

* Regression Analysis *

$$Y = \beta_0 + \beta_1 X + \epsilon$$

Y - response / Dependent / study variable

X - independent / regressor / predictor / explanatory variable

β_0, β_1 - unknown regression coefficients / parameters

ϵ - an observable & we can not control it.

* Cause & Effect

X - cause &

Y - effect

1) X - I.Q. & T - marks

X - cause & Y - effect

X - effect & Y - cause.

2) X - Rainfall & T - yield

Cause - X - Rainfall

Effect - T - Yield

3) X - H.d. & Literates & T - H.d. & Criminals

As X increases, Y also increases because of third variable Z - population. Such correlation is called "Non-sense correlation".

$y = x_1, x_2, \dots, x_k$

The general regression model is,

$$y = f(x_1, x_2, \dots, x_k; \beta_1, \dots, \beta_k) + \epsilon$$

f - well defined f ?

x_1, \dots, x_k - Regressors

β_1, \dots, β_k - Regression coefficients

ϵ - If $\epsilon = 0$, then the model

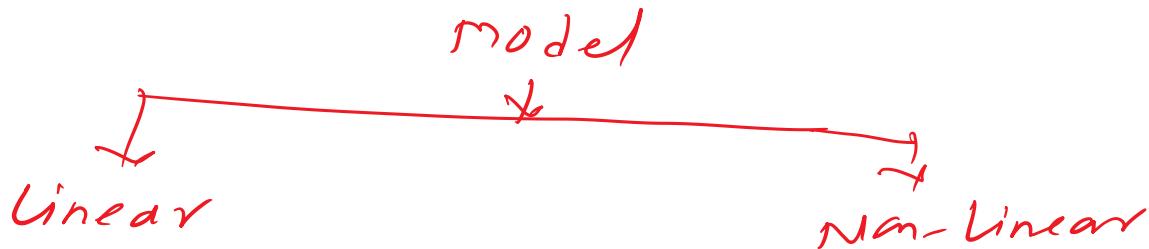
$$y = f(x_1, \dots, x_k, \beta_1, \dots, \beta_k)$$

is a mathematical model

for $\epsilon \neq 0$, model is
statistical model.

*Types of models:-

$$y = f(x_1, \dots, x_k; \beta_1, \dots, \beta_k) + \epsilon$$



* It is linear in parameters,
 β_1, \dots, β_k

* E.g. $y = \beta_1 x_1^2 + \beta_2 \sqrt{x_2} + \beta_3 x_3 + \epsilon$

* It is non-linear in
parameters β_1, \dots, β_k

+ e.g. $y = \beta_1 x_1^2 + \beta_2 \sqrt{x_2}$
 $+ \beta_3 \log x_3 + \epsilon$

$$Y = \beta_1 X_1^2 + \beta_2 \sqrt{X_2} + \beta_3 X_3 + \varepsilon \quad \textcircled{I}$$

$$\frac{\partial Y}{\partial \beta_1} = X_1^2, \quad \frac{\partial Y}{\partial \beta_2} = \sqrt{X_2} \quad \text{and} \quad \frac{\partial Y}{\partial \beta_3} = X_3$$

Hence model \textcircled{I} is linear

$$Y = \beta_1^2 X_1^2 + \beta_2 \sqrt{X_2} + \beta_3 \log X_3 + \varepsilon \quad \textcircled{II}$$

$$\frac{\partial Y}{\partial \beta_1} = 2 \beta_1 X_1^2 - \frac{\partial Y}{\partial \beta_2} = \sqrt{X_2} \quad \text{and} \quad \frac{\partial Y}{\partial \beta_3} = \log X_3$$

Hence, the model is not linear.
i.e. Non-linear.

In general, linear model is defined as,

$$\begin{aligned} Y &= f(x_1, \dots, x_k, \beta_1, \dots, \beta_k) + \varepsilon \\ &= \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon \end{aligned}$$

X & Y - pre-defined variable

Knowledge of model is depend on knowledge of parameters.

Estimation methods:-

- is Maximum Likelihood Estimation method.
- Method of moments
- 3) Least square method, etc.
 do need an knowledge about dist. of y

"Regression" :- To move in the backward direction.

\underline{y} - Yield, \underline{x}_1 - quality & Fertilizer
 \underline{x}_2 - level & irrigation.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

+ Steps in Regression :-

1) statement of the problem/objectives

find the yield.

2) choice of relevant variables:-

- fertilizer, rainfall, irrigation, temp,
 wind speed, etc.

3) Collection of the Data:-

1. temp - °C

2. temp = $\begin{cases} 1 & > 80^\circ C \\ 0 & \text{o.w.} \end{cases}$

17, 19, 31, 32,
0, 0, 1, 1

age = mean - s

y	x_1	x_2	\dots	x_k
y_1	x_{11}	x_{12}	\dots	x_{1k}
y_2	x_{21}	x_{22}	\dots	x_{2k}
\vdots	\vdots	\vdots	\ddots	\vdots
y_n	x_{n1}	x_{n2}	\dots	x_{nk}

4. Specification of the model :-

$$y = \beta_1 x_1^2 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

$$\approx y = \beta_1^2 x_1 + \beta_2^2 x_2 + \beta_3 \log x_3 + \epsilon.$$

$$y = f(x_1, \dots, x_k, \beta_1, \dots, \beta_k) + \epsilon.$$

5. choice of parameter estimation method (model fitting):-

- MLE
- MMSE
- OLS
- Ridge.
- LAP
- etc.

6. Fitting of the model :-

$$\underline{\beta_0, \beta_1, \dots, \beta_k} \quad \underline{\text{OLS}}$$

$$\underline{y = f(x_1, \dots, x_k, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)}$$

$$\{(x_i, y_i) : i=1, 2, \dots, n\}$$

$$y = \beta_0 + \beta_1 x + \epsilon$$

$$\hat{\beta}_0 \text{ & } \hat{\beta}_1$$

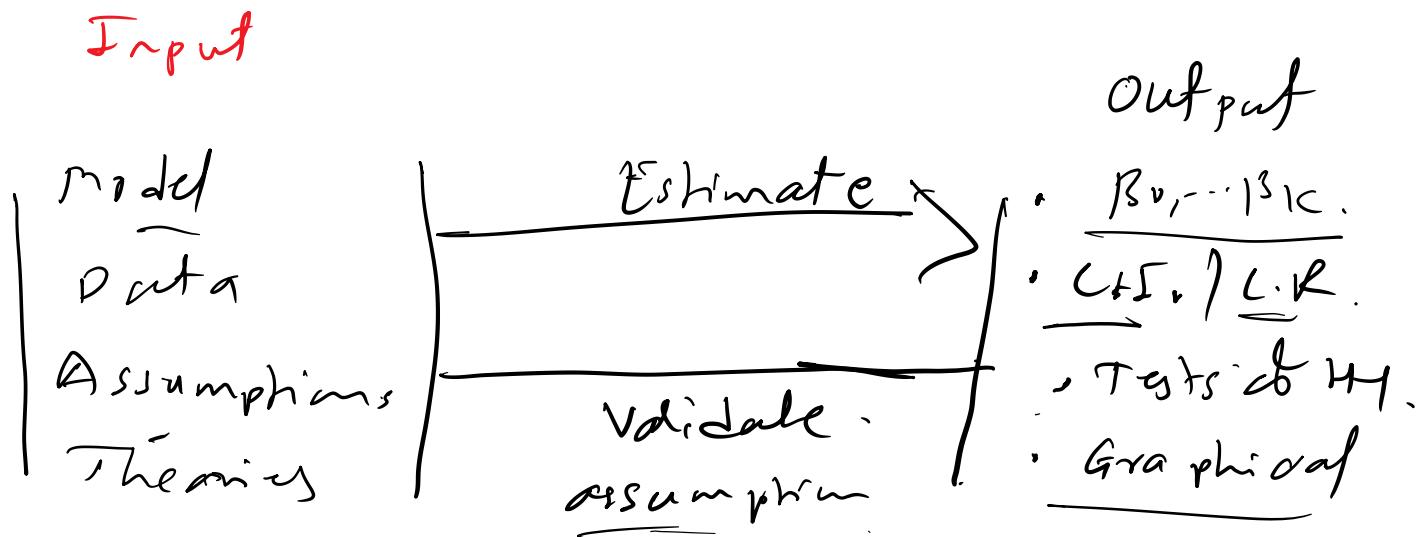
$$\underline{y = \hat{\beta}_0 + \hat{\beta}_1 x}$$

fitting value of y - $x \in \text{dataset}$.

predicted value of y - $\cdot \text{any } x$.

7. model validation :-

- * Assumptions.



* Type of Regression:-

1. Univariate

only one response variable.

2. Multivariate

Two or more response variables.

3. simple

only one regressor

4. multiple

Two or more regressors.

* Simple Linear Reg. Analysis +.

$$\underline{Y} = \beta_0 + \beta_1 X + \underline{\epsilon}$$

Y - response / study / Dependent variable.

X - independent / regressor /

β_0, β_1 - unknown Reg. coeff.

ϵ - error. term.

Assumptions:-

- i) ϵ is r.v. with mean 0 & variance σ^2 .
r.c. $E(\epsilon_i) = 0$ & $V(\epsilon_i) = \sigma^2$ $\forall i = 1, 2, \dots, n$.
- ii) $\text{cov}(\epsilon_i, \epsilon_j) = 0 \Rightarrow \epsilon_i, \epsilon_j$ are uncorrelated.
i.e. no. auto correlation.
- iii) $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$
- iv) model is linear.

If x is not r.v.

$$\begin{aligned} E(y) &= E(\underline{\beta_0 + \beta_1 x + \epsilon}) \quad (\because y = \beta_0 + \beta_1 x + \epsilon) \\ &= \beta_0 + \beta_1 x + 0 \\ &= \beta_0 + \beta_1 x. \end{aligned}$$

& $V(Y) = V(\underline{\alpha + \epsilon}) \quad \alpha = \beta_0 + \beta_1 x.$

$$= V(\epsilon) = \sigma^2$$

$$Y \sim N(\beta_0 + \beta_1 x, \sigma^2)$$

$$\underline{y_i} \sim N(\underline{\beta_0 + \beta_1 x_i}, \sigma^2)$$

y_i are independent but identical.

If \underline{x} is r.v.

$$\underline{E(y|x)} = \beta_0 + \beta_1 x$$

&

$$\underline{V(y|x)} = \sigma^2$$

If β_0, β_1, σ are known then model is completely known.

Now we need to find out unknown parameters.

Method of Estimation :-

1. Least square.

2. MLE

3. MM

4. Ridge... ch.

* Simple Linear Regression *

For the bivariate data,

$\{(x_i, y_i); i=1, 2, \dots, n\}$, the simple linear regression model is

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i=1, \dots, n.$$

y_i - response variable or dependent variable

x_i - regressor or predictor variable or independent variable

β_0, β_1 - regression coefficients

(β_0 - intercept & β_1 - slope)

ϵ_i - Error term

Assumptions:-

$$\Rightarrow E(\varepsilon_i) = 0 \quad \text{and} \quad V(\varepsilon_i) = \sigma^2 \quad i = 1, 2, \dots, n$$

σ^2 is unknown

i) $\text{cov}(\varepsilon_i, \varepsilon_j) = 0 \Rightarrow \varepsilon_i \text{ and } \varepsilon_j \text{ are uncorrelated.}$

i.e. no autocorrelation.

iii) $\varepsilon_i \sim \text{i.i.d. } N(0, \sigma^2)$

iv) model is linear

i) If x is not r.v. (constant)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$\therefore E(y_i) = \beta_0 + \beta_1 x_i \quad (\because E(\varepsilon_i) = 0)$$

& $V(y_i) = \sigma^2 \quad (\because V(\varepsilon_i) = \sigma^2)$

ii) If x is r.v.
i.e. $y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

$$E(y_i | x_i) = \beta_0 + \beta_1 x_i$$

& $V(y_i | x_i) = \sigma^2$

i.e. $y_i | x_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

Example:-

X - Sugarcane & Y - Production of sugar

f) Let simple linear regⁿ model is

$$Y = 10 + 15 X$$

ii) X - weight & Y - height

$$f = 5 + 10 X$$

if $X = 0$

$$\overline{E(Y)} = 5 \quad \text{for } X = 0$$

+ methods for the estimation of unknown parameters
 $(\beta_0, \beta_1, \sigma^2)$:-

+ Least square method :-

→ Direct Regression method :-

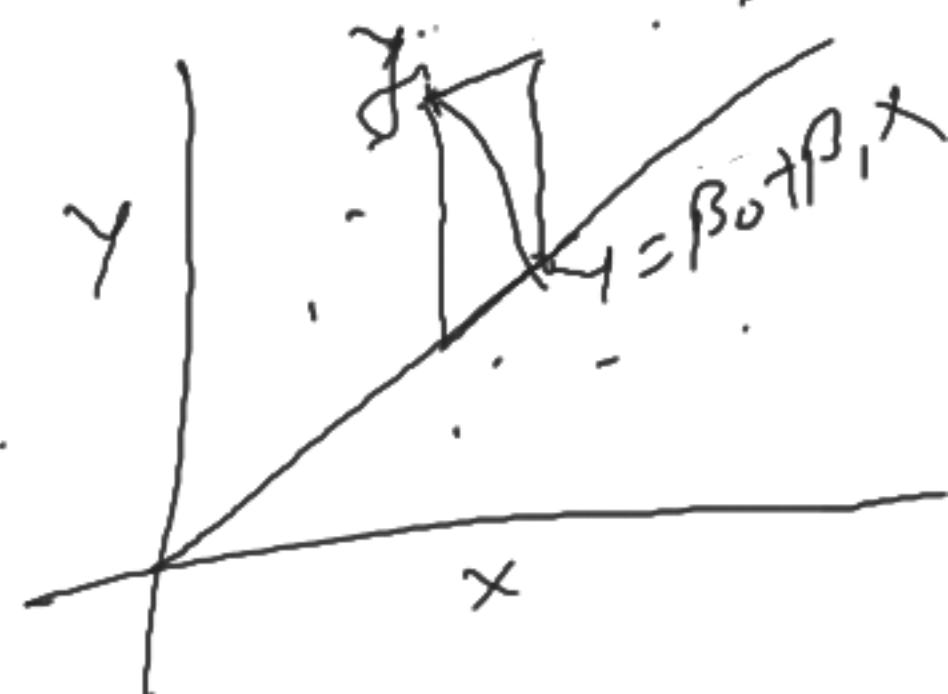
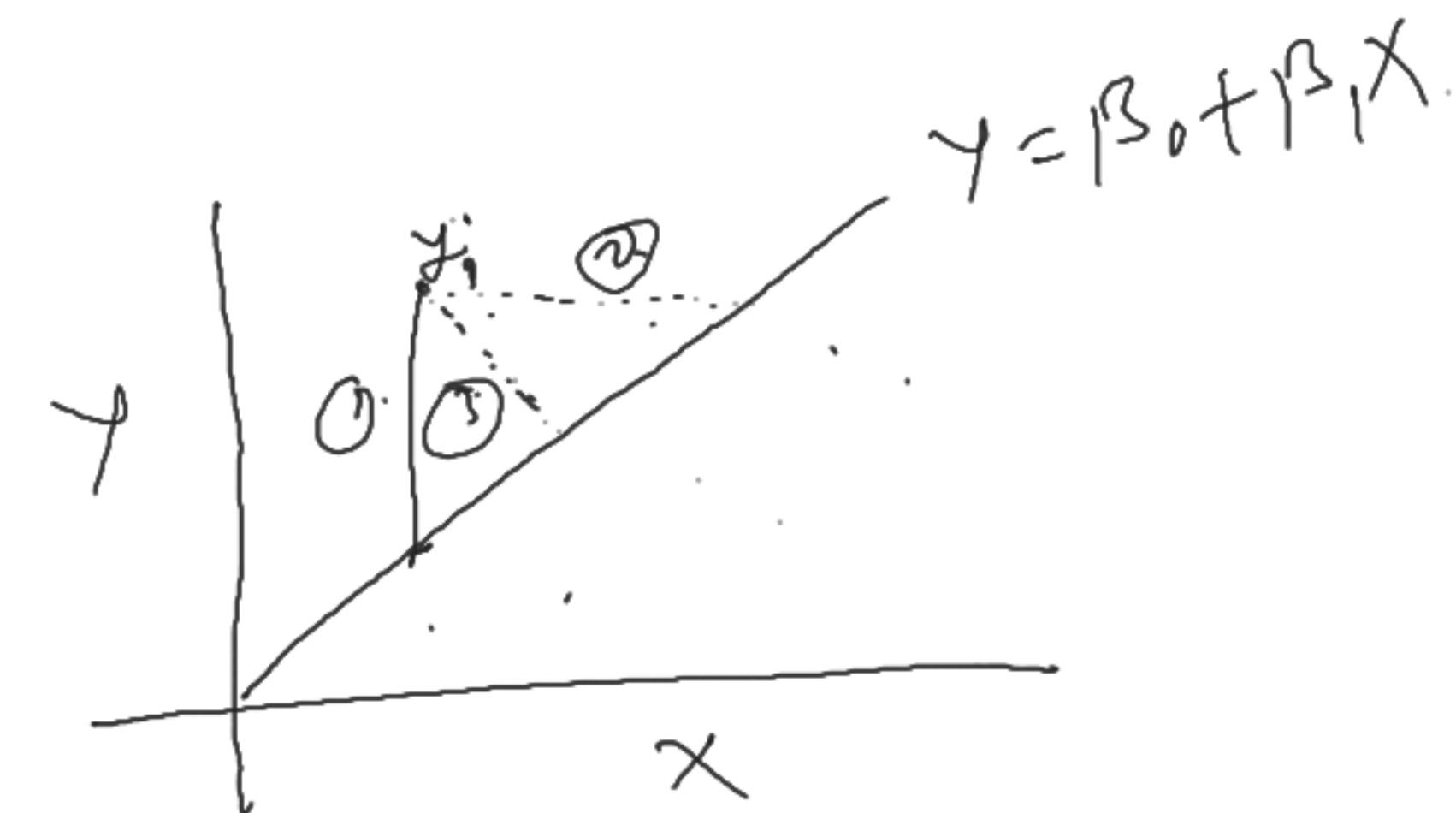
minimize vertical distance

i) Reverse (or Inverse) Reg.

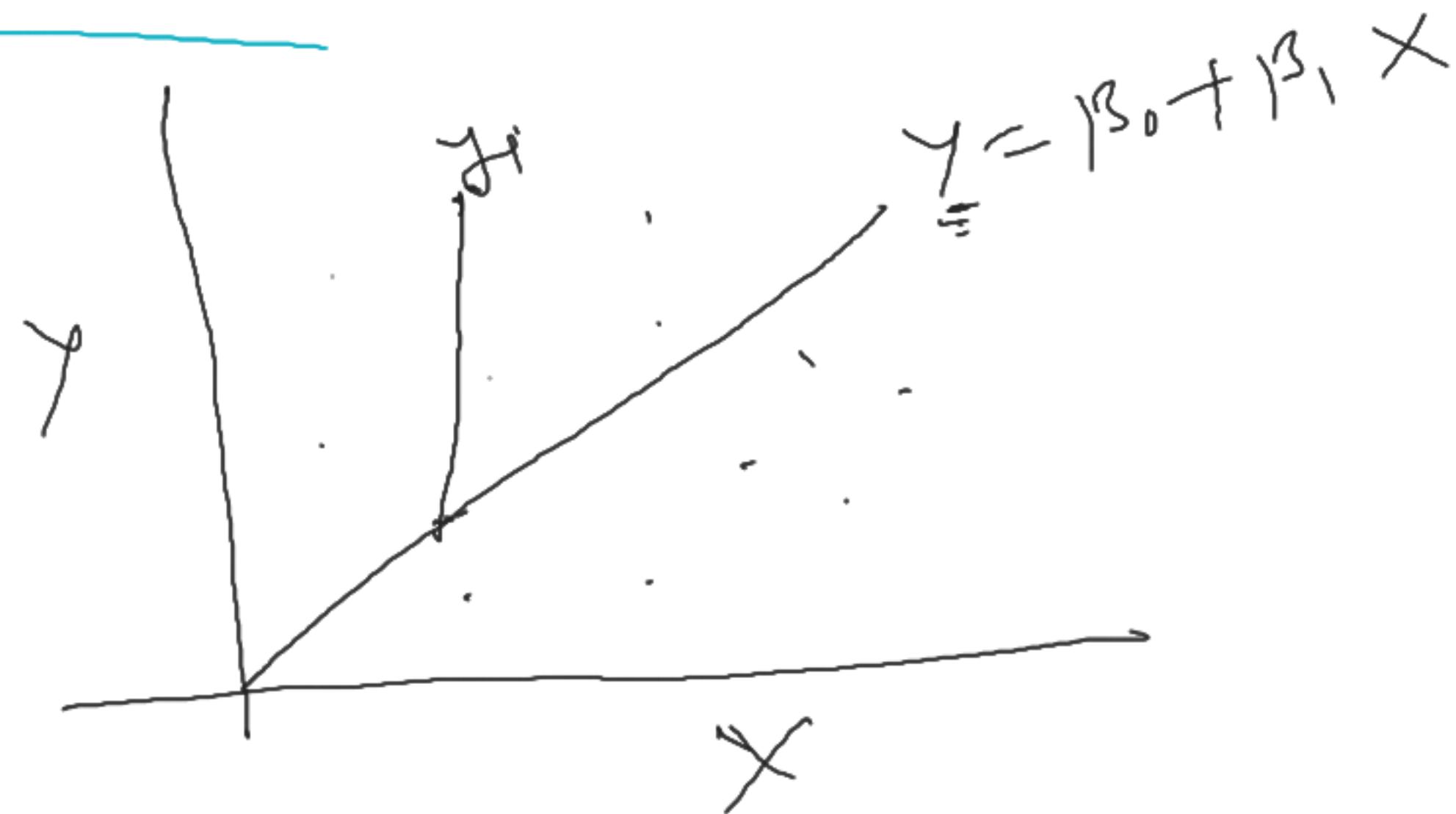
minimize horizontal distance

iii) orthogonal (or major axis) Reg.

iv) Reduced major axis reg.



* Direct Reg. method :-



Least square estimator :- $\underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$

Least absolute deviation :- $\underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n |y_i - \beta_0 - \beta_1 x_i|$

In general, $\underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n g(y_i - \beta_0 - \beta_1 x_i)$, $g(\cdot)$ is function

Suppose

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial S}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i = n \beta_0 + \beta_1 \sum x_i \quad \text{--- } \textcircled{I}$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\Rightarrow \sum_{i=1}^n y_i x_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \quad \text{--- } \textcircled{II}$$

Least square normal eq's.

$$\sum_{i=1}^n y_i = n\beta_0 + \beta_1 \sum_{i=1}^n x_i \quad \text{--- } \textcircled{3}$$

$$g \quad \sum_{i=1}^n y_i x_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \quad \text{--- } \textcircled{4}$$

From eq: $\textcircled{1}$

$$n\beta_0 = \sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i$$

i.e. $\beta_0 = \bar{y} - \beta_1 \bar{x} \quad \text{--- } \textcircled{II}$

From eq's \textcircled{II} & $\textcircled{4}$

$$\sum_{i=1}^n y_i x_i = (\bar{y} - \beta_1 \bar{x}) \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2$$

$$\begin{aligned}\therefore \sum_{i=1}^n y_i x_i &= (\bar{y} - \beta_1 \bar{x}) \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \\ &= (\bar{y} - \beta_1 \bar{x}) n \bar{x} + \beta_1 \sum_{i=1}^n x_i^2 \quad (\because \bar{x} = \frac{\sum x_i}{n})\end{aligned}$$

$$\begin{aligned}&= n \bar{x} \bar{y} - n \beta_1 \bar{x}^2 + \beta_1 \sum_{i=1}^n x_i^2 \\ &= n \bar{x} \bar{y} + \beta_1 \left(\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) \\ \therefore \hat{\beta}_1 &= \frac{\sum y_i x_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \frac{\sum y_i x_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}\end{aligned}$$

Vence

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \text{e} \quad \hat{\beta}_1 = \frac{\sum y_i x_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

$$\text{Let } S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i^2 - 2\bar{x}x_i + \bar{x}^2)$$

$$= \sum x_i^2 - 2\bar{x} \sum x_i + n\bar{x}^2$$

$$= \sum x_i^2 - 2\bar{x} n\bar{x} + n\bar{x}^2$$

$$= \sum x_i^2 - 2n\bar{x}^2 + n\bar{x}^2$$

$$= \sum x_i^2 - n\bar{x}^2$$

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$= \sum x_i y_i - n\bar{x}\bar{y}$$

$$\therefore \text{cov}(x, y) = \frac{S_{xy}}{n} \text{ and } \text{v}(x) = \frac{S_{xx}}{n}.$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

g

$$\hat{\beta}_1 = \frac{\sum y_i x_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \frac{s_{xy}}{s_{xx}}$$

This is a global optimal soln if

$$\frac{\partial S}{\partial \beta_0^2} > 0, \quad \frac{\partial^2 S}{\partial \beta_1^2} > 0$$

g

$$\frac{\partial^2 S}{\partial \beta_0^2} \frac{\partial^2 S}{\partial \beta_1^2} - \left(\frac{\partial^2 S}{\partial \beta_0 \partial \beta_1} \right)^2 > 0$$

Since,

$$S(\beta_0, \beta_1) = \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial S}{\partial \beta_0} = -2 \sum (y_i - \beta_0 - \beta_1 x_i) = -2 \sum y_i + 2n\beta_0 + 2\beta_1 \sum x_i$$

$$\frac{\partial^2 S}{\partial \beta_0^2} = 2n > 0$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum (y_i - \beta_0 - \beta_1 x_i) x_i = -2 \sum y_i x_i + 2\beta_0 \sum x_i + 2\beta_1 \sum x_i^2$$

$$\frac{\partial^2 S}{\partial \beta_1^2} = 2 \sum x_i^2 > 0$$

$$\frac{\partial^2 S}{\partial \beta_0 \partial \beta_1} = 2 \sum x_i$$

$$\frac{\partial^2 S}{\partial \beta_0^2} \quad \frac{\partial^2 S}{\partial \beta_1^2} - \left(\frac{\partial^2 S}{\partial \beta_0 \partial \beta_1} \right)^2$$

$$= 2n \cdot 2 \sum x_i^2 - (2 \sum x_i)^2$$

$$= 4n \sum x_i^2 - 4n^2 \bar{x}^2 \quad (\because \bar{x} = \frac{\sum x_i}{n})$$

$$= 4n \left(\sum x_i^2 - n \bar{x}^2 \right)$$

$$= 4n s_{xx} = 4n \sum (x_i - \bar{x})^2 > 0$$

$(\hat{\beta}_0, \hat{\beta}_1)$ is a global optimum w/?

Alternatively,

Hessian matrix is

$$H = \begin{bmatrix} \frac{\partial^2 S}{\partial \beta_0^2} & \frac{\partial^2 S}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 S}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 S}{\partial \beta_1^2} \end{bmatrix} = \begin{bmatrix} 2n & 2n \bar{x} \\ 2n \bar{x} & 2 \sum x_i^2 \end{bmatrix}$$

$(\hat{\beta}_0, \hat{\beta}_1)$ is a global optimal sol' if H is positive definite matrix.

Note:- A matrix H is p.d. if it's all eigen values are greater than zero.

$$\text{i.e. } |H| > 0$$

$$|H| = \left[4n \sum x_i^2 - 4n^2 \bar{x}^2 \right]$$

$$= 4n \left[\sum x_i^2 - n \bar{x}^2 \right]$$

$$= 4n \sum (x_i - \bar{x})^2$$

$$= 4n S_{xx} > 0$$

Hence, H is positive definite matrix.

$(\hat{\beta}_0, \hat{\beta}_1)$ is a global optimal sol?

Properties of $\hat{\beta}_0, \hat{\beta}_1$:

① Unbiasedness :-

We know that,

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{s_{xx}} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{s_{xx}}$$

$$\sum (x_i - \bar{x}) y_i = \sum x_i y_i - \sum \bar{x} y_i$$

$$= \sum x_i y_i - \bar{x} \sum y_i = \sum x_i y_i - n \bar{x} \bar{y}$$

$$\therefore \hat{\beta}_1 = \frac{\sum (x_i - \bar{x}) y_i}{s_{xx}} = \sum c_i y_i \text{, where } c_i = \frac{x_i - \bar{x}}{s_{xx}}$$

$$\begin{aligned}
 E(\hat{\beta}_1) &= E\left(\sum_{i=1}^n c_i y_i\right) \\
 &= \sum_{i=1}^n c_i E(y_i) \\
 &= \sum_{i=1}^n c_i (\beta_0 + \beta_1 x_i) \quad (\because E(y_i) = \beta_0 + \beta_1 x_i) \\
 &= \beta_0 \sum_{i=1}^n c_i + \beta_1 \sum_{i=1}^n c_i x_i \quad \text{--- } \textcircled{2}
 \end{aligned}$$

where $\bar{c} = \sum_{i=1}^n \frac{(x_i - \bar{x})}{S_{xx}} = \frac{1}{S_{xx}} \sum (x_i - \bar{x}) = 0$

g

$$\begin{aligned}
 \sum c_i x_i &= \sum \frac{(x_i - \bar{x}) x_i}{S_{xx}} = \frac{1}{S_{xx}} \sum (x_i - \bar{x}) x_i \\
 &= \frac{1}{S_{xx}} [\sum x_i^2 - n \bar{x}^2] = 1
 \end{aligned}$$

$$E(\hat{\beta}_1) = \beta_1$$

$\therefore \hat{\beta}_1$ is an unbiased estimator of β_1 .

$$E(\hat{\beta}_0) = E(\bar{y} - \hat{\beta}_1 \bar{x})$$

$$= E(\beta_0 + \beta_1 \bar{x} + \bar{\varepsilon} - \hat{\beta}_1 \bar{x})$$

$$= \beta_0 + \beta_1 \bar{x} + \bar{\varepsilon} - \bar{x} E(\hat{\beta}_1)$$

$$= \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x}$$

$$= \beta_0$$

$\therefore \hat{\beta}_0$ is an unbiased estimator of β_0 .

④ Variance:

$$\text{Var}(\hat{\beta}_1) = \text{Var}(\sum a_i y_i)$$

$$\hat{\beta}_1 = \sum c_i y_i, \quad c_i = \frac{x_i - \bar{x}}{S_{xx}}$$

Since,

$$\begin{aligned}\text{Var}(\sum a_i y_i) &= \sum a_i^2 \text{Var}(y_i) + 2 \sum_{i < j} \sum a_i a_j \text{Cov}(y_i, y_j) \\ &= \sum a_i^2 \text{Var}(y_i) + \sum_{i \neq j} \sum a_i a_j \text{Cov}(y_i, y_j)\end{aligned}$$

If y_i 's, $i = 1, 2, \dots, n$, are independent r.v.'s then.

$$\text{Var}(\sum a_i y_i) = \sum a_i^2 \text{Var}(y_i)$$

$$\therefore \text{Var}(\hat{\beta}_1) = \sum a_i^2 \text{Var}(y_i) = \sum a_i^2 \sigma^2 = \sigma^2 \left\{ \frac{(x_i - \bar{x})^2}{S_{xx}} \right\}$$

($\because y_i$'s are independent)

$$\therefore V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \sum (x_i - \bar{x})^2 = \frac{\sigma^2}{S_{xx}} S_{xx}$$

$$\therefore V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

.

$$V(\hat{\beta}_0) = V(\bar{y} - \hat{\beta}_1 \bar{x})$$

$$= V(\bar{y}) + \bar{x}^2 V(\hat{\beta}_1)$$

$$- 2 \bar{x} \text{cov}(\bar{y}, \hat{\beta}_1)$$

where

$$(\because V(ax+by) = a^2 V(x) + b^2 V(y) + 2ab \text{cov}(x,y))$$

$$V(\bar{y}_i) = \sigma^2, \quad V(\bar{y}) = \sigma^2/n$$

$$V(\hat{\beta}_1 \bar{x}) = \bar{x}^2 V(\hat{\beta}_1) = \bar{x} \frac{\sigma^2}{S_{xx}}$$

$$V(\hat{\beta}_0) = V(\bar{y} - \hat{\beta}_1 \bar{x})$$

$$= V(\bar{y}) + \bar{x}^2 V(\hat{\beta}_1) - 2 \bar{x} C_V(\bar{y}, \hat{\beta}_1) =$$

$$V(\bar{y}) = \frac{\sigma^2}{n},$$

$$\bar{x}^2 V(\hat{\beta}_1) = \frac{\bar{x}^2 \sigma^2}{S_{xx}}$$

$$C_V(\bar{y}, \hat{\beta}_1) = \text{cov}\left(\frac{\sum y_i}{n}, \frac{\sum (x_i - \bar{x}) y_i}{S_{xx}}\right)$$

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

$$\begin{aligned} & (\because \text{cov}(ax, by) \\ & = ab \text{cov}(x, y)) \end{aligned}$$

$$= \frac{\sigma^2}{n S_{xx}} \sum (x_i - \bar{x})$$

$$\begin{aligned} & \text{cov}(x+y, u+v) \\ & = \text{cov}(x, u) + \text{cov}(y, u) \\ & \quad + \text{cov}(x, v) + \text{cov}(y, v) \end{aligned}$$

$$= 0$$

$$\therefore \text{Var}(\hat{\beta}_0) = \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{S_{xx}} - 0$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$

3) $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1)$

$$\begin{aligned}\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) &= \text{Cov}(\bar{y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1) \\ &= \text{Cov}(\bar{y}, \hat{\beta}_1) - \bar{x} \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ &= 0 - \bar{x} \frac{\sigma^2}{S_{xx}} \\ &= -\frac{\sigma^2 \bar{x}}{S_{xx}}\end{aligned}$$

$(\because \text{Cov}(\bar{y}, \hat{\beta}_1) = 0)$

Remark :-

Gauss-Markov thm:-

For the linear regression $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ with assumption I & II, The OLS estimator is a best linear unbiased estimator (BLUE).

H.W. :-

Find the least square estimator of β_0, β_1 for the following regression models.

i) $y_i = \beta_0 + \epsilon_i$

ii) $y_i = \beta_1 x_i + \epsilon_i$

Also, check the properties of least square estimators.

Remarks:

[$\hat{\beta}_{1yx}$ ols of reg'line
on x .

① $\hat{\beta}_{1yx}, \hat{\beta}_{1xy}$ for have same sign.

$\hat{\beta}_{1xy}$ ols of reg'line
 x only

② $\gamma = \sqrt{n} \hat{\beta}_{1xx} \hat{\beta}_{1xy}$

$\gamma = \text{corr}(x, y)$]

③ Both reg. coefficients cannot exceed unity simultaneously.

④ If $\gamma = 1$ ($\text{or } r = -1$) then $\hat{\beta}_{1yx} \hat{\beta}_{1xy} = 1$

⑤ change of origin (or location) does not affect reg. coeff.

⑥ change of scale affect reg. coeff.

⑦ The angle betw. two reg lines y on x & x on y is
 $\theta = \tan^{-1} \left(\frac{1-r^2}{|r|} \cdot \frac{6_x 6_y}{6_x^2 + 6_y^2} \right)$

③

If $\sigma = 1$ or $r = -1$, $\theta = 0$

If $r = 0$, $\theta = \frac{\pi}{2}$

* Residual :-

* \hat{e}_t is diff. bet. y_t & \hat{y}_t .

* $e_t = y_t - \hat{y}_t$

* \hat{e}_t is estimate value of error term.

Error

i) $\epsilon = y_i - \beta_0 - \beta_1 x_i$

ii) ϵ_i 's are independent

& ϵ_i are uncorrelated

iii) ϵ_i 's have same variance.

i.e. constant variance.

iv) $\epsilon_i \sim$ Normal dist

Residual

$$\begin{aligned} r_i &= y_i - \hat{y}_i \\ &= y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \end{aligned}$$

r_i 's are independent & correlated.

Variance is not constant.

[if n is large, they are close to independent & constant variance.]

$\eta \sim$ Normal dist

Properties :-

① $\sum_{i=1}^n \epsilon_i = \sum_{i=1}^n (y_i - \hat{y}_i) = 0$

② $\sum_{i=1}^n \hat{y}_i = \sum_{i=1}^n y_i$

③ $\sum_{i=1}^n x_i \epsilon_i = 0$

④ $\sum \hat{y}_i x_i = 0$

⑤ Regression line always passes through the centroid
of the data (\bar{x}, \bar{y}) .

$$\textcircled{1} \quad \sum \hat{y}_i \cdot r_i = 0$$

$$\left| \begin{array}{l} \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \\ = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i \\ = \bar{y} + \hat{\beta}_1 (x_i - \bar{x}) \end{array} \right.$$

Proof:-

$$\sum \hat{y}_i \cdot r_i = \sum \hat{y}_i (\hat{y}_i - \bar{y})$$

$$= \sum_{i=1}^n (\bar{y} + \hat{\beta}_1 (x_i - \bar{x})) (\hat{y}_i - \bar{y} - \hat{\beta}_1 (x_i - \bar{x}))$$

$$= \sum \bar{y} (\hat{y}_i - \bar{y}) - \sum_{i=1}^n \bar{y} \hat{\beta}_1 (x_i - \bar{x}) + \sum (\hat{y}_i - \bar{y}) \hat{\beta}_1 (x_i - \bar{x})$$

$$- \hat{\beta}_1^2 \sum (x_i - \bar{x})^2$$

$$= 0 - 0 + \hat{\beta}_1 \sum (x_i - \bar{x}) (\hat{y}_i - \bar{y}) - \hat{\beta}_1^2 \sum (x_i - \bar{x})^2$$

$$= \hat{\beta}_1 s_{xy} - \hat{\beta}_1^2 s_{xx}$$

$$= \frac{s_{xy}}{s_{xx}} s_{xy} - \frac{s_{xy}}{s_{xx}} s_{xx} = \frac{s_{xy}^2}{s_{xx}} - \frac{s_{xy}^2}{s_{xx}} = 0$$

* Residual sum of squares:-

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

$$= \sum_{i=1}^n [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

$$= \sum [y_i - (\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i)]^2$$

$$= \sum [y_i - (\bar{y} + \hat{\beta}_1 (x_i - \bar{x}))]^2$$

$$= \sum [(y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x})]^2$$

$$\begin{aligned} &= \sum_{i=1}^n (y_i - \bar{y})^2 + \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\ &\quad - 2 \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \end{aligned}$$

sy

sxp

$$\text{Let } S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2, \quad S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2, \quad S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$\therefore SS_{Reg} = S_{yy} + \hat{\beta}_1 S_{xx} - 2 \hat{\beta}_1 S_{xy}$$

$$= S_{yy} + \hat{\beta}_1 (S_{xx} - 2 S_{xy})$$

$$= S_{yy} + \hat{\beta}_1 \left(\frac{S_{xy}}{S_{xx}} S_{xx} - 2 S_{xy} \right)$$

$$= S_{yy} + \hat{\beta}_1 (S_{xy} - 2 S_{xy})$$

$$= S_{yy} - \hat{\beta}_1 S_{xy} \quad \text{--- (I)}$$

$$= S_{yy} - \hat{\beta}_1 S_{xy} \cdot \frac{S_{xy}}{S_{xx}} \cdot \frac{S_{xx}}{S_{xy}} = S_{yy} - \hat{\beta}_1 S_{xx} \quad \text{--- (II)}$$

* Estimation of σ^2 :

$$\text{Since, } SS_{\text{Res}} = SS_{YY} - \hat{\beta}_1^2 S_{XX}$$

$$E(SS_{\text{Res}}) = E(SS_{YY}) - S_{XX} E(\hat{\beta}_1^2)$$

$$E(SS_{YY}) = E\left(\sum(y_i - \bar{y})^2\right)$$

$$= E(\sum y_i^2 - n\bar{y}^2)$$

$$= E(\sum y_i^2) - n E(\bar{y}^2)$$

$$= \sum E(y_i^2) - n E(\bar{y}^2)$$

$$\left| \begin{array}{l} v(y_i) = E(y_i^2) - E^2(y_i) \\ \therefore \sigma^2 = E(y_i^2) - (\beta_0 + \beta_1 x_i)^2 \\ \therefore E(y_i^2) = \sigma^2 + (\beta_0 + \beta_1 x_i)^2 \\ \therefore v(\bar{y}) = E(\bar{y}^2) - E^2(\bar{y}) \\ \therefore E(\bar{y}^2) = v(\bar{y}) + E^2(\bar{y}) \\ = \frac{\sigma^2}{n} + (\beta_0 + \beta_1 \bar{x})^2 \end{array} \right.$$

$$\begin{aligned}
E(S_{yy}) &= \sum E(y_i^2) - nE(\bar{y}^2) \\
&= \sum_{i=1}^n [\sigma^2 + (\beta_0 + \beta_1 x_i)^2] - n \left[\frac{\sigma^2}{n} + (\beta_0 + \beta_1 \bar{x})^2 \right] \\
&= n\sigma^2 + \sum_{i=1}^n (\beta_0 + \beta_1 x_i)^2 - \sigma^2 - n(\beta_0 + \beta_1 \bar{x})^2 \\
&= (n-1)\sigma^2 + \sum_{i=1}^n (\beta_0 + \beta_1 x_i)^2 - n(\beta_0 + \beta_1 \bar{x})^2 \\
&= (n-1)\sigma^2 + \sum_{i=1}^n (\beta_0^2 + \beta_1^2 x_i^2 + 2\beta_0\beta_1 x_i) \\
&\quad - n(\beta_0^2 + \beta_1^2 \bar{x}^2 + 2\beta_0\beta_1 \bar{x}) \\
&= (n-1)\sigma^2 + \cancel{n\beta_0^2} + \beta_1^2 \sum_{i=1}^n x_i^2 + 2\cancel{\beta_0\beta_1} \sum_{i=1}^n x_i \\
&\quad - \cancel{n\beta_0^2} - n\beta_1^2 \bar{x}^2 - \cancel{2\beta_0\beta_1} \bar{x}
\end{aligned}$$

$$\begin{aligned}
 E(S_{\text{reg}}) &= (n-1)\sigma^2 + \beta_1^2 \sum x_i^2 - n\beta_1^2 s_{\bar{x}}^2 \\
 &= (n-1)\sigma^2 + \beta_1^2 \left(\sum x_i^2 - n\bar{x}^2 \right) \\
 &= (n-1)\sigma^2 + \beta_1^2 s_{xx}
 \end{aligned}$$

Hence.

$$\begin{aligned}
 E(S_{\text{res}}) &= E(S_{\text{reg}}) - s_{xx} E(\hat{\beta}_1^2) \\
 &= (n-1)\sigma^2 + \beta_1^2 s_{xx} - s_{xx} E(\hat{\beta}_1^2)
 \end{aligned}$$

Since, $V(\hat{\beta}_1) = E(\hat{\beta}_1^2) - E^2(\hat{\beta}_1)$

$$\therefore E(\hat{\beta}_1^2) = V(\hat{\beta}_1) + E^2(\hat{\beta}_1) = \frac{\sigma^2}{s_{xx}} + \beta_1^2$$

$$\therefore E(S_{\text{res}}) = (n-1)\sigma^2 + \beta_1^2 s_{xx} - s_{xx} \left(\frac{\sigma^2}{s_{xx}} + \beta_1^2 \right)$$

$$E(SS_{\text{Res}}) = (n-1)\sigma^2 + \cancel{\beta_1^2 S_{xx}} - \underline{\sigma^2 - \cancel{\beta_1^2 S_{xx}}}$$

$$= (n-2)\sigma^2$$

$$\therefore E\left(\frac{SS_{\text{Res}}}{n-2}\right) = \sigma^2$$

Hence, $\frac{SS_{\text{Res}}}{n-2}$ is an unbiased estimator σ^2

$$\hat{\sigma}^2 = \frac{SS_{\text{Res}}}{n-2} = m_{\text{Res}}$$

i.e. residual mean square.

g

$\hat{\sigma}$ - standard error of reg.

Other method :-

Thm:- Suppose $\gamma \sim N_p(\mu, \Sigma)$, A is a symmetric matrix of constants with rank r . If $\lambda = \frac{1}{2} \mu' A \mu$ then

$\gamma' A \gamma \sim \chi^2_{(r, 1)}$ iff $\underline{\underline{A \Sigma}}$ is idempotent matrix.

Cochran's Thm:-

Let x_1, \dots, x_n be a random sample

from $N(0, \sigma^2)$ & $\sum_{i=1}^n x_i^2 = Q_1 + Q_2 + \dots + Q_K$.

where, Q_j is a quadratic form in x_1, \dots, x_n with r_j d.f.
then $\sim \chi^2_{r_j}$ & Q_1, \dots, Q_K are mutually independent & Q_j / σ^2
is $\chi^2_{r_j}$ if $\sum_{i=1}^K r_i = n$.

Using thⁿ₁ & th^m₂, we can prove that,

$$\frac{SS_{\text{Res}}}{\sigma^2} = \frac{n S_{\text{Res}}}{\sigma^2} (n-2) \sim \underline{\chi^2_{n-2}}$$

∴ $E\left(\frac{SS_{\text{Res}}}{\sigma^2}\right) = (n-2)$

∴ $E\left(\frac{SS_{\text{Res}}}{n-2}\right) = \sigma^2$

Other method :-

$$\varepsilon_i \sim N(0, \sigma^2)$$

$$\therefore \frac{\varepsilon_i}{\sigma} \sim N(0, 1)$$

$$\therefore \left(\frac{\varepsilon_i}{\sigma}\right)^2 \sim \chi_1^2$$

$$\therefore \sum_{i=1}^n \frac{(y_i - \beta_0 - \beta_1 x_i)^2}{\sigma^2} \sim \chi_n^2$$

$$\therefore \sum_{i=1}^n \frac{(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{\sigma^2} \sim \chi_{n-2}^2$$

$$\therefore \frac{SS_{Reg}}{\sigma^2} \sim \chi_{n-2}^2$$

$$i.e. E\left(\frac{SS_{Reg}}{n-2}\right) = \sigma^2$$

* Estimate of variance of $\hat{\beta}_0$ & $\hat{\beta}_1$:-

$$\hat{v}(\hat{\beta}_0) = \hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right),$$

$$\text{&} \quad \hat{v}(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{S_{xx}}$$

$$\hat{\sigma}^2 = \frac{SS_{\text{Res}}}{n-2} = S_{\text{Res}}$$

Testing of Hypothesis:-

$$H_0: \beta_1 = \beta_{10} \quad \text{vs} \quad H_1: \beta_1 \neq \beta_{10}$$

i) σ^2 is known

Since, $\epsilon_i \sim N(0, \sigma^2)$,

$$\hat{\beta}_1 \sim \mathcal{N}(\beta_1, \frac{\sigma^2}{s_{xx}})$$

$$Z_{cd} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{s_{xx}}}}$$

under $H_0: \beta_1 = \beta_{10}$

$$Z_{cal} = \frac{\hat{\beta}_1 - \beta_{10}}{\sqrt{\frac{\sigma^2}{s_{xx}}}} \sim \mathcal{N}(0, 1)$$

we reject H_0 if

$$|t_{\text{cal}}| > z_{\alpha/2},$$

otherwise accept H_0 .

ii) σ^2 is unknown :-

$$t_{\text{cal}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{SS_{\text{Res}}}{(n-2)S_{xx}}}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{m S_{\text{Res}}}{S_{xx}}}}$$

($m S_{\text{Res}} = \frac{SS_{\text{Res}}}{n-2}$)

under H_0 ,

$$t_{\text{cal}} = \frac{\hat{\beta}_1 - \beta_{10}}{\sqrt{\frac{m S_{\text{Res}}}{S_{xx}}}} \sim t_{n-2}$$

If x & y are independent r.v.'s & $x \sim N(0, 1)$, $y \sim \chi^2_n$

then,

$$\frac{x}{\sqrt{Y/n}} \sim t_n$$

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2/s_{xx})$$

$$i. x = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2/s_{xx}}} \sim N(0, 1)$$

$$\frac{SS_{res}}{\sigma} \sim \chi^2_{n-2}$$

$$r = \frac{mS_{res}(n-2)}{\sigma} \sim \chi^2_{n-2}$$

$$\frac{x/\sqrt{Y/n}}{\sqrt{\hat{\beta}_1 - \beta_1}} = \frac{\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2/s_{xx}}}}{\sqrt{\frac{mS_{res}(n-2)}{\sigma(n-2)}}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{mS_{res}}{s_{xx}}}}$$

We reject H_0 if $|t_{\text{cal}}| > t_{\alpha/2, n-2}$.
 i.e. we accept H_0 .

$\alpha = 1 - \alpha_s$.

b) Testing of hypothesis for β_0 :-

$$H_0: \beta_0 = \beta_{00} \text{ vs } H_1: \beta_0 \neq \beta_{00}$$

i) σ^2 is known :-

$$z_{\text{cal}} = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim N(0, 1)$$

under H_0 ,

$$z_{\text{cal}} = \frac{\hat{\beta}_0 - \beta_{00}}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim N(0, 1)$$

we reject H_0 if $|t_{\text{cal}}| \geq t_{\alpha/2}$, o.w. accept H_0 .

i) σ^2 unknown,

$$t_{\text{cal}} = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\frac{ms_{\text{res}}}{n-2} \left(\frac{1}{a} + \frac{\bar{x}^2}{s_{xx}} \right)}} \\ = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{ms_{\text{res}} \left(\frac{1}{a} + \frac{\bar{x}^L}{s_{xx}} \right)}} \sim t_{n-2}.$$

under H_0 ,

$$t_{\text{cal}} = \frac{\hat{\beta}_0 - \beta_{00}}{\sqrt{ms_{\text{res}} \left(\frac{1}{a} + \frac{\bar{x}^2}{s_{xx}} \right)}} \sim t_{n-2}.$$

we reject H_0 if $|t_{\text{cal}}| > t_{\alpha/2, n-2}$, o.w. accept H_0 .

Q) Test for σ^2 .

$$H_0: \sigma^2 = \sigma_0^2 \text{ vs } H_1: \sigma^2 \neq \sigma_0^2$$

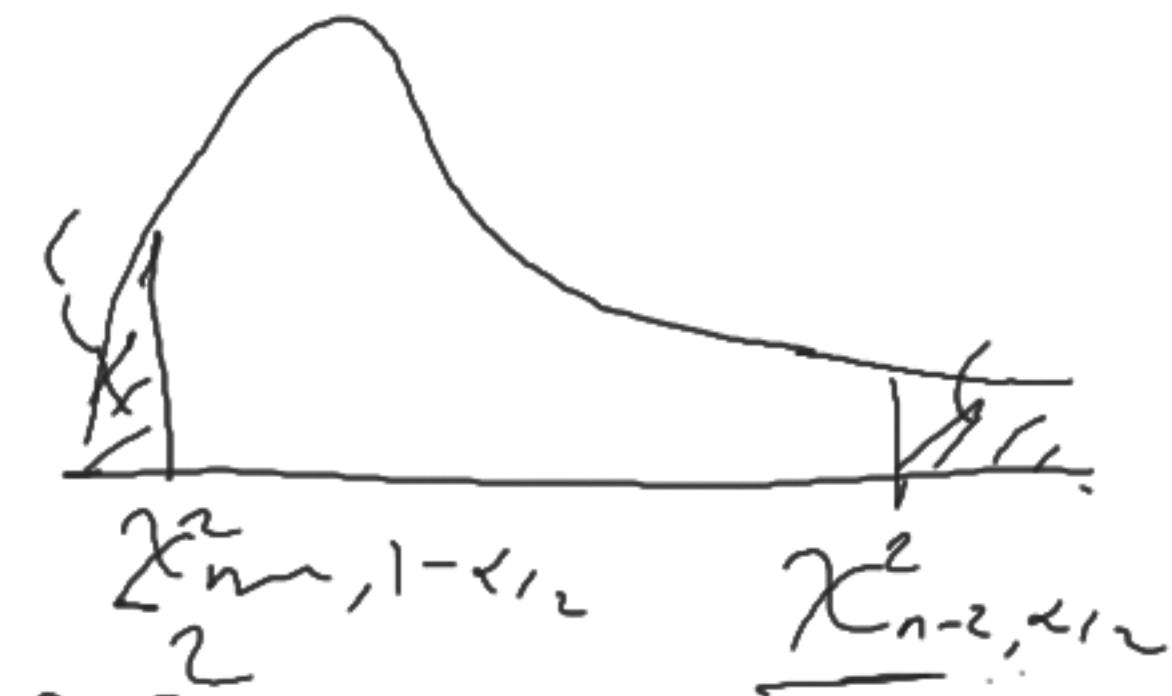
Since, $\frac{SS_{Res}}{\sigma^2} \sim \chi^2_{n-2}$

$$C = \frac{SS_{Res}}{\sigma^2}$$

under H_0 ,

$$C = \frac{SS_{Res}}{\sigma_0^2} \sim \chi^2_{n-2}$$

we reject H_0 if $C < \chi^2_{n-2, 1-\alpha/2}$ or $C > \chi^2_{n-2, \alpha/2}$.
o.w. accept H_0 .



d) linearity:
 $H_0: \beta_1 = 0$ v/s $H_1: \beta_1 \neq 0$
 (no linear rel)
 (linear rel)

i) σ known

$$z_{\text{cal}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2_{\text{Sxx}}}}$$

under H_0

$$z_{\text{cal}} = \frac{\hat{\beta}_1}{\sqrt{\sigma^2_{\text{Sxx}}}} \sim N(0, 1)$$

we reject H_0 if $|z_{\text{cal}}| > z_{\alpha/2}$. i.e. accept H_0 .

ii) σ unknown

$$t_{\text{cal}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{MS_{\text{Res}}}{S_{\text{xx}}}}}$$

under H_0 : $t_{\text{cal}} = \frac{\hat{\beta}_1}{\sqrt{\frac{MS_{\text{Res}}}{S_{\text{xx}}}}} \sim t_{n-2}$

we reject H_0 if $|t_{\text{cal}}| > t_{\alpha/2, n-2}$. i.e. accept H_0 .

* Confidence Interval :-

↪ C.I. for β_1 :-

↪ σ^2 known.

100(1 - α)% C.I. for β_1 is

$$P(-z_{\alpha/2} \leq z_{\text{cal}} \leq z_{\alpha/2}) = 1 - \alpha$$

$$\therefore P\left(-z_{\alpha/2} \leq \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma}{S_{xx}}}} \leq z_{\alpha/2}\right) = 1 - \alpha,$$

$$\therefore P\left(-z_{\alpha/2} \sqrt{\frac{\sigma}{S_{xx}}} \leq \hat{\beta}_1 - \beta_1 \leq z_{\alpha/2} \sqrt{\frac{\sigma}{S_{xx}}}\right) = 1 - \alpha$$

$$P\left(-Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}} \leq \beta_1 - \hat{\beta}_1 \leq Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}}\right) = 1 - \alpha.$$

$$\therefore P\left(\hat{\beta}_1 - Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}} \leq \beta_1 \leq Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}} + \hat{\beta}_1\right) = 1 - \alpha.$$

Now (\rightarrow). C.E. for β_1 is

$$\left(\hat{\beta}_1 - Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}}, \hat{\beta}_1 + Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}}\right)$$

i.e. $\hat{\beta}_1 \pm Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}}$.

ii) σ unknown:

for $(1-\alpha)\%$. C. S. for β_1 is

$$P(-t_{\alpha/2, n-2} \leq t_{\text{cal}} \leq t_{\alpha/2, n-2}) = 1 - \alpha$$

$$\therefore c. P\left(-t_{\alpha/2, n-2} \leq \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{MS_{\text{Res}}}{S_{xx}}}} \leq t_{\alpha/2, n-2}\right) = 1 - \alpha.$$

$$\therefore c. P\left(-t_{\alpha/2, n-2} \sqrt{\frac{MS_{\text{Res}}}{S_{xx}}} \leq \hat{\beta}_1 - \beta_1 \leq t_{\alpha/2, n-2} \sqrt{\frac{MS_{\text{Res}}}{S_{xx}}} \right) = 1 - \alpha.$$

$$\therefore c. P\left(-t_{\alpha/2, n-2} \sqrt{\frac{MS_{\text{Res}}}{S_{xx}}} \leq \beta_1 - \hat{\beta}_1 \leq t_{\alpha/2, n-2} \sqrt{\frac{MS_{\text{Res}}}{S_{xx}}}\right) = 1 - \alpha.$$

C.I. for β_1 is

$$(\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{m s_{\text{res}}}{s_{xx}}}, \quad \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{m s_{\text{res}}}{s_{xx}}})$$

b) C.I. for β_0 :-

i) σ is known :-

(or $(1-\alpha)\%$ C.I. for β_0 is

$$P(-z_{\alpha/2} \leq z_{\text{cal}} \leq z_{\alpha/2}) = 1 - \alpha.$$

$$\text{i.e. } P(-2\leq \beta_0 \leq \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}})}} \leq z_{\alpha/2}) = 1 - \alpha.$$

Hence, $(1-\alpha)\%$ C.E. for β_0 is

$$(\hat{\beta}_0 - z_{\alpha/2} \sqrt{\sigma^2(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}})}, \hat{\beta}_0 + z_{\alpha/2} \sqrt{\sigma^2(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}})})$$

ii) If β_0 is unknown :-

$$(\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MS_{\text{reg}}(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}})}, \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MS_{\text{reg}}(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}})})$$

$\xrightarrow{\text{S.E.}(\hat{\beta}_0)}$

Confidence Interval:-

Parameter.

β_0

σ^2

known

unknown

β_1

known

unknown

C.I

$$\hat{\beta}_0 \pm \underline{z}_{\alpha/2} \frac{s_e(\hat{\beta}_0)}{=}$$

$$\hat{\beta}_0 \pm \underline{t}_{\alpha/2, n-2} s_e(\hat{\beta}_0)$$

$$\hat{\beta}_1 \pm \underline{z}_{\alpha/2} s_e(\hat{\beta}_1)$$

$$\hat{\beta}_1 \pm \underline{t}_{\alpha/2, n-2} s_e(\hat{\beta}_1)$$

c) C.I. for σ^2 :

100(1-\alpha)% C.I. for σ^2 is

$$P\left(\chi^2_{1-\alpha/2, n-2} \leq \frac{SS_{\text{Res}}}{\sigma^2} \leq \chi^2_{\alpha/2, n-2}\right) = 1 - \alpha.$$

$$\therefore C. I. P\left(\frac{1}{\chi^2_{1-\alpha/2, n-2}} > \frac{\sigma^2}{SS_{\text{Res}}} > \frac{1}{\chi^2_{\alpha/2, n-2}}\right) = 1 - \alpha.$$

$$1 - \alpha \cdot P\left(\frac{SS_{\text{Res}}}{\chi^2_{\alpha/2, n-2}} \leq \sigma^2 \leq \frac{SS_{\text{Res}}}{\chi^2_{1-\alpha/2, n-2}}\right) = 1 - \alpha.$$

Hence, C.I. for σ^2 is

$$\left(\frac{SS_{\text{Res}}}{\chi^2_{\alpha/2, n-2}}, \frac{SS_{\text{Res}}}{\chi^2_{1-\alpha/2, n-2}}\right) \text{ or } \left(\frac{n-2M_{\text{Res}}}{\chi^2_{\alpha/2, n-2}}, \frac{n-2M_{\text{Res}}}{\chi^2_{1-\alpha/2, n-2}}\right)$$

d) C.I. für $E(\hat{y}|x_0)$

$$E(\hat{y}|x_0) = \hat{\mu}_{y|x_0} = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

$$\text{v}(\hat{\mu}_{y|x_0}) = \text{v}(\hat{\beta}_0 + \hat{\beta}_1 x_0)$$

$$= \text{v}(\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_0)$$

$$= \text{v}(\bar{y} + \hat{\beta}_1 (x_0 - \bar{x}))$$

$$= \text{v}(\bar{y}) + (x_0 - \bar{x})^2 \text{v}(\hat{\beta}_1) \quad (\stackrel{\text{cov}(\bar{y}, \hat{\beta}_1)}{=} 0)$$

$$= \frac{\sigma^2}{n} + (x_0 - \bar{x})^2 \frac{\sigma^2}{S_{xx}}$$

$$\therefore V(\hat{y}|x_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

or unknown

$$\frac{\hat{y}_x - E(\hat{y}|x_0)}{\sqrt{V_{\text{sees}}\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}\right)}} \sim t_{n-2}$$

$\therefore c \cdot \Sigma \cdot f_s$

$$\hat{y}_{x_0} \pm \sqrt{V_{\text{sees}}\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}\right)}$$

* Analysis of variance :-

We know that,

$$\begin{aligned}x_i &= y_i - \hat{y}_i \\&= y_i - \bar{y} - \hat{y}_i + \bar{y} \\&= (y_i - \bar{y}) - (\hat{y}_i - \bar{y})\end{aligned}$$

$$\begin{aligned}\therefore \sum_{i=1}^n (x_i - \hat{y}_i)^2 &= \sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\&\quad - 2 \sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{y})\end{aligned}$$

where :-

$$\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{y}) = \sum_{i=1}^n (y_i - \bar{y}) \hat{\beta}_1 (x_i - \bar{x})$$

$$\therefore \hat{y}_i - \bar{y} = \hat{\beta}_1 (x_i - \bar{x})$$

$$= \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$= \hat{\beta}_1 S_{xy}$$

$$= \hat{\beta}_1 \cancel{S_{xy}} \frac{S_{x1}}{S_{xx}} \frac{S_{xx}}{\cancel{S_{xy}}} = \hat{\beta}_1 S_{xx}$$

$$= \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

$$= \sum_{i=1}^n [\hat{\beta}_1 (x_i - \bar{x})]^2$$

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

$$(\hat{y}_i - \bar{y} = \hat{\beta}_1 (x_i - \bar{x}))$$

Aence,

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

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

- - -

TSS
or
SS_T
or
SST
Sum of square about mean
or
Corrected sum of squares.

reg sum of squares
SSR
or SS_{Reg}

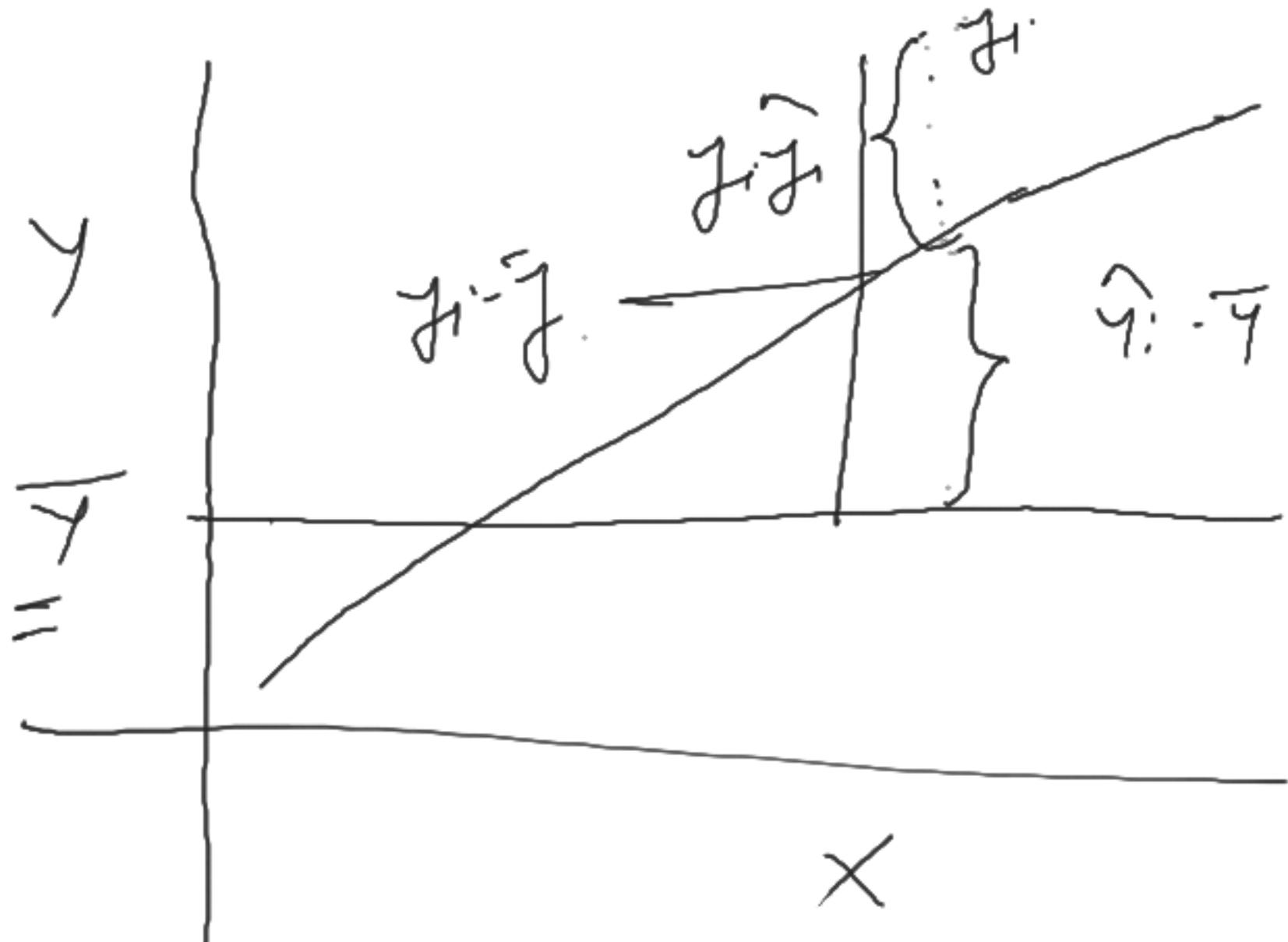
residual sum of squares.
SS_{Res}
SSE
=

Source of variation

$$SS_T = \sum_{i=1}^n (y_i - \bar{y})^2 \quad n-1$$

$$SS_R = \sum (\hat{y}_i - \bar{y})^2 \quad 1$$

$$SS_{Res} = \sum (y_i - \hat{y}_i)^2 \quad n-2$$



Degrees of freedom -

7 chocolates : A, B, C, D, E, F, G

Day chocolate choice

(1)	D	7
(2)	E	6
(3)	B	5
(4)	F	5
(5)	G	3
(6)	C	2
(7)	=	1

D.F.

6

$$x_1, x_2, \dots, x_{10}$$
$$\sum x_i = 100$$
$$(or \bar{x} = 10)$$

$$D.F. = 9$$

Degrees of freedom

= No. of obs. - known relⁿ. betw them
(or no. of constraints or no. of estimated parameters)

Examples:-

i) Data: x_1, \dots, x_{10} & constraint $\sum x_i = 10$

$$D.f. = \text{No. of obs.} - \text{No. of constraint} = 10 - 1 = 9$$

ii) Residual sum of square = $\sum (y_i - \hat{y}_i)^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$

D.f. = No. of obs. - No. of constraint (normal eq's) (No. of parameters estimate)

$$= n - 2$$

ANOVA table

Source of variation	Df	S. S.	mss	F
Reg	1	SS_{Reg}	$mS_{Reg} = \frac{SS_{Reg}}{1}$	$\frac{MS_{Reg}}{MS_E}$
Residual	$n-2$	SS_{Res}	$MS_E: \frac{SS_{Res}}{n-2}$	
Total	$n-1$	SS_T		

Expectation of mean sum of square :-

$$\begin{aligned} E(MS_{Reg}) &= E\left(\sum_{i=1}^n (\hat{y}_i - \bar{y})^2\right) \\ &= E\left(\hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2\right) \\ &= E\left(\hat{\beta}_1^2 s_{xx}\right) \\ &= s_{xx} E(\hat{\beta}_1^2) \quad \left(\because E(\hat{\beta}_1) = E(\hat{\beta}_1) - E^e(\hat{\beta}_1) \right) \\ &= s_{xx} \left(\frac{\sigma^2}{s_{xx}} + \hat{\beta}_1^2 \right) \quad \left(\because \frac{\sigma^2}{s_{xx}} = E(\hat{\beta}_1) - \hat{\beta}_1^2 \right) \\ &= \underline{\sigma^2} + \underline{\hat{\beta}_1^2} s_{xx} \geq 0 \end{aligned}$$

∴

$$E(MSE) = E\left(\frac{SS_{Res}}{n}\right) = \sigma^2 \rightarrow 0$$

$$\frac{(n-2) \text{ MSE}}{\sigma^2} \sim \chi_{n-2}^2 \quad \& \quad \frac{\text{MSE}_{\text{reg}}}{\sigma^2} \sim \chi_1^2$$

Test :-

$$\underline{H_0: \beta_1 = 0} \quad \text{v/s} \quad H_1: \underline{\beta_1 \neq 0}$$

under H_0 :

$$F_{\text{cal}} = \frac{\frac{\text{MSE}_{\text{reg}}}{\sigma^2}}{\frac{\text{MSE}_{\text{unreg}}}{\sigma^2}} =$$

$$= \frac{\text{MSE}_{\text{reg}}}{\text{MSE}_{\text{unreg}}} \sim F_{1, n-2}$$

$$X \sim \chi_{n-2}^2$$

$$Y \sim \chi_{m-1}^2$$

$$F = \frac{X/n_1}{Y/n_2}$$

we reject H_0 if $F_{cal} > F_{\alpha, 1, n-2}$

o.w. accept H_0 .

Coefficient of Determination:-

$$SS_T = SS_{reg} + SS_{res}$$

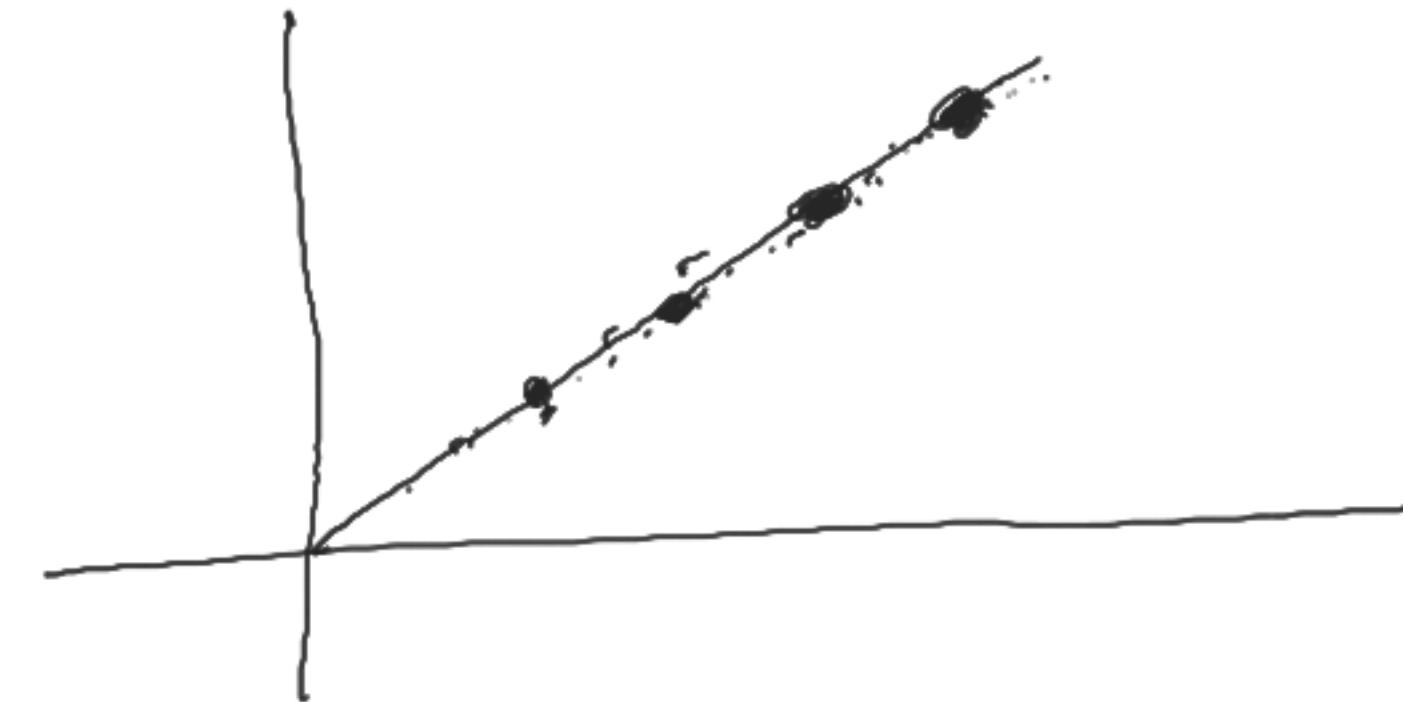
$$R^2 = \frac{SS_{reg}}{SS_T} = \frac{SS_T - SS_{res}}{SS_T} = 1 - \frac{SS_{res}}{SS_T}$$

R^2 measures the propⁿ of variation in response variable that is explained by the regⁿ model.

case 1 : $R^2 = 1$

i.e. $R^2 = \frac{SS_{\text{reg}}}{SS_T} = 1 \Rightarrow SS_{\text{reg}} = SS_T$

i.e. $SS_{\text{res}} = 0$



case 2 : $R^2 = 0$

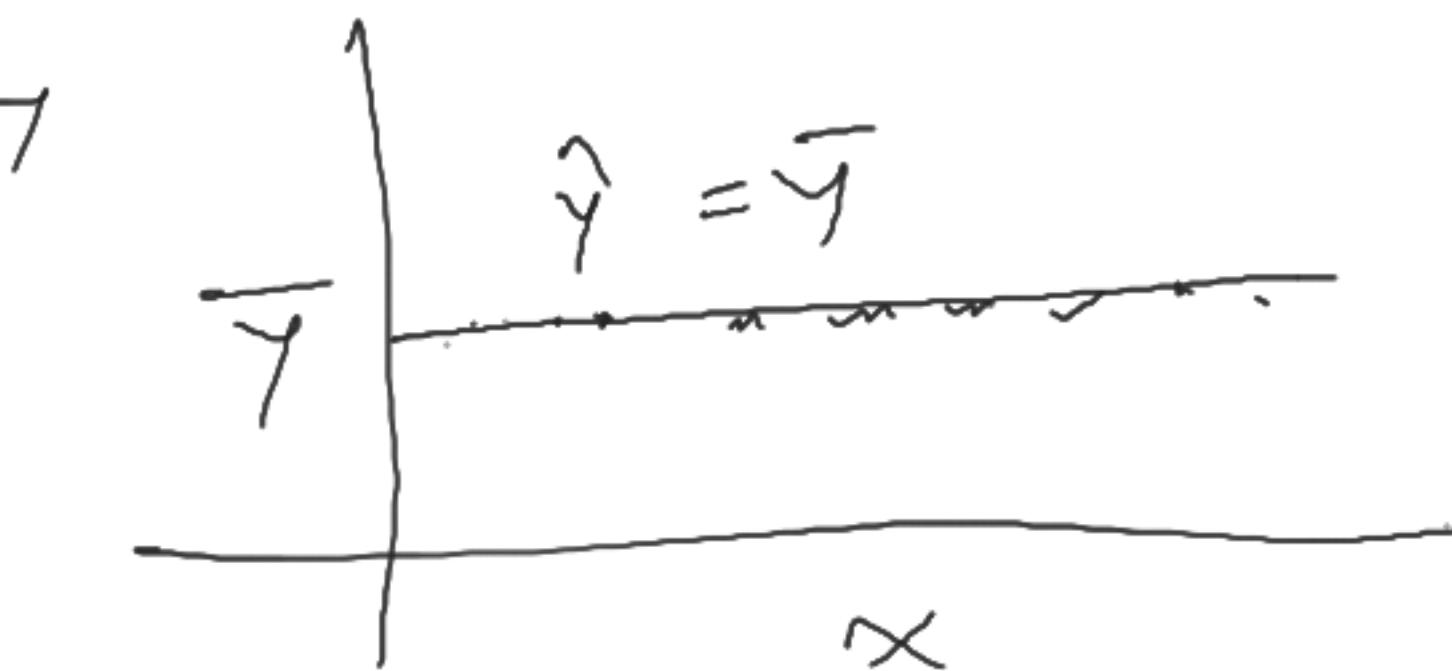
$R^2 = \frac{SS_{\text{reg}}}{SS_T} = 1 - \frac{SS_{\text{res}}}{SS_T} = 0 \Rightarrow SS_{\text{res}} = SS_T$

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

Hence, If $r^2 = 0$

then $x \perp \gamma$ or

uncorrelated.



Remark:-

$$1) R^2 = (\text{corr}(x, \gamma))^2$$

$$2) \text{Since } -1 \leq \text{corr}(x, \gamma) \leq 1 \Rightarrow 0 \leq R^2 \leq 1$$

3)

$$R^2 = \frac{SS_{\text{reg}}}{SS_T}$$

Since, $SS_T = SS_{\text{reg}} + SS_{\text{res}}$

$$\therefore SS_{\text{reg}} \leq SS_T \Rightarrow \frac{SS_{\text{reg}}}{SS_T} \leq 1 \quad \text{--- (I)}$$

if $SS_{\text{reg}} \geq 0$ & $SS_T \geq 0 \Rightarrow \frac{SS_{\text{reg}}}{SS_T} \geq 0 \quad \text{--- (II)}$

$$\therefore 0 \leq \frac{SS_{\text{reg}}}{SS_T} \leq 1$$

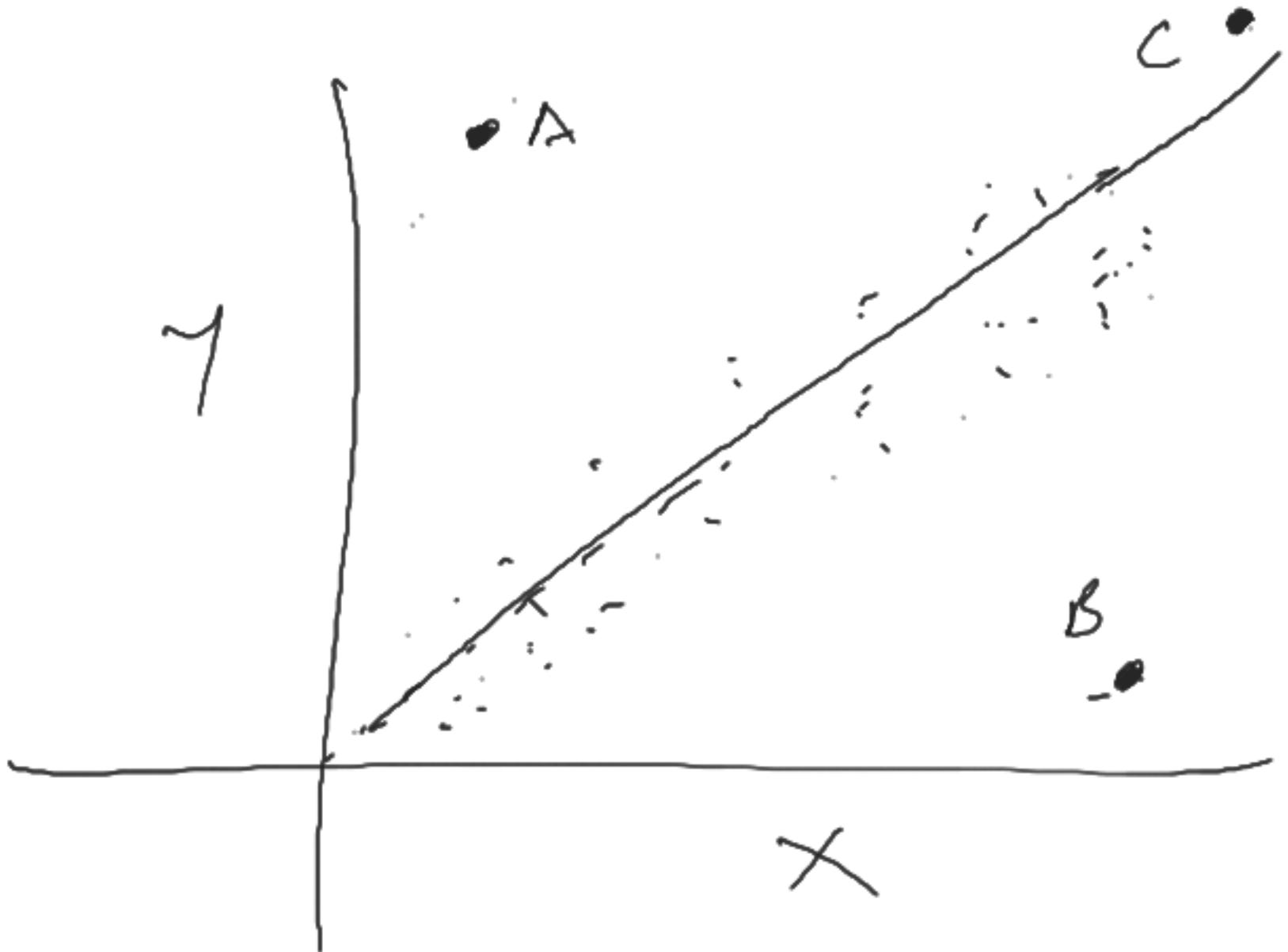
i.e. $0 \leq R^2 \leq 1$

- * Types of outliers:-

- ~ A - vertical outlier

- ~ B - Bad leverage point

- c. Good leverage point



* Model Adequacy checking :-

Assumptions:-

i) $x \text{ & } y$ are linearly related

ii) $E(\epsilon_i) = 0$

iii) $\sqrt{E(\epsilon)} = \sigma$

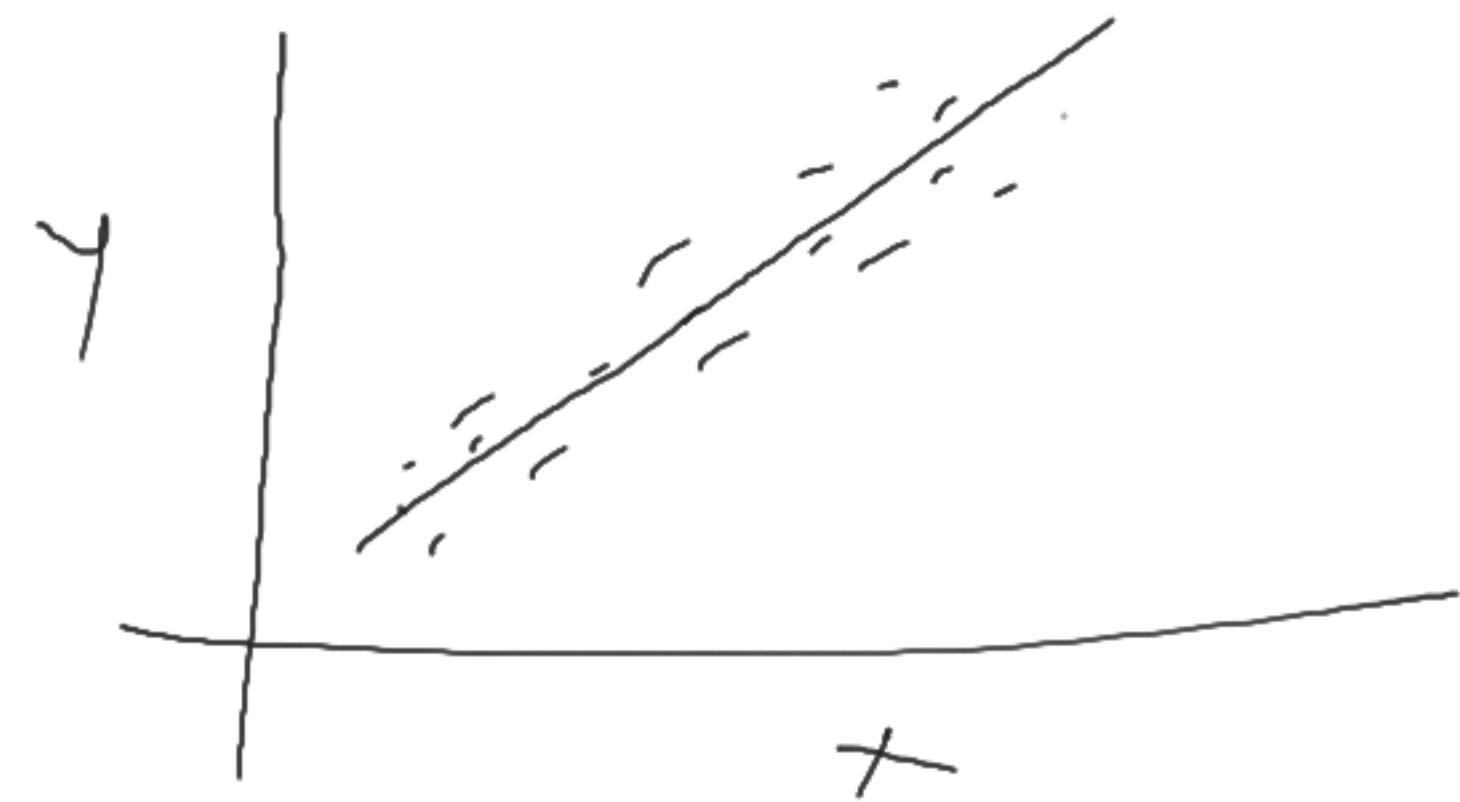
iv) ϵ_i 's are uncorrelated

v) $\epsilon_i \sim N(0, \sigma^2)$

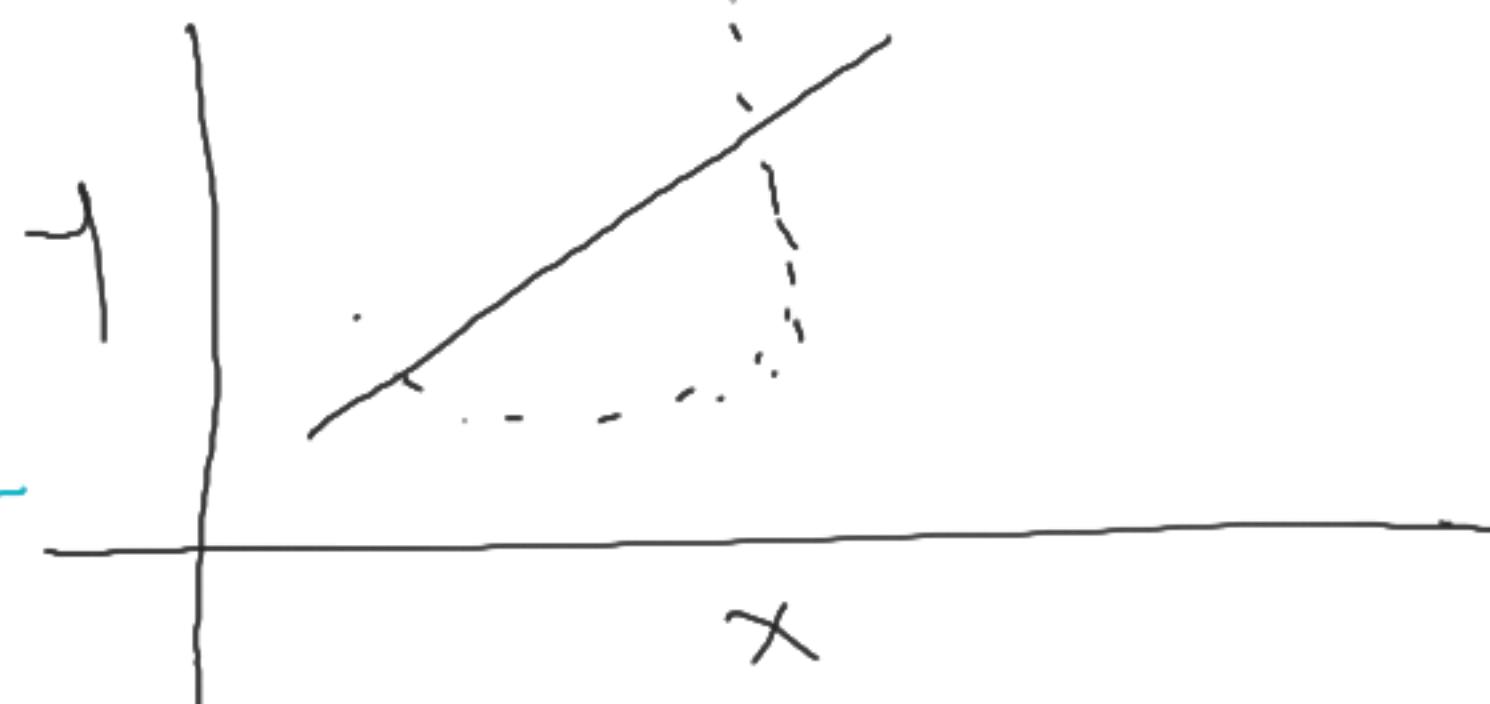


1) Linear relationship bet. x & y .

Linear



Non-
linear



* Residual Analysis:-

ii. Residuals are estimated values of errors.

$$i) E(\epsilon) = 0$$

$$iii) \hat{\sigma}^2 = MSE =$$

$$\frac{\sum (Y - \hat{Y}_i)^2}{n-2} = \frac{SS_{res}}{n-2}$$

+ method for scaling residuals:-

i. standardized residual s_i :-

$$d_i = \frac{y_i - E(y_i)}{\sqrt{\text{MSE}}} =$$

$$r_i = y_i - \hat{y}_i$$
$$\sqrt{\frac{n}{\text{MSE}}}, i=1, 2, \dots, n.$$

* $E(d_i) = 0$

* $V(d_i) \approx 1$

if $|d_i| > 3$, then $i^{\text{th}} \text{ obs}$ is an outlier.
(i.e. $d_i < -3$ or $d_i > 3$)

2) Standardized residual :-

$$z_i^* = \frac{r_i - E(r_i)}{\sqrt{v(r_i)}}, i=1, 2, \dots, n.$$

Since,

$$\sqrt{v(r_i)} = \sigma \left[1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{Sxx} \right) \right]$$

$$r_i^* = \frac{r_i}{\sqrt{MSE \left(1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{Sxx} \right) \right)}}, i=1, 2, \dots, n.$$

Remarks :-

i) if $(x_i - \bar{x})^2$ is large, r_i^* large

ii) if $(x_i - \bar{x})^2$ is small, r_i^* small.

(iii) Let

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}}$$

since, $s_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$

$$0 \leq h_{ii} \leq 1$$

$$0 \leq 1 - h_{ii} \leq 1$$

Hence, $V(\gamma_i) = \text{mse}(1 - h_{ii}) = \text{mse}\left[1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}}\right)\right] \leq \text{mse}$

(iv) for large n , $h_{ii} \rightarrow 0$ & $V(\gamma_i) = \text{mse}$.

$$\gamma_i = \gamma_i^*$$

PRESS | deleted residual :- { Prediction Error Sum of Squares}

We know that,

$$r_i = y_i - \hat{y}_i$$

$$= y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$$

The PRESS residual is defined as,

$$r_{(i)} = y_i - \hat{y}_{(i)}$$

Procedure:-

- (i) delete i^{th} observation.
- (ii) fit the regression line based on remaining $(n-1)$ observations.
- (iii) calculate Fitted value of y_i ($\hat{y}_{(i)}$) using the fitted regression model.
- (iv) calculate $r_{(i)} = y_i - \hat{y}_{(i)}$

Remarks :-

To calculate the PRESS residuals for n observations, n regression lines are required. Instead of this, we can calculate PRESS residuals using the relation betw. $r_{(i)}$ & r_i :

$$r_{(i)} = \frac{r_i}{1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}} \right)}$$

(ii) Since, $0 \leq 1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}} \right) \leq 1$

$$\therefore r_i \leq r_{(i)}$$

(iii) PRESS is used to identify outliers.

Standardized PRESS residual :-

The PRESS residual is,

$$r_{(i)} = \frac{v_i}{1-h_{ii}}, \text{ where } h_{ii} = \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right)$$

$$\therefore E(r_{(i)}) = 0$$

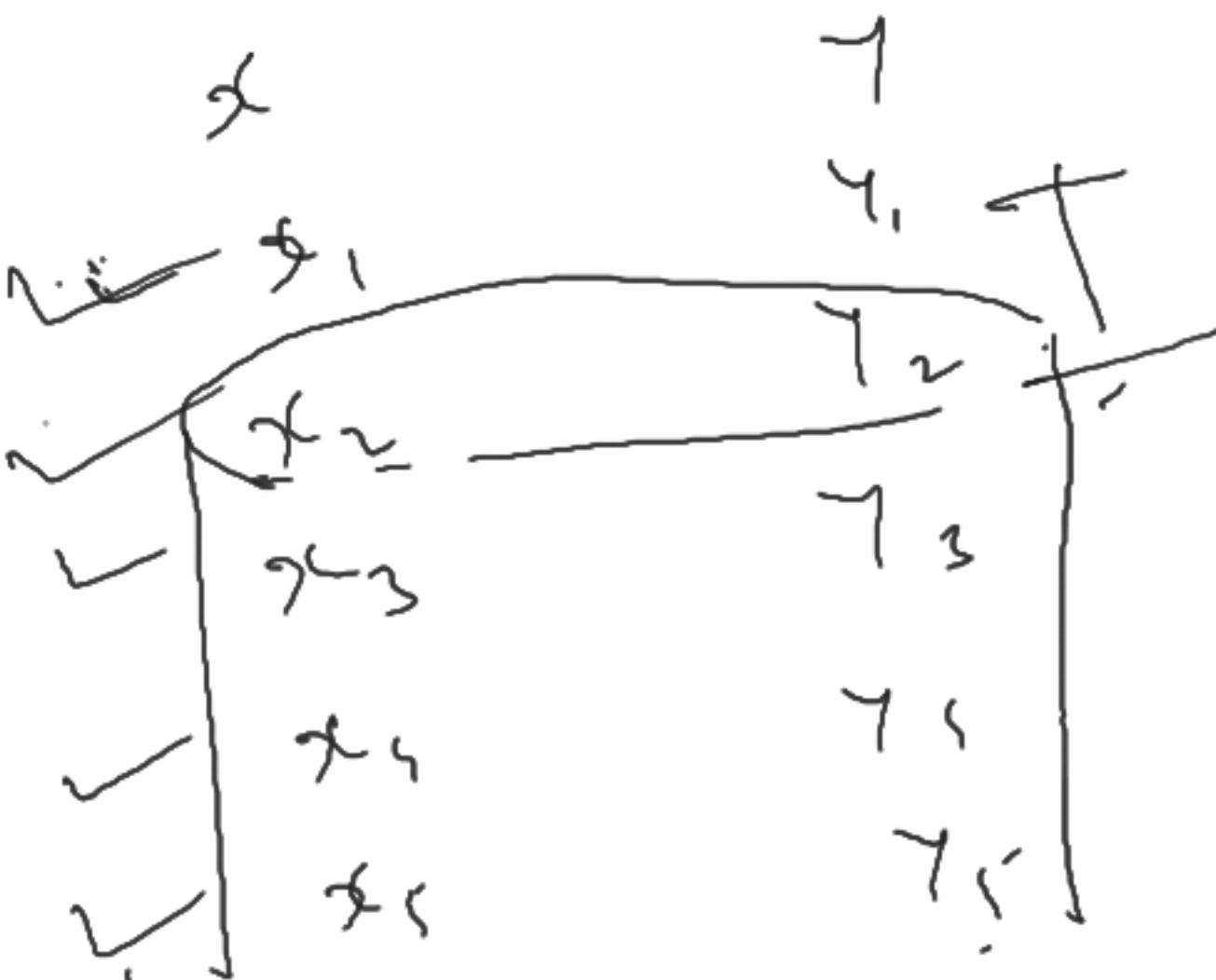
&

$$V(r_{(i)}) = \frac{1}{(1-h_{ii})^2} V(v_i) = \frac{\sigma^2(1-h_{ii})}{(1-h_{ii})^2} = \frac{\sigma^2}{1-h_{ii}}$$

Therefore the standar-dized PRESS residual is

$$\frac{r_{(i)} - E(r_{(i)})}{\sqrt{V(r_{(i)})}} = \frac{r_i / (1-h_{ii})}{\sqrt{\sigma^2 / (1-h_{ii})}} = \frac{r_i}{\sqrt{\sigma^2 (1-h_{ii})}}$$

$$n = 5$$



$$\hat{y}_{(i)} = \frac{\bar{y}_i}{\bar{w}_{(i)}} - \frac{\bar{y}_i - \bar{T}_i}{\bar{r} - \bar{w}_{(i)}}$$

OLS Estimator

$$\hat{\beta}_{00}, \hat{\beta}_{11}$$

$$(\hat{\beta}_{20}, \hat{\beta}_{21})$$

$$(\hat{\beta}_{30}, \hat{\beta}_{31})$$

$$(\hat{\beta}_{40}, \hat{\beta}_{41})$$

$$w_{ii} = \left(k + \frac{(x_i - \bar{x})^2}{S_{xx}} \right)$$

$$\hat{y}_i = \frac{\bar{y}_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)}{(-w_{ii})}$$

$$\underline{\hat{y}_{(1)}} = \bar{y}_1 - \bar{T}_{(1)}$$

$$\hat{y}_{(i)} = \bar{y}_i - (\hat{\beta}_{20} + \hat{\beta}_{11} x_i)$$

$$\bar{y}_i - (\hat{\beta}_{20} + \hat{\beta}_{21} x_i)$$

$$\bar{y}_i - (\hat{\beta}_{30} + \hat{\beta}_{31} x_i)$$

$$\bar{y}_i - (\hat{\beta}_{40} + \hat{\beta}_{41} x_i)$$

R-student residual :-

The R-student residual is defined as,

$$t_i = \frac{r_i}{\sqrt{s_{(i)}^2 (1-h_{ii})}}$$

where,

$$h_{ii} = \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}} \right)$$

&

$$\underline{s_{(i)}^2} = \frac{(n-2) \underline{\text{MSE}} - r_i^2 / (1-h_{ii})}{\underline{n-3}}$$

Note:-

$$h_{ii} \rightarrow 0 \Leftrightarrow n \rightarrow \infty$$

i) Normality of error :-

Let r_1, \dots, r_n are residuals.

Procedure of normal probability plot:-

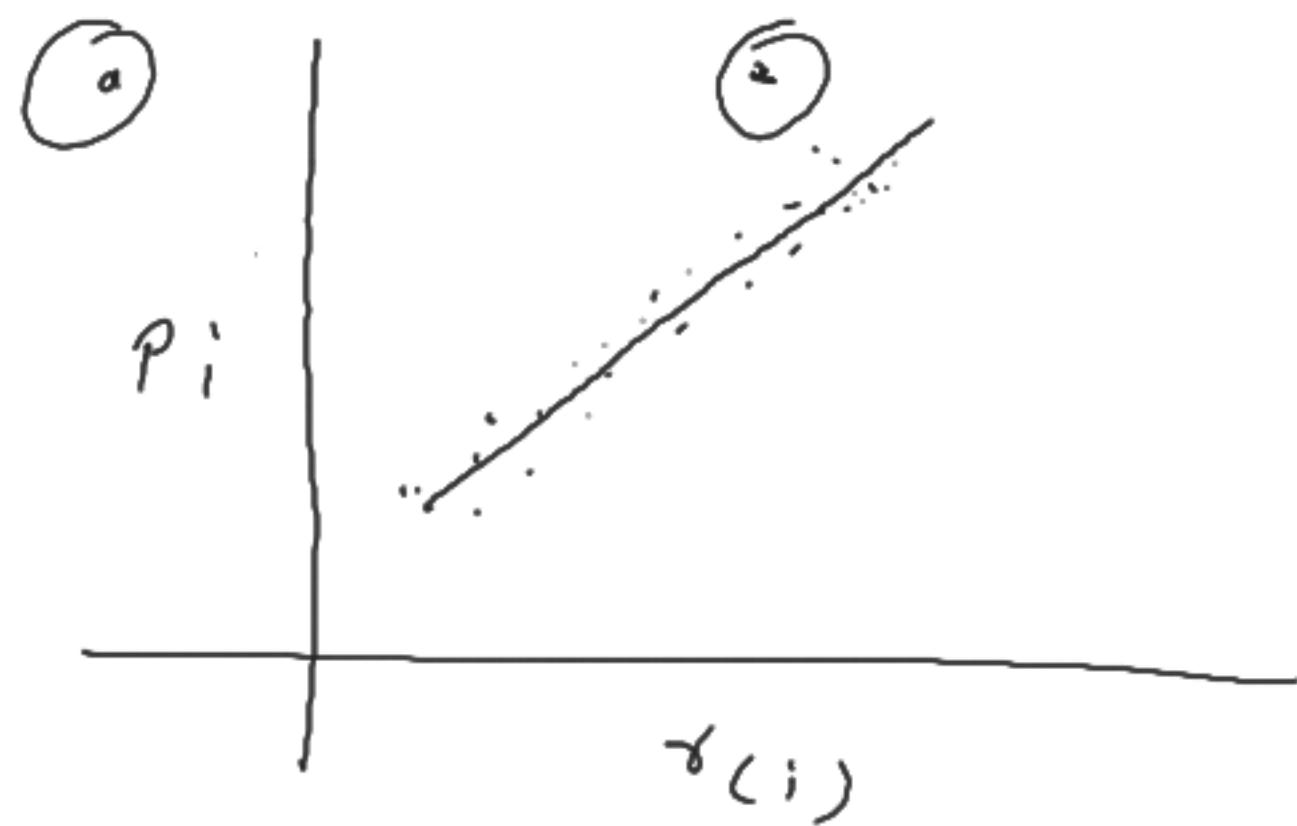
i) Arrange residuals in an increasing order $(r_{(1)}, r_{(2)}, \dots, r_{(n)})$

ii) Calculate cumulative probabilities.

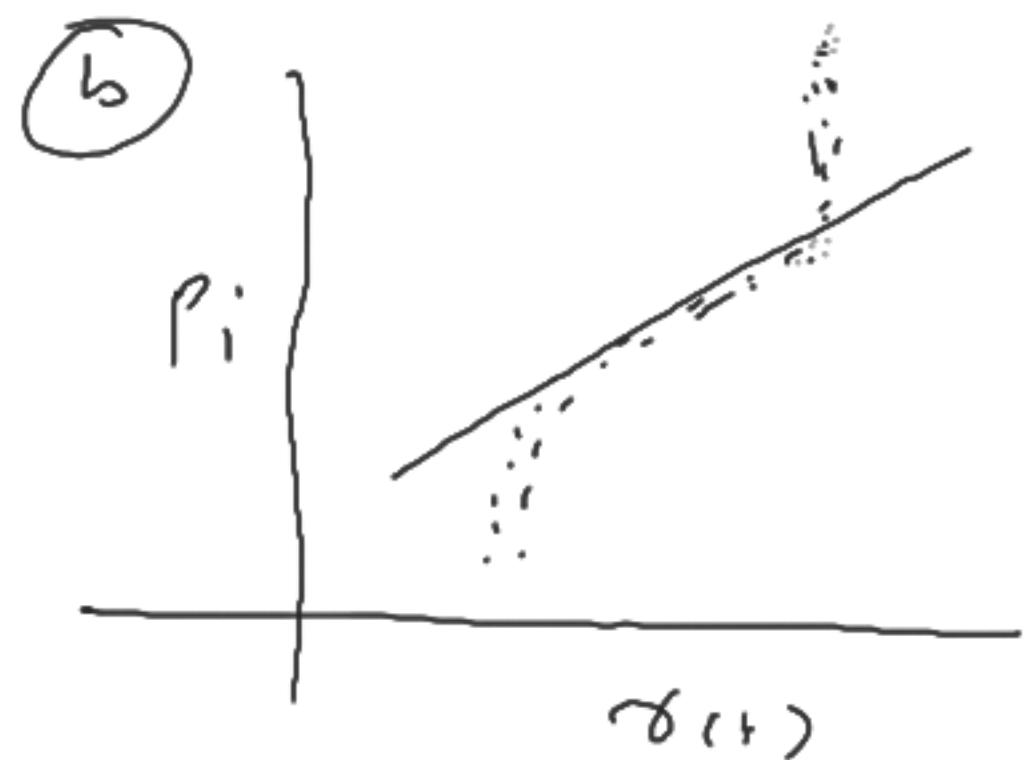
$$P_i = \frac{i - \frac{1}{2}}{n}, \quad i=1, 2, \dots, n.$$

iii) Draw the plot of $(r_{(i)}, P_i)$.

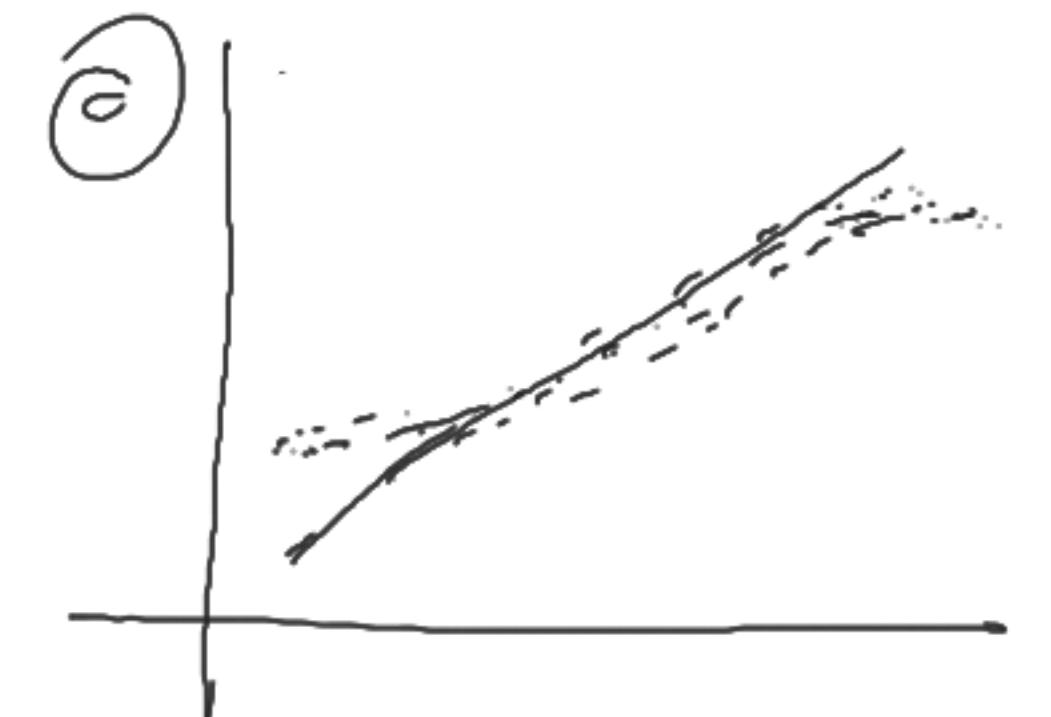
OR. Alternatively, draw the plot of $(v_{(i)}, \phi^{-1}(P_i))$, ($\because v_{(i)} = r_{(i)} - \bar{r}$)
where ϕ^{-1} c.d.f. of $\sim N(0, 1)$



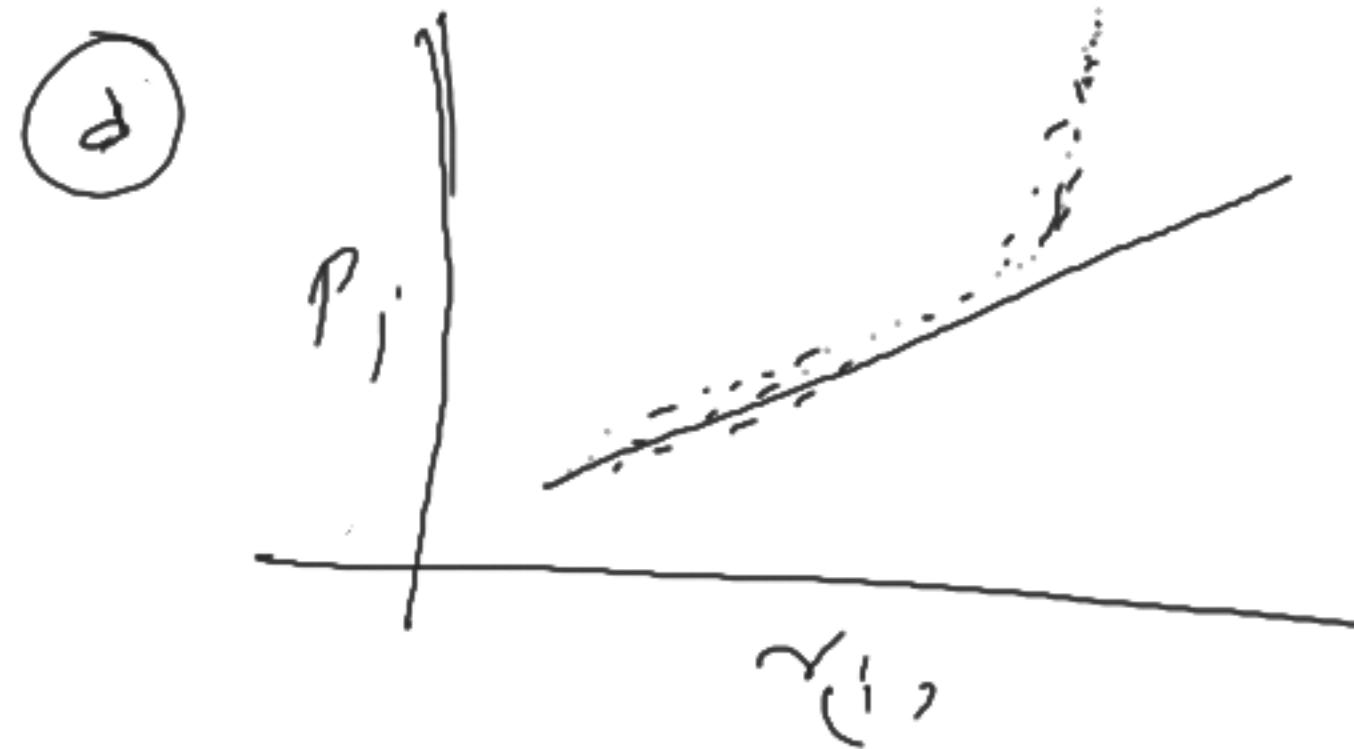
$\gamma_i \sim \text{normal dist}$



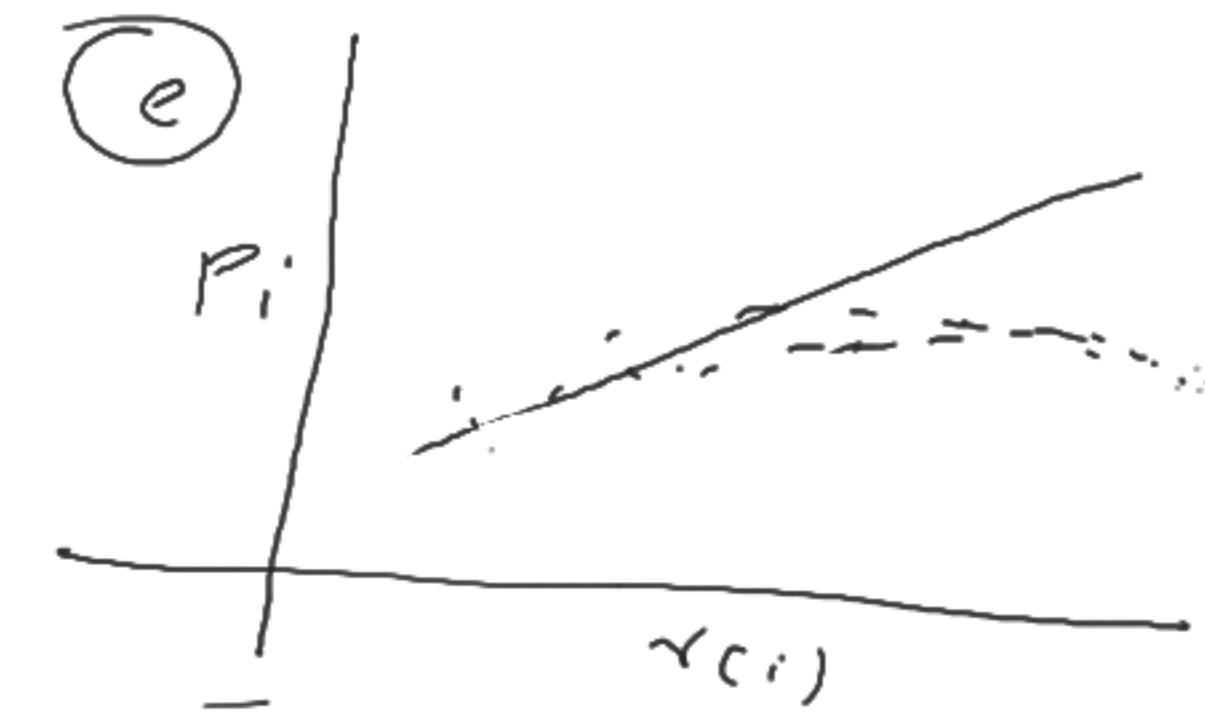
$\gamma_i \sim \text{heavy-tailed dist}$



$\gamma_i \sim \text{light-tailed dist}$

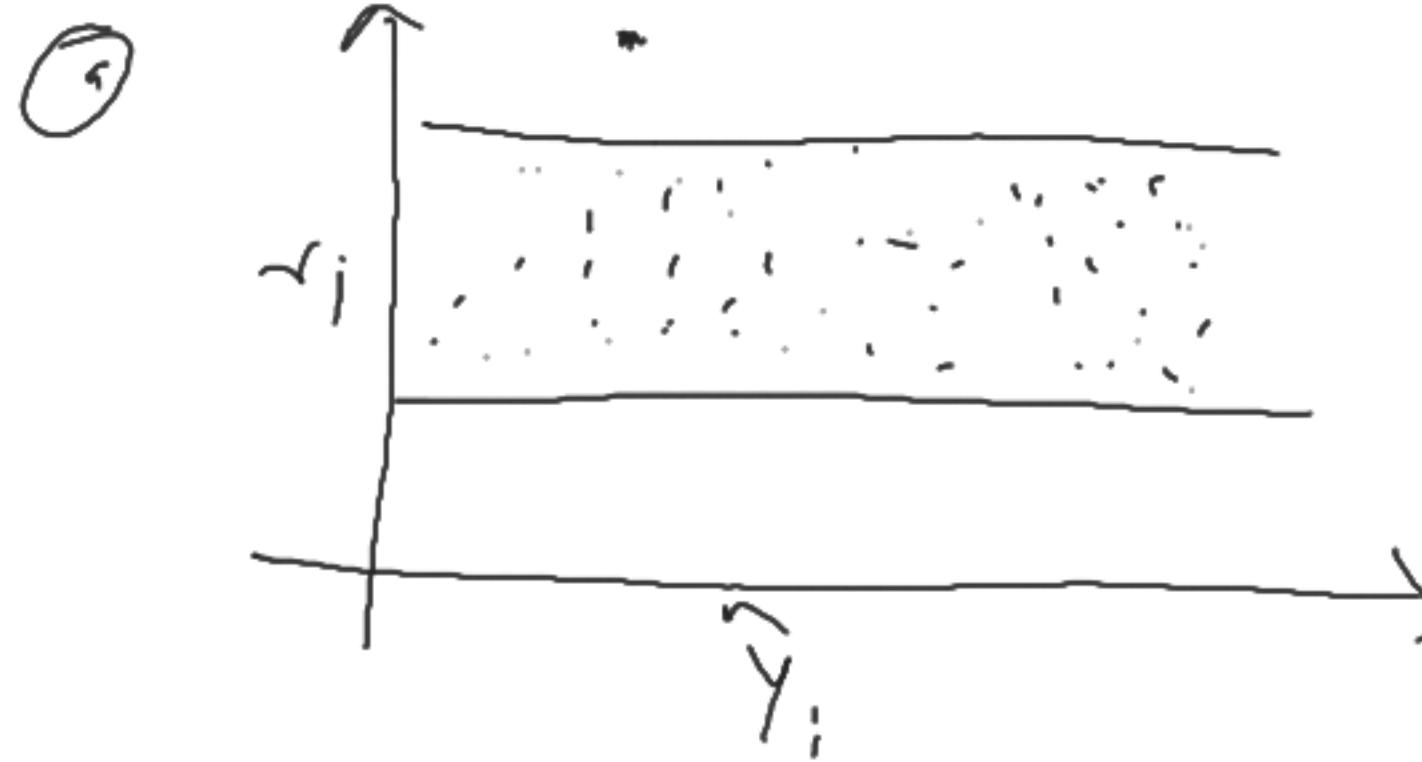


$\gamma_i \sim \text{positively skewed}$

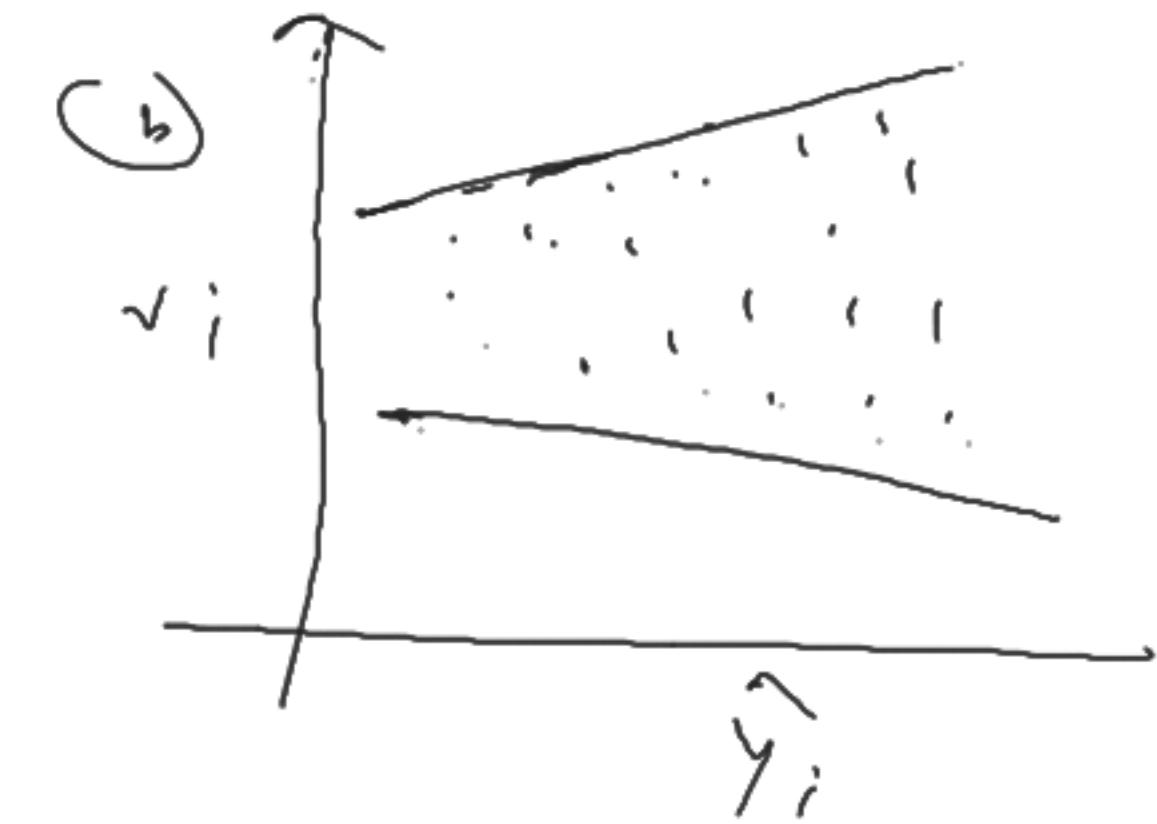


$\gamma_i \sim \text{negatively skewed}$

Constant variance of Error :-



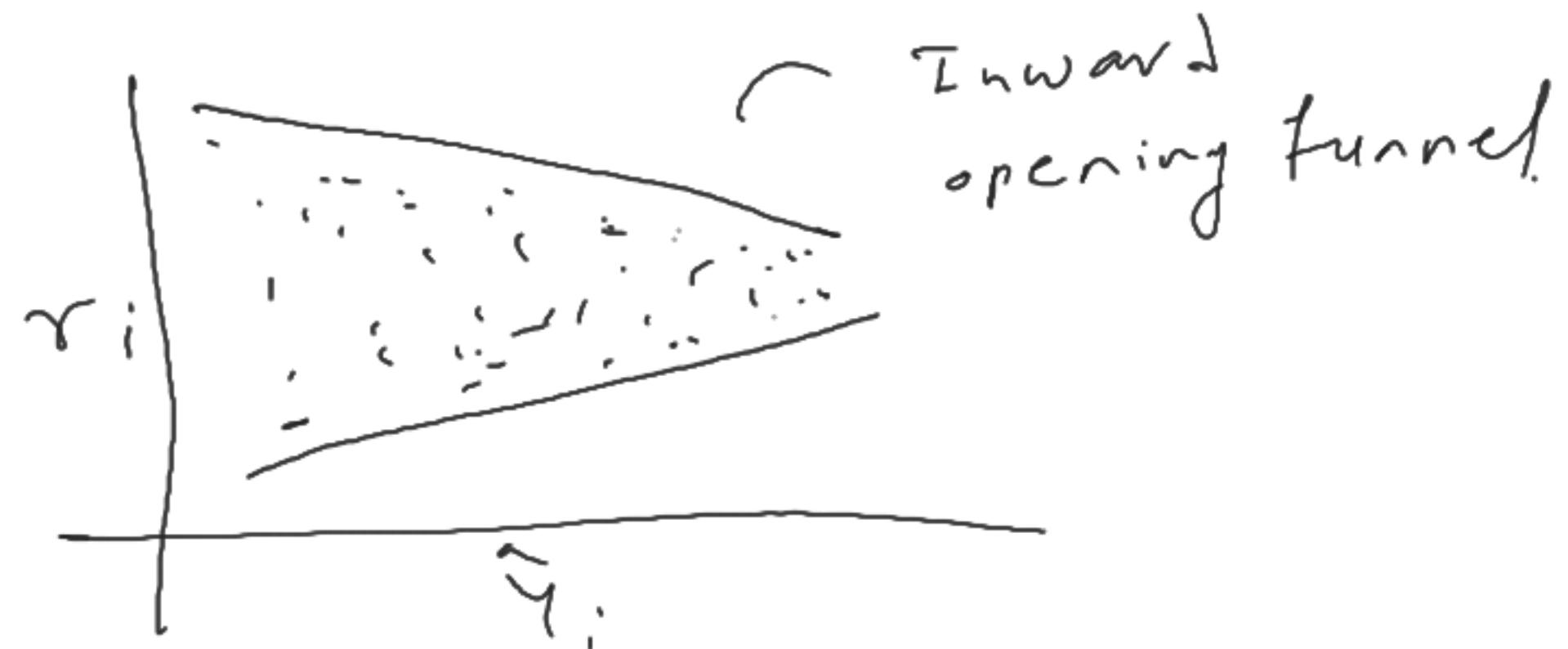
$v(\cdot)$ is constant



$v(\cdot)$ is not constant
 $v(\cdot)$ is increasing
 $f_i \cdot f_j$

outward opening funnel

(c)



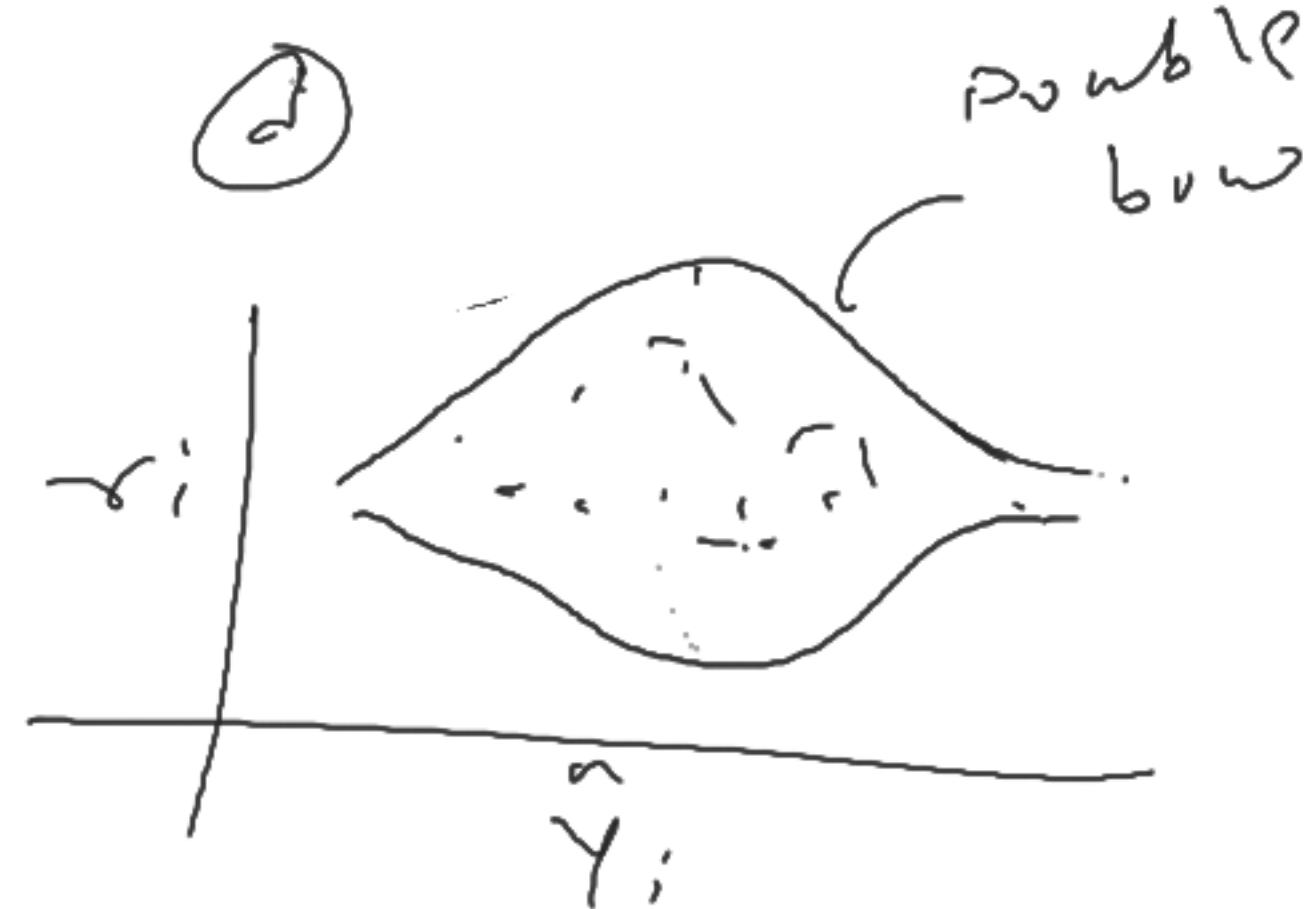
$v(\varepsilon_i)$ is not constant

$v(\varepsilon_i)$ is a decreasing

fnc. of \hat{Y}_i .

Inward
opening funnel

(d)

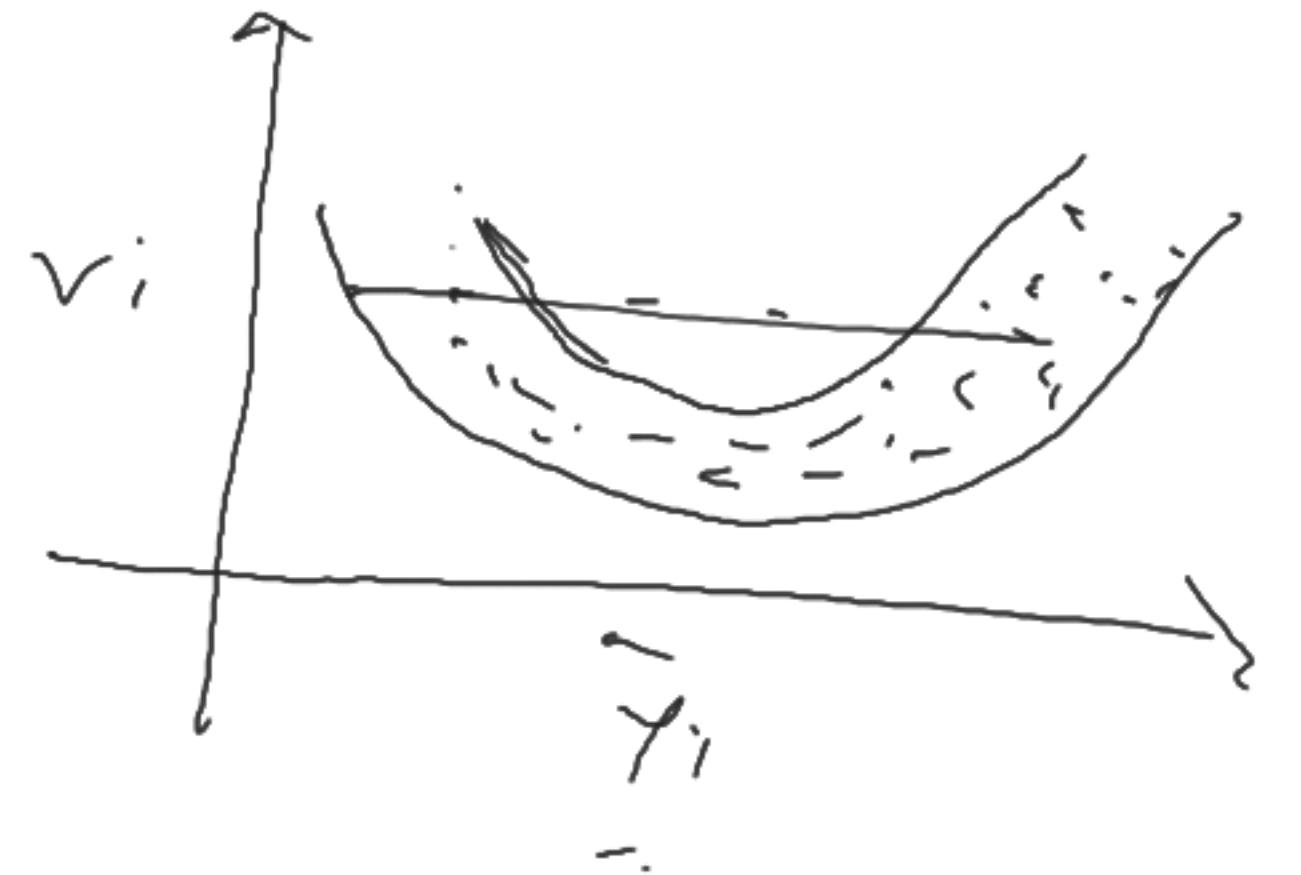


$v(\varepsilon_i)$ is not constant

In this case, y may
follows Binomially
dist

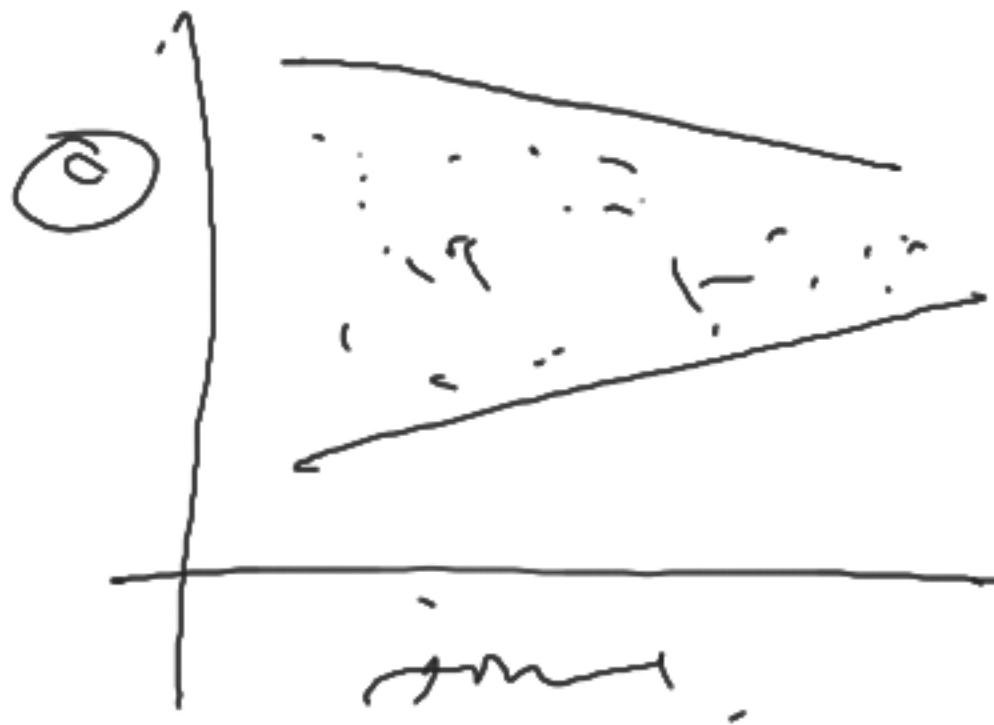
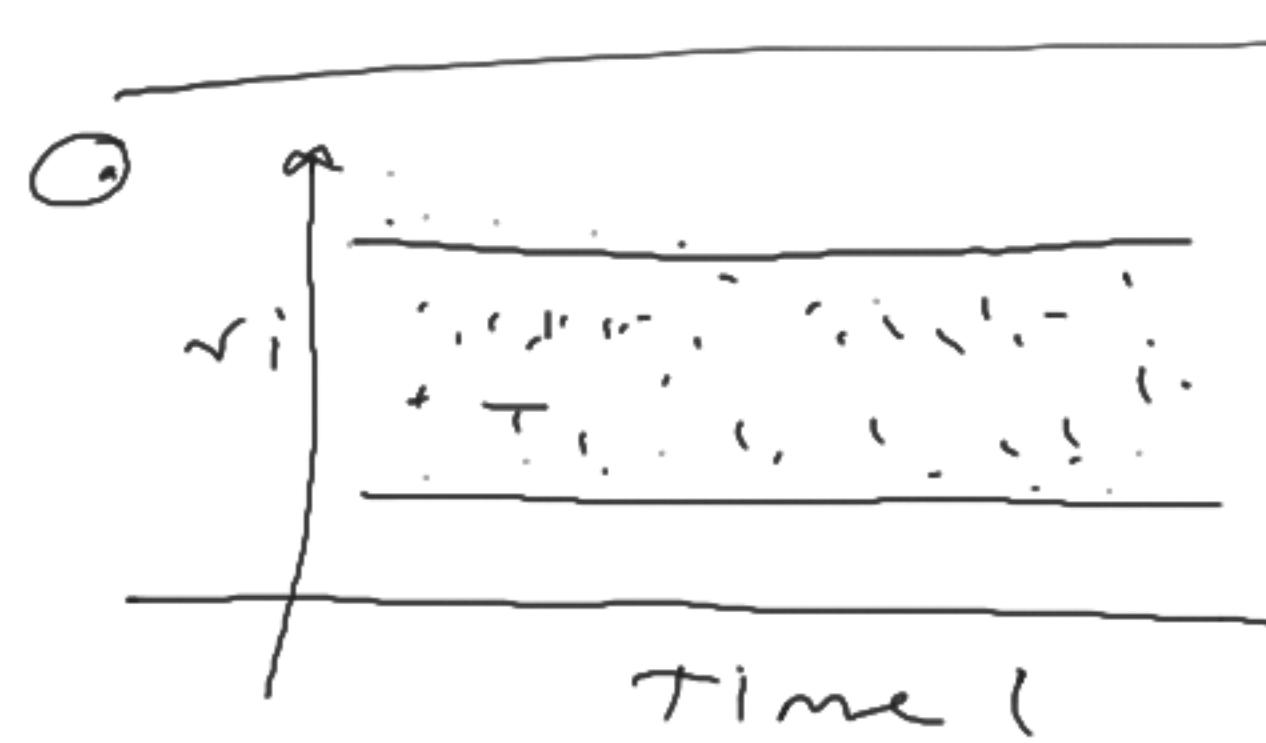
Double
bowl

(e)



Non-linearity

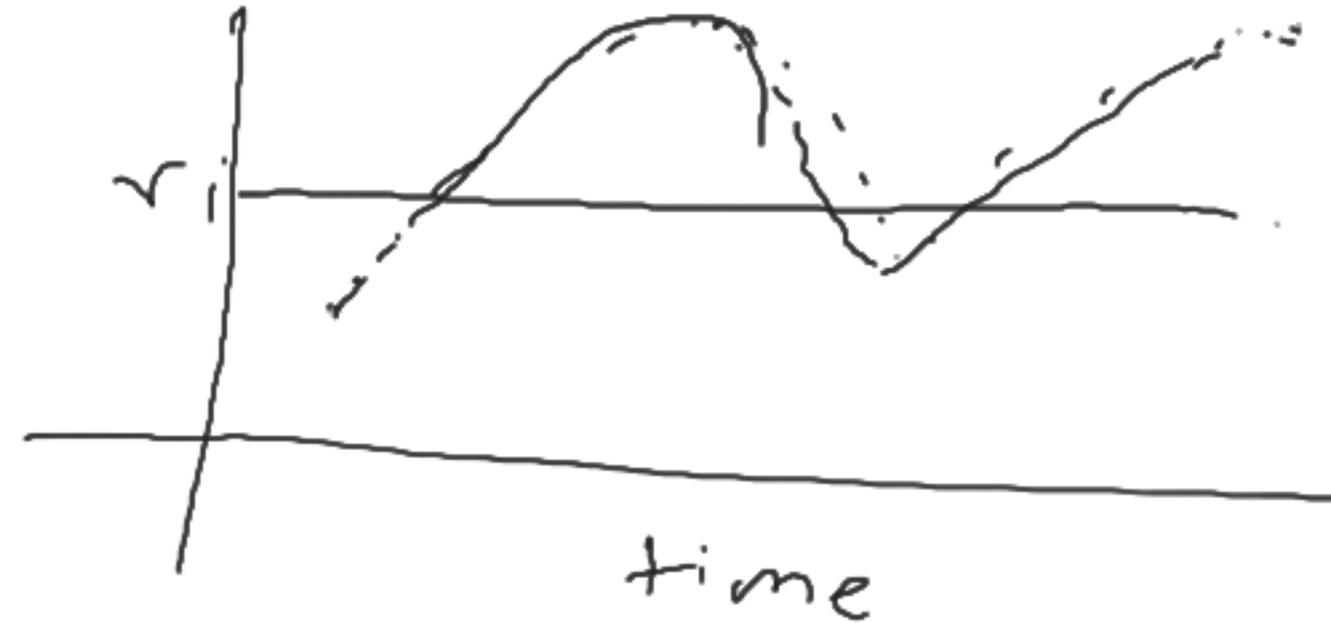
Errors are uncorrelated:-



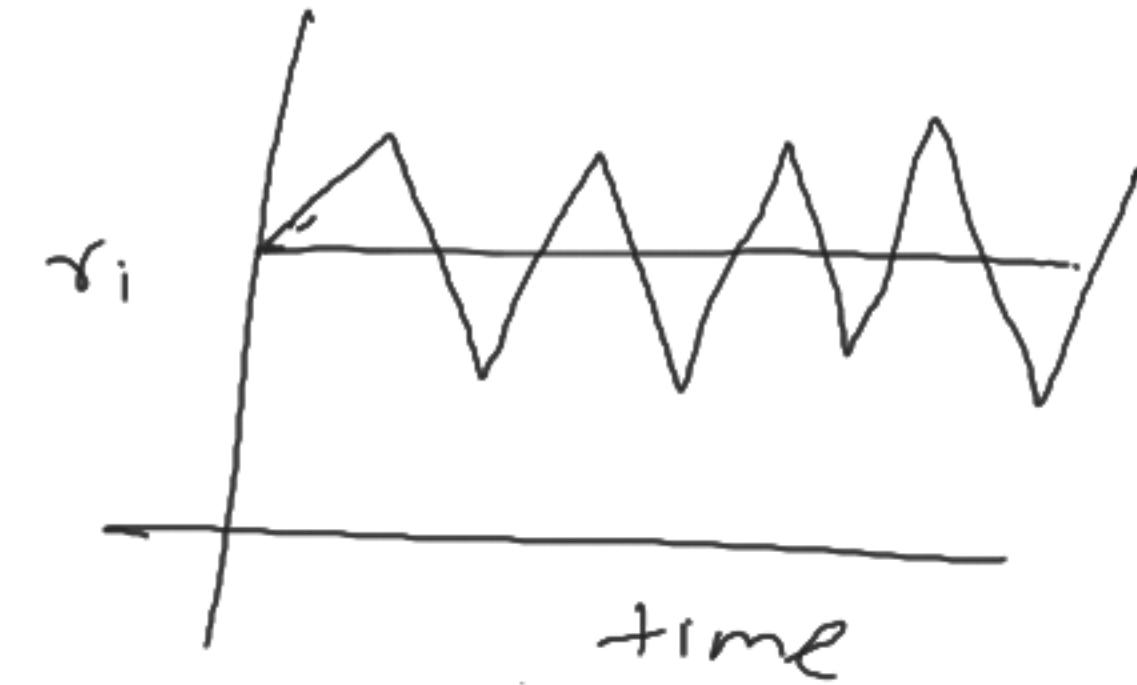
Errors uncorrelated



t



Positive autocorrelation



Negative auto correlation.