

# classmate

# Statistics and Information Theory

$\Rightarrow$  Gaussian Distribution ( $\mu, \sigma^2$ ):

$$\text{PDF: } f_x(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\Rightarrow$  Standard gaussian distribution:

$$\mu = 0 \quad \sigma^2 = 1$$

$$f_x(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$\Rightarrow$  Variance  $\rightarrow$  Variance of data  $\rightarrow$  spread of data across the mean value.



Agriculture

$$y = f(x_1, x_2, x_3)$$

↕  
 Yield of crop (in kg)  
 ↓  
 quantity of fertilizer used.

Model.

$$y = f(x_1, x_2)$$

↕  
 Dosage of Meds to be given  
 ↓  
 Age of patient  
 ↓  
 BP level of patient

$$y = mx + c$$

~~$$\tan \theta = m$$~~

Data sheet 1

$$x = 2, y = 5$$

$$s = 2m + c$$

$$(2, 5)$$

$\frac{1}{+}$  Just a point  
(can't say which  
line it is on)

Not much info

Data sheet 2

$$m = 3, c = 4$$

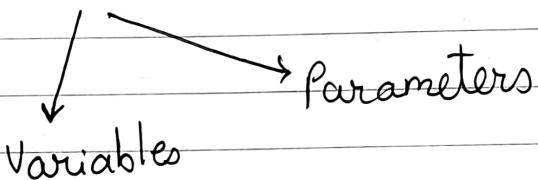
$$y = 3x + 4$$

Gives multiple  
points and  
on a single line.

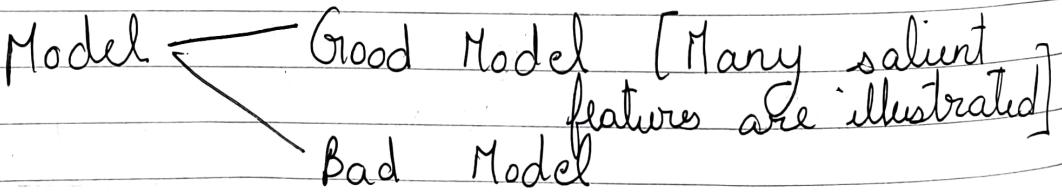
Much info.

$$y = mx + c \rightarrow \begin{array}{l} \text{Variables } (x, y) \\ \text{Parameters } (m, c) \end{array}$$

Model

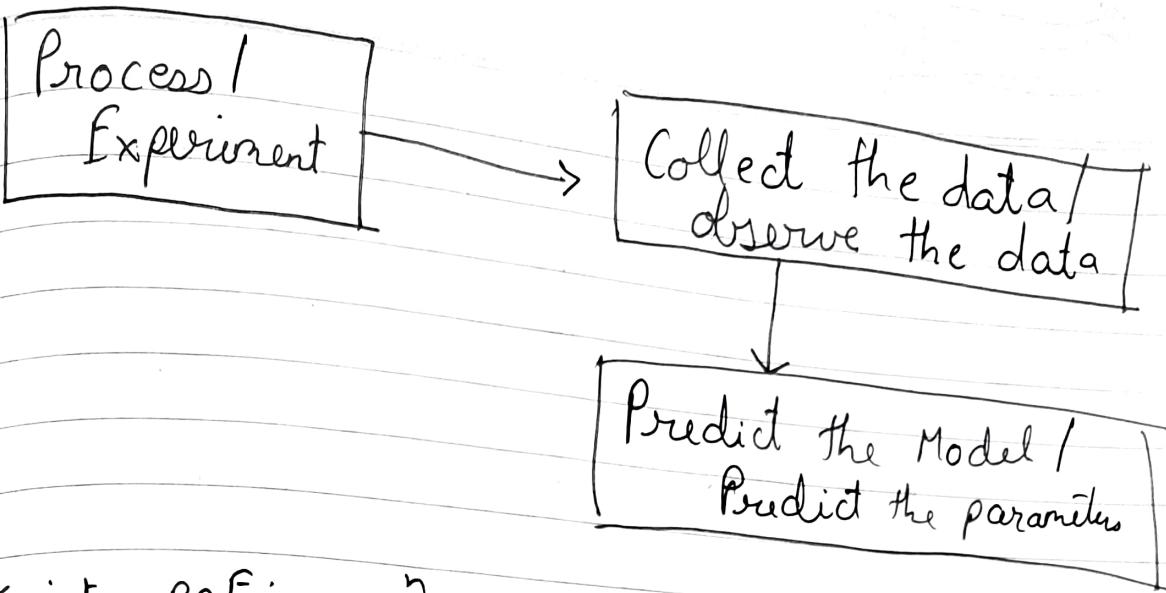


- Parameters play a key role here is state
- Parameters: Component of a model which illustrates complete information of model.



Process/  
Experiment

Collect the data / observe



Joint POF:

$$\frac{1}{(2\pi)^{N/2}} e^{-\frac{1}{2} \sum_{i=1}^n y_i^2}$$

$$y = mx + c \Rightarrow y = \beta_0 + \beta_1 x$$

$\downarrow$  Intercept parameter     $\downarrow$  slope parameter

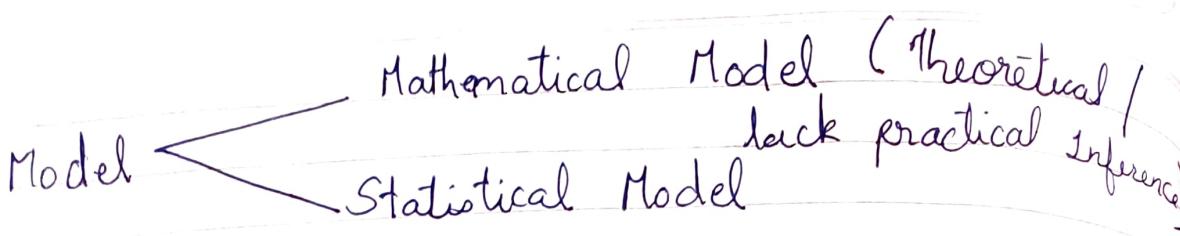
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Generic Model Regression.

$$y = \beta_0 + \beta_1 x_1$$

$\Downarrow$

Data set N  
different values of  
 $y \Rightarrow y$   
 $\beta_0 \Rightarrow \hat{\beta}_0$   
 $\beta_1 \Rightarrow \hat{\beta}_1$



$$y = 3x + 2$$

yield of crop (in kg)      Quantity of field fertilizer

field of loop (in kg)

$$y = 3x + 2$$

Quantity of fertilizer (in kg)

mathematical model :

(Theoretical and lacks practical inferences)

x	y
1	5
2	8
:	:
100	302

Statistical Model :

$x = 1 \text{ kg}$   $y = 5 \text{ kg} \rightarrow$  Not exact depend on various factors such as rainfall, draught etc.

$x = 1 \text{ kg}$   $y = 4.7 \text{ kg}$      $x = 1 \text{ kg}$   $y = 5.1 \text{ kg}$  } Shows the actual yield obtained

according to statistical model, on average when the experiment of crop yield was conducted.

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

By statistical model, on an average when we do an experiment, we can get estimated values  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , which will be close to  $\beta_0$  and  $\beta_1$  respectively.

$$\hat{\beta}_0 = \beta_0 + \epsilon$$

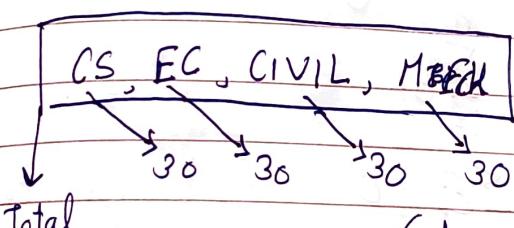
$$\hat{\beta}_1 = \beta_1 + \theta$$

$\theta, \epsilon \rightarrow \text{error}$

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

B.Tech

↓ n Branches



Sample: 120 students  
Population: 400 students

Total students = 400

(Chose 30 from each)

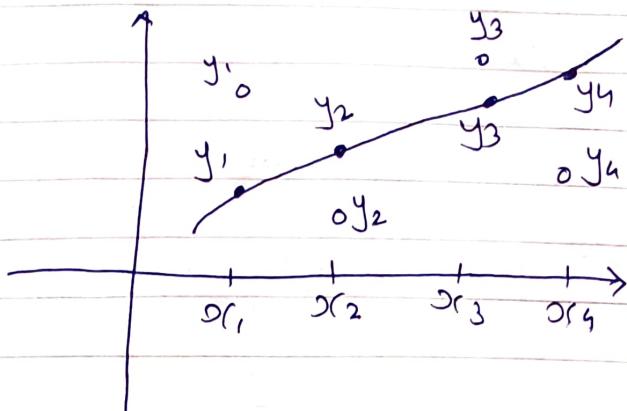
Sample: Subset of population

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

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

n observations:

$y(1), y(2), \dots, y(n) \rightarrow y_1, y_2, \dots, y_n$   
 $x(1), x(2), \dots, x(n) \rightarrow x_1, x_2, \dots, x_n$ .



- → theoretical
- → actual data

Let individual errors for  $n$  observations be  $\epsilon_1, \epsilon_2, \epsilon_3, \dots, \epsilon_n$

Average Error throughout experiment

$$= \frac{1}{n} \sum_{i=1}^n \epsilon_i$$

In some case,  $\sum \epsilon_i = 0$  But you can't say that there were no errors in your experiments as there can be +ve and -ve errors at each observation which can add upto zero.

Hence, we use a more accurate metric of  $\epsilon_i^2$ .

Least squared estimation

$$\sum_{i=1}^n \epsilon_i^2 \rightarrow \text{LSE : Least squared estimation}$$

Sum of squared errors across all observations

$$\sum_{i=1}^n \epsilon_i$$

Sum of errors across all observations

Purpose: To estimate set of data  $(x_i, y_i)$  parameters  $\beta_0, \beta_1$  through

$$\epsilon_i = y_i - \beta_0 - \beta_1 x_i$$

$$\epsilon_i^2 = (y_i - \beta_0 - \beta_1 x_i)^2$$

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

Sample of  $n$  values  
i/p:  $x_1, x_2, \dots, x_n$   
(sample mean)  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$   
o/p:  $y_1, y_2, \dots, y_n$   
 $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

our task is to  
min of  $\sum_{i=1}^n \epsilon_i^2$

$$f(\beta_0, \beta_1) = \sum_{i=1}^n \epsilon_i^2$$

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

to find minima, we find derivative

$$\frac{\partial f(\cdot)}{\partial \beta_0} = 0$$

$$\frac{\partial f(\cdot)}{\partial \beta_1} = 0$$

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

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

$$-2(n\bar{y} - n\beta_0 - n\beta_1 \bar{x}) = 0$$

$$-\bar{y} + \beta_0 + \beta_1 \bar{x} = 0$$

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

$$2(y_i - \beta_0 - \beta_1 x_i)(-x_i) = 0$$

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

$$-2(n\bar{y} - n\beta_0 - n\beta_1 \bar{x}) x_i = 0$$

$$\beta_1 = \frac{\sum_{i=1}^n x_i y_i - \beta_0 n \bar{x}}{\sum_{i=1}^n x_i^2}$$

(you need to know value of  $\beta_0$  to find  $\beta_1$  and vice versa)

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \rightarrow (\bar{y} - \hat{\beta}_0 \bar{x})$$

$$\hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - \bar{y} n \bar{x} + \hat{\beta}_0 n \bar{x}^2$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{y} \bar{x}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}$$

Denominator term:

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

$$= \sum_{i=1}^n x_i^2 + \sum_{i=1}^n \bar{x}^2 - 2\bar{x} \sum_{i=1}^n x_i$$

$$= \sum_{i=1}^n x_i^2 + n \bar{x}^2 - 2n \bar{x}^2$$

$$= \sum_{i=1}^n x_i^2 - n \bar{x}^2$$

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

Numerator term:

$$\sum_{i=1}^n \alpha_i y_i - n\bar{x}\bar{y}$$

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^n [x_i y_i - \bar{y}x_i - \bar{x}y_i + \bar{x}\bar{y}] \\ &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} - n\bar{y}\bar{x} + n\bar{x}\bar{y} \\ &= \sum_{i=1}^n \alpha_i y_i - n\bar{x}\bar{y} \end{aligned}$$

$$\therefore \hat{\beta}_1 = \frac{\sum_{i=1}^n (\alpha_i - \bar{\alpha})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

LSE  $\rightarrow$  Least

Square Estimator

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

Hessian Matrixic

$$H = \begin{bmatrix} \frac{\partial^2 F(\cdot)}{\partial \beta_0^2} & \frac{\partial^2 F(\cdot)}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 F(\cdot)}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 F(\cdot)}{\partial \beta_1^2} \end{bmatrix}$$

$$\begin{aligned} H &= \begin{bmatrix} \sum 2 & \sum 2x_i \\ \sum 2x_i & \sum 2x_i^2 \end{bmatrix} = 2 \begin{bmatrix} 1 & x_i \\ x_i & x_i^2 \end{bmatrix} \\ &= 2 \begin{bmatrix} 2n & 2n\bar{x} \\ 2n\bar{x} & 2 \sum_{i=1}^n \alpha_i^2 \end{bmatrix} \end{aligned}$$

$$\det(H) = 4n \sum_{i=1}^n x_i^2 - 4n^2 \bar{x}^2$$

$$= 4n \left[ \sum_{i=1}^n x_i^2 - n \bar{x}^2 \right]$$

$$= 4n \sum_{i=1}^n (x_i - \bar{x})^2$$

Parameter - fixed unknown value.

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

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

The estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  of the parameters  $\beta_0$  and  $\beta_1$  are linear w.r.t. output values / output observations ( $y_i$ )

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

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

$\underbrace{\qquad\qquad\qquad}_{\textcircled{O}}$

$y_i$  are known as they are input values.

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

$$\text{Let } c_i = \frac{x_i - \bar{x}}{S_{xx}}$$

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i$$

$$c_i = \frac{x_i - \bar{x}}{S_{xx}}$$

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i$$

$$\hat{\beta}_1 = c_1 y_1 + c_2 y_2 + \dots + c_n y_n$$

The estimator  $\hat{\beta}_1$  is linear w.r.t.  $y_i$ .

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

↑  
average sense.

## Unbiased Estimator and Bias of Estimator

The estimator on an average sense should yield the value of the parameter.  
(It is estimating).

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

$$E\{\hat{\beta}_0\} = \beta_0$$

$$\text{Bias}(\hat{\beta}_1; \beta_1) = E(\hat{\beta}_1) - \beta_1$$

$$\text{Bias}(\hat{\beta}_0; \beta_0) = E(\hat{\beta}_0) - \beta_0$$

$$E\{\hat{\beta}_i\} = \beta_i$$

$$\hat{\beta}_i = \sum_{i=1}^n c_i y_i$$

$$E\{\hat{\beta}_i\} = E\left\{ \sum_{i=1}^n c_i y_i \right\}$$

$$= \sum_{i=1}^n c_i E\{y_i\}$$

$$= \sum_{i=1}^n c_i [E(B_0 + \beta_1 x_i + \varepsilon_i)]$$

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

$$= \sum_{i=1}^n c_i \left[ \underbrace{E(B_0)}_{\beta_0} + \underbrace{E(\beta_1 x_i)}_{\beta_1 x_i} + \underbrace{E(\varepsilon_i)}_0 \right]$$

$$E\{\hat{\beta}_i\} = \sum_{i=1}^n c_i [\beta_0 + \beta_1 x_i]$$

$$E\{\hat{\beta}_i\} = \beta_0 \underbrace{\sum_{i=1}^n c_i}_{T_1} + \beta_1 \underbrace{\sum_{i=1}^n c_i x_i}_{T_2}$$

Coming to

$$T_1 : \sum_{i=1}^n c_i = \sum_{i=1}^n \frac{x_i - \bar{x}}{S_{xx}}$$

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

$$= \frac{1}{S_{xx}} \left[ \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} \right]$$

$$\boxed{\sum_{i=1}^n c_i = 0}$$

$$T_1 = 0$$

$$\boxed{n\bar{x}} \quad \boxed{n\bar{\bar{x}}}$$

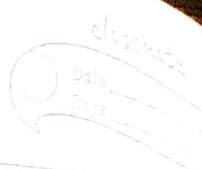
$$\begin{aligned}
 T_2 &= \sum_{i=1}^n c_i x_i = \sum_{i=1}^n \left( \frac{\sigma_i - \bar{\sigma}}{S_{\text{cov}}} \right) x_i \\
 &= \frac{1}{S_{\text{cov}}} \left[ \sum_{i=1}^n (\hat{\sigma}_i)^2 - \bar{x} \bar{x} \right] \\
 &= \frac{1}{S_{\text{cov}}} \left[ \sum_{i=1}^n \sigma_i^2 - n \bar{x}^2 \right] \\
 &= \frac{1}{S_{\text{cov}}} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2 \\
 &= 1
 \end{aligned}$$

$$\begin{aligned}
 \therefore E[\hat{\beta}_1] &= \beta_1(1) \\
 &= \beta_1
 \end{aligned}$$

$$\begin{aligned}
 \text{Bias}(\hat{\beta}_1; \beta_1) &= E(\hat{\beta}_1) - \beta_1 \\
 &= \beta_1 - \beta_1 \\
 &= 0
 \end{aligned}$$

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

$$\begin{aligned}
 E[\hat{\beta}_0] &= E[\bar{y} - \hat{\beta}_1 \bar{x}] \\
 &= E(\bar{y}) - E[\hat{\beta}_1 \bar{x}] \\
 &= E(\bar{y}) - E[\bar{x} \sum_{i=1}^n c_i y_i] \\
 &= E(\bar{y}) - E[\bar{x} \beta_1]
 \end{aligned}$$



$$\begin{aligned}
 \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i \\
 \bar{y} &= \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + \varepsilon_i) \\
 E(\bar{y}) &= E \left[ \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + \varepsilon_i) \right] \\
 &= \frac{1}{n} \sum_{i=1}^n E[\beta_0 + \beta_1 x_i + \varepsilon_i] \\
 &= \frac{1}{n} \sum_{i=1}^n [\beta_0 + \beta_1 x_i + 0] \\
 &= \frac{1}{n} \sum_{i=1}^n \beta_0 + \frac{1}{n} \sum_{i=1}^n \beta_1 x_i \\
 &= \frac{1}{n} (n \beta_0) + \frac{1}{n} (\beta_1 \times n \bar{x}) \\
 &= \beta_0 + \beta_1 \bar{x}
 \end{aligned}$$

$$\begin{aligned}
 E[\hat{\beta}_0] &= E(\bar{y}) - \beta_1 \bar{x} \\
 &= \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} \\
 &= \beta_0
 \end{aligned}$$

$(x_1, y_1) \Rightarrow (x_1, \bar{y}) \dots \dots (x_n, y_n)$

$$\begin{array}{ll}
 \beta_0 & \beta_1 \\
 \downarrow & \downarrow \\
 x & y \\
 \hat{\beta}_0 & \hat{\beta}_1
 \end{array}$$

$$\text{Var}(x-y) = \text{Var}(x) + \text{Var}(-y)$$

$$= \text{Var}(x) + \text{Var}(y)$$

$$- 2\text{Cov}(x, y)$$

discrete

## Variance and Covariance

$x, y$   $\text{Cov}(x, y)$

$$\text{Cov}(x, y) = E[xy] - E[x]E[y]$$

$$= E[(x-\bar{x})(y-\bar{y})]$$

$$\begin{array}{ccc} \downarrow & & \downarrow \\ E[x] & & E(y) \\ \downarrow & & \downarrow \\ \mu_x & & \mu_y \end{array}$$

$$\text{Var}(x+y) = \text{Var}(x) + \text{Var}(y) + 2\text{Cov}(x, y)$$

$$y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\text{non-random terms}} + \underbrace{\varepsilon_i}_{\text{normal distribution}}, \quad \varepsilon_i \stackrel{\text{IID}}{\sim} N(0, \sigma^2)$$

$y_i$  will also follow normal distribution

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

$$\begin{aligned} E(y_i) &= E(\beta_0 + \beta_1 x_i + \varepsilon_i) \\ &= E(\beta_0) + E(\beta_1 x_i) + E(\varepsilon_i) \\ &= \beta_0 + \beta_1 x_i \end{aligned}$$

$$\begin{aligned} E(y_1) &= E(\beta_0 + \beta_1 x_1 + \varepsilon_1) = \beta_0 + \beta_1 x_1, \\ E(y_2) &= E(\beta_0 + \beta_1 x_2 + \varepsilon_2) = \beta_0 + \beta_1 x_2 \end{aligned} \quad \left. \begin{array}{l} E(y_i) = \beta_0 \\ + \beta_1 x_i \end{array} \right\}$$

$$\begin{aligned}\text{Var}(y_i) &= \text{Var}(\beta_0 + \beta_1 x_i + \varepsilon_i) \\ &= \text{Var}(\beta_0) + \text{Var}(\beta_1 x_i) + \text{Var}(\varepsilon_i) \\ &= 0 + 0 + \sigma^2 \\ &= \sigma^2\end{aligned}$$

$$\begin{aligned}\text{Var}(\hat{\beta}_0) \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) &= \text{Var}(\bar{y}) + \text{Var}(\hat{\beta}_1 \bar{x}) - \\ &\quad 2 \text{Cov}(\bar{y}, \hat{\beta}_1 \bar{x}) \\ &= \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) - 2 \text{Cov}(\bar{y}, \hat{\beta}_1 \bar{x})\end{aligned}$$

$$\begin{aligned}\text{Cov}(\bar{y}, \hat{\beta}_1 \bar{x}) &= \bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1) \\ &= \bar{x} E[(\bar{y} - E(\bar{y})) (\hat{\beta}_1 - E(\hat{\beta}_1))] \\ &= \bar{x} E[\underbrace{(\bar{y} - E(\bar{y}))}_{T_1} \underbrace{(\hat{\beta}_1 - \beta_1)}_{T_2}]\end{aligned}$$

$$\begin{aligned}T_1 &= \bar{y} - E(\bar{y}) = \frac{1}{n} \sum_{i=1}^n y_i - E\left(\frac{1}{n} \sum_{i=1}^n y_i\right) \\ &= \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + \varepsilon_i) - \frac{1}{n} \sum_{i=1}^n E(\beta_0 + \beta_1 x_i + \varepsilon_i) \\ &= \frac{1}{n} [n\beta_0 + n\beta_1 \bar{x} + n\bar{\varepsilon}] - \frac{1}{n} [n\beta_0 + n\beta_1 \bar{x} + 0] \\ &= \bar{\varepsilon}\end{aligned}$$

$$\bar{\varepsilon} = \frac{1}{n} \sum_{i=1}^n \varepsilon_i = \frac{1}{n} [\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_n]$$

$$\begin{aligned}E(\bar{\varepsilon}) &= \frac{1}{n} [E(\varepsilon_1) + E(\varepsilon_2) + \dots + E(\varepsilon_n)] \\ &= \frac{1}{n} [0 + 0 + \dots + 0] \\ &= 0 \quad (\because \varepsilon_i \sim N(0, \sigma^2))\end{aligned}$$

$$E(\bar{\varepsilon}) = 0$$

$$\text{Var}(\bar{y}) = \text{Var}\left(\frac{1}{n} \sum_i c_i y_i\right)$$

$$= \frac{1}{n^2} \sum_i c_i^2 \text{Var}(y_i)$$

$$T_2 = \hat{\beta}_1 - \beta_1$$

$$= \frac{1}{n^2} \sum_i c_i^2 \sigma^2 = \frac{\sigma^2}{n}$$

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

$$= \sum_{i=1}^n c_i (\beta_0 + \beta_1 x_i + \varepsilon_i) - \beta_1$$

$$= \beta_0 \sum_{i=1}^n c_i + \beta_1 \sum_{i=1}^n c_i x_i + \sum_{i=1}^n c_i \varepsilon_i - \beta_1$$

$$= \beta_1 (1) + \sum_{i=1}^n c_i \varepsilon_i - \beta_1$$

$$= \sum_{i=1}^n c_i \varepsilon_i$$

$$\text{Cov}(\bar{y}, \hat{\beta}_1, \bar{x}) = \bar{x} E[\bar{\varepsilon} (\sum_{i=1}^n c_i \varepsilon_i)]$$

$$= \bar{x} \bar{\varepsilon} \sum_{i=1}^n c_i E[\varepsilon_i]$$

$$= 0.$$

$$\text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) = \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{n}$$

$$\text{Var}(\sum_{i=1}^n c_i y_i) = \sum_{i=1}^n c_i^2 \text{Var}(y_i) + \cancel{2 \sum_{i=1}^n \sum_{k=1}^n c_i c_k \text{Cov}(y_i, y_k)} = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$

$y_1, y_2 \rightarrow$  are independent.

$$\therefore \text{Var}(\sum_{i=1}^n c_i y_i) = \sum_{i=1}^n c_i^2 \text{Var}(y_i).$$

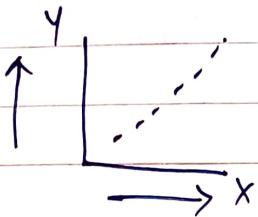
$$\sum c_i^2 = \sum (x_i - \bar{x})^2 = \sum c_i^2 \sigma^2$$

$$= \frac{1}{S_{xx}} \sum (x_i^2 + \bar{x}^2 - 2x_i \bar{x}) = \sigma^2 \sum c_i^2 = \frac{\sigma^2}{S_{xx}}$$

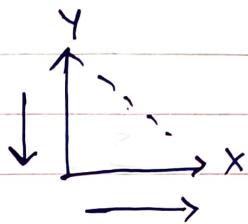
$$= \frac{1}{S_{xx}} (\sum x_i^2 - n \bar{x}^2) = \frac{1}{S_{xx}}$$

- $$\begin{aligned}\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) &= \text{Cov}(\bar{y} - \hat{\beta}_0, \bar{x}, \hat{\beta}_1) \\ &= \text{Cov}(\bar{y}, \hat{\beta}_1) - \text{Cov}(\bar{x}, \hat{\beta}_0, \hat{\beta}_1) \\ &= 0 - \bar{x} \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ &= -\bar{x} \text{Var}(\hat{\beta}_1).\end{aligned}$$

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



$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]$$



- Covariance measures the joint variability of two variables

if  $x \uparrow$  with  $y \uparrow$  then  $\text{Cov}(x, y)$  is +ve.  
 if  $x \downarrow$  with  $y \uparrow$  then  $\text{Cov}(x, y)$  is -ve.  
 if  $x \downarrow$  with  $y \downarrow$  then  $\text{Cov}(x, y)$  is +ve.

$\text{Cov}(x, y)$  gives a linear Relationship

$X:$	-2	-1	0	1	2	$\text{Cov}(x, y) = 0$
$Y:$	4	1	0	1	4	

$y = x^2 \Rightarrow$  relationship does exist b/w  $x$  and  $y$  but it is not linear.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2)$$

IID

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

InID  $\rightarrow$  Independent but not identical distribution

n InID  $\rightarrow$  not independent and not identical distribution

n IID  $\rightarrow$  not independent but identical distribution.

$$\begin{aligned} i=1 &\Rightarrow y_1 = \beta_0 + \beta_1 x_1 + \varepsilon_1 \\ i=2 &\Rightarrow y_2 = \underbrace{\beta_0 + \beta_1 x_2 + \varepsilon_2}_{\text{these are a constant}} \end{aligned}$$

$\hookrightarrow$  these are a constant independent value.

as  $y_i$  has different mean, it is not identically distributed.

Actual Response:  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

Fitted Model/Fitted Response:  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

$$\boxed{r_i = y_i - \hat{y}_i} \rightarrow \underline{\text{residual}}$$

Squared sum of Residual:

$$\begin{aligned} SS_{\text{residual}} &= \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \end{aligned}$$

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

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

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

$$SS_{\text{res}} = S_{yy} + \hat{\beta}_1^2 S_{xx} - 2 \hat{\beta}_1 S_{xy}$$

We know that,  $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$

$$= S_{yy} + \hat{\beta}_1^2 S_{xx} - 2 \frac{S_{xy}}{S_{xx}} S_{xy}$$

$$= S_{yy} + \frac{S_{xy}^2}{S_{xx}^2} S_{xx} - 2 \frac{S_{xy}^2}{S_{xx}^2}$$

$$= S_{yy} - \frac{S_{xy}}{S_{xx}} S_{xy}$$

$$\boxed{SS_{\text{res}} = S_{yy} - \hat{\beta}_1 S_{xy}}$$

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

$$\boxed{SS_{\text{res}} = S_{yy} - \hat{\beta}_1^2 S_{xx}}$$

$$E[SS_{\text{res}}] = (n-2)\sigma^2 \quad || \text{ em nam dhoki didhi}$$

$$\boxed{E \left[ \frac{SS_{\text{res}}}{n-2} \right] = \sigma^2} \Rightarrow SS_{\text{res}} \text{ is unbiased estimator of } \sigma^2$$

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

Maximum Likelihood Estimator (MLE)

n R.V.s  $\Rightarrow X_1, X_2, \dots, X_n$

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

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

$y_1, y_2, \dots, y_n$  are independent.

Joint pdf:  $f(y_1, y_2, \dots, y_n) = f(y_1)f(y_2)\dots f(y_n)$

is product of marginal pdf.

$$= \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_1 - \beta_0 - \beta_1 x_1)^2}{2\sigma^2}} \right) \dots \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_n - \beta_0 - \beta_1 x_n)^2}{2\sigma^2}} \right)$$

$$l(\beta_0, \beta_1, \sigma^2) = \left( \frac{1}{2\pi\sigma^2} \right)^{n/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2}$$

↓ likelihood function

log likelihood function:

$$L(\beta_0, \beta_1, \sigma^2) = \log l(\beta_0, \beta_1, \sigma^2)$$

$$L(\beta_0, \beta_1, \sigma^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

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

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

$$\frac{\partial L(\cdot)}{\partial \sigma^2} = \frac{-n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

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

$$n\bar{y} - n\beta_0 - n\beta_1 \bar{x} = 0$$

$$\therefore \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \dots \textcircled{1}$$

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

$$\sum_{i=1}^n x_i y_i - (\bar{y} - \beta_1 \bar{x}) n \bar{x} - \beta_1 \sum_{i=1}^n x_i^2 = 0$$

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

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

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \dots \textcircled{2}$$

$$\frac{\partial L(\cdot)}{\partial \sigma^2} = 0 \Rightarrow \frac{n}{2\sigma^2} = \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i^2 + \beta_0^2 + \beta_1^2 x_i^2 - 2\beta_0 y_i - 2\beta_1 \beta_0 x_i y_i + 2\beta_0 \beta_1 x_i^2)$$

$$n\sigma^2 = \sum_{i=1}^n y_i^2 + n\beta_0^2 + \beta_1^2 \sum_{i=1}^n x_i^2 - 2n\beta_0 \bar{y} - 2\beta_1 \sum_{i=1}^n x_i y_i + 2n\beta_0 \beta_1 \bar{x}$$

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

Hessian Matrix

$$H = \begin{bmatrix} \frac{\partial^2 L}{\partial \beta_0^2} & \frac{\partial^2 L}{\partial \beta_0 \partial \beta_1} & \frac{\partial^2 L}{\partial \beta_0 \partial \sigma^2} \\ \frac{\partial^2 L}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 L}{\partial \beta_1^2} & \frac{\partial^2 L}{\partial \beta_1 \partial \sigma^2} \\ \frac{\partial^2 L}{\partial \sigma^2 \partial \beta_0} & \frac{\partial^2 L}{\partial \sigma^2 \partial \beta_1} & \frac{\partial^2 L}{\partial \sigma^2 \partial \sigma^2} \end{bmatrix}$$

$$\det(H) < 0$$

$$\text{LSE / MLE} \Rightarrow y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$\varepsilon_i \stackrel{\text{ iid }}{\sim} N(0, \sigma^2)$$

LSE:

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

MLE:

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

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

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

$$\hat{\sigma}^2 = \frac{1}{(n-2)} SS_{\text{residual}}$$

$$\hat{\sigma}^2 = \frac{1}{n} SS_{\text{residual}}$$

$$x_i \sim N(E(x_i), \text{Var}(x_i))$$

$$x_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$$

$$\text{MLE: } \hat{\mu}, \hat{\sigma}^2$$

If these both parameter are unknown

$$\text{Joint PDF: } f(x_1, x_2, \dots, x_n) \\ = f(x_1) f(x_2), \dots, f(x_n)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \dots \\ \dots \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_n - \mu)^2}{2\sigma^2}\right)$$

$$L = \left( \frac{1}{2\pi\sigma^2} \right)^{n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

$$L = \log L$$

$$L = -\frac{n}{2} \log e^{2\pi} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial L}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \Rightarrow \frac{\partial L}{\partial \mu} = 0$$

$$\boxed{\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i}$$

$$\frac{\partial L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial L}{\partial \sigma^2} = 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

# Hypothesis Testing / Confidence Level / Significance level

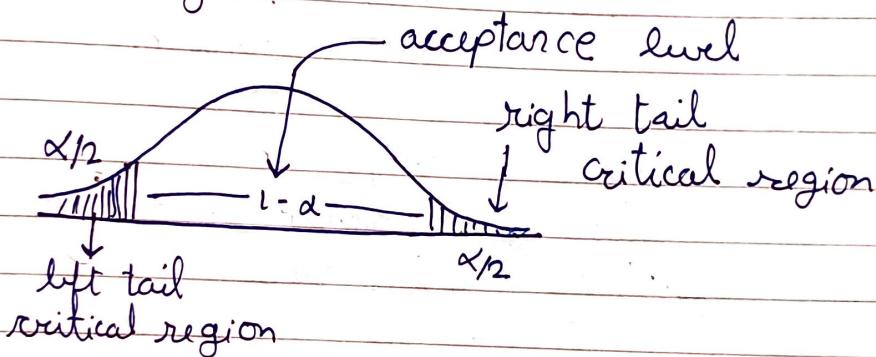
Point estimator  $\rightarrow$  A particular single estimated value of unknown parameter

Interval Estimator  $\rightarrow$  The estimated value of unknown parameter lies in an interval.

Confidence Interval:

PDF is divided into

- (a) acceptance region
- (b) critical region.

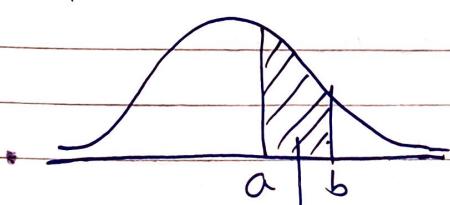


Confidence level ( $1 - \alpha$ )

Significance level ( $\alpha$ )

Significance level =  $1 - \text{Confidence level}$

$$X \sim N(\mu, \sigma^2)$$



gives the probability that r.v.  $X$  lies b/w  $a$  and  $b$

Q. Marks of an exam of 400 students are normally distributed with a population standard deviation of 5.6. A random sample with marks of 40 students is collected which has mean of 32.

Estimate the population mean with a confidence of

- (a) 80%    (b) 90%    (c) 98%

$$N = 400$$

$$n = 40$$

$$\bar{x} = 32$$

$$\sigma = 5.6$$

$$\text{population} \rightarrow \mu, \sigma^2$$

$$\text{sample} \rightarrow \bar{x}, \frac{\sigma^2}{n}$$

mean variance

$$1 - \alpha = 0.8$$

$$\alpha = 0.2$$

$$z_{\alpha/2} = 0.1$$

$$\frac{32 - \mu}{5.6 / \sqrt{40}}$$

$$\text{Population} = 400$$

Choose 4 students from each class  $\Rightarrow$

$$\text{Sample} = 4 \times 10 = 40$$

$$\bar{x} = 32 \text{ for 40 students.}$$

$$E = \underline{z_{\alpha/2} \frac{\sigma}{\sqrt{n}}}$$

According to point estimator, mean for 400 students will also be 32.

(a) Confidence level 80%.

$$1 - \alpha = 0.8 \Rightarrow \alpha = 0.2 \Rightarrow \alpha/2 = 0.1$$

Interval estimator:

(CL) Lower Confidence limit:  $\bar{x} - E$   
 (UCL) Upper Confidence limit:  $\bar{x} + E$

$$E = z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$n \rightarrow$  Sample mean

$$E = Z_{0.1} \left( \frac{S \cdot G}{\sqrt{n}} \right) = 1.13$$

↓  
1.28

$$LCL = 32 - 1.13 = 30.87$$

$$UCL = 32 + 1.13 = 33.13$$

$$\therefore \mu \in [30.87, 33.13]$$

(b) 90% Confidence

$$1-\alpha = 0.9 \Rightarrow \alpha = 0.1 \Rightarrow \alpha/2 = 0.05$$

$$E = Z_{0.05} \left( \frac{S \cdot G}{\sqrt{n}} \right) = 1.46$$

$$LCL = 32 - 1.46 = 30.54$$

$$UCL = 32 + 1.46 = 33.46$$

$$\mu \in [30.54, 33.46]$$

(c) 98% Confidence

$$1-\alpha = 0.98 \Rightarrow \alpha = 0.02 \Rightarrow \alpha/2 = 0.01$$

$$E = Z_{0.01} \left( \frac{S \cdot G}{\sqrt{n}} \right) = 2.06$$

$$\mu \in [29.94, 32.06]$$

$$x_i \sim N(\mu, \sigma^2)$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \sim N(\mu, \sigma^2/n)$$

$$E(\bar{x}) = \mu \quad \text{Var}(\bar{x}) = \sigma^2/n$$

6

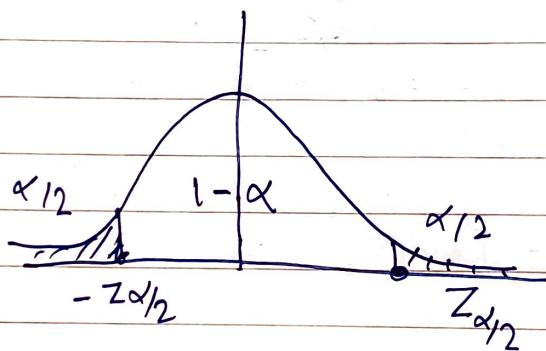
Test statistic  $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$

$$E(Z) = \frac{E(\bar{x}) - E(\mu)}{\sigma/\sqrt{n}} = \frac{\mu - \mu}{\sigma/\sqrt{n}} = 0$$

$$\text{Var}(Z) = \frac{1}{(\sigma/\sqrt{n})^2} [\text{Var}(\bar{x}) + \text{Var}(\mu)]$$

$$= \frac{1}{(\sigma/\sqrt{n})^2} \left( \sigma^2/n + 0 \right)$$

$$= 1$$



$$P[-Z_{\alpha/2} < Z < Z_{\alpha/2}] = 1 - \alpha.$$

$$-Z_{\alpha/2} < \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} < Z_{\alpha/2}$$

$$\underbrace{\bar{x} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}}_{LCL} < \mu < \underbrace{\bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}}_{UCL}$$

Linear  
Regression  
Model:

LRM:

$$\hat{\beta}_1 = \sum_{i=1}^n c_i y_i = \frac{S_{xy}}{S_{xx}}$$

$$\hat{\beta}_1 \sim N(E(\hat{\beta}_1), \text{Var}(\hat{\beta}_1))$$

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

$$c_i = \frac{x_i - \bar{x}}{S_{xx}}$$

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{S_{xx}}\right)$$

Test statistic :  $Z = \frac{\hat{\beta}_1 - \beta_1}{\sigma / \sqrt{S_{xx}}} \sim N(0, 1)$

$$E(Z) = E(\hat{\beta}_1) - E(\beta_1) = \frac{\beta_1 - \beta_1}{\sigma / \sqrt{S_{xx}}} = 0$$

$$P[-Z_{\alpha/2} < Z_1 < Z_{\alpha/2}] = 1 - \alpha$$

$$P\left[-Z_{\alpha/2} < \frac{\hat{\beta}_1 - \beta_1}{\sigma / \sqrt{S_{xx}}} < Z_{\alpha/2}\right] = 1 - \alpha.$$

$$P\left[\hat{\beta}_1 - Z_{\alpha/2} \frac{\sigma}{\sqrt{S_{xx}}} < \beta_1 < \hat{\beta}_1 + Z_{\alpha/2} \frac{\sigma}{\sqrt{S_{xx}}}\right] = 1 - \alpha.$$

$\hat{\beta}_1 \rightarrow$  Point Estimator for  $\beta_1$

$$P\left[\hat{\beta}_1 - Z_{\alpha/2} \frac{\sigma}{\sqrt{S_{xx}}} < \beta_1 < \hat{\beta}_1 + Z_{\alpha/2} \frac{\sigma}{\sqrt{S_{xx}}}\right] = 1 - \alpha.$$

$$\left[\hat{\beta}_1 - E[\hat{\beta}_1], \hat{\beta}_1 + E[\hat{\beta}_1]\right]$$

→ Interval estimators

$$E = Z_{\alpha/2} \frac{\sigma}{\sqrt{S_{xx}}}$$

$$\beta_0 \xrightarrow{\substack{PE \\ IE}} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$[\hat{\beta}_0 - E, \hat{\beta}_0 + E]$$

$$\hat{\beta}_0 \sim N(E(\hat{\beta}_0), \text{Var}(\hat{\beta}_0))$$

$$\hat{\beta}_0 \sim N(\hat{\beta}_0, \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right])$$

Finding the test statistics:

→ should follow standard normal distribution

$$Z_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}}$$

$$= \frac{\hat{\beta}_0 - \beta_0}{\theta}$$

$$\theta^2 = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]$$

$$E[Z_0] = E\left[\frac{\hat{\beta}_0 - \beta_0}{\theta}\right]$$

$$= \frac{1}{\theta} (E[\hat{\beta}_0] - \beta_0)$$

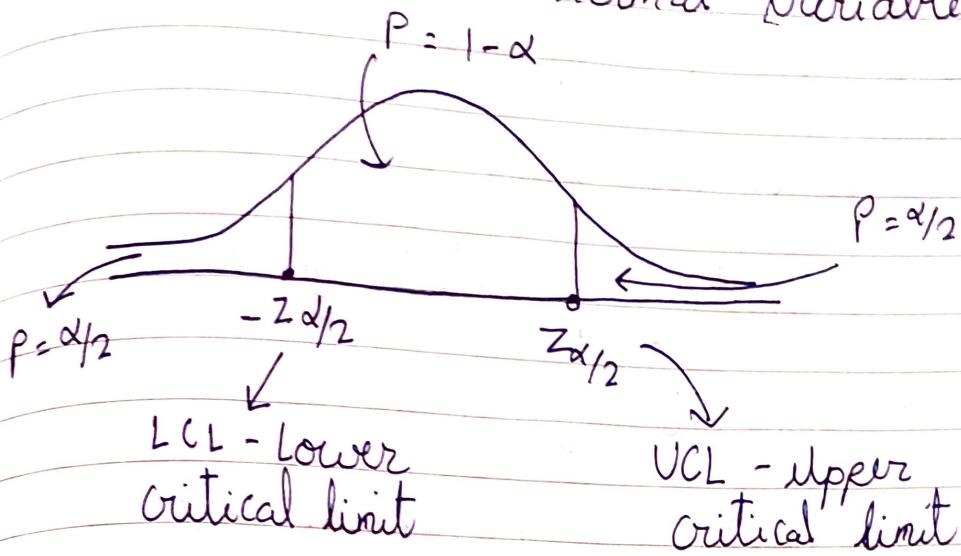
$$= 0$$

$$\text{Var}(Z_0) = \frac{1}{\theta^2} [\text{Var}(\hat{\beta}_0) + \text{Var}(\beta_0)]$$

$$= \frac{1}{\theta^2} (\theta^2 + 0)$$

$$= 1$$

Thus,  $Z_0$  is standard normal variable.



$$P[-Z_{\alpha/2} < Z_0 < Z_{\alpha/2}]$$

$$= 1 - \alpha$$

$$= P\left[-Z_{\alpha/2} < \frac{\hat{\beta}_0 - \beta_0}{\theta} < Z_{\alpha/2}\right]$$

$$= P\left[-Z_{\alpha/2} \cdot \theta < \hat{\beta}_0 - \beta_0 < Z_{\alpha/2} \cdot \theta\right]$$

$$= P\left[\underbrace{\hat{\beta}_0 - Z_{\alpha/2} \cdot \theta}_{\text{Lower Confidence}} < \beta_0 < \underbrace{\hat{\beta}_0 + Z_{\alpha/2} \cdot \theta}_{\text{Upper Confidence}}\right]$$

Lower Confidence  
limit (LCL)

Upper Confidence  
Limit (UCL)

Confidence Interval.

$$CL: 1 - \alpha \rightarrow (1 - \alpha) 100 \%$$

$$SL: 1 - CL \rightarrow (100 - CL) \%$$

$$\begin{array}{ccc} PE & \xrightarrow{\hat{\theta}} & \hat{\theta} \\ \theta & \xrightarrow{\quad} & \end{array}$$

$$IE \xrightarrow{\quad} \hat{\theta} - E \text{ to } \hat{\theta} + E$$

## # Test of Hypothesis:

- Null Hypothesis:

Something presented by  $H_0$

- Alternate Hypothesis:

- Something represented by  $H_1$ .
- $H_0$  and  $H_1$  are complementary to each other
- Alternate hypothesis always contradicts null hypothesis.

Verdict \ truth	Covid - ve Innocent	Covid + ve Guilty
RT - PCR	✓	Type 2 Error Misprediction
Covid - ve Set Free		
Covid + ve Jailed.	Type 1 error False Alarm X False + ve	✓

① One tailed test

② Two tailed test

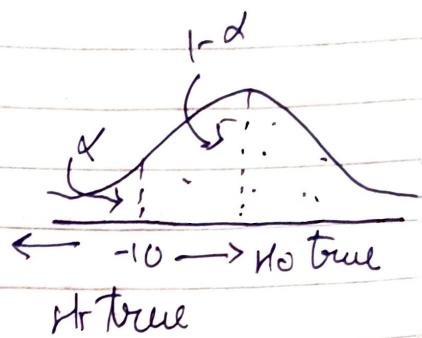


By comparing test statistics value with the critical value. (Obtained from table) we come to the conclusion whether to reject the null hypothesis.  $H_0$  or we fail to reject the null hypothesis.

Left Tail test:

$$\underline{H_0} : H_0 : \theta \geq -10$$

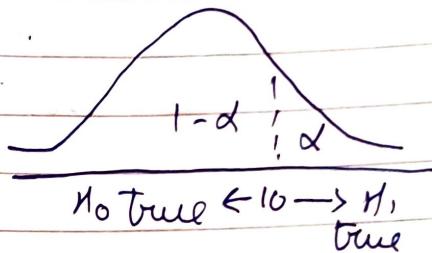
$$H_1 : \theta < -10$$



Right Tail test:

$$H_0 : \theta \leq 10$$

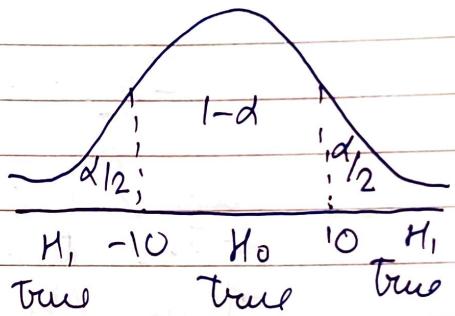
$$H_1 : \theta > 10$$



Two tailed test:

$$H_0 : |\theta| \leq 10$$

$$H_1 : |\theta| > 10$$



Manufacturer is manufacturing 1200 boxes and claims that average weight of a box is 1.84 kg.

Customer randomly chooses 64 boxes and find sample average weight as 1.88 kg. Considering that population standard deviation  $\sigma$  is 0.3 kg, use 95% confidence level and test for hypothesis that true average weight of shipment is 1.84 kg.

$$N = 1200$$

$$\mu = 1.84$$

$$n = 64$$

$$\bar{x} = 1.88$$

$$\sigma = 0.3$$

$$1-\alpha = 0.95$$

$$\alpha = 0.05$$

$$-1.96 < \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} < 1.96$$

-1.69

$$H_0: \mu = 1.84$$

0.0735

$$H_1: \mu \neq 1.84$$

$$1.8065 < \mu < 1.9535$$

$$N = 1200$$

$$n = 64$$

$$\sigma = 0.3$$

$$\bar{x} = 1.82 \text{ kg}$$

$$\sigma^2 = 0.09$$

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$100(1-\alpha) \% = 95\%$$

$$\alpha = 0.05$$

$$\alpha/2 = 0.025$$

$$Z = \frac{1.88 - 1.84}{0.3/\sqrt{64}} = 1.07$$

$$Z_{\alpha/2} = 1.96$$

We can see that  $Z$  lies between  
 $-Z_{\alpha/2}$  to  $Z_{\alpha/2}$

Thus, Null Hypothesis is true.  
 Decision failed to reject Null Hypothesis.

Inconclusive evidence to suggest that avg. weight.

## Q. Estimating the Difference in Mean of two Normal Distribution.

Let  $X_1, X_2, \dots, X_n$  be a sample of size  $n$  from a Normal Population having mean  $\mu$ , and variance  $\sigma^2$ .

Let  $Y_1, Y_2, \dots, Y_m$  be a sample of size  $m$  from a Normal population having mean  $\mu_2$  and variance  $\sigma_2^2$ .

point of concern:  $\mu_1 - \mu_2$

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} \quad x_i \sim N(\mu_1, \sigma_1^2)$$

$$E(\bar{x}) = \mu_1 \quad \text{Var}(\bar{x}) = \frac{\sigma_1^2}{n} \quad \bar{x} \sim N(\mu_1, \sigma_1^2/n)$$

$$\bar{y} = \frac{y_1 + y_2 + \dots + y_m}{m} \quad y_i \sim N(\mu_2, \sigma_2^2)$$

$$E(\bar{y}) = \mu_2 \quad \text{Var}(\bar{y}) = \frac{\sigma_2^2}{m} \quad \bar{y} \sim N(\mu_2, \sigma_2^2/m)$$

$$\mu_1 - \mu_2 \rightarrow \bar{x} - \bar{y} \quad (\text{Point estimator})$$

$$\begin{aligned} &\rightarrow \text{Interval estimator?} \\ &= \text{Point estimator} \pm E \\ &= (\bar{x} - \bar{y}) \pm E \end{aligned}$$

$$\bar{x} - \bar{y} \sim N(E(\bar{x} - \bar{y}), \text{Var}(\bar{x} - \bar{y}))$$

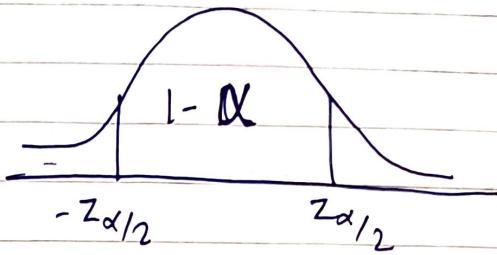
$$\underbrace{\bar{x} - \bar{y}}_{\theta} \sim N(\underbrace{\mu_1 - \mu_2}_{E(\theta)}, \underbrace{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}_{\text{Var}(\theta)})$$

Test statistic:

$$Z_{\text{test}} = \frac{\theta - E(\theta)}{\sqrt{\text{Var}(\theta)}} = \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}}$$

In exam, value of Z multiplier means the value of  $Z_{\alpha/2}$ .

$$P(-Z_{\alpha/2} < Z_{\text{test}} < Z_{\alpha/2}) = 1 - \alpha.$$



$$P\left[-Z_{\alpha/2} < \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{\sqrt{\sigma_1^2/n + \sigma_2^2/m}} < Z_{\alpha/2}\right] = 1 - \alpha$$

$$P\left[(\bar{x} - \bar{y}) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}} < \mu_1 - \mu_2 <\right.$$

$$\left.(\bar{x} - \bar{y}) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}\right] = 1 - \alpha$$

~~wo (1 - \alpha) \therefore two sided confidence interval for the estimator of difference of mean of two distribution.~~

trying  
to find  
interval  
estimator

LRM  $\rightarrow$  linear Regression Model

$\beta_1$  : Slope parameter  $y = \beta_0 + \beta_1 x$

$H_0: \beta_1 = \beta_1^k$  [ $\beta_1^k$ : known constant value]

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

Test statistics :

$$z_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2 / s_{xx}}}$$

$$E(z_1) = \frac{1}{\sqrt{\sigma^2 / s_{xx}}} [E(\hat{\beta}_1) - E(\beta_1)] = 0$$

~~$$\text{Var}(z_1) = \frac{1}{\sigma^2 / s_{xx}}$$~~

$$\begin{aligned} \text{Var}(z_1) &= \frac{1}{\sigma^2 / s_{xx}} (\text{Var}(\hat{\beta}_1) - \text{Var}(\beta_1)) \\ &= \frac{\sigma^2 / s_{xx}}{\sigma^2 / s_{xx}} \\ &= 1 \end{aligned}$$

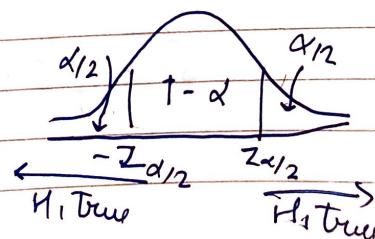
$z_1$  (test statistic) can be used to test hypothesis, significance of slope parameter.

$$H_1: \beta_1 \neq \beta_1^*$$

Decision Rule : When to reject  $H_0$  ?

Ans:  $|z_1| > z_{\alpha/2}$  or  $z_1 < -z_{\alpha/2}$

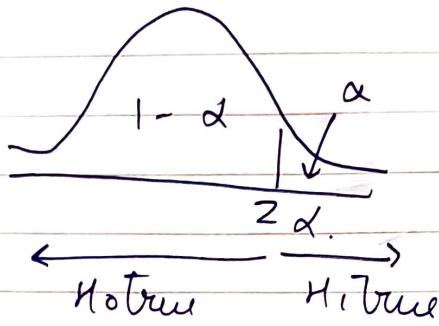
$$|z_1| > z_{\alpha/2}$$



$$H_1: \beta_1 > \beta_1^k$$

Decision Rule:  
Reject  $H_0$

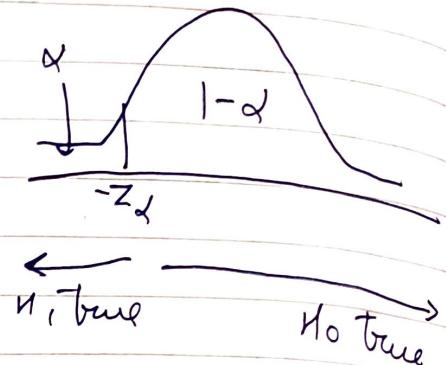
$$Z_1 > Z_\alpha$$



$$H_1: \beta_1 < \beta_1^k$$

Decision Rule:  
Reject  $H_0$

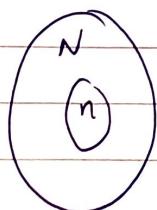
$$Z_1 < -Z_\alpha$$



Normal distribution  
( $\mu, \sigma^2$ )

Population parameter

$$x_i \sim N(\mu, \sigma^2)$$



Population mean

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

population variance.

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Sample is considered  
( $\bar{x}, s^2$ )

$$\begin{aligned}
 E(s^2) &= E \left\{ \frac{1}{n} \sum_{i=1}^n (x_i^2 + \bar{x}^2 - 2x_i\bar{x}) \right\} \\
 &= E \left\{ \frac{1}{n} \left( \sum_{i=1}^n x_i^2 + n\bar{x}^2 - 2n\bar{x} \sum_{i=1}^n x_i \right) \right\} \\
 &= E \left\{ \frac{1}{n} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \right\} \underbrace{\quad}_{2n\bar{x}^2} \\
 &= E \left\{ \frac{1}{n} \sum_{i=1}^n x_i^2 \right\} - E \{ \bar{x}^2 \} \\
 &= \frac{1}{n} \sum_{i=1}^n E \{ x_i^2 \} - E \{ \bar{x}^2 \}.
 \end{aligned}$$

$$\text{Var}[x_i] = E[x_i^2] - (E[x_i])^2$$

$$\sigma^2 = E \{ x_i^2 \} - \mu^2$$

$$E \{ x_i^2 \} = \sigma^2 + \mu^2$$

$$\text{Var}(\bar{x}) = E(\bar{x}^2) - [E(\bar{x})]^2$$

$$\sigma^2/n = E(\bar{x}^2) - \mu^2.$$

$$E(\bar{x}^2) = \frac{\sigma^2}{n} + \mu^2$$

$$\begin{aligned}
 E(s^2) &= \frac{1}{n} \sum_{i=1}^n E \{ x_i^2 \} - E \{ \bar{x}^2 \} \\
 &= \frac{1}{n} \sum_{i=1}^n (\sigma^2 + \mu^2) - \left( \frac{\sigma^2}{n} + \mu^2 \right) \\
 &= \sigma^2 + \mu^2 - \frac{\sigma^2}{n} - \mu^2 \\
 &= \frac{\sigma^2(n-1)}{n}
 \end{aligned}$$

$$\hat{\sigma}^2 = \left(\frac{n}{n-1}\right)s^2$$

If  $\sigma^2$  is unknown, from above expression what is estimator of  $\sigma^2$ ??

$$E\left(\left(\frac{n}{n-1}\right)s^2\right) = \sigma^2$$

Can we say  $s^2$  is an unbiased estimator of  $\sigma^2$ .

If  $s^2$  was unbiased estimator for  $\sigma^2$ , we could have

$\hat{\sigma}^2 = \bar{s}^2$  but here, it's not true  
So,  $s^2$  is biased estimator.

$(\mu, \sigma^2)$

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$\begin{aligned} S_{n-1}^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \\ E\{S^2\} &= \sigma^2 \end{aligned}$$

Confidence interval (CI)  $100[1-\alpha]\%$  for  $\mu$

$$\left[ \bar{x} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right]$$

Interval estimator for population parameter  $\mu$  when population variance  $\sigma^2$  was known.

When  $\sigma^2$  itself is unknown and you need to find Interval estimator for  $\mu$ .

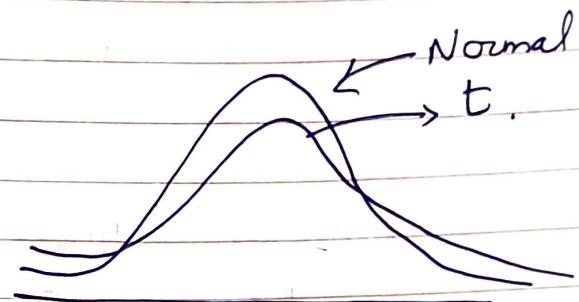
Test statistics:

$$t = \frac{\bar{x} - u}{s/\sqrt{n}} \sim t_{n-1}$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

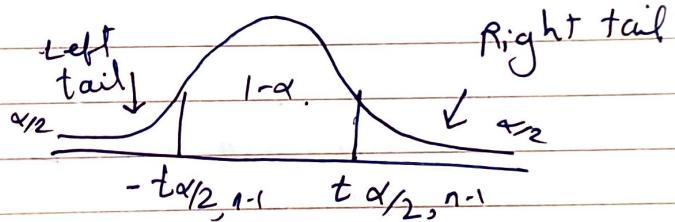
freedom as you used  $s$  instead of  $\sigma$ .

t-distribution:



Degree of freedom  $\uparrow \rightarrow$  t will closely match normal.  
(close margin)

$$t = \frac{\bar{x} - u}{s/\sqrt{n}}$$



Aim: To find interval estimator for unknown population mean parameter  $u$ , when  $\sigma$  is unknown.

100(1- $\alpha$ )% confidence, formulate the interval estimator for  $u$ .

$$P[-t_{\alpha/2, n-1} < t < t_{\alpha/2, n-1}] = 1-\alpha$$

$$P \left[ -t_{\alpha/2, n-1} < \frac{\bar{x} - \mu}{S/\sqrt{n}} < t_{\alpha/2, n-1} \right] = 1 - \alpha$$

$$P \left[ \bar{c} - (t_{\alpha/2, n-1}) \frac{s}{\sqrt{n}} < \mu < \bar{c} + (t_{\alpha/2, n-1}) \frac{s}{\sqrt{n}} \right]$$

$$E = \sum_{j=1}^n t_{\alpha/2} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{x})^2 = \frac{s^2}{n}$$

~~X~~ Intercept Parameters:

$$\beta_0 :$$

$$H_0 : \beta_0 = \beta_0^*$$

$$H_1 : \beta_0 \neq \beta_0^*$$

Interval estimate  
 $(1-\alpha) 100\%$   
 CI

Test statistics  $Z_0 = \frac{\hat{\beta}_0 - \beta_0^*}{\sqrt{\sigma^2 \left( \frac{1}{n} + \frac{\bar{c}^2}{S_{xx}} \right)}}$

standard normal.

$$E(Z_0) = \frac{1}{\theta} [ E(\hat{\beta}_0) - E(\beta_0^*) ] \\ = \frac{1}{\theta} [ \beta_0^* - \beta_0^* ] \\ = 0$$

$$\text{Var}(Z_0) = \frac{1}{\theta^2} (\text{Var}(\hat{\beta}_0) - \text{Var}(\beta_0^*)) \\ = \frac{1}{\theta^2} \cdot \theta^2 = 1$$

When  $\sigma^2$  is KNOWN:

For parameter  $\beta_1$ :

$$CI: \left[ \hat{\beta}_1 - Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}}, \hat{\beta}_1 + Z_{\alpha/2} \sqrt{\frac{\sigma^2}{S_{xx}}} \right]$$

for parameter  $\beta_0$ :

$$CI: \left[ \hat{\beta}_0 - Z_{\alpha/2} \sqrt{\sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}, \hat{\beta}_0 + Z_{\alpha/2} \sqrt{\sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \right]$$

When  $\sigma^2$  is UNKNOWN (so we use t-distribution)

$$t_1 = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}} \sim t(n-2)$$

$$H_0: \beta_1 = \beta_1^K$$

Reject  $H_0$  when  $t_1$

$$\text{We know that: } E\left[\hat{\sum}_{i=1}^n \frac{SS_{res}}{n-2}\right] = \sigma^2$$

$$\text{So, } \hat{\sigma}^2 = \frac{SS_{res}}{n-2}$$

$$\therefore \hat{\sigma} = \sqrt{\frac{SS_{res}}{n-2}}$$

↳ derived earlier.

$$\rightarrow \text{Confidence CI} : \left[ \hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}, \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \right]$$

No, for  $\beta_0$ :

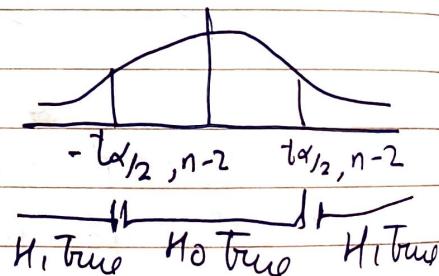
$$t_0 = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim t(n-2)$$

$$\hat{\sigma}^2 = \frac{SS_{res}}{n-2}$$

$$\text{CI} : \left[ \hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}, \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \right]$$

\*  $H_0: \beta_0 = \beta_0^c$   
 $H_1: \beta_0 \neq \beta_0^c$

Reject  $H_0$  when  $|t_0| \geq |t_{\alpha/2}|$



test statistic:  $t_0 = \frac{\hat{\beta}_0 - \beta_0^c}{\sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$

# Chi-Square Distribution:

Consider 'n' dataset:

$$\sum_{i=1}^n \varepsilon_i^2 = \varepsilon_1^2 + \varepsilon_2^2 + \dots + \varepsilon_n^2$$

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

→ Not in standard Normal form

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

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

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$= \sum [(\beta_0 + \beta_1 x_i + \varepsilon_i) - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

$$SS_{res} = \sum_{i=1}^n [(\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)x_i + \varepsilon_i]^2$$

$$\frac{SS_{res}}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n [(\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)x_i + \varepsilon_i]^2$$

$$\frac{\sum_{i=1}^n \varepsilon_i^2}{\sigma^2} = \frac{\varepsilon_1^2}{\sigma^2} + \frac{\varepsilon_2^2}{\sigma^2} + \dots + \frac{\varepsilon_n^2}{\sigma^2}$$

$$\frac{\varepsilon_i^2}{\sigma^2} \sim \chi^2 \text{ follows chi sq. distribution}$$

with 'n' degree of freedom.

$$\frac{SS_{res}}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n [(\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)x_i + \varepsilon_i]^2$$

2 degrees of freedom are lost as we are using estimated value  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .

$$\therefore \frac{SS_{res}}{\sigma^2} \sim \chi^2(n-2)$$

$$H_1: \lambda < \theta$$

(we reject  $H_0$  if  $P(x \leq x_c) = \alpha/2$ )

$$E\left[\frac{SS_{res}}{\sigma^2}\right] = n-2 \Rightarrow E\left[\frac{SS_{res}}{n-2}\right] = \sigma^2$$
$$\therefore \hat{\sigma}^2 = \frac{SS_{res}}{n-2}$$

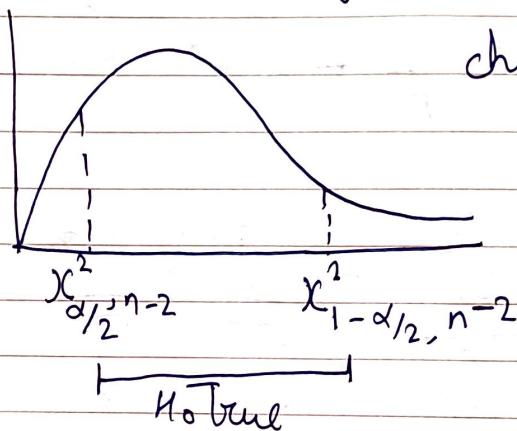
- Test of hypothesis for  $\sigma^2$ .

$$H_0: \sigma^2 = \sigma_0^2$$

$$H_1: \sigma^2 \neq \sigma_0^2$$

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

$$\text{Test Statistic: } C_0 = \frac{SS_{res}}{\sigma^2} \sim \chi^2(n-2)$$



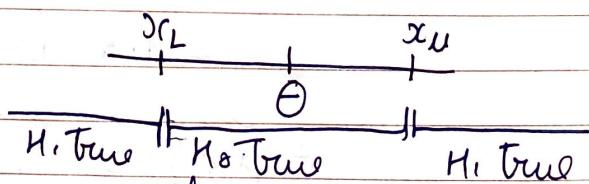
Decision Rule:

Reject  $H_0$  if  $C_0 < x_{\alpha/2, n-2}^2$  or  $C_0 > x_{1-\alpha/2, n-2}^2$

\* Poisson Distribution Hypothesis Testing:  
 $\hookrightarrow f_x(x) = e^{-\lambda} \frac{\lambda^x}{x!} \quad E(x) = \lambda$

$$H_0: \lambda = \theta$$

$$H_1: \lambda \neq \theta$$



(with some energy margin)

We reject  $H_0$  if  $P(x \leq x_c) = \alpha/2$  or  $P(x \geq x_u) = \alpha/2$

$H_1: \lambda > 0$  (We reject  $H_0$  if  $P(x \geq x_u) = \alpha/2$ )

# Information Theory

Information → Engineering concept

$$\begin{array}{l} \text{(1) } T \underset{?}{\underline{H}} \underset{T}{\underline{\text{?}}} \Rightarrow \{E\} \\ \text{(2) } T \underset{?}{\underline{C}} \underset{T}{\underline{\text{?}}} \Rightarrow \{A, U, O\} \text{ (Guess List)} \\ \text{(3) } T \underset{?}{\underline{A}} \underset{T}{\underline{\text{?}}} \Rightarrow \{G, E, H, P, M, F, \dots, Y\} \end{array}$$

The Guess list keeps on increasing.

$$\begin{array}{ll} \text{GO } \underset{?}{\underline{?}} & \text{AL} \\ \text{OL} & \\ \text{AT} & \end{array} \quad \begin{array}{ll} \text{GOA } \underset{?}{\underline{?}} & \text{T} \\ & \text{L} \end{array}$$

Randomness in my choosing is MORE.

Randomness further decreases.

Information is reduction of uncertainty (randomness)

## # Entropy / Information

GNN      RNN      ~~Not~~ Neutral Network

Randomness ↑  
information ↑  
probability ↓

## # Modeling the Uncertainty :

Prime Numbers less than 10  $\Rightarrow \{2, 3, 5, 7, 4\}$   
Uncertainty = 4

$C ? T \Rightarrow \{A, U, O\}$   
Uncertainty = 3

# Measure of Uncertainty ??

Must be modelled mathematically  
 $\log_2( )$

$\mathcal{X}$  : Set of possible answers.

$|\mathcal{X}|$  : No. of possible answers

$$\mathcal{X} = \{2, 3, 5, 7\}$$

$$|\mathcal{X}| = 4$$

$x_i$  :  $i^{\text{th}}$  answer

$P(x_i)$  : Probability that  $i^{\text{th}}$  the  $x_i^{\text{th}}$  answer is correct.

Information :  $I(x_i) = \log_2 \left( \frac{1}{P(x_i)} \right)$

$$I(x_i) = \log_2 \left( \frac{1}{P(x_i)} \right)$$

(event)  $E_1$  : The sun rises from East

$E_2$  : It will be snowfall in Ahmedabad in the month of May/June.

$$P(E_1) = 1$$

$$P(E_2) \approx 0$$

$$I(E_1) = \log_2 \left( \frac{1}{P(E_1)} \right)$$

$$= \log_2 1$$

$$= 0$$

(No information)

$$I(E_2) = \log_2 \left( \frac{1}{P(E_2)} \right)$$

$$= \log_2 (\infty)$$

$$= \infty$$

(Highly informative)

Entropy ( $\mathcal{H}$ )

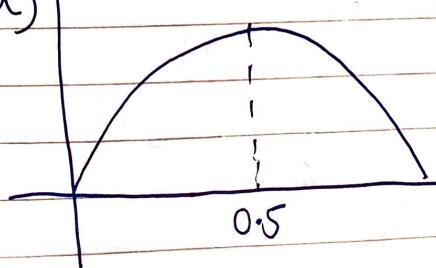
$$H(\mathcal{X}) = \sum_{i=1}^{K \rightarrow |\mathcal{X}|} P(x_i) \log P(x_i)$$

$$H(\mathcal{X}) = \sum_{i=1}^K P(x_i) \cdot \log_2 \left( \frac{1}{P(x_i)} \right)$$

$\Rightarrow$  Uncertainty / Guess-list may NOT be equally likely / Probable.

$\Rightarrow$  Uncertainty can ~~be~~ vary based on knowledge

$$H(\mathcal{X})$$



(When all events are equally likely)

More Entropy

Most uncertainty

2<sup>nd</sup> years



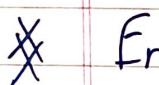
10 students are equally likely to get gold medal.

Prediction: Full of uncertainty.

2019



Uncertainty = 0

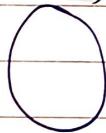


Entropy

Venn Diagram

X: Thaltej traffic

H(X)



Y: traffic in Marin Drive

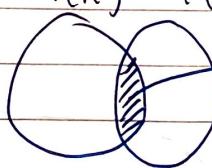
H(Y)



X and Y are independent, so  $H(X)$  has nothing to do with  $H(Y)$ .

X: traffic in Thaltej Metrostation

H(X) H(Y)



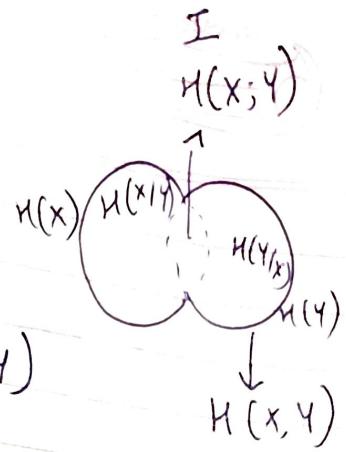
Y: traffic in Palladium Mall.

Mutual Information

There is some dependency between X and Y.

$$I(X) = H(X) - H(X|Y)$$

$$= H(Y) - H(Y|X)$$



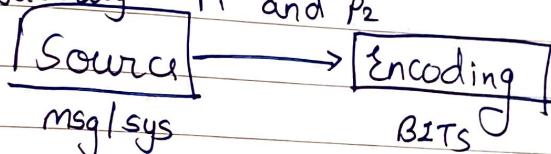
$$I(X) = H(X) + H(Y) - H(X, Y)$$



$$I(X) = H(X, Y) - H(X|Y) - H(Y|X)$$



If we have message  $m_1$  and  $m_2$ , independent with probability  $p_1$  and  $p_2$

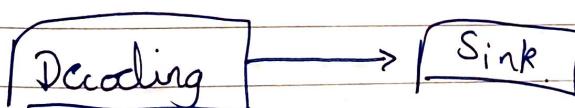


$$P = p_1 \cdot p_2$$

$$I = \log_2 \left( \frac{1}{P} \right) = \log_2 \left( \frac{1}{p_1 \cdot p_2} \right)$$

$$I = \log_2 \left( \frac{1}{p_1} \right) + \log_2 \left( \frac{1}{p_2} \right)$$

$$I = I_1 + I_2$$



- If there are 'm' messages from a source and all the messages are equally likely

$$m = 2^N = 16 \Rightarrow N = 4$$

$$P_i = \frac{1}{m} \quad 1 \leq i \leq 16$$

$$H = \sum_{i=1}^{m=16} P_i \log_2 \left( \frac{1}{P_i} \right)$$

$$= \sum_{i=1}^m \frac{1}{m} \log_2(m)$$

$$= \frac{m}{m} \log_2(m)$$

$$= \log_2(2^N)$$

$$= N$$

$$M = \{m_1, m_2, \dots, m_{16}\}$$

$\downarrow_{16} \quad \downarrow_{16} \quad \downarrow_{16}$

$$H = 4 \text{ bits / Sym}^{\rightarrow \text{msg}}$$

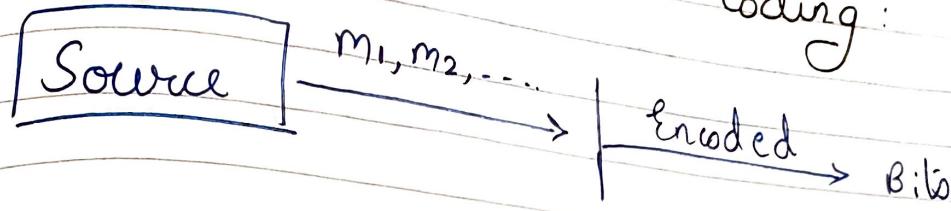
$$I_i = \log_2 \left( \frac{1}{P_i} \right) \quad i = 1, \dots, N$$

$\downarrow$   
Information

$$H = \underbrace{P_1 \cdot I_1 + P_2 \cdot I_2 + \dots + P_N \cdot I_N}_{\text{Average Information}}$$

$$H = \sum_i P_i \log_2 \left( \frac{1}{P_i} \right) = - \sum_i P_i \log_2 P_i$$

# # Encoding Techniques | Source Coding:



## Shannon - Fano Encoding:

- It is an entropy encoding technique for data compression.
- It assigns a code to each symbol/message based on their probability.
- Variable length coding scheme.
- The source symbols are split into two subsets "0" and "1" with probability as close to one-half. Then each of the sub-sets is split into two and so on until there is only one symbol/message left in the set.

Example:

~~Find the codewords~~

Example:

Source  $S = \{S_1, S_2, S_3, S_4\}$

Probability of occurrence of each source message

Message	Prob.	1 <sup>st</sup> stage	2 <sup>nd</sup> stage	3 <sup>rd</sup> stage	Codeword
$S_1$	$\frac{1}{2}$	1	0		0
$S_2$	$\frac{1}{4}$		1	0	10
$S_3$	$\frac{1}{8}$			1	110
$S_4$	$\frac{1}{8}$			1	111

$$\begin{aligned}
 H &= \sum_i p_i \log_2 \frac{1}{p_i} \\
 &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} \\
 &\quad - \frac{1}{8} \log_2 \frac{1}{8} \\
 &= \frac{1}{2} + \frac{1}{2} + 3 \cdot \frac{1}{4} \\
 &= \frac{7}{4}
 \end{aligned}$$

Avg. Codeword length:  $= \sum_{i=1}^4 p_i l_i$

$\hat{H} | L | \bar{L}$

$$\begin{aligned}
 &= P_1 l_1 + P_2 l_2 + P_3 l_3 + P_4 l_4 \\
 &= \left(\frac{1}{2} \times 1\right) + \left(\frac{1}{4} \times 2\right) + \left(\frac{1}{8} \times 3\right)^2 \\
 &= \frac{7}{4}
 \end{aligned}$$

$$n[\text{Efficiency}] = \frac{\text{Entropy}}{\text{Avg. Codeword length}}$$

$$= \frac{7/4}{7/4}$$

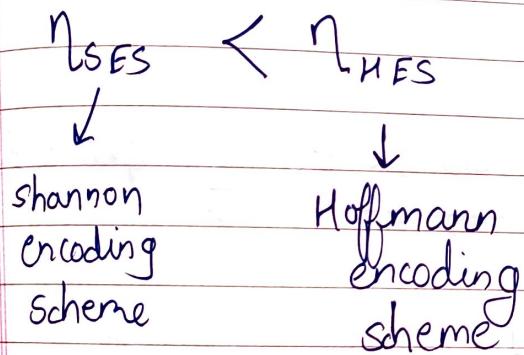
$$= 1$$

discrete  
2nd

$$\text{Redundancy [Re]} = 1-n = 0.$$

Efficiency = 100% when both the subsets are equally probable in the 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup> stage.

$n < 1$  when subset probability are not equal.



# Example:

$$M = \{m_1, m_2, \dots, m_7\}$$

$$P = \{0.4, 0.2, 0.12, 0.08, 0.08, 0.08, 0.04\}$$

By shannon's encoding scheme,

message	Prob	Code word
$m_1$	0.4	10
$m_2$	0.2	0
$m_3$	0.12	1
$m_4$	0.08	0
$m_5$	0.08	1
$m_6$	0.08	0
$m_7$	0.04	1

$$H = 0.4 \log_2 \left( \frac{1}{0.4} \right) + \dots + 0.04 \log_2 \left( \frac{1}{0.04} \right)$$

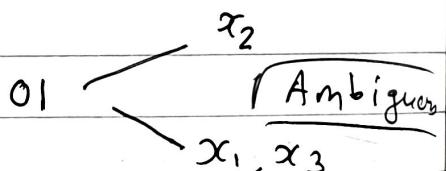
$$\bar{L} = 0.4(1) + 0.2 \times 3 + \dots$$

$$n = \frac{H}{\bar{L}} < 1$$

↓  
There will be redundancy in SEC.

$$X = \{x_1, x_2, x_3, x_4\}$$

$$\begin{matrix} \downarrow & \downarrow & \downarrow & \downarrow \\ 0 & 01 & 1 & 10 \end{matrix}$$



→ Shannon Encoding scheme is prefix-free code

Kraft's Inequality / Theorem:  $\sum_{i=1}^n 2^{-l_i} \leq 1$

If this condition is satisfied, we obtain codes which are prefix-free code.

Aug encoding scheme that does not have prefix subject, is called prefix free code.

Refer shannon example 1<sup>st</sup> question:

{S<sub>1</sub>, S<sub>2</sub>, S<sub>3</sub>, S<sub>4</sub>}.

$$l_1 = 1 \quad l_2 = 2 \quad l_3 = 3 \quad l_4 = 3$$

$$\begin{aligned} \sum_{i=1}^4 2^{-l_i} &= 2^{-1} + 2^{-2} + 2^{-3} + 2^{-3} \\ &= 2^{-1} + 2^{-2} + 2^{-2} \\ &= 2^{-1} + 2^{-1} \\ &= 1 \end{aligned}$$

# Huffman Encoding:

(Shannon - Fano)

↓  
Robert Fano  
Prof. at MIT

Amnon Shashua,

Israeli

Prof. of  
Hebrew University  
for salam.

Intel

# Huffman Coding Scheme:

- ① The source symbols are arranged in order of decreasing probability. Then the two of lowest probabilities are assigned bits 0 and 1. and encode them
- ② Group the last 2 symbols/messages and move the grouped symbol as high as possible.

- ③ Repeat the step (2) until only 2 symbols are left.
- ④ Codeword for each symbol is found by moving in reverse direction.

$$\text{Q. } S = \{S_1, S_2, S_3, S_4, S_5\}$$

$$P = \{0.4, 0.2, 0.2, 0.1, 0.1\}$$

Symbol Prob.

$S_1$	0.4	0.4	0.4	0.4
$S_2$	0.2	0.2	0.2	0.4
$S_3$	0.2	0.2	0.2	0.2
$S_4$	0.1	0.2	0.2	0.2
$S_5$	0.1	0.2	0.2	0.2

$$S_1 \quad 00 \quad l_1 = 2$$

$$S_2 \quad 10 \quad l_2 = 2$$

$$S_3 \quad 11 \quad l_3 = 2$$

$$S_4 \quad 010 \quad l_4 = 3$$

$$S_5 \quad 011 \quad l_5 = 3$$

$$3 \times 2^{-2} + 2 \times 2^{-3}$$

$$\frac{3}{4} + \frac{1}{4}$$

$$H = \sum_{i=1}^5 P_i \log_2 \left( \frac{1}{P_i} \right)$$

$$\bar{L} = \sum_{i=1}^5 P_i l_i$$

prefix  
free  
code.

$$n = \frac{H}{\bar{L}} \quad R_e = 1 - n$$

efficiency (Hoffmann)

better than efficiency (Shannon)

discrete

tribe

face

$$n = \frac{H}{\sum \log_2 M}$$

$$M = 2$$

$$\log_2 2 = 1 \quad (\text{Binary})$$

10  
010  
1