

Name- Suresh Kumar Prajapati  
Class- M.Sc. (Master in data science and  
Applied statistics).

Enroll No- CUSB230222008

Assignment work

Define the coefficient of determination of  $R^2$  and adjusted  $R^2$  for multiple linear regression :-

Let  $R$  be the multiple correlation coefficient between  $y$  and  $x_1, x_2, \dots, x_k$ . Then square of multiple correlation coefficient ( $R^2$ ) is called a coefficient of determination.

The value of  $R^2$  commonly describes how well the sample regression line fits the observed data.

This is also treated as a measure of goodness of fit of the model.

Assuming that the intercept is present in the model as

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_k x_{ik} + u_i$$

then,  $i = 1, 2, \dots, n$ .

$$R^2 = 1 - \frac{e'e}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{SS_{res}}{SS_T} = \frac{SS_{reg}}{SS_T}$$

where,  $SS_{res}$  = sum of squares due to residuals,

$SS_T$  = Total sum of squares.

$SS_{reg}$  = sum of squares due to regression.

$$\therefore e'e = y' [I - X(X'X)^{-1}X']y = y'H'y$$

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

where  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} l'y$  with  $l = (1, 1, \dots, 1)'$

$$y = (y_1, y_2, \dots, y_n)'$$

So,  $R^2 = 1 - \frac{y'H'y}{y'Hy}$

The limits of  $R^2$  are 0 and 1, i.e,

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

$R^2 = 0$  indicates the poorest fit of the model.

$R^2 = 1$  indicates the best fit of the model.

$R^2 = 0.95$  indicates that 95% of the variation in  $y$  is explained by  $R^2$ .

In simple word, the model is 95% good.

Remark:- Any other value of  $R^2$  between 0 and 1 indicates the adequacy of the fitted model.

Adjusted  $R^2$  :-

If more explanatory variables are added to the model, then  $R^2$  increases. In case the variables are irrelevant, then  $R^2$  will still increase and gives an overly optimistic picture.

With a purpose of correction in the overly optimistic picture, adjusted  $R^2$ , denoted as  $\bar{R}^2$  or  $\text{adj } R^2$  is used which is defined as

$$R^2 = 1 - \frac{SS_{res}/(n-k)}{SS_T/(n-1)}$$

$$\boxed{R^2 = 1 - \left( \frac{n-1}{n-k} \right) (1 - R^2)}$$

$(n-k)$  &  $(n-1)$  are degree of freedom associated with the distribution of  $SS_{res}$  and  $SS_T$  respectively

### Limitation:-

- ① If the constant term is absent in the model, then  $R^2$  can not be defined. In such cases  $R^2$  can be negative. Some ad-hoc measures based on  $R^2$  for regression line through origin have been proposed in the literature.

Reason that why  $R^2$  is valid only in linear models with intercept term:-

$$\begin{aligned} TSS &= \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= SS_{reg} + SS_{res} - 2\bar{y}\mathbf{l}'\mathbf{e} \quad (\text{because } \mathbf{X}'\mathbf{e} = 0) \end{aligned}$$

Fisher Cochran theorem requires

$TSS = SS_{reg} + SS_{res}$  to hold true in the context of analysis of variance and further to define by  $R^2$ . true only if  $\mathbf{e}'\mathbf{e} = 0$ . i.e.  $\mathbf{e}'\mathbf{e} = \mathbf{e}'(\mathbf{y} - \hat{\mathbf{y}}) = 0$  which is possible only when there is an intercept term in the model.

Some key points:-

- ① A simple linear regression model with intercept term  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$  ( $i = 1, 2, \dots, n$ )

where the parameters  $\beta_0$  and  $\beta_1$  are estimated

$$\text{at. } b_0 = \bar{y} - \beta_1 \bar{x}, \quad b_1 = \frac{S_{xy}}{S_{xx}}$$

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

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

$$l'e = \sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - \hat{y}_i)$$

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

$$= \sum_{i=1}^n (y_i - \bar{y} + b_1 \bar{x} - b_1 x_i)$$

$$= \sum_{i=1}^n [(y_i - \bar{y}) - b_1 (x_i - \bar{x})]$$

$$= \sum_{i=1}^n (y_i - \bar{y}) - b_1 \sum_{i=1}^n (x_i - \bar{x})$$

$$= 0$$

In case of multiple regression.

$y = \beta_0 l + X\beta + \epsilon$  where parameters  $\beta_0$  and  $\beta$  are estimated as  $\hat{\beta}_0 = \bar{y} - b \bar{x}$  &  $b = (X'X)^{-1} X' y$  respectively

$$l'e = l'(y - \hat{y})$$

$$= l'(y - \hat{\beta}_0 - Xb)$$

$$= l'(y - \bar{y} + \bar{X}b - Xb)$$

$\sum (y - \bar{y}) + k(x - \bar{x})/b = 0$

②  $R^2$  is sensitive to extreme values, so  $R^2$  lacks robustness.

③  $R^2$  always increase with an increase in the number of explanatory variable in the model. This indicates that the model is getting better, which is not really correct.

Relationship of analysis of variance test and coefficient of determination :-

Assuming the  $\beta_1$  to be an intercept term, then for  $H_0: \beta_2 : \beta_3 = \dots = \beta_k = 0$ , then F-statistic in analysis of variance test is.

$$F = \frac{MS_{reg}}{MS_{res}} = \frac{(n-k) SS_{reg}}{(k-1) SS_{res}} = \frac{(n-k) SS_{reg}}{(k-1)(SS_T - SS_{reg})}$$

$$F = \frac{(n-k) R^2}{(k-1)(1-R^2)}$$

∴  $R^2$  is the coefficient of determination.

∴ F and  $R^2$  are closely related.

when  $R^2 = 0$ , then  $F = 0$

when  $R^2 = 1$ , then  $F = \infty$

So both F and  $R^2$  vary directly.

If F is highly significant, it implies that we can reject  $H_0$  i.e. y is linearly related to  $X$ 's.

Prove that OLS is the best linear unbiased estimator of beta for M.R. :- Gauss-Markov Theorem:-

Proof:- The OLS of  $\beta$  is

$$b = (X'X)^{-1}X'y$$

which is linear function of y. Consider the arbitrary linear estimator

$$b^* = a'y$$

of linear parametric function  $\beta'$  where the elements of  $a$  are arbitrary constants.

Then for  $b^*$

$$E(b^*) = E(a'y) = a'X\beta$$

and so  $b^*$  is an unbiased estimator of  $\beta'$  when

$$E(b^*) = a'X\beta = \beta'$$

$$\Rightarrow a'X = \beta$$

Since we wish to consider only those estimators that are linear and unbiased, so we restrict ourselves to those estimators for which

$$a'X = \beta$$

Further,

$$\text{Var}(\alpha'y) = \alpha' \text{Var}(y) \alpha = \sigma^2 \alpha' \alpha$$

$$\text{Var}(J'b) = J \text{Var}(b) J$$

$$= \sigma^2 \alpha' \gamma(X'X)^{-1} X' \alpha$$

Consider,

$$\text{Var}(\alpha'y) - \text{Var}(J'b) = \sigma^2 [\mathbf{I} - \gamma(X'X)^{-1}\gamma] \alpha$$

$$= \sigma^2 \alpha' [\mathbf{I} - H] \alpha$$

$\because (\mathbf{I} - H)$  is a true semi-definite matrix, so,

$$\text{Var}(\alpha'y) - \text{Var}(J'b) \geq 0$$

This reveals that if  $b^*$  is any linear unbiased estimator then its variance must be no smaller than that of  $b$ .

Consequently  $b$  is the best linear unbiased estimator, where 'best' refers to the fact that  $b$  is efficient with the class of linear unbiased estimators.

### Likelihood Ratio test :-

The likelihood ratio test statistic is:

$$\lambda = \frac{\max L(\beta, \sigma^2 | y, X)}{\max L(\beta, \sigma^2 | y, X, R\beta = r)} = \frac{L(\lambda)}{L(w)}$$

where  $\Lambda$  is the whole parametric space  
and  $w$  is the sample space.

If both the likelihoods are maximized, one constrained, and the other unconstrained then the value of the unconstrained will not be smaller than the value of constrained.

Hence  $\lambda \geq 1$ .

First, we discuss the likelihood ratio test for a more straightforward case when

$$R = I_k \text{ and } r = \beta_0 \text{ i.e. } \beta = \beta_0$$

This will give us a better and detailed understanding of the minor details and then we generalized it for  $R\beta = r$ , in general.

Likelihood Ratio test for  $H_0: \beta = \beta_0$

Let the null hypothesis related to  $k \times 1$  vector  $\beta$  is

$$H_0: \beta = \beta_0$$

where  $\beta_0$  is specified by the investigator.

The element of  $\beta_0$  can take on any value including zero. The concerned alternative hypothesis is.

$$H_1: \beta \neq \beta_0$$

$$\therefore \epsilon \sim N(0, \sigma^2 I) \text{ in } y = X\beta + \epsilon \text{ s.o.}$$

$y \sim N(X\beta, \sigma^2 I)$ . Thus the whole parameteric space and sample space are  $\mathcal{R}$  and  $\omega$  respectively given by

$$\mathcal{R} = \left\{ (\beta, \sigma^2) : -\infty < \beta_i < \infty, \sigma^2 > 0 \quad i=1,2,\dots,k \right\}$$

$$\mathcal{W} : \left\{ (\beta, \sigma^2) : \beta = \beta_0, \sigma^2 > 0 \right\}$$

The unconstrained likelihood under  $\mathcal{R}$

$$L(\beta, \sigma^2 | y, X) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp \left[ -\frac{1}{2\sigma^2} (y - X\beta)' (y - X\beta) \right]$$

This is maximized over  $\mathcal{R}$  when,

$$\hat{\beta} = (X'X)^{-1} X' y$$

$$\hat{\sigma}^2 = \frac{1}{n} (y - X\hat{\beta})'(y - X\hat{\beta})$$

where  $\hat{\beta}$  and  $\hat{\sigma}^2$  are maximum likelihood estimates of  $\beta$  and  $\sigma^2$  which are values maximizing the likelihood function

$$\hat{L}(\mathcal{R}) = \max L(\beta, \sigma^2 | y, X)$$

$$= \frac{1}{\left[ \frac{2\pi}{n} (y - X\hat{\beta})'(y - X\hat{\beta}) \right]^{\frac{n}{2}}} \exp \left[ -\frac{(y - X\hat{\beta})'(y - X\hat{\beta})}{2\frac{(y - X\hat{\beta})'(y - X\hat{\beta})}{n}} \right]$$

$$= \frac{n^{\frac{n}{2}} \exp \left( -\frac{n}{2} \right)}{\left( \frac{2\pi}{n} \right)^{\frac{n}{2}} [(y - X\hat{\beta})'(y - X\hat{\beta})]^{\frac{n}{2}}}$$

The constrained likelihood under  $\mathcal{W}$  is

$$\hat{L}(\mathcal{W}) = \max L(\beta, \sigma^2 | y, X, \beta = \beta_0)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp \left[ -\frac{1}{2\sigma^2} (y - X\beta_0)' (y - X\beta_0) \right]$$

$$\hat{L}(\omega) = \frac{n^{n/2} \exp(-\frac{1}{2})}{(2\pi)^{n/2} [(y - X\beta_0)'(y - X\beta_0)]^{n/2}}$$

The likelihood ratio is:

$$\begin{aligned}\frac{\hat{L}(R)}{\hat{L}(\omega)} &= \left( \frac{\frac{n^{n/2} \exp(-\frac{1}{2})}{(2\pi)^{n/2} [(y - X\tilde{\beta})'(y - X\tilde{\beta})]^{n/2}}}{\frac{n^{n/2} \exp(-\frac{1}{2})}{(2\pi)^{n/2} [(y - X\tilde{\beta}_0)'(y - X\tilde{\beta}_0)]^{n/2}}} \right) \\ &= \frac{[(y - X\beta_0)'(y - X\beta_0)]^{n/2}}{[(y - X\tilde{\beta}_0)'(y - X\tilde{\beta}_0)]^{n/2}} \\ &= \left( \frac{\hat{\sigma}_w^2}{\hat{\sigma}^2} \right)^{n/2} = (1)^{n/2}\end{aligned}$$

where

$$1 = \frac{(y - X\beta_0)'(y - X\beta_0)}{(y - X\hat{\beta})'(y - X\hat{\beta})}$$

is the ratio of the quadratic forms.

Now we simplify the numerator in 1 as follows:

$$\begin{aligned}(y - X\beta_0)'(y - X\beta_0) &= [(y - X\hat{\beta}) + (\hat{\beta} - \beta_0)X]' \\ &\quad [(y - X\hat{\beta}) + X(\hat{\beta} - \beta_0)]\end{aligned}$$

$$(y - x\hat{\beta})'(y - x\hat{\beta}) + 2y'[\mathbf{I} - \mathbf{x}(\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}']\mathbf{x}(\hat{\beta} - \beta_0) +$$

$$(\hat{\beta} - \beta_0)'x'x(\hat{\beta} - \beta_0)$$

$$= (y - x\hat{\beta})'(y - x\hat{\beta}) + (\hat{\beta} - \beta_0)'x'x(\hat{\beta} - \beta_0)$$

Thus

$$1 = \frac{(y - x\hat{\beta})'(y - x\tilde{\beta}) + (\hat{\beta} - \beta_0)'x'x(\hat{\beta} - \beta_0)}{(y - x\hat{\beta})'(y - x\hat{\beta})}$$

$$1 = 1 + \frac{(\hat{\beta} - \beta_0)'x'x(\hat{\beta} - \beta_0)}{(y - x\hat{\beta})'(y - x\hat{\beta})}$$

$$1 - 1 = 1_0 = \frac{(\hat{\beta} - \beta_0)'x'x(\hat{\beta} - \beta_0)}{(y - x\hat{\beta})'(y - x\hat{\beta})}$$

where

$$0 \leq 1_0 < \infty$$

## Confidence interval Estimation:-

The confidence intervals in a multiple regression model can be constructed for individual regression coefficients as well as jointly.

### confidence interval on the individual regression coefficient :-

Assuming,  $\epsilon_i$ 's are identically and independently distributed following  $N(0, \sigma^2)$  in  $y = X\beta + \epsilon$ ,

we have,

$$y \sim N(X\beta, \sigma^2 I)$$

$$\hat{b} \sim N(\beta, \sigma^2 (X'X)^{-1})$$

Thus the marginal distribution of any regression coefficient estimate

$$\hat{b}_j \sim N(\beta_j, \sigma^2 c_{jj})$$

where  $c_{jj}$  is the  $j^{\text{th}}$  diagonal element of  $(X'X)^{-1}$

Thus

$$t_j = \frac{\hat{b}_j - \beta_j}{\sqrt{\sigma^2 c_{jj}}} \sim t_{(n-k)} \text{ under } H_0, j=1, 2, \dots$$

$$\text{where } \sigma^2 = \frac{SS_{\text{res}}}{n-k} = \frac{y'y - b'X'y}{n-k}$$

so the  $100(1-\alpha)\%$  confidence interval for  $\beta_j$  ( $j=1, 2, \dots, k$ ) is obtained as follow.

$$P \left[ -t_{\alpha/2, n-k} \leq \frac{\hat{b}_j - \beta_j}{\sqrt{\sigma^2 c_{jj}}} \leq t_{\alpha/2, n-k} \right] = 1 - \alpha$$

$$P \left[ b_j - t_{\alpha/2, n-k} \sqrt{s_{\text{reg}}^2 g_{jj}} \leq \beta_j \leq b_j + t_{\alpha/2, n-k} \sqrt{s_{\text{reg}}^2 g_{jj}} \right] = 1-\alpha$$

Thus the confidence interval is

$$(b_j - t_{\alpha/2, n-k} \sqrt{s_{\text{reg}}^2 g_{jj}}, b_j + t_{\alpha/2, n-k} \sqrt{s_{\text{reg}}^2 g_{jj}})$$

Simultaneous confidence intervals on regression coefficients :-

A set of confidence intervals that are true simultaneously with probability  $(1-\alpha)$  are called simultaneous or joint confidence intervals.

It is relatively easy to define a joint confidence region for  $\beta$  in multiple regression model.

Since

$$\frac{(b-\beta)' X' X (b-\beta)}{k M S_{\text{reg}}} \sim F_{k, n-k}$$

$$\Rightarrow P \left[ \frac{(b-\beta)' X' X (b-\beta)}{k M S_{\text{reg}}} \leq F_{\alpha}(k, n-k) \right] = 1-\alpha$$

So, a  $100(1-\alpha)\%$  joint confidence region for all parameters in  $\beta$  is

$$\frac{(b-\beta)' X' X (b-\beta)}{k M S_{\text{reg}}} \leq F_{\alpha}(k, n-k)$$

which describes an elliptically shaped region.