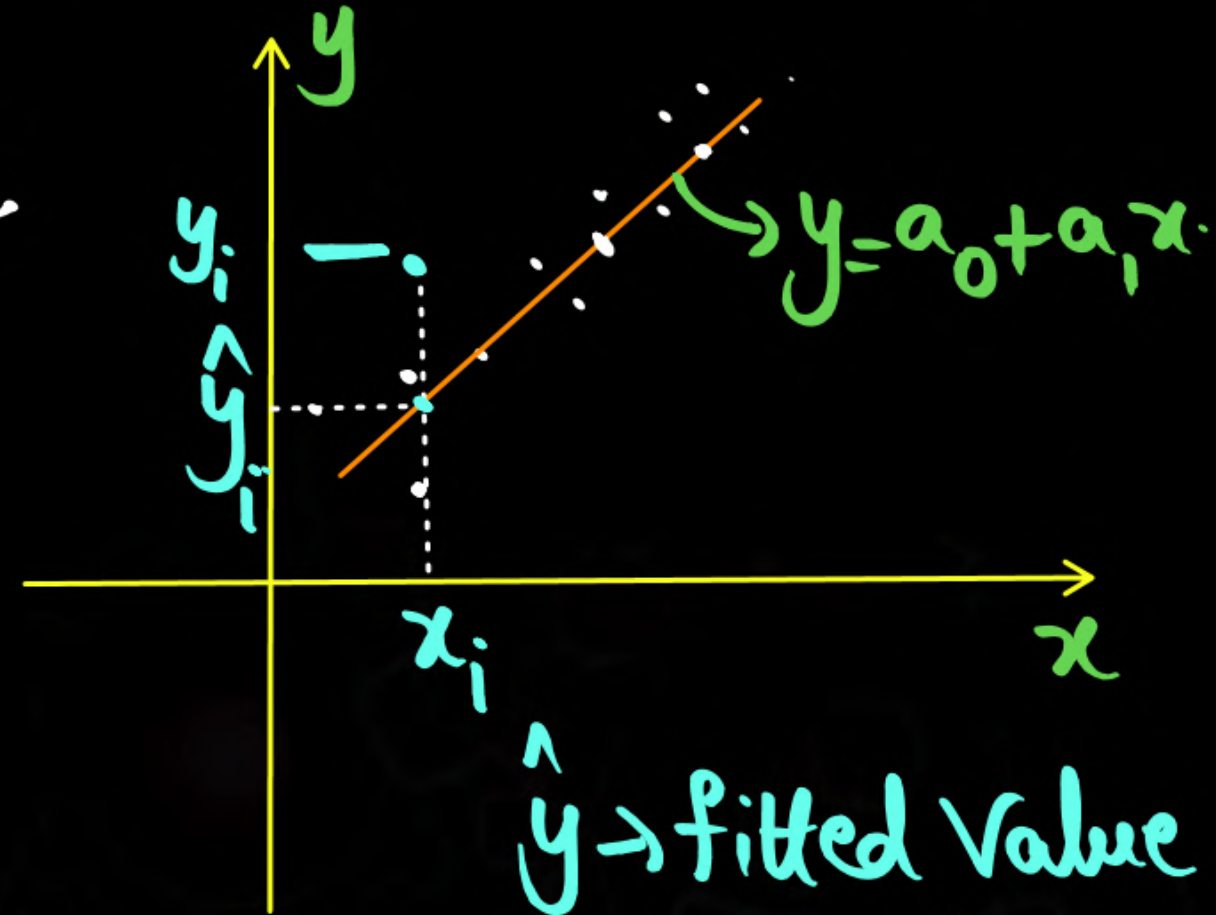| Distribution: | Statistical Attributes | |
|---|---|---|
| 1. Binomial | $n, P$ | $\longrightarrow 2$ |
| 2. Poisson | $\lambda$ | $\longrightarrow 1$ |
| 3. Gaussian/Normal | $\mu, \sigma$ | $\longrightarrow 2$ |
| 4. Exponential | $\lambda$ | $\longrightarrow 1$ |

# Linear Regression

Regression $\rightarrow$ fitting of a curve b/w a set of Variables.

Let the Data Points be $(x_1, y_1), (x_2, y_2), (x_3, y_3),$
$\ldots \ldots (x_N, y_N)$

$y_i \rightarrow$ Actual Data Point Value.

$\hat{y}_i \rightarrow$ fitted Value.

$$\text{Error} = e_i = y_i - \hat{y}_i = y_i - (a_0 + a_1 x_i)$$

$\rightarrow y = a_0 + a_1 x$

$\hat{y} \rightarrow$ fitted Value at $x_i$

$\Rightarrow \hat{y}_i = a_0 + a_1 x_i$

Sum of Squares of errors: $\sum\limits_{i=1}^{N} e_i^2$

∴ For the fit be a good fit, Sum of Squares of errors should be least Possible.

$\Rightarrow \sum\limits_{i=1}^{N} e_i^2$ should be least. (Least Square Regression) (LSR).

$\Rightarrow \sum\limits_{i=1}^{N} (y_i - a_0 - a_1 x_i)^2 \longrightarrow$ Should be least

$\Rightarrow f(a_0, a_1) = \sum\limits_{i=1}^{N} (y_i - a_0 - a_1 x_i)^2$

$$f(a_0, a_1) = \sum_{i=1}^{N} (y_i^2 + a_0^2 + a_1^2 \cdot x_i^2 - 2a_0 y_i + 2a_0 a_1 x_i - 2y_i a_1 x_i)$$

For $f(a_0, a_1)$ to be minimum,

$$\frac{\partial f}{\partial a_0} = 0 \quad \text{and} \quad \frac{\partial f}{\partial a_1} = 0$$

$$\frac{\partial f}{\partial a_0} = 0 \Rightarrow \sum_{i=1}^{N} (2a_0 - 2y_i + 2a_1 x_i) = 0 \Rightarrow \sum_{i=1}^{N} a_0 - \sum_{i=1}^{N} y_i + a_1 \sum_{i=1}^{N} x_i = 0$$

$$\frac{\partial f}{\partial a_1} = 0 \Rightarrow \sum_{i=1}^{N} (2a_1 x_i^2 + 2a_0 x_i - 2x_i y_i) = 0 \Rightarrow \sum_{i=1}^{N} a_1 x_i^2 + \sum_{i=1}^{N} a_0 \cdot x_i - \sum_{i=1}^{N} x_i y_i = 0$$

$$\Rightarrow a_0 \cdot N + a_1 \cdot \sum_{i=1}^{N} x_i = \sum_{i=1}^{N} y_i \longrightarrow \text{①}$$

$$a_0 \cdot \sum_{i=1}^{N} x_i + a_1 \cdot \sum_{i=1}^{N} x_i^2 = \sum_{i=1}^{N} x_i y_i \longrightarrow \text{②}$$

$$\sum_{i=1}^{N} 1 = 1 + 1 + \cdots + 1$$
$$= \underline{N}$$

Solving above two Equations

$$\text{①} \times \sum_{i=1}^{N} x_i - \text{②} \times N$$

$$\Rightarrow a_1 \left\{ \left( \sum_{i=1}^{N} x_i \right)^2 - N \cdot \sum_{i=1}^{N} x_i^2 \right\} = \left( \sum_{i=1}^{N} x_i \right) \left( \sum_{i=1}^{N} y_i \right) - N \cdot \sum_{i=1}^{N} x_i y_i$$

$$\Rightarrow a_1 = \frac{N \cdot \left( \sum_{i=1}^{N} x_i y_i \right) - \left( \sum_{i=1}^{N} x_i \cdot \sum_{i=1}^{N} y_i \right)}{N \cdot \sum_{i=1}^{N} x_i^2 - \left( \sum_{i=1}^{N} x_i \right)^2}$$

Since $y = a_0 + a_1 x \longrightarrow a_1$ is slope of the line.

$$\therefore a_1 = \frac{N \cdot \Sigma xy - (\Sigma x)(\Sigma y)}{N \cdot \Sigma x^2 - (\Sigma x)^2}$$

Substituting $a_1$ in ①

$$\Rightarrow a_0 \cdot N + \left\{ \frac{N \cdot \sum_{i=1}^{N} x_i y_i - \left(\sum_{i=1}^{N} x_i\right)\left(\sum_{i=1}^{N} y_i\right)}{N \cdot \sum_{i=1}^{N} x_i^2 - \left(\sum_{i=1}^{N} x_i\right)^2} \right\} \sum_{i=1}^{N} x_i = \sum_{i=1}^{N} y_i$$

$$\Rightarrow a_0 \cdot N + \frac{N \cdot \sum_{i=1}^{N} x_i y_i \cdot \sum_{i=1}^{N} x_i - \left(\sum_{i=1}^{N} x_i\right)^2 \left(\sum_{i=1}^{N} y_i\right)}{N \cdot \sum_{i=1}^{N} x_i^2 - \left(\sum_{i=1}^{N} x_i\right)^2} = \sum_{i=1}^{N} y_i$$

$$\Rightarrow N \cdot a_0 \cdot \sum_{i=1}^{N} x_i^2 - N \cdot a_0 \left(\sum_{i=1}^{N} x_i\right)^2 + N \sum_{i=1}^{N} x_i y_i \cdot \sum_{i=1}^{N} x_i - \sum_{i=1}^{N} y_i \left(\sum_{i=1}^{N} x_i\right)^2$$

$$= N \cdot \sum_{i=1}^{N} x_i^2 \cdot \sum_{i=1}^{N} y_i - \left(\sum_{i=1}^{N} x_i\right)^2 \left(\sum_{i=1}^{N} y_i\right)$$

$$\Rightarrow a_0 \left( N \cdot \sum_{i=1}^{N} x_i^2 - \left( \sum_{i=1}^{N} x_i \right)^2 \right) = \sum_{i=1}^{N} x_i^2 \cdot \sum_{i=1}^{N} y_i - \sum_{i=1}^{N} (x_i y_i) \cdot \sum_{i=1}^{N} x_i$$

$$\Rightarrow a_0 = \frac{\displaystyle\sum_{i=1}^{N} x_i^2 \cdot \sum_{i=1}^{N} y_i - \sum_{i=1}^{N} (x_i y_i) \cdot \sum_{i=1}^{N} x_i}{N \cdot \displaystyle\sum_{i=1}^{N} x_i^2 - \left( \sum_{i=1}^{N} x_i \right)^2}$$

$\therefore$ For a linear fit, $y = a_0 + \boxed{a_1} x$ → Slope of the line

$$a_0 = \frac{(\sum x^2)(\sum y) - (\sum xy) \cdot (\sum x)}{n(\sum x^2) - (\sum x)^2} \quad ; \quad a_1 = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

Goodness of the fit: $(R^2)$. $R^2$ lies in 0 to 1.

$R^2 \rightarrow$ Coefficient of determination.

$$R^2 = 1 - \frac{SS|_{Res}}{SS|_{Tot}}$$

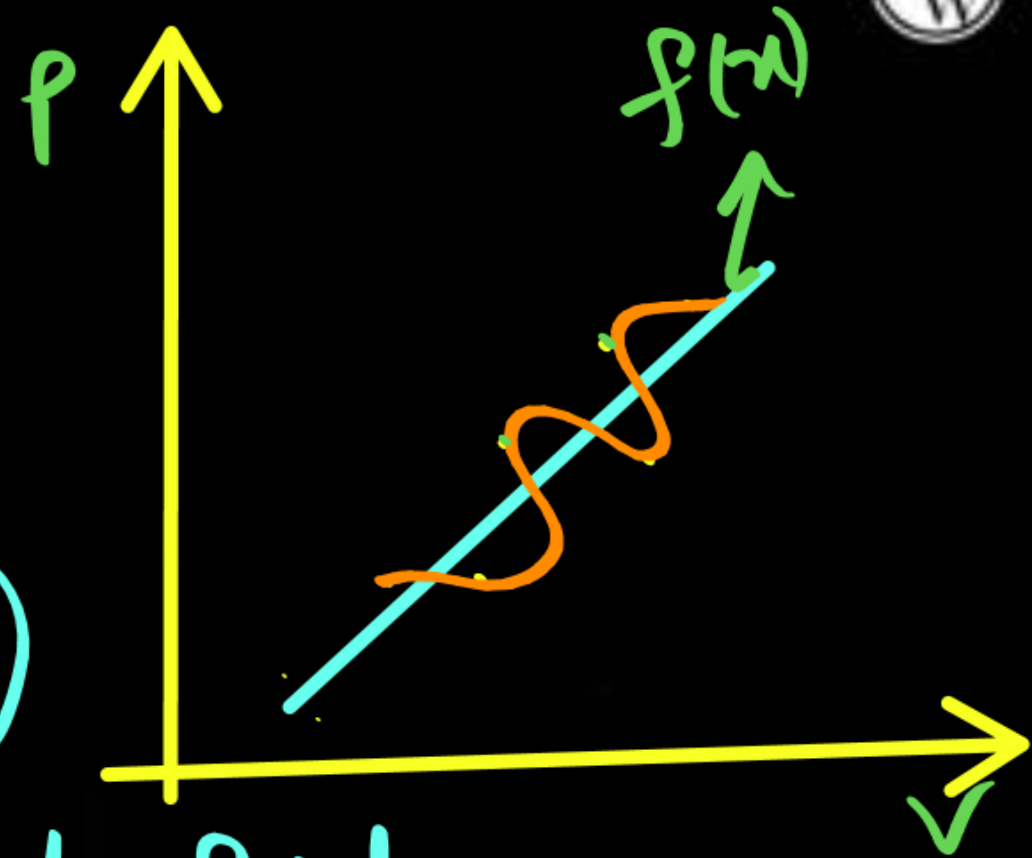$$\Rightarrow R^2 = 1 - \frac{\sum\limits_{i=1}^{N}(y_i - f_i)^2}{\sum\limits_{i=1}^{N}(y_i - \bar{y})^2}$$

Residue
$\uparrow$
error $= (y_i - f_i)$

$y_i \rightarrow$ Actual Data Point

$f_i \rightarrow$ fitted value.

$$\bar{y} = \sum\limits_{i=1}^{n} y_i$$

$(y_i - \bar{y}) \rightarrow$ Deviation

# Correlation Coefficient

$$\text{Correlation coefficient} = \lambda = \frac{Cov(x,y)}{n \cdot \sigma_x \cdot \sigma_y}$$

$$Cov(x,y) = E\left((x-\bar{x})(y-\bar{y})\right)$$

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

$$\sigma_x = \sqrt{\frac{1}{N} \cdot \sum (x_i - \bar{x})^2} \;;\; \sigma_y = \sqrt{\frac{1}{N} \cdot \sum (y_i - \bar{y})^2}$$

$$y = a + bx. \rightarrow \text{Linear fit}$$

$$\Sigma y = a\Sigma 1 + b\cdot\Sigma x.$$

$$\Rightarrow \frac{1}{N}\cdot\Sigma y = \frac{a}{N}\cdot\Sigma 1 + \frac{b}{N}\cdot\Sigma x$$

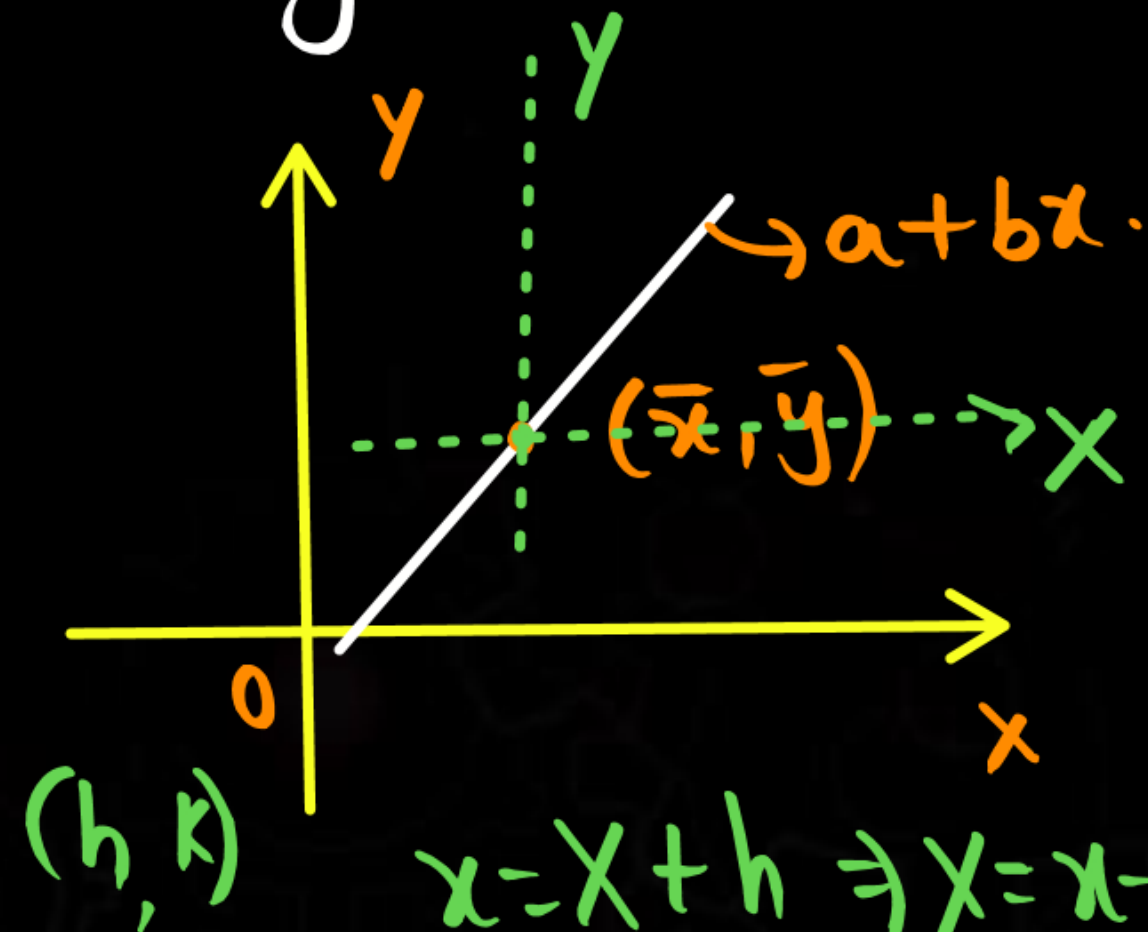$$\boxed{\Rightarrow \bar{y} = a + b\cdot\bar{x}}$$

$\rightarrow$ For 'N' data Points

$$(x_1, y_1), (x_2, y_2), \cdots (x_n, y_n),$$

the Means $(\bar{x}, \bar{y})$ lie on the fit.

since $\Sigma y = a + b\cdot\Sigma x$

$$\Rightarrow \Sigma xy = a\cdot\Sigma x + b\cdot\Sigma x^2.$$



$\rightarrow a + bx.$

$(\bar{x}, \bar{y})$

$(h, k)$

$$x = X + h \Rightarrow X = x - \bar{x}$$
$$y = Y + k \Rightarrow Y = y - \bar{y}$$

$$\Rightarrow \Sigma(x-\bar{x})(Y-\bar{y}) = a \cdot \Sigma(x-\bar{x})^{\nearrow 0} + b \cdot \Sigma(x-\bar{x})^2$$

(sum of deviations is 0).

$$\Rightarrow b = \text{Slope} = \frac{\Sigma(x-\bar{x})(Y-\bar{y})}{\Sigma(x-\bar{x})^2}$$

$$\sigma_x^2 = \frac{1}{n} \cdot \Sigma(x-\bar{x})^2$$

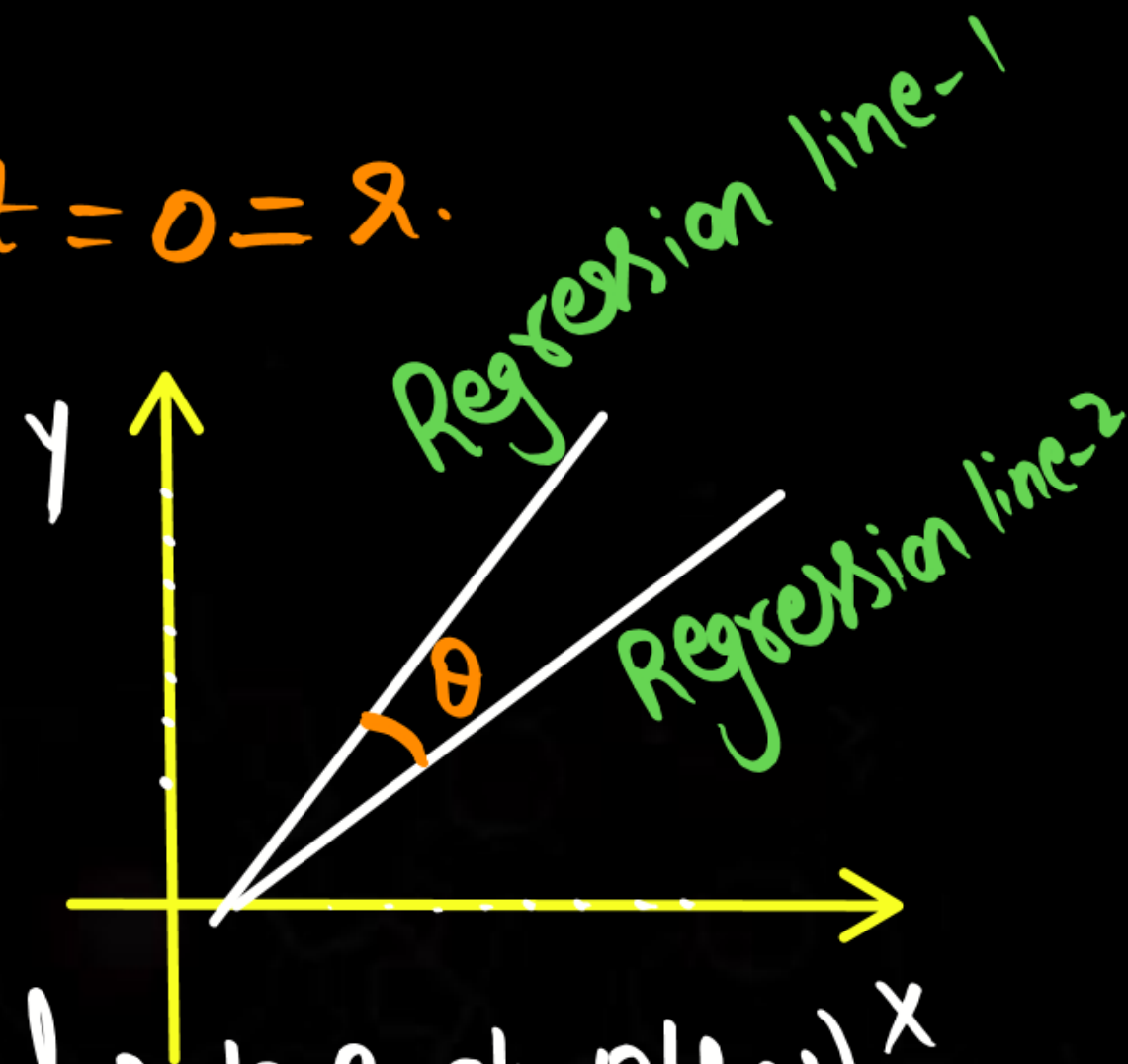$$\Rightarrow b = \text{Slope} = \frac{\Sigma(x-\bar{x})(Y-\bar{y})}{n \cdot \sigma_x^2}$$

$$\Rightarrow b = \frac{\lambda \cdot n \cdot \sigma_x \cdot \sigma_y}{n \cdot \sigma_x^2}$$

$$\boxed{\Rightarrow \text{Slope} = \lambda \cdot \frac{\sigma_y}{\sigma_x}}$$

For two independent Variables,
$$cov(X, Y) = 0$$

$$\Rightarrow \text{Correlation coefficient} = 0 = r.$$

Angle b/w the Correlation lines

$$\boxed{Tan\theta = \frac{1 - r^2}{r} \cdot \frac{\sigma_x \cdot \sigma_y}{\sigma_x^2 + \sigma_y^2}}$$

Regression line-1

Regression line-2

Regression line-2

If $r = 0 \Rightarrow \theta = \frac{\pi}{2}$ (lines are perpendicular to each other)

If $r = 1 \Rightarrow \theta = 0$ ; (lines are parallel or coincident)