

- Application of statistical methods to measure economic relationship.
- Economic data were classified into 5 - But 3 will be discussed
- ① Cross sectional data → large time period
 - ② Time series data → differentiate b/w the 3
 - ③ Panel data

Model II
linear in
parameters
OLS Assumption

- Linearity / Multicollinearity
- Normality
- Independence of obs
- Homoscedasticity (some variance across groups)

Ts.: Monthly sales data - $\begin{cases} \text{Yearly} \\ \text{Quarterly} \\ \text{Monthly} \\ \text{Weekly} \end{cases}$ } basic.

Panel data: Combines existence of cross sectional data and Time series -
More than one variable More than one time period.

Fixed Model / Pool model / Random model /

Carry test on the data to know which model to be used. depending on the result of that test

$$Y_i = \beta_0 + \beta_1 X_{i1} + \epsilon_i \quad i=1, 2, 3, \dots, n$$

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_K X_{iK} \quad i=1, 2, \dots, n$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad X = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1K} \\ X_{21} & X_{22} & \dots & X_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{nK} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_K \end{bmatrix}, \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$Y = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_K X_{iK}$$

$$\text{But } \beta_1 = \frac{\partial Y_i}{\partial X_{i1}}, \quad \beta_2 = \frac{\partial Y_i}{\partial X_{i2}}, \quad \beta_3 = \frac{\partial Y_i}{\partial X_{i3}}$$

Rank of matrix

Number of linearly independent columns. X_i doesn't affect other X_{ij} . If model is not fully ranked, then inverse of X cannot be defined. Full rank matrix - obs are not independent of one another.

No column is linearly dependent on one another.

Homoscedasticity of error
Error variances are constant

$X_1 \quad X_2 \quad \dots \quad X_n$

ϵ_i

$$E(\epsilon_i) = 0$$

$$y = X\beta + \epsilon$$

$\sim \sim \sim \sim$

Rank = no of non zero columns in the matrix

Full Rank \Rightarrow

$$V(\epsilon_i) = \sigma_i^2$$

$$= E[(\epsilon - E(\epsilon))(\epsilon - E(\epsilon))^T]$$
$$= E[\epsilon\epsilon^T - E(\epsilon)\epsilon^T - E(\epsilon)\epsilon + E(\epsilon)\epsilon^T]$$

$$= E(\epsilon\epsilon^T)$$

$$\Rightarrow \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \begin{bmatrix} E(\epsilon_1\epsilon_1) & E(\epsilon_1\epsilon_2) & \dots & E(\epsilon_1\epsilon_n) \\ E(\epsilon_2\epsilon_1) & E(\epsilon_2\epsilon_2) & \dots & E(\epsilon_2\epsilon_n) \\ \vdots & \vdots & \ddots & \vdots \\ E(\epsilon_n\epsilon_1) & E(\epsilon_n\epsilon_2) & \dots & E(\epsilon_n\epsilon_n) \end{bmatrix}$$

$$V(\epsilon_i)$$

scalar $\in \mathbb{R}^n$. $E(\epsilon_{i1}) = E(\epsilon_{i2}) = \dots = E(\epsilon_{in})$

$$E(\epsilon_{i1}) = \text{Var}(\epsilon_i) = \sigma_i^2$$

$$\text{generally } E(\epsilon_i\epsilon_i) = E(\epsilon_i^2) = \sigma_e^2$$

- variance b/w errors are zero —ols/MR assumption.

$$\mathbb{E}(g_{i,j}^{\text{eff}}) = 0$$

$$F(e_{i+1}) = \neg (e_i \rightarrow E(e_i)) \quad (e_j \rightarrow E(e_j)).$$

$$E(e_i e_j) = 0 \quad \forall i \neq j \quad (9x-1)(9x-1) = 0$$

If $i=j$ then we get the variance of x_i

$$V_{(cc)} = \begin{pmatrix} \sigma_1^2 & 0 & 0 & \dots \\ 0 & \sigma_2^2 & 0 & \dots \\ 0 & 0 & \sigma_3^2 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} = \sigma^2 I_n$$

note of changes has σ^2 from σ_1^2 plus other

Normality assumptions. — Assumption 5

Since $y = xp + c$

$$F(x) = x^{\beta}$$

$$\tilde{v}(y) = v(x_B + e)$$

$$= \sqrt{x\beta} + \sqrt{r}$$

$$N(y) = 0 + V(x)$$

The Sampling distribution of \bar{Y} is

$$y \sim N(x_\beta, \sigma_e^2)$$

Least Square Estimation

$$SSE = (y - X\beta)' (y - X\beta)$$
$$= (y - Xb)' (y - Xb)$$

$$= y'y - b'X'y - y'Xb + b'X'Xb$$

$X'y$ | Xb $X'X$ | $\beta X'$

$X'X$ | βX $X'X$ | βX

$= Xb$ scalar

$X'X$ scalar

$$= y'y - 2b'Xy + b'X'Xb$$

Differentiate partially w.r.t. b and equate to zero
Clearing the

$\Rightarrow \text{differentiate w.r.t. } b$

$$\Rightarrow 2(X'y) - 2X'b = 0$$

$$\Rightarrow X'b = 0$$

$$\Rightarrow (X'X)b = 0$$

$$\Rightarrow (X'X)^{-1} (X'X)b = 0$$

$$\Rightarrow b = 0$$

$$\Rightarrow b = 0$$

If X' is non singular

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

Define Econometrics

- Econometrics is the field of economics in which statistical methods are developed and applied to estimate the economic relationships, testing economic theories and evaluating and implementing government and business policy. i.e. policies implemented by private industry, government and supranational organizations are studied.
- ## Define Economic Models Vs Economic Model
- ### Econometric Models

- Econometric Models consists of 3 things.

I - A set of equations describing the behaviour of variables involved. The equations are derived from the econometric model and are of 2 parts

- ① observed variables
- ② disturbances

II Statement about errors in the in the observed values of the variables.

III Specification of the probability distribution of the disturbances or error.

- Econometric Models specifies Statistical relationship b/w the various economic quantities pertaining to a particular economic phenomenon.

e.g. Monthly ~~expenditure~~ by customers is linearly dependent on their income in the previous month.

$$C_t = a + b_1 I_{t-1} + \epsilon_t$$

↓
Customer spending
in a month

↓ Constant or Intercept
↓ Slope
↓ Income
↓ Previous month's labor force
↓ Error term showing the extent to which model cannot fully explain the consumption

↓ right side shows work and leisure

↓ Objectives

(1) To estimate parameters and fit the model

(2) To predict for future observations.

Econometric Models consists of a set of equations designed to provide a quantitative explanation of the behaviour of the economic variables.

Type of Econometric Models:

- Linear Regression
- Generalized Linear Model
- ARIMA
- Probit
- Logit
- Tobit

Vector Auto Regression (VAR)

VAR model is a multivariate time series model that relates current values of a variable with past values of itself and other variables in the system.

Recent ex: NSE

Historical data → Forecasting with ARIMA

Economic Models

- Economic model is a set of assumptions that describes behaviour of an economy or phenomenon under study
- Economic model consists of mathematical equations that describes various relationships.

Utility maximization is made use of here., individual makes choices to maximise their well being which gives power framework for creating economic model.

- Utility maximization lead to a set of demand equations
- In demand equations, quantity demanded depends on
① price of goods, ② price of substitute and complementary good
③ the consumer's income & individual ④ tastes that affects taste.

$$Y \sim X_1 + X_2 + X_3 + X_4$$

$$Y = GDP \quad Y = \text{food expenditure}$$

$X_1 = \text{Income}$

$X_2 = \text{education}$

$X_3 = \text{family size}$

$X_4 = \text{age}$

Another definition of Econometrics

- Econometrics include the study of methods for selecting the models, estimating the model and making inference on the model

How Econometrics —
select model
estimate model
make inferences

- ~~cross~~ ~~out~~ ~~structure~~ of economic data
- ① Cross sectional data
 ② Time series data
 ③ Panel data OR longitudinal data
 ④ Others are clustered data AND spatial data
- (3) structures
- Cross
Sectional
data
Time
Series
data
Panel
data
Longitudinal
data
Clustered
data
Spatial
data

Cross sectional Data:

- Consists of sample of individuals, households, firms, cities, states, countries or a variety of other units taken at a given point in time
- Cross sectional data have one observation per individual
- It can be assumed that the CSO have been obtained by random sampling underlying the population.

Example

CSO on 528 random sampled persons in 1976 with variables					
class	Wage	Education level	Expenditure	Sex	Marital status
1	34	11	2	1	1
2	3.24	12	22	0	0
528	3.88	14	5	1	0

Further notes - not notes manipulated to make it look good

Time series. Data of wage is very useful when

- Set of observations on a Variable or Several Variables
- which are collected over time.
- * Example: ① GDP per year
 ② CPI
 ③ Money supply
 ④ Stock prices
 ⑤ Automobile sales figures etc.

- In time series data, we have lags - Past events can influence future events
- With lags and influences in the data; it is more difficult to analyse a time series data. because observations are strongly related to recent histories.

② Data frequency at which data is collected

The frequency may be Weekly, Monthly, Quarterly, ~~Daily~~ etc

Weekly - Daily \Rightarrow Stock prices

Monthly Weekly \Rightarrow Money Supply (US)

Monthly \Rightarrow Inflation, Unemployment rates

Quarterly (every 3 months) \Rightarrow GDP

Yearly \Rightarrow Infant Mortality rate

Example time series on Minimum Wage effect in a certain country

Panel data or Longitudinal data:

Panel data combines both cross-sectional data and time series data

- Data are collected on time-series basis; for each cross-sectional member in the data set

- Panel data consists of set of individuals, persons, households, corporations etc which are measured repeatedly over time

Example:

① Wage, education and employment data for a set of individuals followed over a 10-year period

② Commection of investment, financial information about the same set of firms (eg Banks) over a 5-year period

Assumptions of panel data modelling

✓ Individuals are mutually independent of one another BUT
given individuals observations are mutually dependent

Example ③ over 3 years & starting behaviour is static

3-year panel data set on Crime and related statistics
for 120 cities in a country
Crime variables may be Murder, Rape, Theft etc.

$$\text{dependent variable} = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

static behaviour = α

exogenous variables = β

dependent variable is labour start rate which 2.10 after three yrs.

dependent variable is labour start rate ① static behaviour is α

determinants not related to the world ②

determinants not related to themselves ③

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_n X_{nt}$$

④

$$\text{if } \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_n X_{nt} = b$$

then b is called the intercept. Intercept is 3.20

and $\alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_n X_{nt} - 3.20$ is called R^2

$$b + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_n X_{nt} - 3.20 + 3.20 = b$$

perhaps we can do this for all observations and then b

first \hat{y}_1 for first observation \hat{y}_2 for second observation \hat{y}_3 for third observation \hat{y}_4 for fourth observation \hat{y}_5 for fifth observation \hat{y}_6 for sixth observation \hat{y}_7 for seventh observation \hat{y}_8 for eighth observation \hat{y}_9 for ninth observation \hat{y}_{10} for tenth observation

Basic Multiple Regression Model

- Regression Models are used to study relationship b/w variables
- This model can be econometric model which seeks to

explain a dependent variable Y in terms of some N, X

- OR seeking to study how DV varies with changes in IV

$$Y = f(X_1, \beta_1, \dots, X_K, \beta_K) + e$$

Y = Regressed DV | response Var (explained variable)

$X = X_1, X_2, X_3, X_4, \dots, X_K$ = Regressors IV | explanatory Var

β = Parameter Vector

e = random error or disturbance

- We will use OLS which states that model is linear in parameters

f in general may be ① Linear in Variables

② Linear or Nonlinear in parameters

③ Parametric or Non parametric.

egn ① can be expressed as

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_K X_K + e \quad \text{--- (2)}$$

here e is Unobserved factors

If model is linear in parameters we have

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_K X_{iK} + e_i$$

β_s and e are Unobservable which are generated

- ① Deterministic component

~~stochastic component~~

β_i are marginal effects of X_i with other factors held constant

- Aim:
- To estimate unknown parameters
 - To test hypothesis about parameters
 - To predict values on y outside sample

Assumption of Classical ^{Linear} Regression Model (5) Assumptions

- ① Linearity
- ③ Errors have zero mean
- ⑤ Normality of errors
- ② Full rank of X
- ④ Homoscedasticity {Spherical errors}

Linearity: Model is linear in parameters. This allows us to write the model in matrix form

$$Y_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_K X_{iK} + \epsilon_i$$

Linearity in parameters helps us to write the model in Matrix form.

$$Y = X\beta + \epsilon$$

$$\begin{matrix} Y \\ \vdots \\ Y_n \end{matrix} = \begin{matrix} n \times 1 \\ \text{row } X \text{ columns} \\ \vdots \\ n \times K \text{ columns} \end{matrix} \begin{matrix} K \times 1 \\ \text{row } X^T \text{ columns} \end{matrix} \beta = \begin{matrix} n \times 1 \\ \text{row } \beta \text{ columns} \end{matrix} + \begin{matrix} n \times 1 \\ \text{row } \epsilon \text{ columns} \end{matrix}$$

$$Y_i = X_i \beta + \epsilon_i$$

② Full Rank

- There are no exact linear dependencies among the columns of X
- If there is linear dependency, then one or more regressor is redundant

Rank = no of linearly independent columns in the matrix X

Since X is $n \times K$ matrix, then $\text{rank}(X) = K$ if $n > K$

$\text{rank}(A) \leq \min(\text{rows, columns})$ which indicates that $n > K$

③ Errors have zero mean
 If it is assumed in the population that $E(e_i) = 0$, $\forall i=1, 2, 3, \dots, n$
 e is $n \times 1$ matrix

$$E(e) = E\begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = 0$$

④ Homoscedasticity. [constancy of Variance of e_i]

- In the population, the errors are generated by a process where Variance, σ^2_e is constant. And that e_i is Uncorrelated with each other.

$$\text{Var}(e_i) = \sigma^2_e$$

- constancy of variance = assumption of homoscedasticity.
- e_i are uncorrelated.

$$\text{Cov}(e_i, e_j) = 0 \quad i \neq j \quad \{ \text{No autocorrelation} \}$$

• Covariance Matrix for random vector e can be generalized

$$\begin{aligned} \text{Var}(e) &= E\{(e - E(e))^2\} = E\{(e_i - E(e_i))(e_j - E(e_j))'\} \\ &= E\{ee' - E(e)e' - E(e')e + E(e)E(e')\} \\ &= E\{ee'\} - E(eEe') - E(Eee') + E(eEe') \end{aligned}$$

but $E(e) = E(e') = 0$, then

$$\text{Var}(e) = E(ee') - E(0) - E(0)' + E(0) = E(ee')$$

$$\text{Var}(e) \geq E(ee')$$

$$\text{Var}(e) = E(ee')$$

$$\text{Now } E(ee') = \text{Var}(e) = \frac{1}{n} \begin{vmatrix} e_1 & e_2 & \cdots & e_n \\ e_1 & e_1e_1 & e_1e_2 & \cdots & e_1e_n \\ e_2 & e_2e_1 & e_2e_2 & \cdots & e_2e_n \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ e_n & e_ne_1 & e_ne_2 & \cdots & e_ne_n \end{vmatrix}$$

$$E(e_i e_j) = E(e_i^2) = E(e_i - E(e_i))^2$$

But $E(e_i) = 0$

$$E(e_i e_j) = E((e_i - 0)^2) = E(e_i^2) = \text{Var}(e_i) = \sigma^2$$

(Because e_i is an independent error term)

$$E(e_i e_j) = E((e_i - E(e_i))(e_j - E(e_j)))$$

$$= E(e_i - 0)(e_j - 0)$$

$$= E(e_i)e_j$$

$$= \text{Cov}(e_i, e_j) = 0$$

No autocorrelation errors

$$\text{Var}(e) = \begin{Bmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & \dots & 0 \\ 0 & 0 & \dots & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \dots & \sigma_n^2 \end{Bmatrix} = \sigma^2 I_n$$

$I_n = \text{Identity Matrix}$

$\text{Var}(e)$ is a scalar matrix

③ Normality of errors

Errors e_i are normally distributed as $e_i \sim N(0, \sigma^2 I_n)$

$$\text{Since } y = X\beta + e$$

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

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

$$E(e) = 0$$

$$E(y) = X\beta + 0$$

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

Its variance

$$V(y) = V(X\beta + e)$$

$$V(y) = V(X\beta) + V(e)$$

$$V(y) = 0 + V(e)$$

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

$\therefore y$ has sampling distn

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

Least squares Estimation of parameters β_0 , β_1

from $y = X\beta + \epsilon$ $\epsilon_i \sim N(0, \sigma^2)$ in
We use Least square strategy because we have $(n+1)$ parameters.

Our Unobserved error ($\epsilon_i = y_i - \hat{y}_i \beta$)

equivalently for estimate $\hat{\epsilon}_i = y_i - \hat{y}_i \hat{\beta}$

$$\hat{y}_i = X\hat{\beta} + \epsilon$$

$$\hat{y}_i = X\hat{\beta} \quad \text{or } X\beta$$

$$\text{then } y_i = \hat{y}_i + \epsilon$$

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

To use least square, i.e. we minimize ϵ .

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

$$\hat{\epsilon}_i^2 = (y_i - \hat{y}_i)^2$$

$$\hat{\epsilon}_i^2 = (y_i - \hat{y}_i)(y_i - \hat{y}_i)$$

$$\text{Since } \hat{y} = X\beta$$

$(\hat{\epsilon}_i^2, 0)$ has to be initialized since it is not zero.

$$(s + qx)v = (p)v$$

$$s + qx = k$$

$$(sv + qx)v = (pv)$$

$$sv + qx = (p)v$$

$$(sv + 0) = (pv)$$

$$0 = (sv)$$

$$sv = (pv)$$

$$sv = (pv)$$

which implies $v = 0$

$$sv = (pv)$$

$$(s, qx) \mapsto p$$

$$Y = X\beta + \epsilon$$

$$\epsilon = Y - \hat{Y} \quad \text{But } \hat{Y} = X\beta$$

$$\epsilon_i = Y_i - \hat{Y}_i$$

$$\epsilon_i = (Y - X\beta)^T$$

$$\epsilon_i^2 = (Y - X\beta)^T (Y - X\beta)$$

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

$$= Y^T Y - Y^T X\beta - X^T \beta^T Y + X^T \beta^T X\beta$$

$$= Y^T Y - Y^T X\beta - X^T \beta^T Y + X^T X\beta^T \beta$$

$$= Y^T Y - 2\beta^T X^T Y + \beta^T X^T X\beta$$

diff wrt β

$$\frac{\partial \epsilon^2}{\partial \beta_j} = \frac{\partial}{\partial \beta_j} (Y^T Y - 2\beta^T X^T Y + \beta^T X^T X\beta)$$

$$= 0 - 2X^T Y + 2X^T X\beta$$

$$= -2X^T Y + 2X^T X\beta$$

To minimise, $\frac{\partial \epsilon^2}{\partial \beta} = 0$

$$0 = -2X^T Y + 2X^T X\beta$$

$$-2X^T Y = 2X^T X\beta$$

$$X^T Y = X^T X\beta$$

$$\beta = X^T Y$$

$$\beta = \frac{X^T Y}{X^T X}$$

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

$$\text{Variance of error } S_e^2 = \frac{\sum y - b \sum xy}{n-2} \quad \text{OR} \quad \frac{\sum yy}{n-2}$$

$$S_{yy} = \sum y^2 - \frac{(\sum y)^2}{n}$$

$$S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} \quad S_{yy} = SST$$

$$b = \frac{\sum xy}{\sum x^2}$$

$$S^2 = \frac{SSE}{n-2} = \frac{SST - SSR}{n-2}$$

$$S_e^2 = \frac{\sum y - \hat{\beta}_1 \sum x_1 y}{n-r-1}$$

$r = \text{no of variables}$

$x_1 x_2$

$$S_{yy} = \sum y^2 - \frac{(\sum y)^2}{n} \quad \text{But } \bar{y} = \frac{\sum y}{n} ; n\bar{y} = \sum y \text{ then } \frac{(\sum y)^2}{n} = (\bar{y})^2 = \frac{n\bar{y}^2}{n} = \bar{y}^2$$

We deduct $n\bar{y}^2$ from $\hat{\beta}_1 \sum x_1 y$ i.e. $\hat{\beta}_1 \sum x_1 y - n\bar{y}^2$ when cal S_e^2

$$\text{Since } \hat{\beta}_1 = (x_1 x_1)^{-1} (x_1 y)$$

$$v(\hat{\beta}_1) = \sigma^2 (x_1 x_1)^{-1} \quad \text{where } \sigma^2 = S_e^2 \text{ calculated}$$

Testing parameter β_i

$$t_{\text{cal}} = \frac{\hat{\beta}_i - 0}{S_e(\hat{\beta}_i)} \quad \text{with } H_0: \hat{\beta}_i = 0 \quad \text{Reject } H_0 \text{ if } |t_{\text{cal}}| > t_{\text{tab}}$$

$$t_{\text{tab}} = t_{1-\alpha/2, n-r-1}$$

$$(I \Rightarrow \hat{\beta}_i \pm t_{1-\alpha/2, n-r-1} S_e(\hat{\beta}_i))$$

$$R^2 = 1 - \frac{SSE}{SST}$$

$$R^2 = 1 - \frac{(n-1)}{(n-r)} \left(\frac{SSE}{SST} \right)$$

$$R^2 = 1 - \frac{(n-1)}{(n-r-1)} \left(\frac{SSE}{SST} \right)$$

ANOVA Table

SV	df	SS
Reg on x_1	1	$\hat{\beta}_1 S_{yy}$
Reg on x_2	1	$\hat{\beta}_2 S_{yy}$
Error	$n-3$ no of variables	SSE
Total	$n-1$	$SST = S_{yy} = \sum y^2$

$$S_{yy} = \sum y^2 - \frac{\sum y^2}{n}$$

$$S_{yy} = \sum x_1 y - \sum x_2 y$$

$$SSE = \hat{\beta}_1 \sum x_1 y + \hat{\beta}_2 \sum x_2 y$$

SV	df	SS	MJ	F
Reg	p-1			
Error	$n-p$			
Total	$n-1$			

here
 $n = \text{no of all variables}$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 \quad \text{--- (1)}$$

from (1) multiply through by Σ

$$\Sigma Y = \Sigma \beta_0 + \beta_1 \Sigma X_1 + \beta_2 \Sigma X_2$$

$$\Sigma Y = n\beta_0 + \beta_1 \Sigma X_1 + \beta_2 \Sigma X_2 \quad \text{--- (2)}$$

from (1) multiply through by ΣX_1

$$\Sigma XY = \beta_0 \Sigma X_1 + \beta_1 \Sigma X_1^2 + \beta_2 \Sigma X_1 X_2 \quad \text{--- (3)}$$

from (1) multiply through by ΣX_2^2

$$\Sigma X_2 Y = \beta_0 \Sigma X_2 + \beta_1 \Sigma X_1 X_2 + \beta_2 \Sigma X_2^2 \quad \text{--- (4)}$$

combine (2), (3), (4)

$$\Sigma Y = n\beta_0 + \beta_1 \Sigma X_1 + \beta_2 \Sigma X_2 \quad \text{--- (5)}$$

$$\Sigma XY = \beta_0 \Sigma X_1 + \beta_1 \Sigma X_1^2 + \beta_2 \Sigma X_1 X_2 \quad \text{--- (6)}$$

$$\Sigma X_2 Y = \beta_0 \Sigma X_2 + \beta_1 \Sigma X_1 X_2 + \beta_2 \Sigma X_2^2 \quad \text{--- (7)}$$

$$\Sigma X_2 Y = \beta_0 \Sigma X_2 + \beta_1 \Sigma X_1 X_2 + \beta_2 \Sigma X_2^2 \quad \text{--- (8)}$$

$$\begin{pmatrix} \Sigma Y \\ \Sigma XY \\ \Sigma X_2 Y \end{pmatrix} = \begin{pmatrix} n & \Sigma X_1 & \Sigma X_2 \\ \Sigma X_1 & \Sigma X_1^2 & \Sigma X_1 X_2 \\ \Sigma X_2 & \Sigma X_1 X_2 & \Sigma X_2^2 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}$$

$$\begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} n & \Sigma X_1 & \Sigma X_2 \\ \Sigma X_1 & \Sigma X_1^2 & \Sigma X_1 X_2 \\ \Sigma X_2 & \Sigma X_1 X_2 & \Sigma X_2^2 \end{pmatrix}^{-1} \begin{pmatrix} \Sigma Y \\ \Sigma XY \\ \Sigma X_2 Y \end{pmatrix}$$

$$\hat{\beta} = X\beta$$

$$Y = X\beta + \epsilon$$

establishing this allows for estimation (1)

quantifying bias and variance for residuals (2)

standardizing the distribution of unbiased

estimators are called unbiased estimators (3)

if two different

- ① Multicollinearity
- ② Heteroscedasticity
- ③ Autocorrelation

MULTICOLLINEARITY

- Presence of perfect linear or nearly perfect linear relationship among some or all the predictor variables of the regression model.
- The parameters becomes Indeterminate and OLS because either of the predictor variables are perfectly already correlated.
- There is no problem of multicollinearity if regression equation has correlation coefficient ($x_i x_j$) for the i 's is zero. Then each parameter β_i can be estimated by SLR of y on x_i .
- Multicollinearity can be present, But the issue is its magnitude or severity.
- Multicollinearity is common / inherent in most economic relationships.

Consequences of Multicollinearity

- OLS is still BLUE {Best Linear Unbiased Estimate} But have large variances & covariances; Making precision estimation difficult.
- If there is perfect correlation;
 - ① Estimates of coefficients are Indeterminate
 - ② Std error of estimates are infinitely large
- Under confidence interval, t ratio of one or more coefficients tends to be statistically significant.
- R^2 is high, even if t-ratio of one or more coefficients is statistically insignificant. {There must be agreement between t-ratio and R^2 }

- OLS estimators and their standard errors can be very sensitive to ~~large~~ small changes in the data
- Small changes in the data can change values of OLS estimators and their standard errors

Detection of Multicollinearity

(5)

- VIF • High R^2 but few significant t-ratios

- High pair-wise correlation among regressors

- Through partial correlations

Auxiliary regression procedure $F_j = \frac{R_j^2}{(1-R_j^2)(n-k)}$

Rule: If $F_j > F_{(k-1)(n-k)}$, we reject H_0

H_0 : X_j is not collinear

H_1 : X_j is collinear

OR Use Klein Rule of thumb: if $R^2 < 0.5$

H_0 : No multicollinearity

H_1 : Multicollinearity exists

Rule: If $R_j^2 > R^2$ (overall from Reg of all X_i), reject H_0

Auxiliary Regression is the regression of the predictor IV on regressors.

We regress X_j on remaining X_i and compute R_j^2

Each of these regressions is known as Auxiliary Regression.

Partial correlation: measuring the association b/w 2 variables while controlling the effect of one or more additional variables

Considering X_1 X_2 X_3

Partial correlation of X_1 and X_2 is done by keeping X_3 fixed

X_2 and X_3 is $\checkmark - \checkmark$ X_1 fixed

X_1 and X_3 $- - \checkmark \checkmark$ X_2 fixed,

VIF (Variance Inflation Factor)

- To measure how variables are associated with each other
- VIF reflects how the variance of an estimator is inflated by presence of Multicollinearity
- VIF is used to decide/detect Multicollinearity between variables

For X_1 and X_2 ,

$$VIF = \frac{1}{1 - R^2_{X_1, X_2}}$$

- High correlation between X_1 and X_2 suggest high VIF

R^2_j = coeff of determination of reg of X_j on remaining regressors

If R^2_j increases towards Unity (i.e 1) collinearity of X_j with other regressors increases; the VIF also increases and can even be finite

- If $VIF > 10$, X_j is highly collinear
- The higher the value of VIF, the more collinear the X_j
- e.g if $R^2_j > 0.90$

High values of VIF are usually considered to indicate multicollinearity between regressors

Example
 Response predictor $X_6 = 6$
 $n = 70$. R^2 values for the Auxiliary Regression on X_6 are given
 and also given is R^2 overall for Reg of γ on all X_5

Test for multicollinearity Using ① F-test procedure
 ② Klein's rule.

X_j	R_j^2	R^2_{overall}	F_j	F_{tab}	$F_{\text{Conclusion}}$	Klein Conclusion
X_1	0.982	0.623	698.311	2.362	> *	> *
X_2	0.803	0.623	520.746	2.362	> *	> *
X_3	0.651	0.623	23.876	2.362	> *	> *
X_4	0.314	0.623	5.8589	2.362	> *	< **
X_5	0.213	0.623	3.4643	2.362	> *	< **
X_6	0.124	-	1.8119	2.362	< **	< **

Hypothesis: F-procedure Klein rule
 $H_0: X_j \text{ is not collinear}$ $H_0: \text{No multicollinearity}$
 $H_1: X_j \text{ is collinear}$ $H_1: \text{Multicollinearity exists}$

Test statistic:

$$F_j = \frac{R_j^2 / (k-1)}{(1-R_j^2) / (n-k)} \quad \text{for } j=1, 2, 3, \dots, 6 \quad R^2_{\text{overall}} = 0.623$$

$$F_1 = \frac{R^2 (k-1)}{(1-R^2) (n-k)} = \frac{(0.982)(6-1)}{(1-0.982)(70-6)} = \frac{(0.982)/5}{0.018/64} = \text{not}$$

$$F_1 = \frac{0.1964}{0.00028125} = \frac{698.311}{18.0} = 38.8$$

(not significant)

$$F_2 = \frac{0.803/5}{(1-0.803)/64} = \frac{0.1606}{0.003078} = 52.1746$$

$$F_3 = \frac{0.651/5}{(1-0.651)/64} = \frac{0.1302}{0.005453} = 23.8762$$

$$F_4 = \frac{0.314/5}{(1-0.314)/64} = \frac{0.0628}{0.01072} = 5.8589$$

$$F_5 = \frac{0.213/5}{(1-0.213)/64} = \frac{0.0426}{0.012297} = 3.4643$$

$$F_6 = \frac{(0.124)/5}{(1-0.124)/64} = \frac{0.0248}{0.01369} = 1.8119$$

DR: If $F_j > F_{(k-1)(n-k)}$

$$F_{(6-1)(70-6)} = F_{(5, 64)} = 2.362$$

Alternatively, as $F_{5, 64}$ is not readily calculable in F-table.

$$F_{5, 64} = x \quad F_{5, 60} = 2.37 \quad F_{5, 720} = 2.29$$

$$\text{Then } \frac{2.37 - 2.29}{60 - 64} = \frac{2.37 - x}{60 - 64} = \frac{(1-x) / 4}{(4-6)(7-1)} = 7$$

$$\frac{0.08}{F_{60}} = \frac{2.37 - x}{4}$$

$$0.0013333 \times \frac{2.37 - x}{4}$$

$$0.005333 = 2.73 - 7x / 4 = 7$$

$$7x = 2.37 - 0.005333$$

$$x = \underline{\underline{2.365}} \quad \text{[As our F-table]}$$

Comments:

F-procedure: x_1, x_2, x_3, x_4 and x_5 are all collinear as we reject the null hypothesis. Only x_6 is not collinear as we fail to reject the null hypothesis.

Klein rule: x_1, x_2 and x_3 are collinear while

x_4, x_5 and x_6 are not collinear as

e.g. collinear $\rightarrow R_i^2, R_2^2, R_3^2 > R^2_{\text{overall}}$ and linear

Not collinear $\rightarrow R_4^2, R_5^2, R_6^2 \leq R^2_{\text{overall}}$

transformation to remove collinearity is done

collinearity is removed

base and not heterogeneous after transformation

check model using four plots (residuals, etc) and check

studentized residuals & influential points

interpretation of regression coefficients

standard deviation of error term is not constant

heteroscedasticity is present

homoscedasticity is assumed



Solutions to Multicollinearity

- ① Increase the sample size:
gather more observations. It is possible that in another sample of greater size multicollinearity may not be serious as seen in 1st sample taken.
Since covariances are inversely proportional to sample size, multicollinearity can be reduced.-
This is true when ① multicollinearity characterises the original sample and not the population.
- ② if multicollinearity is due to error of measurement
- ③ - Transformation of Variables
First difference & Ratio transformation can be used.
- ④ - Combine cross sectional data and time series data
- ⑤ - Use Factor Analysis or Principal component Analysis
- ⑥ - Introduce additional equations to express the relationship between multicollinear X_i . This can reduce multicollinearity.
Hint: We now have simultaneous equation system instead of one equation.

$$\text{Given } Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t \quad \text{--- (1)}$$

First differencing:

$$Y_{t-1} = \beta_0 + \beta_1 X_{1,t-1} + \beta_2 X_{2,t-2} + U_{t-1} \quad \text{--- (2)}$$

(1) minus (2)

$$Y_t - Y_{t-1} = \beta_1 + \beta_1 X_{1t} - \beta_1 X_{1,t-1} + \beta_2 X_{2t} - \beta_2 X_{2,t-2} + U_t - U_{t-1}$$

$$Y_t - Y_{t-1} = \beta_1(X_{1t} - X_{1,t-1}) + \beta_2(X_{2t} - X_{2,t-2}) + U_t - U_{t-1}$$

But let $V_t = U_t - U_{t-1}$, then

$$Y_t - Y_{t-1} = \beta_1(X_{1t} - X_{1,t-1}) + \beta_2(X_{2t} - X_{2,t-2}) + V_t \quad *$$

Eqn * is reg model of the differences, not that of the original variables.

NB: X_1, X_2 have high correlation doesn't mean their difference will have high correlation.

Ratio Transformation

$$\text{From } Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t \quad \text{--- (1) --- I}$$

Divide (1) by one of the X_i s say X_2 i.e. divide only X_{2t}

$$\frac{Y_t}{X_{2t}} = \frac{\beta_0}{X_{2t}} + \frac{\beta_1 X_{1t}}{X_{2t}} + \frac{\beta_2 X_{2t}}{X_{2t}} + \frac{U_t}{X_{2t}} \quad \text{--- II}$$

$$\frac{Y_t}{X_{2t}} = \frac{\beta_0}{X_{2t}} + \frac{\beta_1 X_{1t}}{X_{2t}} + \beta_2 + \frac{U_t}{X_{2t}}$$

NB: Transformation of error term may not satisfy some of OLS assumptions again.

OLS assumption

- Linearity
- Normality ϵ_i
- Homoscedasticity of error $\text{Var}(\epsilon_i) = \sigma^2$
- Full rank

Notables on Multicollinearity

If sole purpose of analysis is prediction, then multicollinearity is not a serious problem.

Reason: The higher the R^2 , the better the predictions.

II If purpose of ^{regression} analysis is reliable estimation of parameters, multicollinearity should not be ignored.

Biased estimates would have large standard errors; and this suggests multicollinearity.

VIF
Frobust
klein

Test of Multicollinearity.

Farrar and Glauber Set of Tests -

- 3 tests

I - Chi-square - detects existence & severity of multicollinearity.

II - F-test - locating variables that are multicollinear.

III - t-test - determine variables responsible for appearance of multicollinear variables

- I Chi-square
- Test whether the regressors are pairwise orthogonal or not
 - orthogonality mean $X_i X_j = 0 \quad \forall i \neq j$
 - We use standardized determinant, D
 - D is of order $(K+1)$ by $(K+1)$

$$D = \begin{vmatrix} 1 & r_{X_1 X_2} & \dots & r_{X_1 X_K} \\ r_{X_2 X_1} & 1 & r_{X_2 X_3} & \dots & r_{X_2 X_K} \\ \vdots & & & \ddots & \\ r_{X_K X_1} & r_{X_K X_2} & \dots & r_{X_K X_{K-1}} & 1 - r_{X_K X_K} \end{vmatrix}$$

$r_{X_1 X_2} = r_{X_2 X_1} = r_{X_2 X_3} = \dots = r_{X_K X_K} = 1$
 $r_{X_2 X_1} = r_{X_2 X_3} = \dots = r_{X_K X_1} = r_{X_K X_2}$

not at $K=3$ regressor

$$\frac{\text{std det}}{D} = \begin{vmatrix} 1 & r_{X_1 X_2} & r_{X_1 X_3} \\ r_{X_2 X_1} & 1 & r_{X_2 X_3} \\ r_{X_3 X_1} & r_{X_3 X_2} & 1 \end{vmatrix}$$

$r_{X_1 X_2} = r_{X_2 X_1}$
 $r_{X_2 X_3} = r_{X_3 X_2}$
 $r_{X_3 X_1} = r_{X_1 X_3}$

for $K=2$ regressor, But X_i s are perfectly linearly correlated.

$$\text{then } D = \begin{vmatrix} 1 & r_{X_1 X_2} \\ r_{X_2 X_1} & 1 \end{vmatrix} \Rightarrow 1 + r_{X_1 X_2}^2 = 1 + 1 = 2$$

then

$$D = \begin{vmatrix} 1 & 1 \\ 1 & 1 \end{vmatrix} \Rightarrow D = |1| - |1| = 1 - 1 = 0$$

for $K=2$ regressor, But X_i s are Orthogonal, i.e. $X_i X_j = 0$

$$D = \begin{vmatrix} 1 & r_{X_1 X_2} \\ r_{X_2 X_1} & 1 \end{vmatrix} \Rightarrow 1 + r_{X_1 X_2} r_{X_2 X_1} = 1 + 0 = 1$$

$$D = \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} \Rightarrow D = |1| - |0| = 1 - 0 = 1$$

Standardized det = 1

Hypothesis

H_0 : X_i s are orthogonal { No multicollinearity }

H_1 : X_i s are not orthogonal { Multicollinearity exists }

Test statistic χ^2

$$\chi^2_{\text{cal}} = - \left\{ n - 1 - \frac{1}{6} (2k + 5) \right\} \log D$$

$$\chi^2_{\text{tab}} = \chi^2_{\frac{1}{2}(K(K-1))}, \alpha$$

$K = \text{no of Regressors}$
OR
 no of predictors

DR

If $\chi^2_{\text{cal}} > \chi^2_{\text{tab}}$ reject H_0

Notables on D { Standardized determinant }

- If D lies btw 0 and 1, there is some degree of multicollinearity
- If D is close to 0 (zero) — then stronger the multicollinearity
- If D is close to 1 — the weaker the multicollinearity.

$$D = 1 - r = |x_1 - x_2| = 0 \leq 1$$

Example:

A cross-sectional sample of 10 families expenditure on clothing (y), disposable income (x_1), family size (x_2) and index for clothing (x_3) revealed the following

$$n = 10$$

$$K = 3 \text{ regressors}$$

$$\gamma = 20.3 + 2.1x_1 + 1.5x_2 - 0.2x_3 \quad \text{For Chi-square}$$

$$r_{x_1 x_2} = 0.82 \quad r_{x_1 x_3} = 0.70 \quad r_{x_2 x_3} = 0.43$$

$$R^2_{x_1, x_2 x_3} = 0.69 \quad R^2_{x_1, x_3} = 0.32 \quad R^2_{x_2, x_3} = 0.59 \quad F_{\text{test}}$$

$$r_{x_1 x_2 \cdot x_3} = 0.71 \quad r_{x_1 x_3 \cdot x_2} = 0.91 \quad r_{x_2 x_3 \cdot x_1} = 0.36$$

① Perform Forward and Glauber set of tests for multicollinearity at $\alpha = 0.05$

Solution

Chi-square

Hypothesis

H_0 : x_i s are orthogonal {No multicollinearity}

H_1 : x_i s are not orthogonal {Multicollinearity exists}

Test statistic

$$\chi^2_{\text{cal}} = - \left\{ n - 1 - \frac{1}{6} (2k + 5) \right\} \log_e D$$

DR If $\chi^2_{\text{cal}} > \chi^2_{\text{tab}}$, reject H_0

$$\chi^2_{\text{tab}} = \chi^2_{\frac{1}{2} k(k-1), \alpha}$$

$k = 3$ regressors/predictors

$$D = \begin{vmatrix} 1 & r_{X_1 X_2} & r_{X_1 X_3} \\ r_{X_1 X_2} & 1 & r_{X_2 X_3} \\ r_{X_1 X_3} & r_{X_2 X_3} & 1 \end{vmatrix}$$

$$D = \begin{vmatrix} 1 & -0.82 & 0.70 \\ -0.82 & 1 & 0.43 \\ 0.70 & 0.43 & 1 \end{vmatrix}$$

$$D = 1 \begin{vmatrix} 1 & 0.43 & -0.82 \\ 0.43 & 1 & 0.70 \\ -0.82 & 0.70 & 1 \end{vmatrix} + 0.70 \begin{vmatrix} 0.82 & 1 & 0.70 \\ 0.70 & 1 & 0.43 \\ 0.43 & 0.70 & 1 \end{vmatrix}$$

$$D = 1 \{1 - 0.1849\} - 0.82 \{0.82 - 0.301\} + 0.70 \{0.3526 - 0.701\}$$

$$D = 1 (0.8151) - 0.82 (0.519) + 0.70 (-0.3474)$$

$$D = 0.8151 - 0.42558 - 0.24318$$

$$D = \cancel{0.1533} \quad 0.14634$$

$$\begin{aligned} R_{\text{cal}}^2 &= - \left\{ n - 1 - \frac{1}{6} (2k+5) \right\} \log_e D \\ &= - \left\{ n - 1 - \frac{1}{6} (2k+5) \right\} \ln D \\ &= - \left\{ 10 - 1 - \frac{1}{6} (23+5) \right\} \ln 0.1463 \\ &= - \left\{ 9 - \frac{1}{6} (28) \right\} \ln 0.1463 \quad \mathcal{E} = -1 \end{aligned}$$

$$\chi^2_{\text{cal}} = - \left\{ 9 - \frac{11}{6} \right\} \ln 0.1463$$

$$= - \left\{ \frac{54-11}{6} \right\} \ln 0.1463$$

$$\chi^2_{\text{obs}} = - \left\{ \frac{43}{6} \right\} \ln 0.1463$$

$$= - [7.1667] \ln 0.1463$$

$$= -7.1667 \times -1.922$$

$$= 13.77511407$$

$$\approx 13.7751$$

$$\chi^2_{\text{tab}} = \chi^2 \cdot \frac{1}{2} k(k-1), \quad \times$$

$$= \chi^2 \frac{1}{2} 3(3-1) , 0.05$$

$$= \chi^2 \frac{1}{2} 3(2)_{0.05} = \chi^2 \frac{6}{2}, 0.05$$

$$= \chi^2_{3, 0.05}$$

{ Check 3 under 0.05 }
straight up

$$\chi^2_{\text{tab}} = 7.815$$

DR: If $\chi^2_{\text{cal}} > \chi^2_{\text{tab}}$, reject H_0

$13.7751 > 7.815$, we reject H_0 at 0.05 SL

Conclusion: there is multicollinearity, i.e. X_i s are not orthogonal. Multicollinearity exists.

Test II : Feature test $\{ H_1 - P \} = \text{false}$

We use R^2 from auxiliary regressions. To locate variables which are collinear.

Hypothesis

$$H_0: R_{x_j \cdot x_1 x_2 \dots x_k}^2 = 0 \quad \{ x_j \text{ is not multicollinear} \}$$

$$H_1: R_{x_j \cdot x_1 x_2 \dots x_k}^2 \neq 0 \quad \{ x_j \text{ is multicollinear} \}$$

Test statistics:

$$F_j = \frac{R_{x_j \cdot x_1 x_2 \dots x_k}^2 / (k-1)}{(1 - R_{x_j \cdot x_1 x_2 \dots x_k}^2) / (n-k)} \sim F_{(k-1), n-k}$$

$$F_{\text{tab}} = F_{(k-1)(n-k)}$$

DR:

If $F_j > F_{\text{tab}}$, reject H_0 .

Example given suppose

$n=10$ $k=3$ regressors. Locate variable that are collinear

$$\text{find } R^2 \\ R_{x_j \cdot x_1 x_2 x_3}^2 = 0.69 \quad R_{x_j \cdot x_1 x_3}^2 = 0.32 \quad R_{x_j \cdot x_2 x_3}^2 = 0.5$$

$$H_0: R_{x_j \cdot x_1 x_2 \dots x_k}^2 \leq 0 \quad \{ x_j \text{ is not multicollinear} \}$$

$$H_1: R_{x_j \cdot x_1 x_2 \dots x_k}^2 \neq 0 \quad \{ x_j \text{ is multicollinear} \}$$

three planes intersecting at one point

Test statistic

$$F_j = \frac{R_{X_j}^2 \cdot X_1 X_2 \dots X_K / (k-1)}{(1 - R_{X_j}^2 \cdot X_1 X_2 \dots X_K) / (n-k)}$$

$$\Rightarrow R_{X_1}^2 \cdot X_2 X_3 = 0.69 \quad n=10 \quad k=3$$

$$F_{X_1 \text{ cal}} = \frac{0.69 / (3-1)}{(1 - 0.69) / (10-3)} = \frac{(0.69) / 2}{(0.31) / 7} = \frac{0.345}{0.044} \approx 7.84$$

$$R_{X_2 \cdot X_1 X_3}^2 = 0.32$$

$$F_{X_2 \text{ cal}} = \frac{0.32 / 2}{(1 - 0.32) / 7} = \frac{0.16}{0.68 / 7} = \frac{0.16}{0.0971} \approx 1.65$$

~~$$R_{X_3 \cdot X_1 X_2}^2 = 0.36$$~~

~~$$F_{X_3 \text{ cal}} = \frac{0.36 / 2}{(1 - 0.36) / 7} = \frac{0.18}{0.64 / 7} = \frac{0.18}{0.09143} \approx 1.97$$~~

~~$$R_{X_3 \cdot X_1 X_2}^2 = 0.59$$~~

~~$$F_{X_3 \text{ cal}} = \frac{0.59 / 2}{(1 - 0.59) / 7} = \frac{0.295}{0.41 / 7} = \frac{0.295}{0.05857} \approx 5.03$$~~

$$F_{\text{tab}} = F_{(k-1)(n-k)}^{0.05} = F_{(3-1)(10-3)}^{0.05} = F_{2(7)}^{0.05} = 4.74$$

$$F_{\text{tab}} = F_{2,7,0.05} = 4.74$$

DR: If $F_{cal} \geq F_{tab}$, we reject H_0

• $7.84 > 4.74$, we reject H_0

• $1.65 < 4.74$, we fail to reject H_0

• $5.03 > 4.74$, we reject H_0

Conclusion:

X_1 is multicollinear

X_2 is not multicollinear

X_3 is multicollinear

Test III: t-test

* We detect variable variables that causes the multicollinearity

* We use partial correlation coefficients and test for their standardized significance using t-test.

Hypothesis:

$H_0: \gamma_{X_i X_j \cdot X_1 X_2 \dots X_k} = 0$ { X_i and X_j are not responsible for multicollinearity}

$H_1: \gamma_{X_i X_j \cdot X_1 X_2 \dots X_k} \neq 0$ { X_i and X_j are responsible for multicollinearity}

Test statistic:

$$t_{cal} = \frac{(\gamma_{X_i X_j \cdot X_1 X_2 \dots X_k}) \sqrt{n-k}}{\sqrt{1 - \gamma_{X_i X_j \cdot X_1 X_2 \dots X_k}^2}}$$

$\sim t_{(n-k)}/2$

DR: If $|t_{cal}| > t_{tab}$ reject Ho at $\alpha = 0.8\%$

$$t_{tab} = t_{n-k}, \frac{\alpha}{2}$$

Example

$$n=10, k=3$$

define variable.

Given

$$r_{x_1 x_2 \cdot x_3} = 0.71 \quad r_{x_1 x_3 \cdot x_2} = 0.91 \quad r_{x_2 x_3 \cdot x_1} = 0.36$$

Hypothesis?

$$H_0: r_{x_i x_j \cdot x_1 x_2 \dots x_k} = 0$$

$\{x_i \text{ AND } x_j \text{ are not responsible for the}\}$
 $\text{multicollinearity -}$

$$H_1: r_{x_i x_j \cdot x_1 x_2 \dots x_k} \neq 0 \quad \{x_i \text{ AND } x_j \text{ are responsible for the}\}$$

 multicollinearity

Test statistic

$$t_{cal} = \frac{r_{x_i x_j \cdot x_1 x_2 \dots x_k} \sqrt{n-k}}{\sqrt{1 - r^2}} \sim t_{(n-k), \frac{\alpha}{2}}$$

$$t_{cal} = \frac{r \sqrt{n-k}}{\sqrt{1 - r^2}}$$

$$t_{tab} = t_{(n-k), \frac{\alpha}{2}} = t_{(10-3), 0.05} = t_{7, 0.025} = 2.865$$

$$t_{cal} = 2.365$$

$$t_{\text{cal}1} = \bar{t}_1 = \frac{r_{x_1 x_2 \cdot x_3} \sqrt{n-1}}{\sqrt{1 - r_{x_1 x_2 \cdot x_3}^2}}$$

$$= \frac{(0.71) \sqrt{10-3}}{\sqrt{1 - (0.71)^2}} = \frac{(0.71)(\sqrt{7})}{\sqrt{1 - 0.5041}} = \frac{1.8785}{\sqrt{0.4959}}$$

$$= \frac{1.8785}{0.7042} = \frac{2.667566032}{\simeq 2.67}$$

$$t_{\text{cal}2} = \bar{t}_2 = \frac{r_{x_1 x_3 \cdot x_2} \sqrt{n-1}}{\sqrt{1 - r_{x_1 x_3 \cdot x_2}^2}}$$

$$= \frac{0.91 \sqrt{10-3}}{\sqrt{1 - (0.91)^2}} = \frac{0.91(\sqrt{7})}{\sqrt{1 - 0.8281}} = \frac{0.91 \times 2.645}{\sqrt{0.1719}}$$

$$= \frac{2.407678}{0.4146082488} = \frac{5.807115528}{\simeq 5.81}$$

$$t_{\text{cal}3} = \bar{t}_3 = \frac{r_{x_2 x_3 \cdot x_1} \sqrt{n-1}}{\sqrt{1 - r_{x_2 x_3 \cdot x_1}^2}}$$

$$= \frac{0.36 \sqrt{7}}{\sqrt{1 - (0.36)^2}} = \frac{0.36(2.645)}{\sqrt{1 - 0.1296}}$$

$$= \frac{0.952488}{\sqrt{0.18704}} = \frac{0.952488}{0.9329523032} = \frac{1.0209365}{\simeq 1.02}$$

$$t_{tab} = t_{(n-1), \frac{\alpha}{2}} = t_{7, 0.025} = 2.365$$

D12: If $t_{cal} > t_{tab}$, reject H_0

$-2.67 > 2.365$, we reject H_0 .

$-5.81 > 2.365$, we reject H_0

$-1.02 < 2.365$, we fail to reject H_0

Conclusion:

X_1 AND X_2 are intercorrelated. X_1 and X_2 are responsible for multicollinearity

X_1 and X_3 are intercorrelated, X_1 and X_3 are responsible for multicollinearity

X_2 and X_3 are not intercorrelated, X_2 AND X_3 are not responsible for multicollinearity

$$t_{cal} = t_{(n-1) \frac{\alpha}{2}} = t_{7, 0.025} = 2.365$$

DR: If $t_{cal} > t_{tab}$, reject H_0

. $2.67 > 2.365$, we reject H_0 , X_1

. $5.81 > 2.365$, we reject H_0

. $-0.2 < 2.365$, we fail to reject H_0

Conclusion?

X_1 AND X_2 are intercorrelated. X_1 and X_2 are responsible for multicollinearity

X_1 and X_3 are intercorrelated. X_1 and X_3 are responsible for multicollinearity

X_2 and X_3 are not intercorrelated, X_2 AND X_3 are not responsible for multicollinearity

Auto correlation refers to degree of correlation between different observations in the data.
Some variables are different observations in the data.

Autocorrelation

$$\begin{cases} E(u_i u_j) = 0 \\ \text{Cov}(u_i u_j | X_i X_j) \neq 0 \end{cases}$$

- Assumption of Auto Correlation
- This is about error terms.
- error terms are independent i.e. error term relating to any observation is influenced by the error term relating to any other observation.
- i.e. one error term of observation doesn't influence error term of other observation
- Successive values of error terms are independent.

Example

- ① Effect of labour strike on output over fuel price hike should affect output over fuel price hike. Should not affect output in the following month. {i.e. successive values of u_i are independent}
- ② Cross-sectional data: Effect of an increase in one family's income on its consumption is not expected to affect consumption expenditure of another family.

Serial correlation is also known as auto-correlation.

It occurs when error of one time period a is correlated with error for a subsequent time period b.

$$\text{i.e. } E(u_i u_j) \neq 0 \text{ as well as } \text{cov}(u_i u_j | x_{itj}) \neq 0$$

• Autocorrelation is the correlation between members of series of observations ordered in time {time series data} or ordered in space {cross-sectional data}

• Auto correlation is the relationship between successive values of same variable. BUT NOT relationship between two or more variables.

+ Auto correlation occurs well in time series. It can be found in time series data.

- Auto correlation occurs least in cross sectional data except when sample is random.

• Auto correlation ranges btw -1 and +1

Sources of autocorrelation ④

① Omitted predictor/explanatory Variables

② Mis-specification of Model ③ Mis-specification of true error term + random

④ Interpolations in observations

{Estimation of a value with 2 known value in sequence of values}

5 Consequences of Autocorrelation.

- 1 OLS estimates are not BLUE
- 2 OLS estimates remain Unbiased
- 3 OLS estimates have larger Variances
- 4 Underestimated Variance of error term
- 5 Prediction based on OLS estimates may not be efc

Solutions to Autocorrelation.

→ Solution depends on whether autocorrelation is from misspecification or it is a case of pure autocorrelation.
If it is a case of ~~exist~~ omitted explanatory variable the solution/remedy is to include such variables.
Otherwise we resort to ~~the~~ transformation of original ~~model~~ model leading to method of Generalized Least square (GLS).

In GLS, we use knowledge of & usually not known. It is usually obtained through iterative - process.

Autoregression of order P

$$U_t = \phi_1 U_{t-1} + \phi_2 U_{t-2} + \dots + \phi_p U_{t-p} + V_t$$

where $E(V_t) = 0$ $E(V^2) = \sigma_v^2$ $E(V_i V_j) = 0 \quad \forall i, j$

If $\phi = 1$, then we have first order autoregressive AR(1)

$$U_t = \phi U_{t-1} + V_t$$

AR(1) is mostly used for autocorrelation problem

AC(1) is ~~also~~ first order Autocorrelation.

Test for Autocorrelation.

- Von Neumann test
- Durbin-Watson test ✓
- Run test
- Breusch-Godfrey (BG) test

Detection of Autocorrelation

- Plot residuals e_t vs e_{t-1} residuals
- plot residuals e_t vs time (t)

Durbin Watson Test

- Suitable for first order autoregressive

AR(1)

$$U_t = \rho U_{t-1} + V_t$$

$H_0: \rho = 0 \quad \{ U_t \text{ (errors) are not auto correlated} \}$

$H_1: \rho \neq 0 \quad \{ U_t \text{ (errors) are auto correlated} \}$

$E(V_t)$
$E(U_t V_t)$
$\text{cov}(V_t, V_t)$
$V(Y_t)$

$$dW = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}$$

$$dW = \frac{\sum_{t=2}^n (e_t^2 - 2e_t e_{t-1} + e_{t-1}^2)}{\sum_{t=1}^n e_t^2}$$

Opening brackets $\sum_{t=2}^n$

$$f_w = \sum_{t=2}^n e_t^2 - 2 \sum_{t=2}^n e_t e_{t-1} + \sum_{t=2}^n e_{t-1}^2$$

$$\sum_{t=1}^n e_t^2$$

$$f_w = \sum_{t=2}^n e_t^2 + \sum_{t=2}^n e_{t-1}^2 - 2 \sum_{t=2}^n e_t e_{t-1}$$

$$\sum_{t=1}^n e_t^2$$

But recall that

$$\sum_{t=1}^T \hat{e}_t^2 \approx \sum_{t=2}^T \hat{e}_{t-1}^2 \text{ and } \sum_{t=2}^T \hat{e}_t \hat{e}_{t-1} \text{ are}$$

approximately same, then it becomes,

$$f_w = \frac{\sum \hat{e}_t^2 + \sum \hat{e}_{t-1}^2 - 2 \sum \hat{e}_t \hat{e}_{t-1}}{\sum \hat{e}_t^2}$$

$$f_w = \frac{2 \sum \hat{e}_t^2 - 2 \sum \hat{e}_t \hat{e}_{t-1}}{\sum \hat{e}_t^2}$$

$$f_w = \frac{2 \left(\sum \hat{e}_t^2 - \sum \hat{e}_t \hat{e}_{t-1} \right)}{\sum \hat{e}_t^2}$$

$$f_w = \frac{2 \left(\sum \hat{e}_t^2 - \sum \hat{e}_t \hat{e}_{t-1} \right)}{\sum \hat{e}_t^2}$$

$$f_w = 2 \left(\frac{\sum \hat{e}_t^2 - \sum \hat{e}_t \hat{e}_{t-1}}{\sum \hat{e}_t^2} \right)$$

$$f_w = 2 \left(1 - \frac{\sum \hat{e}_t \hat{e}_{t-1}}{\sum \hat{e}_t^2} \right)$$

$$f_w = 2 \left(1 - \rho \right)$$

$$\text{Where } \rho = \frac{\sum e_t e_{t-1}}{\sum e_t^2}$$

$$dw = 2(1 - \rho)$$

For correlation, $d = 2$

+ve correlation; $0 < \rho < 1 \quad 0 < d < 2$
-ve correlation; $-1 < \rho < 0 \quad 2 < d < 4$

i.e. d lies betw 0 and 4

i.e. $0 \leq d \leq 4$

DW assumptions

⑥

- Regression Model Includes Intercept term
- Regressors are truly exogenous
- U_t are from first order AR(1) autoregressive
 - U_t is normally distributed
- Reg Model doesn't have lagged values of dependent variable (Y)
- No missing obs in the data

Limitations/problems of DW

- ① If there are lagged ~~variable~~ values of dependent variable, we cannot use DW.
- ② There is inconclusive region.
Inconclusiveness of results may arise
- ③ Can only be used for first order auto-regressive AR(1)
i.e It ~~can~~ cannot be used for higher order auto correlation.