## Session - 53

### multicollinearity :- (MC)

↳ it occurs when two or more independent variable ( i/p column) in a multiple reg column are highly correlated (0.8, 0.9).

**Problem with mc -**

$$y_q = \beta_0 + \beta_1 (cgpa + \beta_2 iq$$

$\beta_1 \Rightarrow$ If we keep iq constant and increase cgpa by 1 then how much lpq will get ↑

but if cgpa & iq are correlated than if we CGPA ↑ then iq ↑ so $\beta_0, \beta_1, \beta_2$ become unreliable

• If we are doing inferences than now mc is bad as $\beta_1, \beta_2, \beta_0$ are not reliable

• If we are doing prediction than no problem

$$X_1 \quad X_2 \quad | \quad Y$$

$X_1 \& X_2$ are related

$$X_1 = a_0 + a_1 X_2 + \eta$$
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$
$$\beta_0 + \beta_1 a_0 + \beta_1 a_1 X_2 + \eta \beta_1 + \beta_2 X_2 + \varepsilon$$

$$Y = (\beta_0 + \beta_1 a_0) + (\beta_1 a_1 + \beta_2) X_2 + (\eta \beta_1 + \varepsilon)$$

Since Y is only f of $X_2$ so no problem

→ Generally in correct$^n$ -

$$X_1 = a_0 + a_1 X_2 + error$$

due to this error we can't find exact value of $X_1$.

but if $X_1 = a_0 + a_1 X_2$

means no error, than this situation is called perfect multicollinearity

· lets assume perfect MC -

$$\begin{bmatrix} CGPA & Percent & lpa \\ 8 & 80 & 3 \\ 6 & 60 & 4 \end{bmatrix}$$

$$lpa = \beta_1 + \beta_1 cgpa + \beta_2 percent + error$$

we want to find $\beta$ -

$$\beta = (X^T X^*)^T X^T Y$$

↳ while solving it we have to find $det(X^T X$
which is zero. means its a singular
matrix means we can't calculate inverse
So we can't find $\beta$.

Hence in perfect MC we can't find $\beta$. and most of data is not perfect MC

· what exactly happens in MC?
→ Difficult to tell which indep. variable has most significant affect on dep. column

→ inflated Standard error (big SE)

↳ unstable & unreliable estimates

(If we change small values in i/p then lots of change)

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

$$Var(\beta) = SE(\beta) \quad \sigma^2 (x^T x)^{-1}$$

$$= \sigma^2 \begin{bmatrix} \underbrace{a} & & \\ & \underbrace{b} & \\ & & \underbrace{c} \end{bmatrix}$$

b→var. covar. matrix

$a \Rightarrow var(\beta_0)$        $SE(\beta_0) = \sqrt{var(\beta_0)}$

$b = var(\beta_1)$

$c = var(\beta_2)$        $SE(\beta) = \sqrt{dig(\sigma^2 (x^T x)^{-1})}$

$\Rightarrow$ If Strong mc then $det(x^T x)$ will be very small due to which inverse will be very high hence Standard Error will also be high.

→ Types of multicollinearity -

i) Structural mc -
   is it arises due to the way in which
   variables are defined

$$y \overset{\curvearrowleft}{\underset{x^2}{x'}} \begin{array}{c|c} x & y \\ \hline x' & x^2 & y \end{array}$$

due to polynomial reg. Structural mc arises as $x^0$, $x'$, $x^2$ are dep correlated

2) Data-driven mc -
   is Indep. variables are highly correlated
   due to specific data being analysed

more area of flat ⇒ more washroom.

- How to detect MC —

1) Correlation

$$X_1 = a_1 X_2 + a_0 + Error$$

Correlation b/w $X_1$ & $X_2$

check values $> 0.8$ or $0.9$

2) Variance Inflation factor (VIF)

A  B  C  |  O

1) we make A, B as input & C as output and apply L.R & find $R^2$ score
2) A, C as input & B as o/p
3) B, C " " ; A " "

$$VIF = \frac{1}{1 - R^2}$$  ) if VIF $> 5$ or $10$

then M.C exist

else not

3) Condition Number —

matrix $X^T X$

$$Con. No = \frac{Largest\ Eigen\ value}{Smallest\ Eigen\ value}$$

Cond$^n$ No. tells about ill conditioning of matrix. means if small change in error cause large change in $Sq^a$

Cond$^n$ no $> 30$ means, <u>MC</u>

· How to remove MC —

→ Collect more data

→ Remove one of the highly correlated variable

→ Combine correlated variable

→ Use partial least Squares (PLS) repression