

Multiple Linear Regression

Look up

Fundamental Assumptions.

①

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \epsilon_i \quad \text{--- (1)}$$

β_0 = Expected value of y when $x_i = 0$

β_i = Change in y for unit change in x_i

ϵ_i = Random var with zero mean, std. dev. = σ_i

Rewriting in matrix Form.

$$y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_k x_{1k} + \epsilon_1$$

\vdots

$$y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_k x_{nk} + \epsilon_n$$

or

$$\begin{matrix} y \\ [n \times 1] \end{matrix} = \begin{matrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \\ n \times 1 \end{matrix} \begin{matrix} \beta \\ [(k+1) \times 1] \end{matrix} = \begin{matrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \end{matrix}$$

$$\begin{matrix} [X] \\ n \times (k+1) \end{matrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & \dots & x_{nk} \end{bmatrix}$$

Thus we have

$$\begin{matrix} y \\ [n \times 1] \end{matrix} = \begin{matrix} X & B \\ \downarrow & \searrow \\ n \times (k+1) & (k+1) \times 1 \end{matrix}$$

Least Square Estimator

We have.

2 errors as: $\hat{y} = x\beta$ $(y - \hat{y})^2 = e^2$

$$= [y - x\beta]' [y - x\beta]$$

Thus $e^2 = [y'y + (x\beta)'x\beta - (x\beta)'y - y'x\beta]$

$$e^2 = [y'y + \beta'x'x\beta - \beta'x'y - y'x\beta]$$

Now for minimising squared error

$$\frac{\partial e^2}{\partial \beta} = 0 \Rightarrow 2x'x\beta - (x'y + y'x) = 0$$

Same scalar Product

Now using $x'y = y'x \therefore x'x\beta = x'y$

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

Soln for $\hat{\beta} = (x'x)^{-1} x'y$

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

$$\hat{y} = Hy$$

$H =$ hat matrix

$$H = x(x'x)^{-1}x'$$

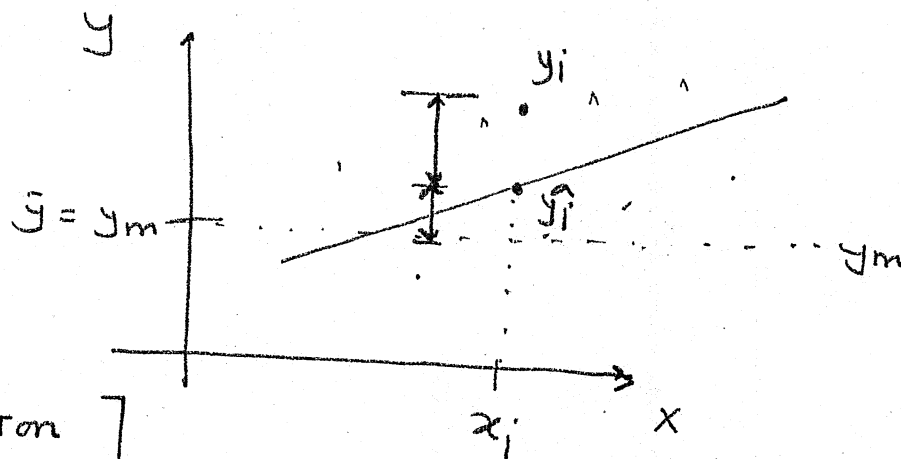
$$y_i = \hat{y}_i + e$$

$$e = y_i - \hat{y}_i = y_i - Hy_i$$

$$e = (I - H)y$$

R^2 and adjusted R^2

(3)



Explained Variation
Unexplained Variation
Total Variation

$$\text{Total Error} = y_i - \bar{y}$$

$$\text{Error / unexplained} = y_i - \hat{y}_i$$

$$\text{Explained} = \hat{y}_i - \bar{y}$$

$$\text{Total error} = y_i - \bar{y} = \text{unexplained } (y_i - \hat{y}_i) + \text{Explained } (\hat{y}_i - \bar{y})$$

Square both sides

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

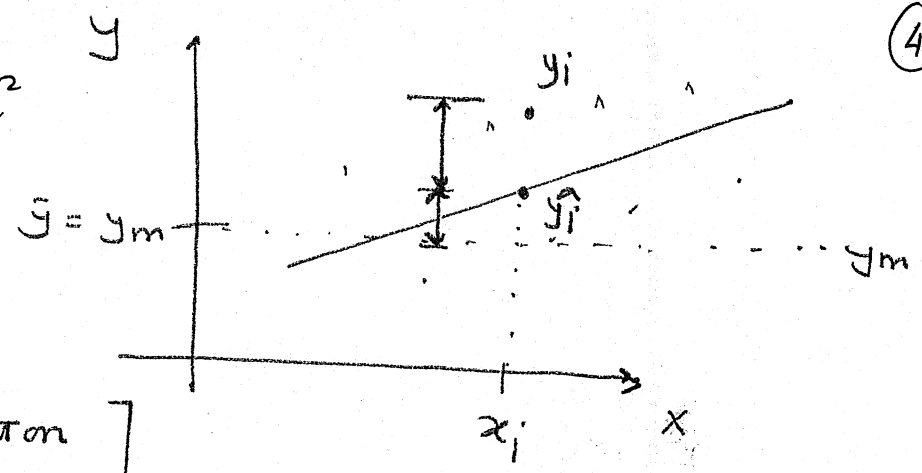
$$\text{on simplifying} = \underbrace{(y_i - \hat{y}_i)^2}_{\text{SSE}} + \underbrace{(\hat{y}_i - \bar{y})^2}_{\text{SSR}}$$

$$\text{Tot S.S} = \text{SSE} + \text{SSR}$$

$$R^2 = \frac{\text{Explained Variation}}{\text{Total Variation}}$$

$$= \frac{\text{SSR}}{\text{SSE} + \text{SSR}}$$

~~Regression~~
 R^2 and adjusted R^2



Explained Variation
Unexplained Variation
Total Variation

$$\begin{aligned}\text{Total Error} &= y_i - \bar{y} \\ \text{Error / unexplained} &= y_i - \hat{y}_i \\ \text{Explained} &= \hat{y}_i - \bar{y}\end{aligned}$$

$$\begin{aligned}\text{Total error} &= y_i - \bar{y} = \text{unexplained } (y_i - \hat{y}_i) \\ &\quad + \text{Explained } (\hat{y}_i - \bar{y})\end{aligned}$$

Square both sides

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

$$\text{on simplifying} = \underbrace{(y_i - \hat{y}_i)^2}_{\text{SSE}} + \underbrace{(\hat{y}_i - \bar{y})^2}_{\text{SSR}}$$

[SST]

$$\text{Tot S.S} = \text{SSE} + \text{SSR}$$

$$R^2 = \frac{\text{Explained Variation}}{\text{Total Variation}} = \frac{\text{SSR}}{\text{SSE} + \text{SSR}}$$

$$\text{Adjusted } R^2 = \frac{SS_{\text{Reg}}}{SST}$$

$$R^2_{\text{adjusted}} = 1 - \frac{SSE / dfe}{SST / df_{\text{total}}}$$

$$(n-p-1) \quad (n-1)$$

$$\text{Total} = F + R$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$(n-1) \quad (k-1) \quad (k-p-1)$$

$$R = p+1$$

F Test of Linear Regression.

Null $H_0: \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$

$H_1: \beta_j \neq 0$ For at least one value of j

$$F_{\text{stat}} = \frac{MSG_{\text{of Reg}}}{MSG_{\text{of Error}}} = \frac{SSR / dfr}{SSE / dfe}$$

Example:

$n = 35$ observations.

$R = 9$ independent variables

$$P = R + 1 = 10 \text{ parameters.}$$

$$SSE = 134$$

$$SSR = 289$$

Test to 0.05

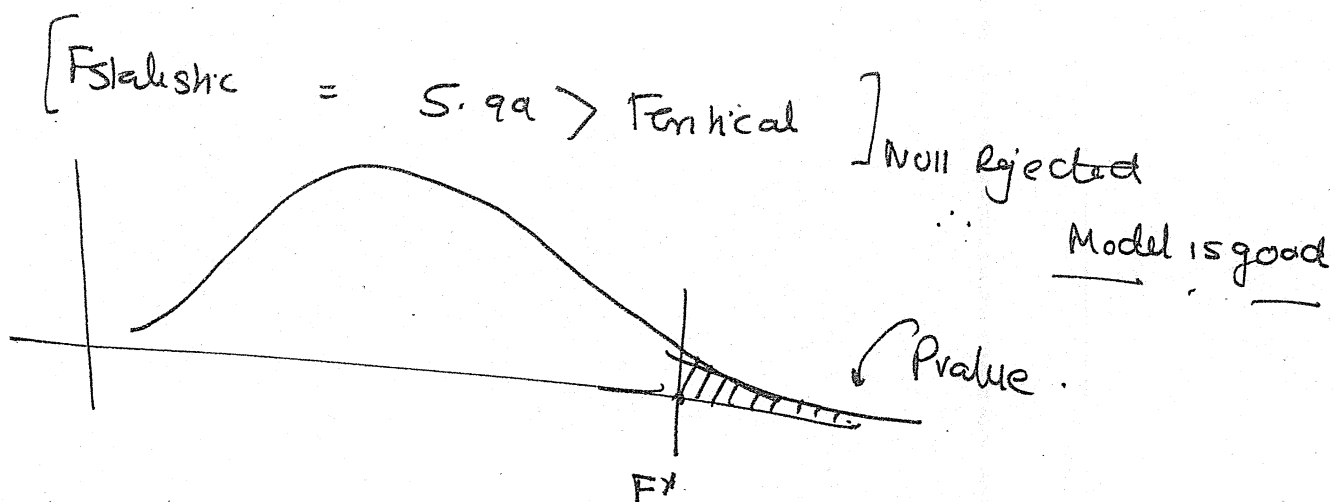
$$F_{\text{Statistic}} = \frac{MSR}{MSE} = \frac{SSR / dfr}{SSE / dfe}$$

$$= \frac{(289) / 9}{134 / 25} \quad F_{\text{stat}} = 5.99$$

Look up F table $(9, 25)$
 $\alpha = 0.05$

(6)

$$F_{\text{critical}} (9, 25) \quad \alpha = 0.05 = 2.28$$



Now check using p values:

$P_{\text{crit}} = 0.05$

$F_{\text{stat}} = 5.99 \quad (9, 25)$

We see from table for $F_{\text{stat}} = 5.99, (9, 25)$
 P_{value} will be less than 0.01

$\therefore P_{\text{less}} < P_{\text{crit}}$ Null is rejected

Significance of Individual Coefficients:

Methods \rightarrow (1) t test of individual coeff
 (2) Partial F test.

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

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

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

$$\sqrt{\frac{\sum (x_i - e_i)^2}{n-1}} = \sqrt{\hat{\sigma}^2} \quad \hat{\sigma}^2 = \text{Error Mean Square}$$

Regression Estimate

$$\hat{y} = x \hat{\beta}$$

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

Var. Covar.

$$\text{Var} = \frac{\sum (y_i - \bar{y})^2}{(n-1)}$$

$$\text{Covar} = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{n-1}$$

$$C = \hat{\sigma}^2 (X^T X)^{-1} \rightarrow \text{Var / Covar Matrix.}$$

$\hat{\sigma}^2$ - Error Mean Square.

$$SE(\beta_j) = \sqrt{C_{jj}}$$

Test statistic, [Significance of individual coefficients]

$$\alpha = 0.05$$

$$T_0 = \frac{\hat{\beta}_j}{SE(\beta_j)}$$

$$H_0: \beta_j = 0$$

$$H_1: \beta_j \neq 0$$

$$\text{if } T_0 > t_{\alpha/2, n-2}$$

$$\text{or } T_0 < -t_{\alpha/2, n-2}$$

also called Partial test.

Partial F test

(8)

$$y = \underbrace{\beta_0 + \beta_1 x_1 + \dots}_{(k+1)-q} + \underbrace{\beta_k x_k}_{q}$$

$(k+1)$

$\Theta_1 \rightarrow$ Contains first $(k+1)-q$ Coeff.

$\Theta_2 \rightarrow$ Contains last q Coeff.

then $\Theta_1 = [\beta_0, \beta_1, \dots, \beta_{k-q}]$

$$\Theta_2 = [\beta_{k-q+1}, \dots, \beta_k]$$

Null $H_0: \Theta_2 = 0$

$H_1: \Theta_2 \neq 0.$

Example: Consider $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e$

Partial F test for β_1

$$F_0 = \frac{SSR(\beta_1/\beta_2)/q}{MSE} \quad \text{as only 1 param.}$$

$\sqrt{q} \quad q=1$

$$= SSR(\beta_0, \beta_1, \beta_2) - SSR(\beta_0, \beta_2)$$

How to quickly compute $SSR(\beta_0, \beta_2)$

$$SSR(\beta_0, \beta_2) = y' \left[H_{\beta_0, \beta_2} - \left(\frac{1}{n} \right) J \right] y$$

Hat Matrix

We have for OLS $\beta = (X'X)^{-1} X'y$

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

$$\hat{y} = HY$$

$H =$ hat matrix

$$C = \hat{\sigma}^2 (X'X)^{-1}$$

↳ error Mean Square.

$$S.E(\beta_j) = \sqrt{C_{jj}}$$

$$t_{\text{test } \beta_j} = \frac{\beta_j}{SE(\beta_j)}$$

Properties of the hat Matrix.

1) Symmetric $H = H^T$

$$(AB)^T = B^T A^T$$

$$H = X(X^T X)^{-1} X^T$$

$$H^T = \left[X(X^T X)^{-1} X^T \right]^T = X \left[X(X^T X)^{-1} \right]^T = X \left[(X^T X)^{-1} \right]^T X^T$$

We know $X^T X$ is symmetric

$$\therefore \left[(X^T X)^{-1} \right]^T = (X^T X)^{-1} \quad \therefore H^T = H \quad \left[\begin{array}{l} X(X^T X)^{-1} X^T = H \\ X^T X \text{ is symmetric.} \end{array} \right]$$

2) $(I-H)$ is idempotent

Means $(I-H) = (I-H)(I-H)$

$$(I-H)(I-H) = I(I-H) - H(I-H)$$

$$= I - H - H + H \cdot H$$

$$= I - H - H + H = I - H$$

$[H \cdot H = H]$ We use this

$$\left| \begin{array}{l} X(X^T X)^{-1} X^T \cdot X(X^T X)^{-1} X^T \\ = X(X^T X)^{-1} X^T = H \end{array} \right|$$

$$\text{Var}(e^1)$$

$$e^1 = y - \hat{y} \quad y \text{ and } \hat{y} = Hy$$

$$\hat{y} = Hy$$

$$\therefore e^1 = y - Hy = (I - H)y$$

$$\text{Var}(e^1) = \text{Var}[(I - H)y]$$

1 For MLR

$$\beta_j \pm t_{\alpha/2, n - (k+1)} \sqrt{e_{jj}}$$

R Studentised Residual

$$\text{DFFIT}_j = \frac{\hat{y}_j - \hat{y}_j(e_j)}{e_j}$$

Change in $y_{\text{pred}, j}$ obs due to deletion of j th obs.

Residual Analysis

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

Standardised Residuals

$$d_i = \frac{e_i}{\sqrt{\text{MSE}}}$$

Studentised Residual

$$r_i = \frac{e_i}{\sqrt{\text{MSE}(1 - h_{ii})}}$$

Find outliers

$$\text{dof} = n - (k+1)$$

t test, 2 tailed

$$t_{\alpha/2, \text{dof}}, -t_{\alpha/2, \text{dof}}$$

Std Form

$$= \frac{\hat{y}_j - \hat{y}_j(-j)}{\sigma_{e(-j)} \sqrt{h_{jj}}}$$

$$\sigma_{e(-j)} \sqrt{h_{jj}}$$

If Large

j th obs was

influential.

Cook's Distance 'D' [Outlier detection]

$$D_i = \frac{r_i^2}{(k+1)} \frac{h_{ii}}{(1 - h_{ii})}$$

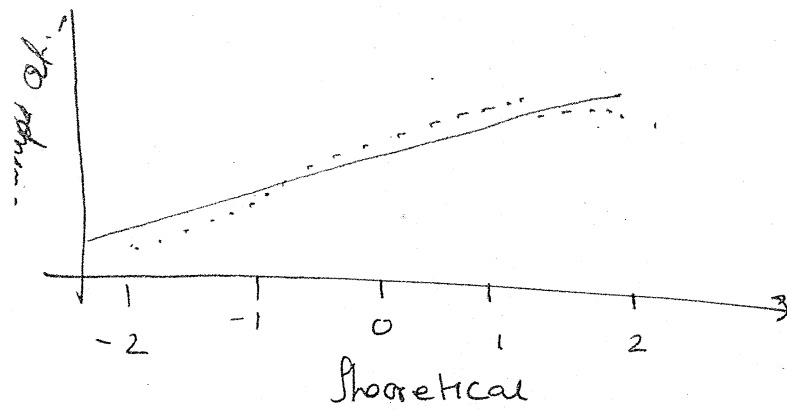
$$r_i = \frac{e_i}{\sqrt{\text{MSE}(1 - h_{ii})}}$$

If $D_i > 50$ th Percentile in F table

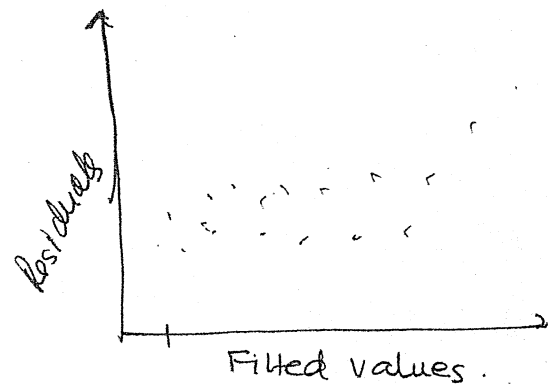
i th observation is influential

Residual plots . the QQ norm.

(11)



Residual Plot.



→ at Lower values of y
 e_i is Less

→ $e_i \uparrow$ at higher values

→ Variance of e_i
is higher at higher
values of y_i

Multi collinearity: VIF R_i^2 = Regression → when i th var is treated as dep. regressed on $(k-1)$ other indep. vars.

When Linear relationship involves more than two features.

$VIF = \frac{1}{1-R_i^2}$ Variance inflation factor.

Thumb Rule

$VIF > 4$

Suspect

$VIF > 10$

Strong multicollinearity.

Feature Selection

1] Check variance of the feature, Remove features which have Low variation.

2] Forward Selection.

Forward Selection

AIC OR Adjusted R^2

(12)

- 1] Start with Null Model.
- 2] Take one feature Rem model
Repeat with each feature \rightarrow Rem 'k' times
- 3] 1st Feature is set (one with Lowest AIC Value)

4] \vdots $AIC = 2k - 2\ln(L)$

Continue adding features till
adding another feature does not
improve the model significantly.

$$1 + (k + (k-1) + \dots + 1) = 1 + \sum_{i=1}^k i = 1 + \frac{k(k+1)}{2}$$

Note: order of Magnitude of computation
is of the scale of k^2

which is less than 2^k which would come
from full grid search.

usages of the hat Matrix.

(13)

We have. $\hat{y} = Hy \quad \therefore y - \hat{y} = (I - H)y$

The Variance operator

$$\text{Var}(\hat{e}) = \text{Var}[(I - H)y] = (I - H) \text{Var}(y) (I - H)'$$

Digression

$$\text{Var}(y) = \begin{bmatrix} \text{Var } y_1 & \text{Cov}(y_1, y_2) & & \\ & \text{Var } y_2 & \dots & \\ & & \ddots & \\ & & & \text{Var } y_N \end{bmatrix}$$

Now $\text{Var}(ax + by) = a^2 \text{Var}(x) + b^2 \text{Var}(y) + 2ab \text{Cov}(x, y)$

Now $y_i = x_i \beta + e_i$

$$\text{Var}(y_i) = \text{Var}(x_i \beta) + \text{Var}(e_i) + 2 \text{Cov}(x_i \beta, e_i)$$

Now $\text{Var}(\beta) = 0$

$$\text{Cov}(x_i \beta, e_i) = 0 \quad \therefore x_i \text{ is not corr. with } e_i$$

$$\therefore \text{Var}(y_i) = \text{Var}(e_i) = \sigma_{e_i}^2$$

Now Homoskedasticity assumption:

$$\text{Var } y = \begin{bmatrix} \sigma_e^2 & & & \\ & \sigma_e^2 & & \\ & & \ddots & \\ & & & \sigma_e^2 \end{bmatrix} = \sigma_e^2 [I]$$

Thus $\text{Var}(\hat{e}) = \sigma_e^2 (I - H)(I - H)'$ as $\text{Var}(y) = \sigma_e^2 I$

$$\therefore \text{Var}(\hat{e}) = \sigma_e^2 (I - H)$$

$$\text{eg. } \left[\text{Var}(\hat{e}_{10}) = \sigma_e^2 (1 - h_{10,10}) \right]$$

Brief Recap:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & x_{11} & \dots & x_{k1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \dots & x_{kn} \end{bmatrix}_{[n \times (k+1)]} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}_{[(k+1) \times 1]} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}_{n \times 1}$$

$$y = X\beta + e \quad \therefore e = y - X\hat{\beta}$$

$$\therefore SSE = (y - X\hat{\beta})'(y - X\hat{\beta})$$

Now $e'e = y'y - 2\hat{\beta}'X'y + \hat{\beta}'X'X\hat{\beta}$

$$\frac{\partial e'e}{\partial \beta} = -2X'y + 2X'X\hat{\beta} = 0 \quad] \text{ Normal equation}$$

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

Normal Eqn's.

$$y = X\hat{\beta} + e \quad \downarrow$$

$$(X'X)\hat{\beta} = X'(X\hat{\beta} + e)$$

$$\boxed{X'X\hat{\beta}} = \boxed{X'X\hat{\beta}} + X'e$$

$$\Rightarrow X'e = 0 \quad \downarrow$$

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{k1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \dots & x_{kn} \end{bmatrix}_{n \times (k+1)} \quad \xrightarrow{\text{n obs}} \quad X' = \begin{bmatrix} 1 & x_{11} & x_{21} & \dots & x_{k1} \\ 1 & x_{12} & x_{22} & \dots & x_{k2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & x_{2n} & \dots & x_{kn} \end{bmatrix}_{(k+1) \times n} \quad e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}_{n \times 1}$$

\Rightarrow Multiplying 2nd Row out

$$\underbrace{X_{11}} e_1 + \underbrace{X_{12}} e_2 + \dots + \underbrace{X_{1n}} e_n = 0$$

$$\hookrightarrow \underbrace{X'_{i1}} e_1 + \underbrace{X'_{i2}} e_2 + \dots + \underbrace{X'_{in}} e_n = 0$$

\Rightarrow Each Regressor has zero correlation with residuals.

First Row Multiplied out:

$$\sum e_i = 0$$

Predicted value of 'y'

$$\hat{y} = X \hat{\beta} \quad \underbrace{y'e = (X \hat{\beta})' e = \beta' \underbrace{X'e}_{=0}} = 0$$

\hat{y} Predicted is uncorrelated with residuals.

Gauss Markov Assumptions:

1] $y = X\beta + e$ y, X are Linearly related.

2] X is a $n \times k$ Matrix of full rank

\rightarrow Cols of X have no perfect multicollinearity.

\rightarrow Cols of X are linearly independent

$$E(y) = X\beta$$

$$\Omega = \sigma^2 I \rightarrow \text{homoskedasticity}$$

X] May be fixed or random but must be generated by a mechanism that is unrelated to e

$$e \sim N[0, \sigma^2 I]$$

Gauss Markov Theorem:

OLS is the best linear unbiased and efficient estimator.

Proof: $\hat{\beta}$ is the unbiased estimator of β

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y \quad Y = X\beta + e$$

$$\text{Thus } \hat{\beta}_{OLS} = (X^T X)^{-1} X^T (X\beta + e)$$

$$= \beta + (X^T X)^{-1} X^T e$$

$$E[\hat{\beta}] = E[\beta] + (X^T X)^{-1} X^T E(e)$$

but $E(e) = 0$

$$\therefore E[\hat{\beta}] = \beta$$

Proof $\hat{\beta}$ is a linear estimator of β .

$$\hat{\beta}_{OLS} = \beta + (X^T X)^{-1} X^T e$$

$$\text{or } \hat{\beta} = \beta + Ae$$

CROSS VALIDATION

- > A Model validation Technique for assessing how the results of a statistical analysis will generalize to an independent data set.
- > Involves partitioning a sample of data into complementary subsets, performing the analysis on one subset (training), validating the analysis on the other (validation set).

'k' fold Cross validation

- Original sample is randomly partitioned in ' k ' equal sized sub samples.
- One sub sample is retained as validation set, the other $k-1$ sub samples form the training set.
- The cross validation process is repeated ' k ' times.
- Each ' k ' sub samples are exactly used once.
- The ' k ' results can then be averaged to produce a single estimation.

Adv] All observations are used for both training and validation.

↳ Can be used for feature selection.

↳ How?

e.g.] Say 20 protein expression levels are being used to test whether a cancer patient will respond to a drug.

→ Which subset of features give the best in-sample error rates?

Limitations: Cross validation only yields meaningful results if the validation set and training set are drawn from the same population and only if human biases are controlled