

Jan 24, 2025

Review of data and statistics

Classification of data

1. Cross-sectional / Time Series / Panel Data

(Combined form of cross-sectional & Time series)

Cross-sectional: Collecting the data in a single frame of time from large number of respondents.

e.g: Employee satisfaction survey

→ conducting a survey only once

Structure of data:

Respondent	Variable 1	Variable 2	...	Variable K
Household	Income	Expenditure		Net worth
Time	Age	Gender		Education level

Purposes of conducting cross-sectional study

- * Identifying the current status of study area
 - ↓
 - central values, graphical representation
- * Comparison of data by analytical domain.
 - ↓
 - provinces, age group, gender etc
- * To identify the relationship among variables.
- * Prediction and estimation.
- * To make plans, policies and strategies.

Time Series

Collection of data over a period of time

Used in macroeconomics

Tranquillity

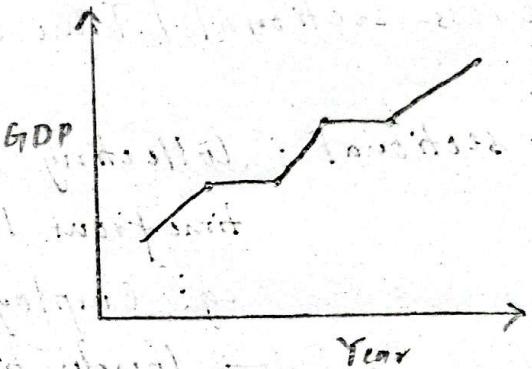
Year GDP

Year vs Remittance

Year vs Profit

- * Analysis of past trends
- * Forecasting for the future
- * Construction of business cycle
- * Studying the seasonal impact

Regression in Time Series



Nepal driving
Standard Survey
by World Bank

2. Categorical vs Numerical

Categorical: Any question having alphabetic/grouping response

Gender	Frequency	Percentage
Female	1	
Male	2 ✓	

Numerical: Any question having a numeric response

Salary \$1200

Min
Max
Mean
Standard Deviation

3. According to number of variables

Univariate data

Canonical correlation

Bivariate data

Multivariate data

dependence
technique

Dependent variable → one
Independent variable →
multiple regression

Jan 25, 2025

Population : Totality of all items under the study
(study area — coverage)

Census: Complete enumeration

Collecting the information / data from each and every element of population

Parameter: Result / outcome of population data

Characteristics of population data

Population size → N

Population mean → μ

Population proportion → π (or P)

Population variance → σ^2

Population correlation → ρ

Population regression → β_i
co-efficient

Sample: Subset of the population

Sampling: Collecting the information / data from sample.

Estimator / Statistic: Result or outcome of sample data.

Characteristics of sample size.

Sample size → n

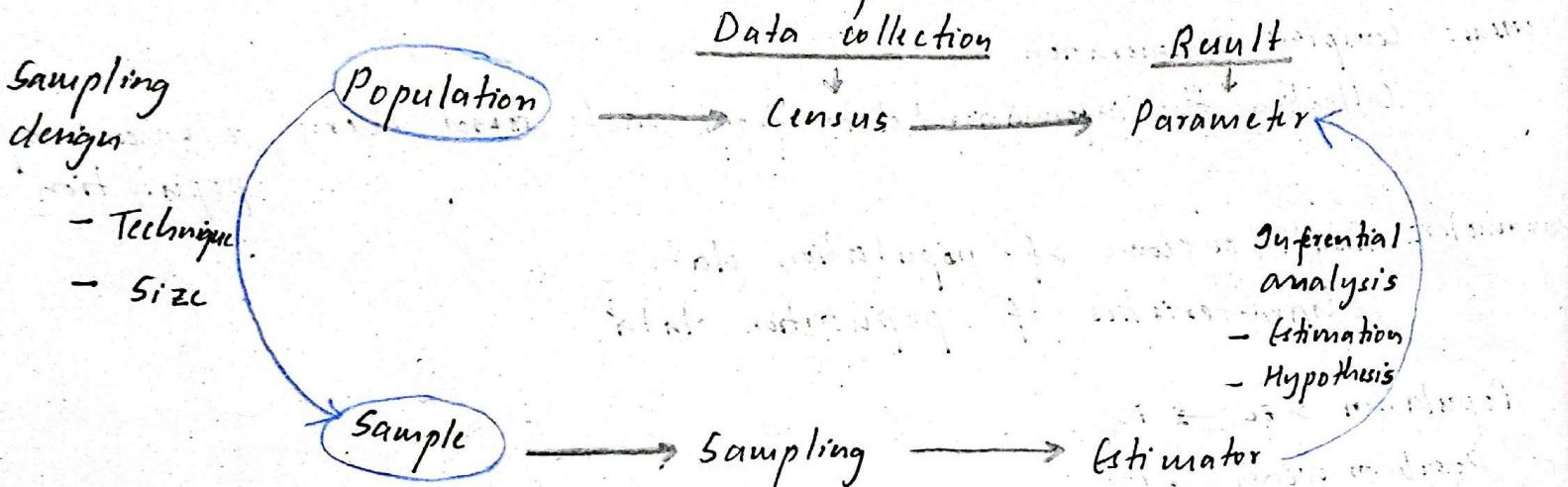
Sample mean → \bar{X} ($\hat{\mu}$)

dist. of all
population
items →
frame

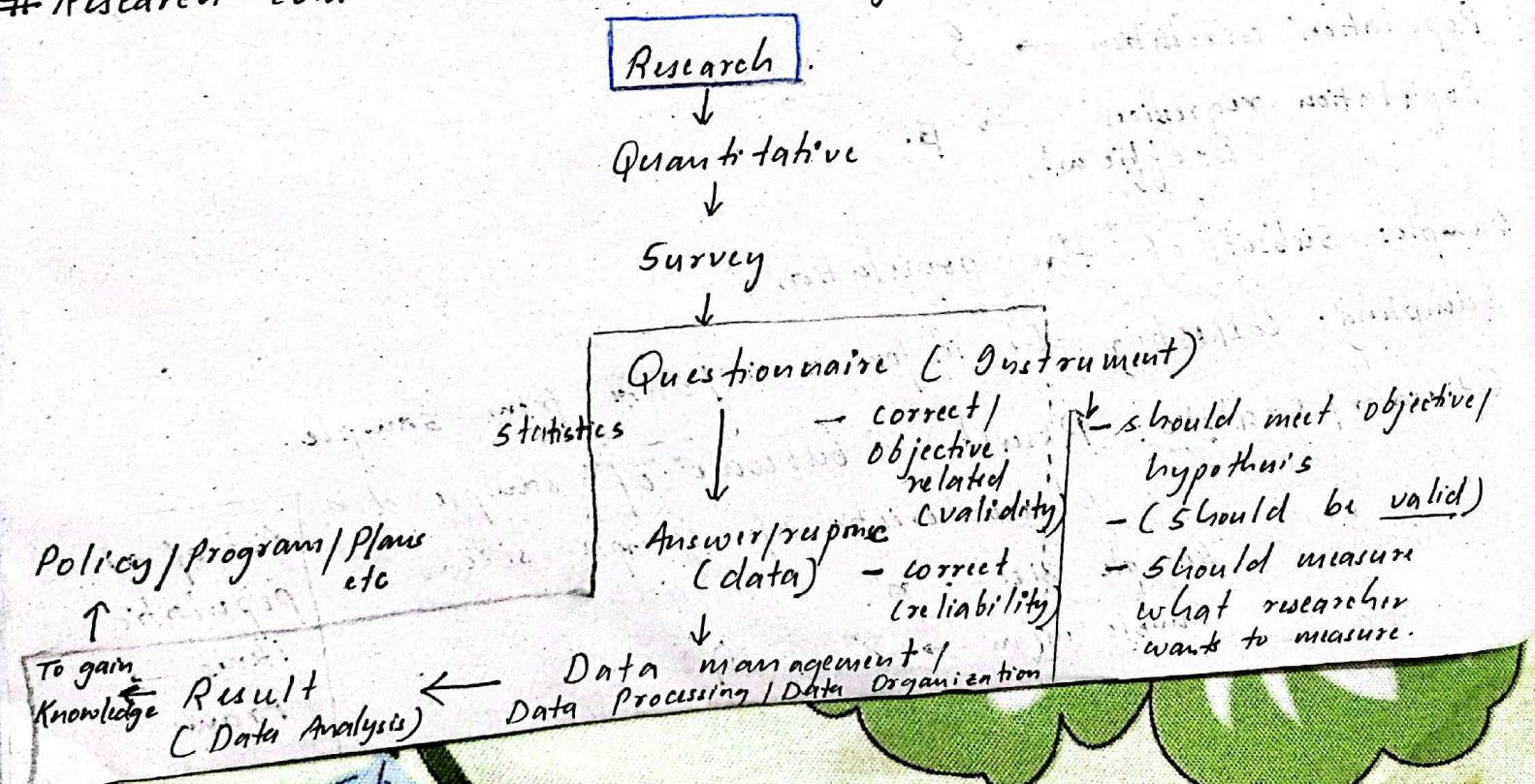
$$\begin{aligned}
 \text{Sample proportion} &= p \quad (\hat{p}) \\
 \text{Sample variance} &\rightarrow s^2 \quad (\hat{s}^2) \\
 \text{Sample correlation} &\rightarrow r \quad (\hat{r}) \\
 \text{Sample regression coefficient} &\rightarrow b_i \quad (\hat{\beta}_i)
 \end{aligned}$$

Sampling Error: Difference between Estimator and Parameter
 \downarrow
 $= |\text{Estimator} - \text{Parameter}|$

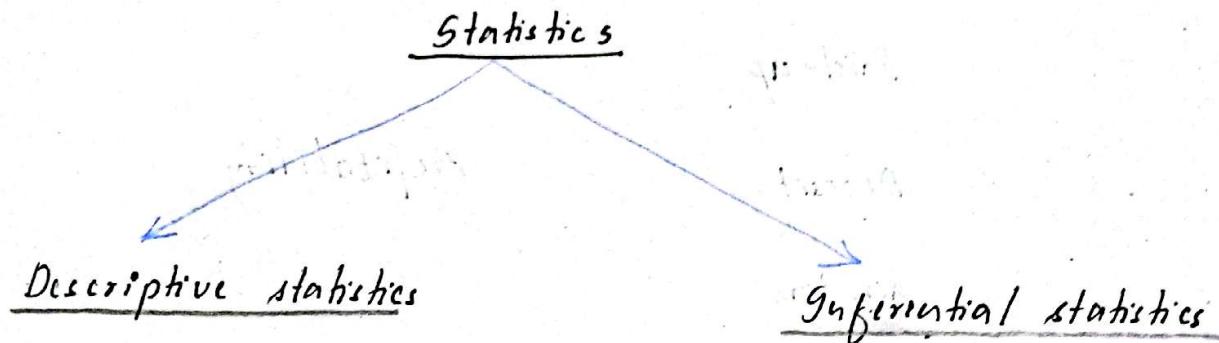
Used to determine the sample size.



Research can be done in two ways:



Statistics is the collection, organization and analysis of data.



Describing the status

- Table / chart
- central values
- Frequency / %
- Cross-Tabulation

Techniques of giving conclusion about the pop. parameter on the basis of sample data / statistic / estimator.

- consists of
- Estimation
 - Hypothesis

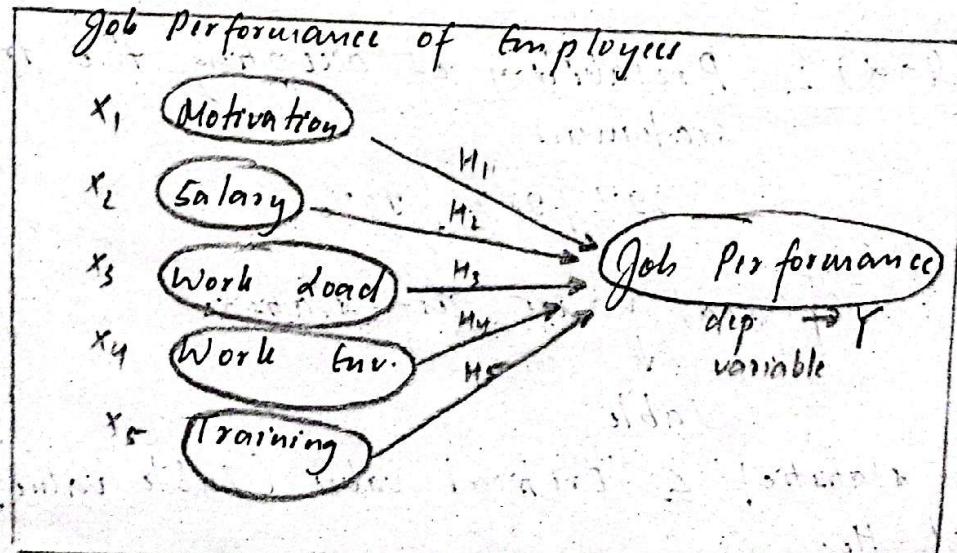
Hypothesis

Pre-assumption / statement / guess about the population parameter.
(outcome / result of the research)

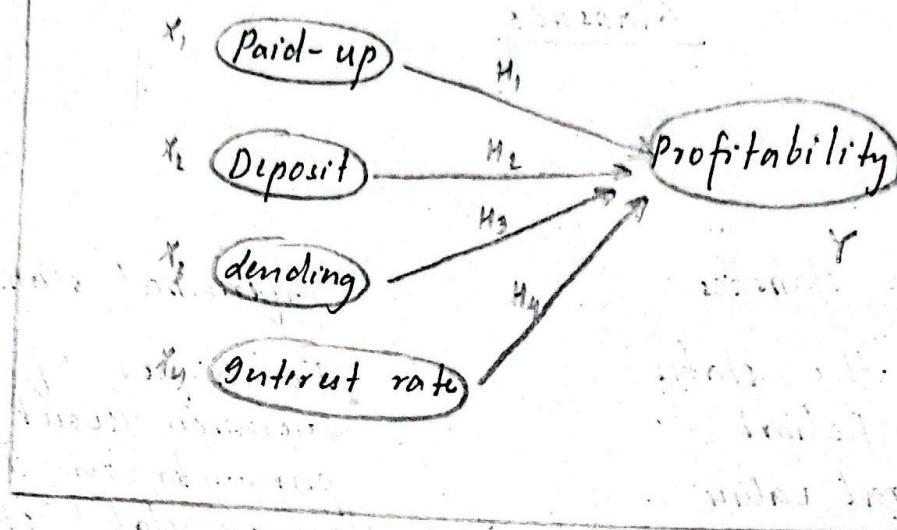
→ Set on the basis of Past literature

→ Conduction of Qualitative research / Pilot study

Null hypothesis: Statement of no difference or no relationship.
 (H_0)
=, \geq , \leq (equality)



Banks Profitability



Alternative hypothesis (H_1): Statement of relationship / difference.

$\neq, >, <$
↓
(Two tailed test) (One Tailed Test)

$n \rightarrow N$

level of Significance (α): Probability of rejecting true statement.
Estimator \rightarrow Parameter
1% / 5% / 10%

Confidence level ($1 - \alpha$): Probability of accepting the true statement

99% / 95% / 90%

Decision Criteria: Critical value / Test statistic

↓

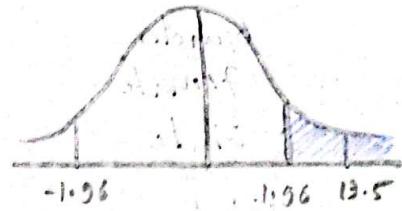
Table

If $|test\text{-statistic}| \leq \text{critical value (Table value)}$
do not reject H_0 .

p-value approach

if $p\text{-value} < \alpha$ (reject H_0)

if $p\text{-value} \geq \alpha$ (do not reject H_0)



Test statistic

$$t = 13.5$$

Table (α) = 5%

Example

Q1 Gender

Female - 1

Male - 2

Nominal

Q2 Post □

Assistant

1 } Ordinal

Officer

2 } Cross-
Sectional

Manager

3 }

Q3 Income 34000

Ratio

Q4 Expenditure 15000

Q1	Q2	Q3	Q4
1	2	34000	15000

Q1 → Frequency %.

Q2 → Frequency %.

Q3 → Min, Max, Mean, SD

Q4 →

Testing or comparing

Q1, Q2 → χ^2 Test

Q1, Q3 → ANOVA Test

Q3, Q4 → Z-test, t-test etc

Reference Book
Econometrics
 - Wooldridge

Data analysis:

I. Univariate Analysis

Case a: Analysis of categorical data

Frequency / Percentage analysis
 Graph / Visualization: Bar diagram or pie chart

* Distribution by Gender

Gender	Frequency	%
Female		
Male		
Total		

Case 6: Analysis of numerical data

Min / Max / Mean / St. deviation

Test: One sample t-test

Graph: Histogram / Box-plots

II. Bivariate Data Analysis

Analysis of two variables at a time

Case a: If both are categorical / grouping

Cross tabulation between two variables with frequency and appropriate %

Test → Chi-Square Test
(χ^2 Test)

Gender	Assistant	Officer	Manager	
Female	70 (0.7)	20 (0.2)	10 (0.1)	100
Male	100 (0.5)	50 (0.25)	50 (0.25)	200
	170	70	60	300

Test: χ^2 -Test

Case b: if one is categorical and another is numerical

Title: # Comparison of numerical by categorical

Gender	Avg	Min	Max	SD
Female				
Male				
Total				

Comparison between two means \Rightarrow Independent sample t-test

Title: # Comparison of income by Post

Post	Avg	Min	Max	SD
Art				
Officer				
Manager				
Total				

Comparison b/w three means \Rightarrow ANOVA (F-test)

NOTE:

If categorical variable consists of exactly two sub-groups, we use independent sample t-test, else for more than two subgroups, we use F-test.

Case c: if both are numerical

Correlation : measures the relationship between variables.

Regression : Should separate dependent and independent variables.

Feb 1, 2025

Developing an equation

$$Y = \beta_0 + \beta_1 x + \epsilon \quad (\text{Population})$$

$$Y = b_0 + b_1 x + e \quad (\text{Sample})$$

$Y \rightarrow$ Observed value (Actual value)

$\hat{Y} \rightarrow$ Predicted value / Estimated value

$Y - \hat{Y} =$ error / residual

Objectives

① Prediction of dependent variables by using information of independent variables.

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2$$

② Impact analysis

Impact of independent variables on dependent variables. (Slope / Coefficient)

③ Computation of amount / percentage change in dependent variable per unit / % change in independent variables. [meaning of slope]

Income	Expenditure
X	Y

$$\hat{Y} = b_0 + b_1 x$$

$$b_1 = \frac{\Delta Y}{\Delta X} = \frac{\text{amount change in } Y}{\text{percent change in } X}$$

$$\ln \hat{Y} = b_0 + b_1 \ln X$$

$$b_1 = \frac{\Delta Y/Y}{\Delta X/X} = \frac{\% \text{ change in } Y}{\% \text{ change in } X}$$

↓

(Elasticity)

$$\ln \hat{Y} = b_0 + b_1 X$$

$$b_1 = \frac{\Delta Y/Y}{\Delta X} = \frac{\% \text{ change}}{\text{unit change}}$$

Year	GDP
X	Y
	$\ln \hat{Y}$

$$\ln \hat{Y} = b_0 + b_1 X$$

~ Per annum

- ④ Computation of growth rate/amount per time period.
 - ⑤ Proportion of ~~explained~~ variation in Y (dependent variable) that is explained by independent variables.

[Coefficient of determination $\rightarrow R^2$]

Price	Sales	Volume
27.00	81.00	40.00
37.00	53.00	30.00
47.00	33.00	20.00

$b_1 = -7.5$

$R^2 = 0.82$

• १०५

$$R^2 = 0.82$$

\Rightarrow 82% variation in sales volume is explained by price.

- ⑥ To make policies/ programs/ plans/ decisions.

Simple Regression Analysis

dependent variable $\rightarrow Y$
 independent variables $\rightarrow X$

Equation

$$Y = \beta_0 + \beta_1 X + \epsilon \quad (\text{Population})$$

$$y = b_0 + b_1 X + e \quad (\text{Sample})$$

By using least squares method.

$$b_1 = \frac{n \sum XY - (\sum X)(\sum Y)}{n \sum X^2 - (\sum X)^2}$$

$$= \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

$$b_0 = \frac{\sum Y}{n} - b_1 \frac{\sum X}{n}$$

$$\hat{Y} = b_0 + b_1 X$$

Example

X (GDP) (in '000)	Y (GDP) (in '000)	XY	X ²	Y ²	$\hat{Y} = 2.206 + 0.572X$	(Y - \hat{Y})	(Y - \hat{Y}) ²
15	10	150	225	100	10.786	-0.786	0.6177
18	12	216	324	144	12.502	-0.502	0.2520
20	12	240	400	144	13.646	-1.646	2.7093
20	17	340	400	289	13.646	3.354	11.2493
25	17	425	625	289	16.506	0.494	0.2490
29	18	522	841	324	18.794	-0.794	0.6304
127	86	1893	2815	1290			
					$\sum = 0.12$	$\sum (Y - \hat{Y})^2 = 15.701$	≈ 0

$e = Y - \hat{Y} = \text{Observed value} - \text{Estimated value}$

$$b_1 = \frac{6 \times 1893 - 127 \times 86}{6 \times 1815 - 1127^2}$$

$$= \underline{0.572}$$

$$b_0 = \frac{\Sigma Y}{n} - b_1 \frac{\Sigma X}{n}$$

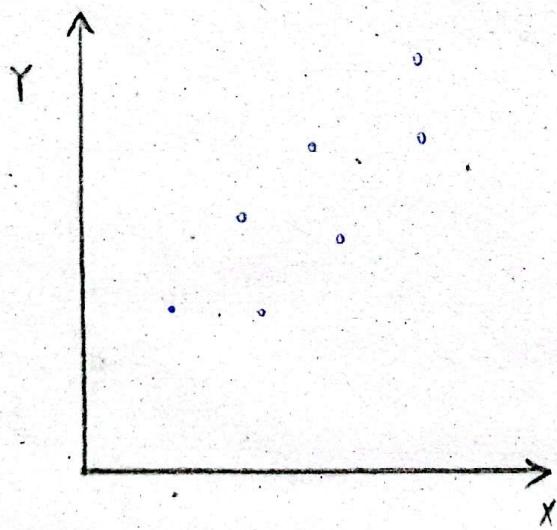
$$= \frac{86}{6} - 0.572 \times \left(\frac{127}{6} \right)$$

$$= \underline{2.206}$$

The estimating equation

$$\hat{y} = \underline{2.206} + 0.572 X$$

\downarrow \downarrow
 b_0 b_1



$$* E(e) = 0$$

* Signs should not be in a pattern (\Rightarrow Autocorrelation)

$$* \text{Standard error of estimate } (\hat{y}) = \sqrt{\frac{\sum e^2}{n-2}}$$

$$SE(\hat{y}) = S_{YX} = se(\hat{y})$$

$$\sqrt{\frac{\sum (Y - \bar{Y})^2}{n-p-1}}$$

$p \rightarrow$ no. of independent variables

$$= \sqrt{MSE}$$

* It measures the average variation/deviation of observed values of dependent variable (Y) around its fitted equation (\hat{Y}).

* We can use this for model selection.

$$Y = b_0 + b_1 x_1 + b_2 x_2^2 + b_3 x$$

linear

Feb 7, 2025

Multiple regression analysis

Example

Important Packages

car,
tidyverse,
lmtest

Sales Volume	Price	no. of stores	no. of ads
₹ 100	₹ 10	50	50

→ RStudio

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

where, β_0 = intercept

β_1 = parameter associated with X_1 ,

β_2 = parameter associated with X_2 and so on.

For sample data it is written as

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p + e$$

(VVI)

Assumptions of CLRM (Classical Linear Regression Model)

- Zero mean of ϵ_i $E(\epsilon_i) = 0$ for each i .
- Homoscedasticity $Var(\epsilon_i) = \sigma^2$ constant
- Non autocorrelation $Cov(\epsilon_i, \epsilon_j) = 0$ where $\epsilon_i \neq \epsilon_j$
- Normality: ϵ_i is normally distributed
- Non stochastic x_s , the values of x -variables are same in repeated samples.

- Zero covariance between ϵ_i and X variables.

$$\text{Cov}(\epsilon_i, X_{1i}) = \text{Cov}(\epsilon_i, X_{2i}) = 0$$

- No exact relationship exists between X variables i.e. X s are not correlated (no multicollinearity).

Model specifications and Assumption (in vector and matrix form)

- The general population regression model involving the dependent variable y_i and the independent variables $x_{1i}, x_{2i}, x_{3i}, \dots, x_{Ki}$ is specified as:

y_i	x_{1i}	x_{2i}	x_{3i}	\dots	x_{Ki}
y_1					
y_2					
\vdots					
y_n					

The multiple regression equation is written as:

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_K x_{Ki} + \epsilon_i$$

$$y_n = \beta_0 + \beta_1 x_{1n} + \dots + \beta_K x_{Kn} + \epsilon_n$$

The system of linear equations can be written as:

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{K1} \\ 1 & X_{12} & & & \vdots \\ \vdots & \vdots & & & \vdots \\ 1 & X_{1n} & & \dots & X_{Kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_K \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

That is

$$Y = X\beta + \epsilon \quad (1)$$

Assumptions

1) Zero mean of ϵ

That is

$$E(\epsilon) = E \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} = \begin{bmatrix} E(\epsilon_1) \\ E(\epsilon_2) \\ \vdots \\ E(\epsilon_n) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

2) Constant variance of ϵ

$$\text{Var}(\epsilon) = E(\epsilon\epsilon')$$

$$\begin{aligned} &= E \cdot \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \begin{bmatrix} \epsilon_1 & \epsilon_2 & \dots & \epsilon_n \end{bmatrix} \\ &= E \begin{bmatrix} \epsilon_1^2 & \epsilon_1\epsilon_2 & \dots & \epsilon_1\epsilon_n \\ \epsilon_2\epsilon_1 & \epsilon_2^2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \epsilon_n\epsilon_1 & \dots & \dots & \epsilon_n^2 \end{bmatrix} = \begin{bmatrix} E(\epsilon_1^2) & E(\epsilon_1\epsilon_2) & \dots & E(\epsilon_1\epsilon_n) \\ E(\epsilon_2\epsilon_1) & E(\epsilon_2^2) & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ E(\epsilon_n\epsilon_1) & \dots & \dots & E(\epsilon_n^2) \end{bmatrix} \end{aligned}$$

$$= \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \ddots & \ddots & \sigma^2 \\ 0 & \dots & \dots & \sigma^2 \end{pmatrix}$$

$$\sigma^2 \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & & \ddots & \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

$\therefore \text{Var}(\epsilon) = \sigma^2 I$ where I is $n \times n$ identity matrix.

3. Non stochastic x_s : This implies that all explanatory variables are non stochastic and hence, independent of ϵ_s .

OLS Estimation (derive OLS for multiple regression model)
given the population regression model

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

The sample regression model is

$$Y = X\hat{\beta} + e$$

Here, e is an estimate of ϵ .

$$\text{And } e = Y - X\hat{\beta}$$

The least squares estimators are obtained by minimizing the sum of squares and which is

$$\begin{aligned}\sum e_i^2 &= e_1^2 + e_2^2 + \dots + e_n^2 \\&= (e_1 \ e_2 \ \dots \ e_n)^t \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} \\&= e^t e \\&= (Y - X\hat{\beta})^t (Y - X\hat{\beta}) \\&= (Y^t - \hat{\beta}^t X^t) (Y - X\hat{\beta}) \\&= Y^t Y - Y^t X\hat{\beta} - \hat{\beta}^t X^t Y + \hat{\beta}^t X^t X\hat{\beta} \\&= Y^t Y - (\hat{\beta}^t X^t Y)^t - (\hat{\beta}^t X^t Y) + \hat{\beta}^t X^t X\hat{\beta}\end{aligned}$$

$$\therefore \sum e_i^2 = Y^t Y - 2\hat{\beta}^t X^t Y + \hat{\beta}^t X^t X\hat{\beta} \quad (1)$$

for least squares

$$\frac{\partial \sum G^2}{\partial \hat{\beta}} = 0$$

$$\Rightarrow \frac{\partial (Y^t Y - 2 \hat{\beta}^t X^t Y + \hat{\beta}^t X^t X \hat{\beta})}{\partial \hat{\beta}} = 0$$

$$\Rightarrow -2X^t Y + 2X^t X \hat{\beta} = 0$$

$$\Rightarrow (X^t X) \hat{\beta} = X^t Y \quad \text{--- (2)}$$

The equations contained in (2) are called OLS normal equations in the context of the general linear model.

Therefore,

$$\hat{\beta} = (X^t X)^{-1} (X^t Y) \quad \text{--- (3)}$$

The vector contains estimators for all unknown parameters.

Software output and interpretation

* Format of ANOVA Table

Source	Dof	SS	MSS	F-Value	P-Value
Regression	k	SSR	$MSE = SSR/k$	$F = \frac{MSE}{MSE}$	
Residual or Error	$n-k-1$	SSE	$MSE = SSE/(n-k-1)$		
Total	$n-1$	TSS or SST			

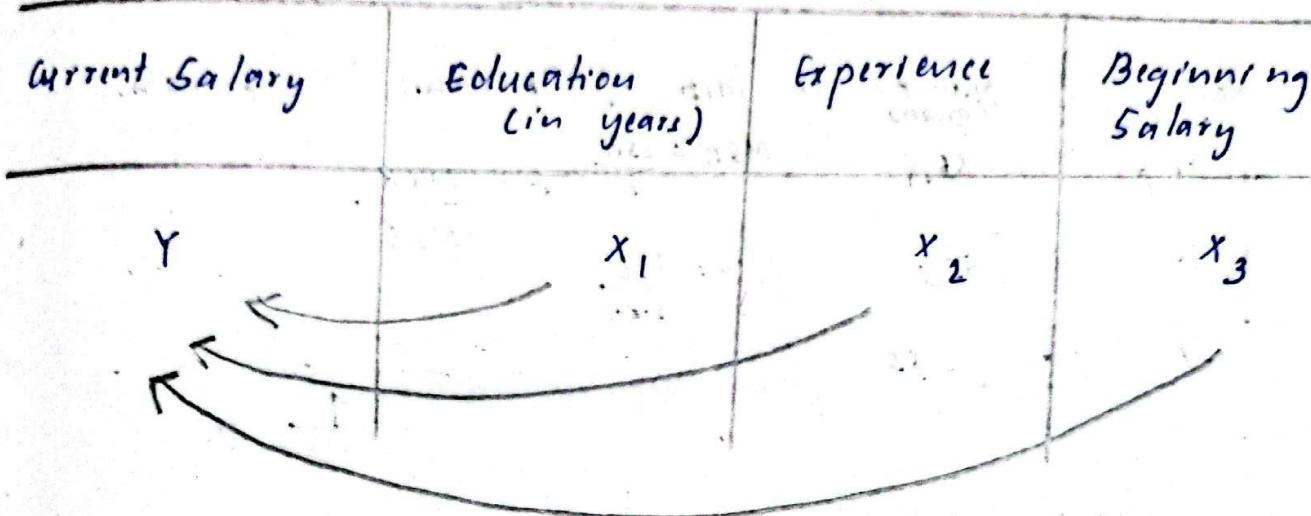
SST: Total sum of squares

n = Sample size, k = no. of independent variables

SSR = explained sum of squares SSE = residual sum of squares

F- Test: Overall Test

t- test: Individual Test



Individual Test

$$H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0 \quad (\text{Overall Test})$$

There is no significant impact of all independent variables on dependent variable.

H_1 : There is a significant impact of at least one independent variable on dependent variable.

$$p\text{-value} < \alpha(0.05)$$



Reject H_0

$$p\text{-value} > \alpha(0.05)$$



Accept H_0

Feb 8, 2025

Format of ANOVA Table

Source	dof	Sum of squares	MSB	F-value	p-value
Regression	k ($\approx p$)	SSR	$MSR = \frac{SSR}{k}$	$F = \frac{MSR}{MSE}$	
Error	$n-k-1$	SSE	$MSE = \frac{SSE}{n-k-1}$		
Total	$n-1$	TSS			

↑ Overall Test

Format of Coefficient Table

Predictor	Regression coeff.	st. error	t-value	p-value
Constant	b_0	s_{b_0}	$t = \frac{b_i}{s_{b_i}}$	
x_1	b_1	s_{b_1}		
x_2	b_2	s_{b_2}		
\vdots	\vdots	\vdots		
x_k	b_k	s_{b_k}		

↑ Individual Test

(VVI)

From ANOVA table and Coefficient table

① Developing an equation and prediction

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$$

② Computing multiple coefficient of determination

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

TSS = Total Sum of Squares (Total Variation)

SSR = Sum of squares due to regression (Explained)

Variation)

$SSE = \text{Sum of squares due to error}$

$$TSS = SSR + SSE$$

Thus,

$$R^2 = \frac{SSR}{TSS}$$

$$= 1 - \frac{SSE}{TSS}$$

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

Meaning: Suppose $R^2 = 0.81$, $k = 5$

\Rightarrow 81% of variation in dependent variable is explained by 5 independent variables.

③ Adj R^2

→ Used for model selection

$Y x_1$	$Y x_1 x_2$	$Y x_1 x_2 x_3$
R^2	R^2	R^2

↑
AIC
↓
VIF

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

AIC
VIF

④ Standard error of estimate

$$SE(\hat{y}) = s_{yn} = \sqrt{MSE}$$

It measures the average variation of observed values of dependent variables around its fitted equation.

⑤ Overall test / Goodness-of-fit

Hypothesis testing for all regression co-efficients simultaneously

Null hypothesis (H_0): $\beta_1 = \beta_2 = \dots = \beta_k = 0$

There is no significant impact of all independent (explanatory) variables on dependent variables.

Alternative hypothesis (H_1): at least one $\beta_i \neq 0$

There is a significant impact of at least one independent variable on dependent variable.

Test-statistic: $F = \frac{MSR}{MSE}$

$$\text{Calc } F = F$$

$$\text{Tab } F = F_{K, n-K-1} \text{ at } \alpha\%$$

$K = \text{dof for numerator}$

$n-K-1 = \text{dof for denominator}$

Decision: If calc $F \leq \text{Tab. } F$, we do not reject H_0 .

If calc $F > \text{Tab. } F$, we reject H_0 .

Errors
Sampling error
Non-sampling error
Standard error

p-value approach (Usually for software)

If p-value $< \alpha$, we reject H_0

If p-value $\geq \alpha$, we do not reject H_0 .

Individual test (t-test)

[Co-efficient table]

Null Hypothesis (H_0): $\beta_i = 0$. (i ranges from 1 to K)

There is no significant impact of Sales | Price | Ads
 x_i on y .

Alternative Hypothesis (H_1): $\beta_i \neq 0$

There is a significant impact of an independent variable on dependent variable.

Test Statistic:

$$t = \frac{b_i}{s b_i}$$

Decision: $| \text{cal } t | = |t|$

$\text{Tab } t = t_{n-k-1, \alpha/2}$.

If $| \text{cal } t | \leq \text{Tab } t$, we do not reject H_0 .

If $| \text{cal } t | > \text{Tab } t$, we reject H_0 .

Software

p-value $< \alpha \Rightarrow$ reject H_0

p-value $\geq \alpha \Rightarrow$ do not reject H_0 .

Confidence Interval Estimation
(Coeff table) $\leftarrow b_i \pm t_{n-k-1, \alpha} \xrightarrow{\text{Table} \rightarrow \text{statistical}} s b_i \rightarrow$ Coefficient table

Example

By using following ANOVA table obtained from 30 observations.

Source	df	SS	MSS	F-value
Regression	4	300	$MSR = \frac{300}{4} = 75$	$F = 75$
Error	25	200	$MSE = \frac{200}{25} = 8$	$= 9.375$
Total	29	500	$TSS = 500$	

$$TSS = 500$$

Complete

- Compute the given ANOVA table
- Compute standard error of estimate & interpret its meaning.
- Compute multiple coefficient of determination & interpret its meaning.
- Compute Adj R²
- Set up null and alternative hypothesis and carry out F-Test.

Solution

$$(b) SE(\hat{Y}) = \sqrt{MSE} = \sqrt{8} = 2.828$$

The avg. variation of observed values of dependent variables around fitted eqn is $\boxed{2.828}$

$$(c) R^2 = \frac{300}{500} = 0.60 \quad (60\% \text{ of variation in dependent variable is explained by } 4 \text{ independent variables})$$

$$(d) \text{Adj } R^2 = 1 - \left\{ (1 - R^2) * \frac{n-1}{n-k-1} \right\}$$

$$= 1 - \left(0.4 * \frac{29}{25} \right) = \underline{\underline{0.536}}$$

$$① H_0: \beta_1 = \beta_2 = \dots = \beta_K = 0$$

$$H_1: \beta_1 = \beta_2 = \dots = \beta_K \neq 0$$

Critical value: $F_{tab} = F_{4, 25, 0.05} = 2.759$

Decision: Since $F_{calc}(9.375) > F_{tab}(2.759)$, we
reject H_0 .

[Mean Vector

Correlation matrix

Dispersion

Factor analysis, PCA, Clustering

>Data reduction

February 14, 2025

Degrees of freedom = No. of obs - no. of unknown parameters
(no. of restrictions)

1) Adj R^2

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

2) Standard Error

$$S_{yx} = \sqrt{MSE}$$

4. Testing significance of individual regression coefficient

$$H_0: \beta_i = 0$$

$$H_1: \beta_i \neq 0$$

From software, we get the t-value

$$\frac{VVI}{Q} \frac{C/W}{}$$

Solution
 $n = 24$

i)

Predictor	Regression co-efficient b_i	Standard Error s_{bi}	t	p-value
Constant	6356.617	838.701	7.578	
x_1	0.56038	0.15811	3.5442	0.000
x_2	-31.2077	8.95905	-3.4833	0.000
x_3	-327.503	149.169	-2.1955	0.001
x_4	-113.895	16.2604	-7.0044	0.000
x_5	-621.458	147.828	-4.2039	0.000

ANOVA

Source	SS	df	MSS	F	p-value
Regression	12204000	5	2440800	$F = 20.28$	0.000
Residual	2166000	18	120333.333		
Total	14370000	23			

$$F_{5,18,0.05} = 2.77$$

Decision: Since $F_{\text{calc}} (20.28) > F_{\text{tab}} (2.77)$, we
reject H_0 .

iii) Here given $\alpha = 0.05$

$$H_0: \beta_3 = 0$$

$$H_1: \beta_3 \neq 0$$

If p-value
not
given
p < 0.05

~~test~~

Here, p-value (0.001) < $\alpha (0.05)$, we reject null hypothesis.

Decision: Thus, β_3 is statistically significant at 5% significance level.

$$H_0: \beta_3 = 0$$

$$H_1: \beta_3 \neq 0$$

Test statistic: $t = \frac{b_3}{s b_3} = -2.195$

Cal $|t| = |t| = 2.195$

Test $t = t_{n-k-1}$ at α

$= t_{18}$ at 5%.

iii) 99% CI

For \bar{x}_1 ,

Regression co-efficient = $b_1 \pm t_{n-k-1}, \alpha / 2 s b_1$ [Two Tailed]

For square feet heated space,

$$b_1 \pm t_{18, 0.01} s b_1 = \left(1 - \alpha = 99\%, \alpha = 1\%\right)$$

$$= 0.56038 \pm 2.878 * 0.15811$$

$$= 0.56038 \pm 0.45504$$

$$= 0.10534 - 1.01542$$

$$= [0.1053 \leq \beta_1 \leq 1.0154]$$

↓
Parameter

iv) Standard Error of Estimate

$$\begin{aligned} SE(\hat{y}) &= S_{yx} = \sqrt{MSG} \\ &= \sqrt{120333.333} \\ &= 346.890 \end{aligned}$$

Interpretation: The average variation of observed values of dependent variable around its fitted equation is $\boxed{346.890}$.

v) $R^2 = \frac{SSR}{TSS} = \frac{12204000}{14370000}$

$$= 0.849 \quad (\text{KWh/month})$$

Int: 84.9% of variation in dependent variable is explained by

$$\text{Adj } R^2 = 1 - \left\{ (1 - R^2) * \frac{n-1}{n-k-1} \right\} \quad \begin{matrix} \text{five independent} \\ \text{variables} \\ (x_1, x_2, \dots, x_5) \end{matrix}$$

$$= 1 - \left\{ (1 - 0.849) * \frac{23}{18} \right\} y$$

$$= 1 - 0.1929$$

Interpretation: After adjusting for given degrees of freedom.

vi)

$$Y = 6356.17 + 0.56038x_1 + -31.2077x_2 - 327.503x_3 \\ - 113.895x_4 - 621.458x_5$$

For $x_1 = 1295, x_2 = 18, x_3 = 5, x_4 = 3$ and $x_5 = 1,$

$\hat{Y} = 6356.17 + 0.56038 \times 1295 - 31.2077 \times 18 - 327.503 \times 35 \\ - 113.895 \times 3 - 621.458 \times 1 \\ = \boxed{3919.4655}$

vii)

$$H_0: \beta_1 = \beta_2 = \dots = \beta_5 = 0$$

At least one

$$H_1: \beta_i \neq 0 \text{ where } i \in \{1, 2, 3, \dots, 5\}$$

Critical value:

Calculated F

$$F_{calculated} = \frac{MSR}{MSE}$$

$$= 20.28$$

$$\text{Critical value: } F_{k, n-k-1, \alpha=0.05} = 2.77$$

Decision: Since $F_{calc} (20.28) > F_{tab} (2.77),$ we reject

null hypothesis. Thus, there is a significant impact of at least one independent variable on

dependent variable

Feb 15, 2025

= BLUE Properties

(Best linear Unbiased Estimate)

- Unbiasedness

- linearity

- Best ness (minimum variance)

- Consistency

Regression Coefficient $\rightarrow \beta_i \rightarrow$ Parameter
(Slope)

$\rightarrow b_i \rightarrow$ Estimator

BLUE Properties

(1) The estimators are linear, that is, they are linear functions of the dependent variables. Linear estimators are easy to understand and deal with compared to nonlinear estimators.

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_n x_n$$

linear

(2) The estimators are unbiased, that is, in repeated applications of the method, on average, the estimators are equal to their true values.

Example $N=3$
Variance $= 0.2$

2	4
6	

$n=2$

$(2, 4)$

$(4, 6)$

$(2, 6)$

\bar{x}

3

5

4

$\bar{x} = \frac{3+5+4}{3} = 4$

$\sigma^2/n = 3 = 1$

Different values
for sample mean.

$N = 35$

$$10 \rightarrow 35 \\ C_{10}$$

$\hat{\beta}_1$
$\hat{\beta}_2$
\vdots
$\hat{\beta}_r$

Thus, $E(\bar{x}) = \mu$

mean = β_0

* Distribution of all possible sample means and sample variance is known as sampling distribution.

③ In the class of linear unbiased estimators, OLS estimators have minimum variance. As a result, the true parameter values can be estimated with least possible uncertainty; an unbiased estimator with the least variance is called an efficient estimator.

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{\epsilon}$$

And we know that:

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

Now substituting $\underline{Y} = \underline{X}\underline{\beta} + \underline{\epsilon}$ in the equation, we get,

$$\begin{aligned} \hat{\beta} &= (\underline{X}^T \underline{X})^{-1} \underline{X}^T (\underline{X}\underline{\beta} + \underline{\epsilon}) \\ &= (\underline{X}^T \underline{X})^{-1} (\underline{X}^T \underline{\beta}) + (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{\epsilon} \end{aligned}$$

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

Taking expectation on both sides:

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

$$\therefore E(\hat{\beta}) = E(\beta) = \beta$$

This proves that $\hat{\beta}$ is an unbiased estimator of β .

Bestness

$$\hat{\beta} = \beta + (X^T X)^{-1} X^T \epsilon \quad \text{--- ①}$$

$$\text{or, } \hat{\beta} - \beta = (X^T X)^{-1} X^T \epsilon$$

The variance of $\hat{\beta}$ is

$$\begin{aligned} \text{Var}(\hat{\beta}) &= E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T] \\ &= E[(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}] \\ &= (X^T X)^{-1} X^T E(\epsilon \epsilon^T) X (X^T X)^{-1} \\ &= (X^T X)^{-1} X^T \sigma^2 X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} (X^T X) (X^T X)^{-1} \\ &= \sigma^2 \end{aligned}$$

Measuring goodness-of-fit in multiple regression analysis

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

Adj R^2 penalizes for adding additional explanatory variables in the model so that Adj R^2 values of models having

- different number of explanatory variables become comparable.
- Other model selection criteria
- Akaike Information Criteria (AIC)

$$AIC = -2 * \log(L) + 2 * k$$
- Schwarz Bayesian Criteria (SBC) (BIC)

$$BIC = -2 * \log(L) + k * \ln(n)$$
- Hannan-Quinn Criteria (HQC)

$$HQ = -2 * \log(L) + 2 * k * \ln(\ln(n))$$

we look for minimum values of these indices while selecting a model.

Multicollinearity

- Multicollinearity refers to a situation in which there is high inter correlations among the explanatory variables of a multiple regression model. It arises only in the context of multiple regressions.
- ANOVA
- F-value \leftarrow significant

t	p

$|y| \geq 0.05$ (Insignificant)

- A good number of the estimated coefficients can be statistically insignificant. This occurs when variances and hence standard errors of the estimated coefficients are large. This is possible when there is little variation in explanatory variables or high inter correlations among the explanatory variables or both.

February 21, 2025

Confidence Interval Estimate and Prediction Interval
 Estimate for mean and individual response of Y.
 Example: Simple Regression Equation

X (Income) in '000	Y (Expenditure) in '000	XY	X^2	Y^2
20	15	300	400	225
24	17	408	576	289
25	18	450	625	324
30	18	540	900	324
40	20	800	1600	400
$\Sigma X = 139$		$\Sigma XY = 2498$	$\Sigma X^2 = 4101$	$\Sigma Y^2 = 1562$

$$b_1 = \frac{n \sum XY - (\sum X)(\sum Y)}{n \sum X^2 - (\sum X)^2} = \frac{5 \times 2498 - 139 \times 88}{5 \times 4101 - (139)^2}$$

$$= 0.2151$$

$$b_0 = \frac{\sum Y}{n} = \frac{b_1 \sum X}{n} = \frac{88}{5} - \frac{0.2151 \times 139}{5}$$

$$= 11.623$$

The estimating equation/ regression equation / best line/ line of best fit is

$$\hat{Y} = 11.623 + 0.2151 X$$

$\downarrow \quad \downarrow$
 $b_0 \quad b_1$

when $x = 35$: (in '000)

$$\begin{aligned}\hat{Y} &= 11.623 + 0.2151 \times 35 \\ &= 19.1515 \text{ (in '000)}\end{aligned}$$

hat matrix element (h_{ii})

$$\begin{aligned}h_{ii} &= \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x - \bar{x})^2} \\ &= \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum x^2 - \frac{n}{n} \bar{x}^2}\end{aligned}$$

$$\bar{x} = \frac{\sum x}{n} = \frac{139}{5} = \underline{\underline{27.8}}$$

$$\begin{aligned}\text{Now, } h_{ii} &= \frac{1}{5} + \frac{(35 - 27.8)^2}{4101 - 5 \times (27.8)^2} \\ &= \underline{\underline{0.4189}}\end{aligned}$$

Confidence interval estimate for mean response of y
Always two-tailed

$$\hat{Y} \pm t_{n-2, \alpha} s_{yx} \sqrt{h_i}$$

$$LL \leq M_{yx} \leq UL$$

Prediction interval estimate for individual response of

$$\hat{Y} \pm t_{n-2, \alpha} s_{YX} \sqrt{1+h_i}$$
$$LL \leq Y_i \leq UL$$

$$s_{YX} = \sqrt{\frac{\sum (Y_i - \bar{Y})^2}{n-2}}$$

$$= \sqrt{\frac{\sum Y^2 - b_0 \sum Y - b_1 \sum X}{n-2}}$$

Now, $s_{YX} = \sqrt{\frac{1562 - 11.623 \times 88 - 0.2151 \times 2498}{5-2}}$

$$= 1.8110684 \quad 0.8378$$

for mean response of \bar{Y} ,

Now, $\bar{Y} \pm t_{n-2, \alpha} s_{YX} \sqrt{h_i}$

$$19.1515 \pm t_{3, 0.05} * \frac{0.8378}{\sqrt{0.4189}}$$

$$= 19.1515 \pm 3.182 * 0.5422$$

$$= 19.1515 \pm 1.7254$$

$$= 17.4261 - 20.8769$$

$$\therefore 17.4261 \leq Y_i \leq 20.8769$$

for individual response of Y ,

$$\hat{Y} \pm t_{n-2, \alpha} s_{YX} \sqrt{1+h_i}$$
$$= 19.1515 \pm 3.182 * 0.8378 * \sqrt{1+0.4189}$$

$$= 19.1515 \pm 3.1757$$

$$\therefore 16.3242 \leq Y_i \leq 22.3256$$

In Multiple regression analysis,

Confidence interval estimate for mean response of Y is

$$\hat{Y} \pm t_{n-p-1, \alpha} s_{yx} \sqrt{h_i}$$

$$LL \leq U_{yx} \leq UL$$

Prediction interval for individual response of Y is

$$\hat{Y} \pm t_{n-p-1, \alpha} s_{yx} \sqrt{1+h_i}$$

where, p = no. of explanatory / independent variables.

EXAM

EXAMPLE

Suppose, $b_0 = 33$, $b_1 = 2.15$, $b_2 = 7.15$, $b_3 = 0.008$.

$$s_{b_0} = 10.5, s_{b_1} = 0.18, s_{b_2} = 1.28, s_{b_3} = 0.001$$

$$n = 30, \sum (Y - \bar{Y})^2 = 20.18, h_i = 0.131$$

- Fit a regression line
- Predict Y when $x_1 = 10, x_2 = 5, x_3 = 7$ and also construct 95% confidence interval estimate for mean response of Y .
- Also construct 95% prediction interval estimate for individual response of Y .
- Test the significance of impact of x_2 on Y .
- Construct 95% CI estimate for regression coefficient of x_2 .

Solution

$$\text{Here, } p = 3$$

a) Regression line,

$$Y = 33 + 2.15x_1 + 7.15x_2 + 0.008x_3$$

b)

$$\text{when } x_1 = 10, x_2 = 5, x_3 = 7,$$

$$\begin{aligned}\hat{Y} &= 33 + 2.15 \times 10 + 7.15 \times 5 + 0.008 \times 7 \\ &= 90.306\end{aligned}$$

$$s_{yx} = \sqrt{\frac{\sum (Y - \bar{Y})^2}{n-p-1}} = \sqrt{\frac{20.18}{30-3-1}} = 0.880$$

Now, CI estimate is:

$$\hat{Y} \pm t_{n-p-1, \alpha/2} s_{yx} \sqrt{h_i}$$

$$= 90.306 \pm t_{26, 0.05} * 0.880 * \sqrt{0.131}$$

$$= 90.306 \pm 2.056 * 0.880 * \sqrt{0.131}$$

$$= 90.306 \pm 0.654$$

$$\therefore 89.652 \leq Y_{yx} \leq 90.960$$

c) PI estimate,

$$\hat{Y} \pm t_{n-p-1, \alpha/2} s_{yx} \sqrt{1+h_i}$$

$$= 90.306 \pm 2.056 \times 0.880 \times \sqrt{1 + 0.131}$$

$$= 90.306 \pm 1.924$$

$$\therefore 88.382 \leq Y_i \leq 92.23$$

d) $H_0: \beta_2 = 0$

$$H_1: \beta_2 \neq 0$$

$$\text{Cal } t, |t| = \frac{b_2}{s_{b_2}} = \frac{7.15}{1.28} = 5.585$$

Tablet test $t = t_{n-p-1, \alpha}$

$$= t_{26, 0.05}$$

$$= 2.056$$

Decision: Since $\text{cal } t > \text{Table } t$, we reject H_0 .

Thus, there is a significant impact of x_2 on y .

e) For x_2 ,
Regression co-efficient = $b_i \pm t_{n-p-1, \alpha} \cdot s_{b_i}$

$$\text{for } i = 2,$$

$$\begin{aligned} \text{Regression coefficient} &= b_2 \pm t_{26, 0.05} \cdot s_{b_2} \\ &= 7.15 \pm 2.056 \times 1.28 \\ &= 7.15 \pm 2.63 \end{aligned}$$

$$4.52 \leq \beta_2 \leq 9.781$$

$$\therefore 4.52 \leq \beta_2 \leq 9.781$$

March 8, 2025

Multicollinearity

Multicollinearity refers to a situation where there are high inter correlations among the explanatory variables of a multiple regression model. Multicollinearity problem arises only in the context of multiple regressions, it is considered as a problem because when the explanatory variables are highly correlated, most of their variation is common so that there is little variation unique to each variable.

Consequences of multicollinearity

- Perfect multicollinearity
- Under perfect multicollinearity it is not possible to compute the values of OLS estimates.
- Imperfect multicollinearity

B L U.E
 ↴ ↓
 But linear Unbiased Estimate

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p$$

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + \underbrace{b_3 x_3}_{\text{Excluded}} + b_4 x_4 + b_5 x_5$$

Imperfect multicollinearity

- Correlation between two independent variables → High
- Most of the independent variables → Insignificant

F-test → reject

t-test: see multi-variate statistics for more information
 If there is no difference in correlation with previous relationships →
 do not reject null hypothesis (Insignificant). If it is significant →
 then it is a significant relationship with other variables.
 This will be significant in multivariate analysis.

#Detection of multicollinearity

1. Correlation matrix → Gives idea; Not utilized for final decision
2. F-statistic and t-statistic
3. Klein's rule of thumb
4. Variance Inflation Factor (VIF)
5. Tolerance method (TOL)
 - $VIF = 1$ (no multicollinearity)
 - $1 < VIF < 5$ (less multicollinearity)
 - $5 \leq VIF \leq 10$ (Moderate multicollinearity)
 - $VIF > 10$ (High multicollinearity)

Two-way Randomized

Correlation matrix →

	Y	x_1	x_2	x_3	x_4
Y	1				
x_1		1			
x_2			1		
x_3				1	
x_4					1

Remedial Technique

- 1) By increasing sample size
- 2) Transformation of variables → Using literature review.
Eg: ds transformation, inverse transformation etc.
- 3) Using extraneous estimates
- 4) Dropping the variables.
 - We have to drop the variable having highest VIF (>10) at a time. Again run the regression of remaining variables check whether there is $VIF > 10$ or not and repeat the same process until we get $VIF < 10$.

Cobb - Douglas Production

- To illustrate the CD function, we use Cobb data on output (as measured by value added, in thousands of dollars), labor input (worker hours, in thousands), and capital input (capital expenditure, in thousands of data dollars) for the US manufacturing sector.
- The data is cross-sectional, covering 50 states and Washington, DC, for the year 2005.

log(Labor)

0.47

The interpretation of the coefficient of log(Labor) of about 0.47 is that if we increase the labor input by 1% on average, output goes up by about 0.47%.

#AUTOCORRELATION & HETEROGENEITY

Autocorrelation

- One of the assumptions of the Classical Linear Regression Model (CLRM) is that the disturbance term of the model is independent.

$$\text{Cov}(\epsilon_t, \epsilon_s) = E(\epsilon_t, \epsilon_s) = 0 \text{ for } t \neq s$$

This feature of regression disturbance is known as serial independence or non autocorrelation. It implies that the value of disturbance term in one period is not correlated with its value in another period.

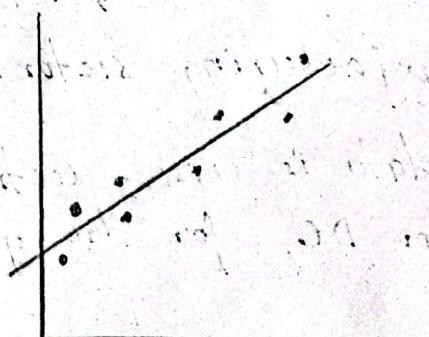
In time series, the disturbance term at period t may be related with the disturbance term at $t-1, t-2, \dots$, and $t+1, t+2, \dots$ and so on. In that case, $\text{Cov}(\epsilon_t, \epsilon_s) \neq 0$ for $t \neq s$ and we say that the disturbances are autocorrelated.

$$Y - \hat{Y} = e$$

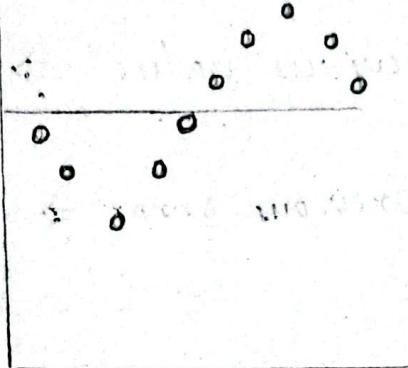
$$E(e_i) = 0$$

$$\text{Var}(e_i) = \text{constant}$$

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



Errors



Pattern Visible



No pattern

Random

Important \Rightarrow Assumptions + BLUE Properties

> Bestness: Minimum Variance

\rightarrow Impact on Bestness

> Impact on test /

Hypothesis testing

Y	x_1	x_2	x_3	\dots	x_p	\hat{Y}	$Y - \hat{Y}$
							+
							+
							-
							-
							+
							+
							+
							+
							-
							-
							-
							-
							-

This is
preferable
pattern
observed

However, in real data, multicollinearity and autocorrelation are preferred / useful.

* Autocorrelation is mostly seen in time-series data.

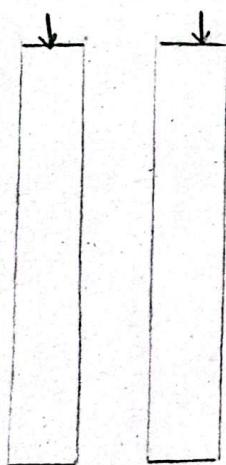
data → dependent on previous value \Rightarrow Autoregression

error → dependent on previous error \Rightarrow Autocorrelation

$$\frac{Y - \hat{Y}}{\text{---}}$$


ARIMA — Autoregression Integrated Moving Average
(Dependency on Previous Error)

dag-1 Correlation



Equilibrium

$$Y_t = \frac{x_t + x_{t-1} + x_{t-2}}{\text{Distributed lag}}$$

ARDL Model

$$Y_t = \underbrace{Y_{t-1} + Y_{t-2} + x_t + x_{t-1}}_{\text{ARDL model}}$$

Specification of Autocorrelation Relationships # Mostly used in Time Series

$$\epsilon_t = \beta \epsilon_{t-1} + u_t \quad (1)$$

↓
White noise

$$E(u_t) = 0$$

$$\text{Var}(u_t) = E(u_t^2) = \sigma_u^2 \text{ for all } t.$$

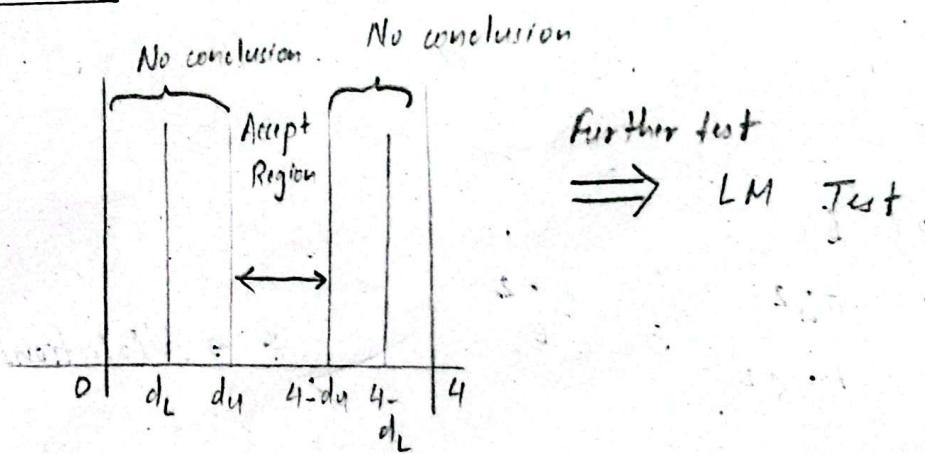
u_t is normally distributed

$$E(u_t u_{t-1}) = 0$$

↓

No autocorrelation

Durbin-Watson Test → Autocorrelation Test



$$\epsilon_t = \beta \epsilon_{t-1} + u_t$$

$$\epsilon_{t-1} = \beta \epsilon_{t-2} + u_{t-1}$$

⋮

$$\epsilon_t = u_t + \beta u_{t-1} + \beta^2 u_{t-2} + \beta^3 u_{t-3} + \dots \quad (2)$$

This shows that under the first-order autoregressive scheme the

effect of past disturbance wears off gradually as $|s| < 1$.

where, s represents the correlation between ϵ_t and ϵ_{t-1} (first order autocorrelation), s^2 is correlation coefficient between ϵ_t and ϵ_{t-2} (second order autocorrelation) and so on.

- $s = 0 \Rightarrow$ no autocorrelation
- strength of autocorrelation becomes high as s approaches to unity (+1).
- strength of autocorrelation becomes high again as s approaches to -1.

Mean

$$E(\epsilon_t) = 0$$

*Variance of ϵ_t

$$\text{Var}(\epsilon_t) = \frac{\sigma_u^2}{1-s^2} = \sigma_\epsilon^2 \quad \rightarrow \text{Violations}$$

*Covariance of ϵ_t and ϵ_{t-1}

$$\text{Cov}(\epsilon_t, \epsilon_{t-1}) = s \sigma_\epsilon^2$$

Assumptions
+

BLUE Properties

March 14, 2025

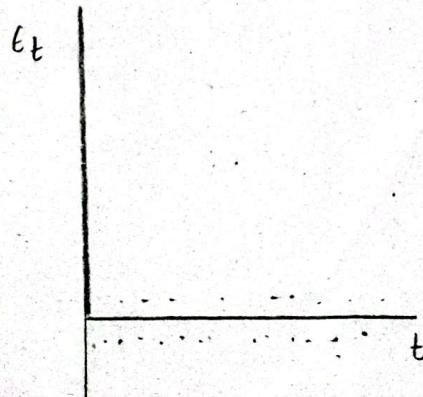
D) Consequences of autocorrelation are:

- i) The OLS estimators are still unbiased and consistent.
- ii) The OLS estimators are no longer minimum variance or best estimators. Hence, they are not efficient and BLUE.
- iii) If we disregard the problem of autocorrelation and believe that all assumptions are valid, following problems will arise:
 - * The estimated variance of disturbance term will be under estimate of its true variance.
 - * The standard error of the estimated slope coefficient will be much smaller if it is computed with usual OLS formula. This will provide spurious (wrong) impression about the statistical significance.
 - * The usual t and F test will become invalid.

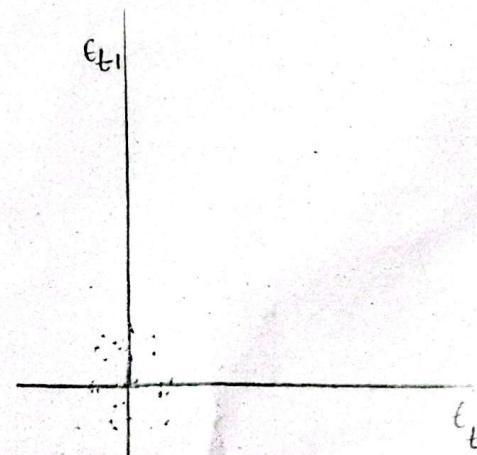
Detection (of autocorrelation)

- Graphical Method (Residual Plot)
 - Graph of ϵ_t against time
 - Graph of ϵ_t against ϵ_{t-1}

i) $\delta = 0$



ϵ_{t-1}



(IMP)

Durbin-Watson Test (D-W Test)

1 Numerical

This is the simplest and most widely used test for autocorrelation. It is based on the following assumptions:

- a. The regression model includes a constant or intercept term.
- b. We are examining presence of first order autocorrelation.
- c. The regression model does not include a lagged dependent variable as an explanatory variable. Now, to understand how the test is performed, consider the model:

$$Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_K X_{Kt} + \epsilon_t$$

where $\epsilon_t = \rho \epsilon_{t-1} + v_t \quad | \rho | < 1$

Null hypothesis (H_0): $\rho = 0$

or, $H_0: \rho \geq 0$

or, $H_0: \rho \leq 0$

There is no significant autocorrelation.

Alternative hypothesis (H_1): $\rho \neq 0$

or, $H_1: \rho < 0$

or, $H_1: \rho > 0$

There is a significant autocorrelation.

Test-statistic

$$D-W \text{ statistic} = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}$$

for tabulated value,

Given, n , α ,

k' (or p) = number of independent variables

= 1 (if it is not given)

Example

$$\alpha = 0.05$$

$$n = 10$$

$$k = 1$$

$$\alpha = 0.05$$

n	k = 1		k = 2		k = 3	
	dL	dU	dL	dU	dL	dU
10	0.88	1.32				

Decision Criteria

		Do NOT REJECT		Reject	
		NO C O N L L	NO C O N C L	Rejection	
0		U S 1 0 N	U S 1 0 N	4	
	+ve	1.32 (dL)	2.68 (4-dU)		
	-ve	0.58 (dL)	3.12 (4-dL)		

Example: By using following information

Period	Error / Residual / e_i	e_{i-1}	$(e_i - e_{i-1})^2$	e_i^2
1	+5	—	—	25
2	+4	+5	1	16
3	+3	+4	1	9
4	+2	+3	1	4
5	+1	+2	1	1
6	-1	+1	4	1
7	-2	-1	1	4
8	-3	-2	1	9
9	-4	-3	1	16
10	-5	-4	1	25

$$\sum_{i=2}^n (e_i - e_{i-1})^2 = 12$$

$$\sum e_i^2 = 110$$

At $\alpha = 0.05$, test the presence of autocorrelation.

$$\sum_{i=1}^{10} e_i^2 = 110$$

Null Hypothesis (H_0): $\rho = 0$ i.e. There is no significant autocorrelation.

Alternative Hypothesis (H_1): $\rho \neq 0$ i.e. There is significant autocorrelation.

Tut-statistic:

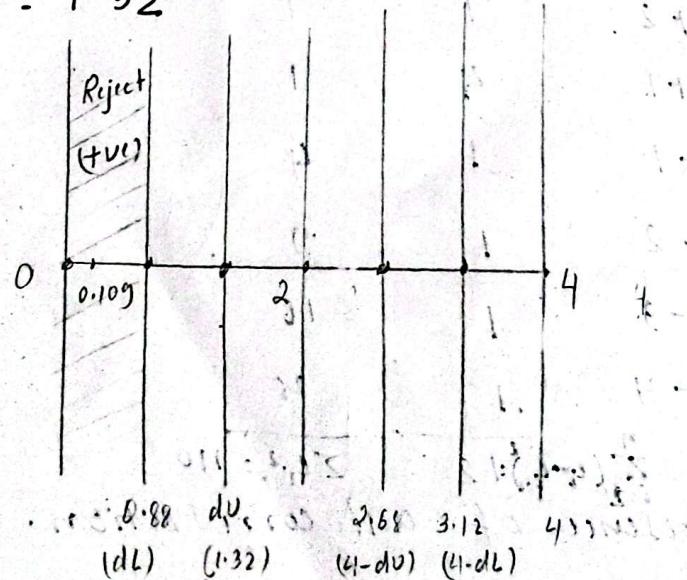
$$\begin{aligned} \text{D-W Statistic} &= \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2} \\ &= \frac{12}{110} \\ &= 0.109 \end{aligned}$$

From Table

$$n = 10, d = 0.05, k = 1$$

$$d_L = 0.88$$

$$d_U = 1.32$$



~~Deny~~ Hence, $0 < DW = 0.109 < d_L (0.88)$, we reject H_0

i.e. There is a significant autocorrelation.

Limitations

- 1) g_t can not be used for testing higher order autocorrelation.
- 2) This test is biased.

Breusch-Godfrey Lagrange Multiplier Test (LM test)

Null Hypothesis (H_0): $s_1 = s_2 = s_3 = \dots = s_p = 0$ [no autocorrelation]

Alternative Hypothesis (H_1): At least one s_i is not zero

$$s_i \neq 0 ; i \in \{1, 2, \dots, p\}$$

(IMP)

5 steps

- 1) Estimate eqn (1) by OLS and estimate $\hat{\epsilon}_t = e_t$
- 2) Run the auxiliary regression of e_t on $x_{1t}, x_{2t}, \dots, x_{kt}, e_{t-1}, e_{t-2}, \dots, e_p$

3) To test the validity of above null hypothesis compute LM statistic as follows.

$$LM = (n-p) R^2$$

where R^2 is coefficient of determination of auxiliary equation.

n is number of observations

p is order of autocorrelation.

The LM statistic follows a Chi-Square distribution with degrees of freedom p .

March 15, 2025

(EXAM)

HETEROSKEDASTICITY

For the sample of two variable model

$$Y_i = \alpha + \beta X_i + \epsilon_i \quad (1)$$

We assumed the variances of disturbance term ϵ_i is constant for all observations

$$\text{i.e. } \text{Var}(\epsilon_i) = E(\epsilon_i^2) = \sigma^2 \text{ (constant for all } i) \quad (2)$$

This feature of disturbances terms of the regression model is known as homoskedasticity.

However, it is quite common in regression analysis to have cases where the variance of disturbance term becomes variable rather than ~~the~~ remaining constant.

In this situation the disturbance is said to be heteroskedasticity.

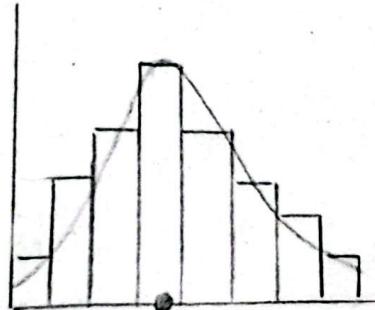
$$\text{i.e. } \text{Var}(\epsilon_i) = E(\epsilon_i^2) = \sigma_i^2 \quad \text{--- (3)}$$

- which means that the variance of disturbance term can change for every different observation in the sample $i=1, 2, \dots, n.$

Marketing Research
- Narash Malhotra

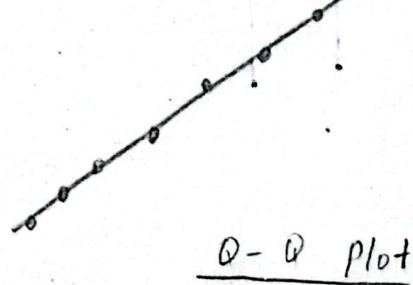
$$Y | x_1, \dots, x_k | \hat{Y} | \overset{\epsilon}{\textcircled{Y - \hat{Y}}}$$

\downarrow
Normality Test



Histogram

Mean:
Median:
Mode



Q-Q Plot

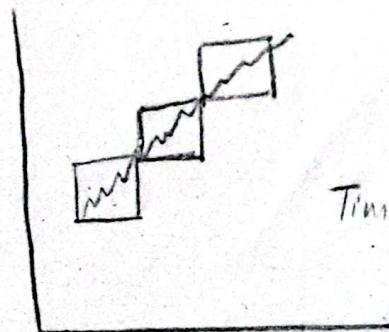
$$\epsilon_i \sim N(0, \sigma^2)$$



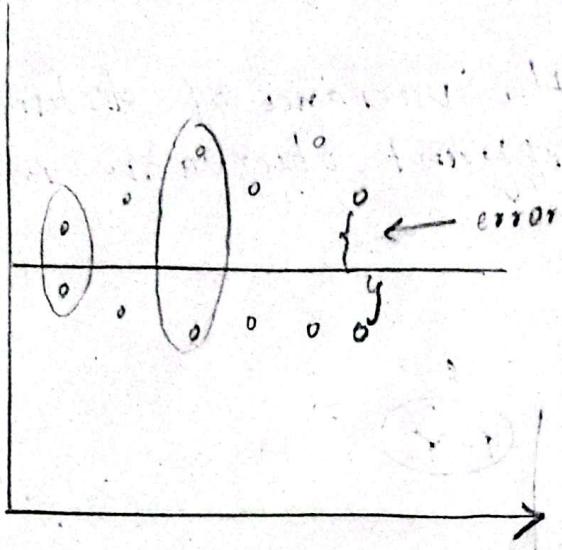
Variance should be constant

\downarrow
Violation

leads to heteroscedasticity.



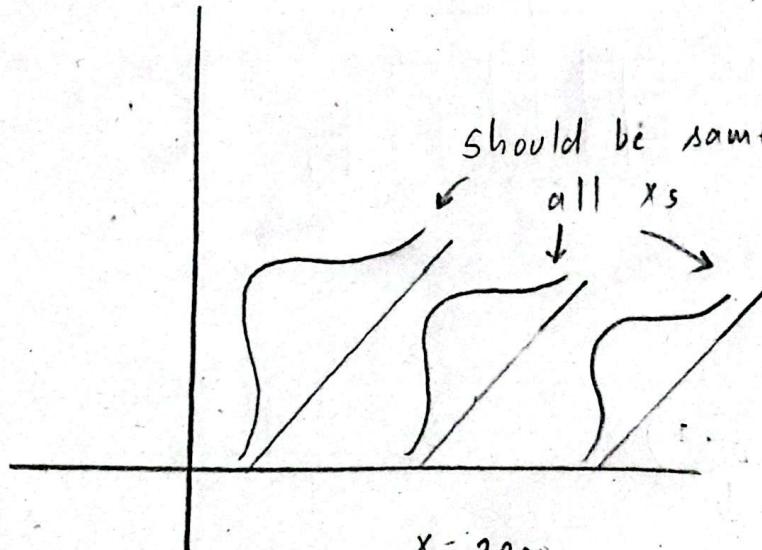
Time Series Data



Income Expenditure

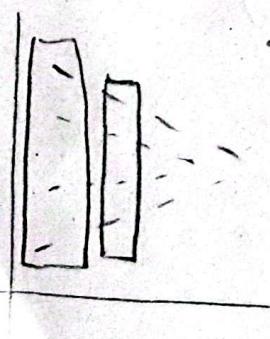
X	Y	\hat{Y}	e
20			
20			
20			
20			

$$\hat{Y} = \beta_0 + \beta_1 X$$



$\text{Var}(\epsilon_i | y) = \sigma^2 \rightarrow \text{Homoskedasticity}$

$\text{Var}(\epsilon_i^2) = \sigma_i^2 \rightarrow \text{Heteroskedasticity}$



dike Flashlight
(Heteroskedasticity)

Sources of Heteroskedasticity

- i) When we are dealing with micro-economic or cross-section data, we are very likely to have a heteroskedasticity problem.
- ii) Presence of outliers.
- iii) If some relevant variables have been mistakenly omitted.
- iv) Inclusion of explanatory variables in the model whose distributions are skewed.
- v) Heteroskedasticity may also arise due to incorrect data transformation.

Consequences

Unbiasedness

$$E(\hat{\beta}) = \beta$$

This shows that $\hat{\beta}$ remains unbiased when the disturbance term of the model ϵ_i is heteroskedasticity.

Bestness

Thus $\text{Var}(\hat{\beta})$ if heteroskedasticity $> \text{Var}(\hat{\beta})$ if homoskedasticity.
So, there is no longer minimum variance and hence not best estimator.

$\hat{\beta}$ is unbiased but not the best.

Consistency

$\hat{\beta}$ is consistent when the disturbance term is heteroskedastic.

- The OLS estimators continue to remain unbiased and consistent under heteroskedasticity.
- Heteroskedasticity increases the variances of the distributions of estimator of β thereby turning the OLS estimators inefficient (not best).

Heteroskedasticity also affects the variance of OLS and their standard error. In fact, the presence of heteroskedasticity causes the OLS method to underestimate the variances and hence standard error of the estimators. As a consequence, we have higher than expected values of t and F statistic.

As the OLS estimators are unbiased under heteroskedasticity

Detection Techniques

Graphical method:

Plotting square of residuals (e_i^2) against explanatory variables (X_i) to which it is suspected the disturbance variance is related. Since e_i^2 is unknown, its proxy measure $e_i^2 = e_i \cdot e_i$

Breusch-Pagan - Godfrey Test

(FAM)

+ steps

They developed a Lagrange Multiplier (LM) test to examine the presence of heteroskedasticity in data.

Considering the model.

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_K X_{Ki} + \epsilon_i \quad (1)$$

And suppose that

$$\text{Var}(\epsilon_i) = \sigma_i^2 = f(r_0 + r_1 Z_{1i} + \dots + r_r Z_{ri})$$

Goldfeld-Quandt Test

Goldfeld-Quandt Test (1965) proposed a test of heteroskedasticity that may be applied when one of the explanatory variables is suspected to be the heteroskedasticity culprit.

- We assume $\sigma_i^2 \propto x_i$ and also $\epsilon_i \sim \text{Normal distribution}$.
The hypothesis are:

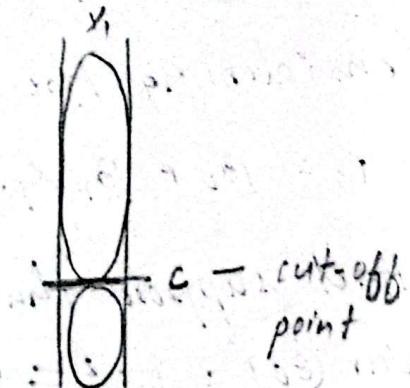
Null Hypothesis (H_0): $\text{Var}(\epsilon_i | X_i) = \sigma^2$, a constant (Homoskedasticity)

Alternative Hypothesis (H_1): $\text{Var}(\epsilon_i | X_i) = \sigma_i^2$, a variable
(Heteroskedasticity)

Steps

- Identify the variable to which the variance of disturbance term is suspected to be related.

- 2) Sort the raw data in ascending order (starting with lowest and going to be highest) of the values of x_i .
- 3) Cut out some central observations (c), breaking data into



Demerits

- 1) Determination of appropriate value of c.
- 2) Difficulty in identifying x-variable

Remedial Techniques

- Log Transformation
- Weighted Least Squares Method

Generalized Least Squares (GLS)

Consider the model:

$$Y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + \epsilon_i \quad \text{--- (1)}$$

And suppose that there is heteroskedasticity.

Packages

- car \rightarrow qq plot
- lmtest \rightarrow for BP Test

April 4, 2025

[Regression Analysis]

Dummy Variable

$$b_1 = \frac{n \sum XY - \cancel{\sum X} \sum Y}{n \sum X^2 - (\sum X)^2}$$

$$= \frac{8 \times 141 - 40 \times 236}{8 \times 4 - (4)^2}$$

Ex 1

Salary (000)	Gender	XY	X ²	
30	1	30	1	$= 11.5$
18	0	0	0	$b_0 = \bar{Y} - b_1 \bar{X}$
20	0	0	0	$= \frac{236}{8} - 11.5 \frac{4}{8}$
34	1	34	1	$= 29.5 - 5.75$
40	1	40	1	$= 23.75$
25	0	0	0	
32	0	0	0	
37	1	37	1	$\hat{Y} = 23.75 + 11.5X$
$\bar{Y} = 29.5$				
$\bar{X} = 0.5$				
$\sum Y = 236$		$\sum X = 4$	$\sum XY = 141$	$\sum X^2 = 4$

Ex 2

Y Salary (000)	Gender (X)	XY	X ²	$b_1 = \frac{8 \times 95 - 4 \times 236}{8 \times 4 - (4)^2}$
30	0	0	0	$= -11.5$
18	1	18	1	$b_0 = \bar{Y} - b_1 \bar{X}$
20	1	20	1	$= 29.5 - (-11.5)$
34	0	0	0	$\frac{4}{8} \times 4$
40	0	0	0	$= 29.5 + 5.75$
25	1	25	1	$= 35.25$
32	1	32	1	
37	0	0	0	$\sum X^2 = 4$

$$\boxed{\hat{Y} = 35.25 - 11.5X}$$

for 1st model,

$$\hat{Y} = 23.75 + 11.5 X$$

when $X=1$,

$$\hat{Y} = 23.75 + 11.5(1) = \underline{35.25} ('000)$$

Average salary
of male.

when $X=0$,

$$\hat{Y} = 23.75 + 11.5(0) = \underline{23.75} ('000).$$

Average salary
of female.

Binary number, Dichotomous variable

Dummy Variable: i) Any variable having two options, denoting them by 1 or 0, is known as dummy variable.
ii) We can include dummy variable as an independent variable in regression analysis.

In second model,

$$\hat{Y} = 35.25 - 11.5 X$$

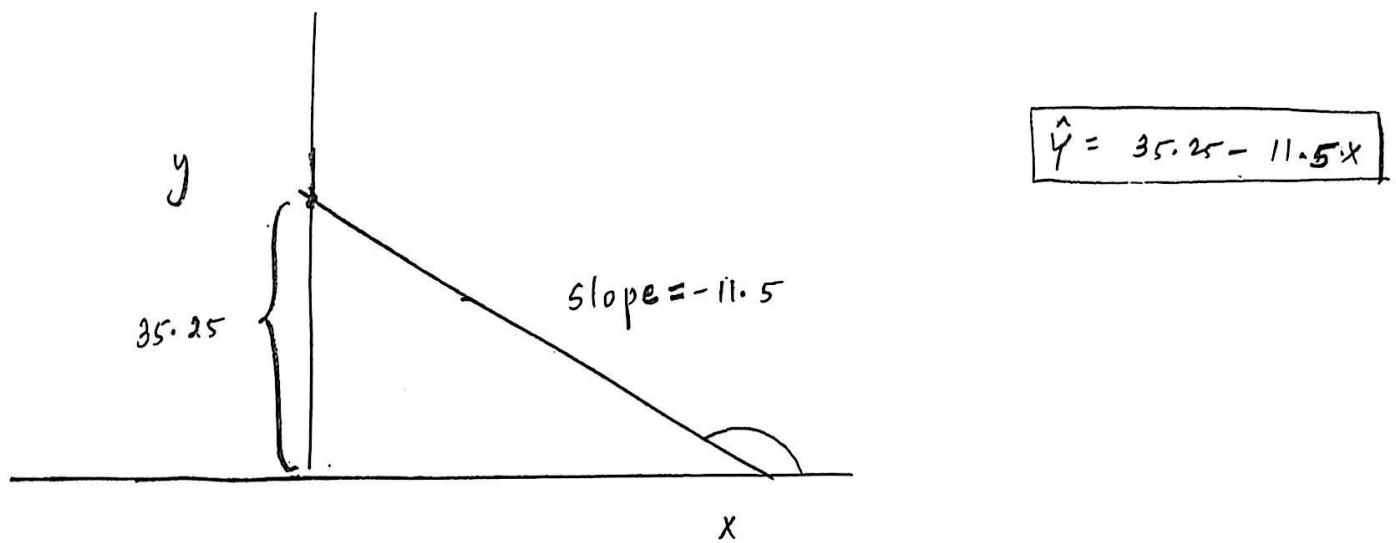
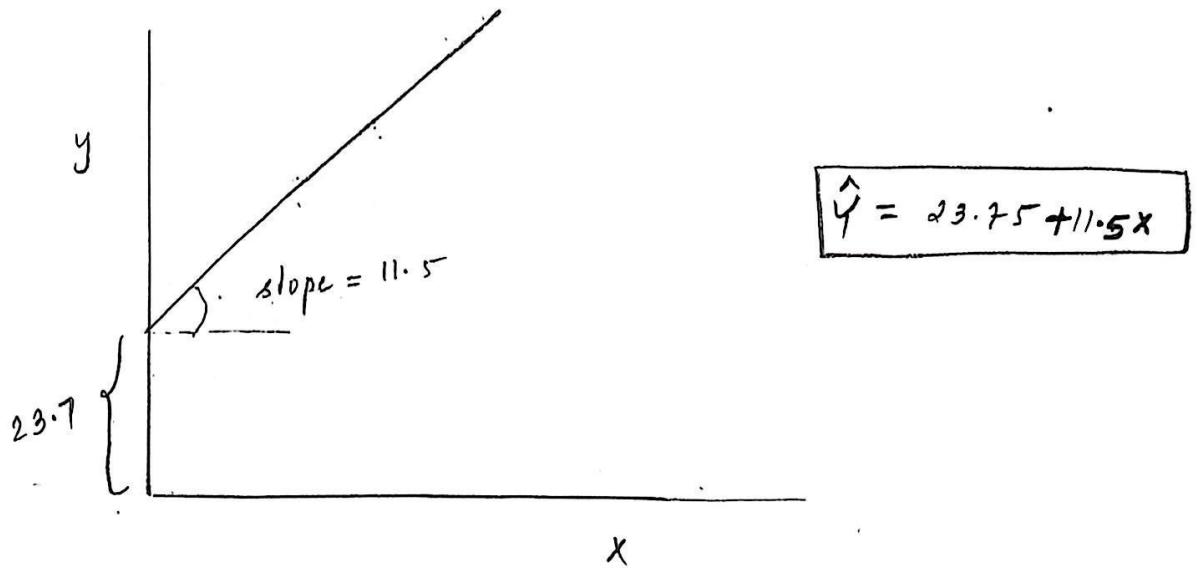
when $X=1$,

$$\hat{Y} = 35.25 - 11.5$$

$$= 23.75$$

when $X=0$,

$$\hat{Y} = 35.25 - 11.5(0) = 35.25$$



$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

$t = \frac{\hat{\beta}_1}{s_{\hat{\beta}_1}}$ reject \Rightarrow There is a significant impact of gender on salary.

\Rightarrow Salary discrimination between male and female.

Gg: I am satisfied with the mileage.

Toyota - 1

Kia - 0

While encoding,

0 → Reference category

Q. Salary _____

Post

Assistant 1

Officer 2

Manager 3

Salary	Post
Doesn't give correct results!!!	1
	2
	3

(Y) Salary	x_1 Assistant	x_2 Officer	x_3 Manager
1	1	0	0
0	0	1	0
1	0	0	1

$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3$ we can not get the result / outcome.
 (Perfect Multicollinearity)
 ↓

Dummy variable trap.

$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 \rightarrow$ (left variable is reference category compared against Manager.)

If categorical variable consists of k number of options / groups / categories, we have to make $(k-1)$ dumy variables.

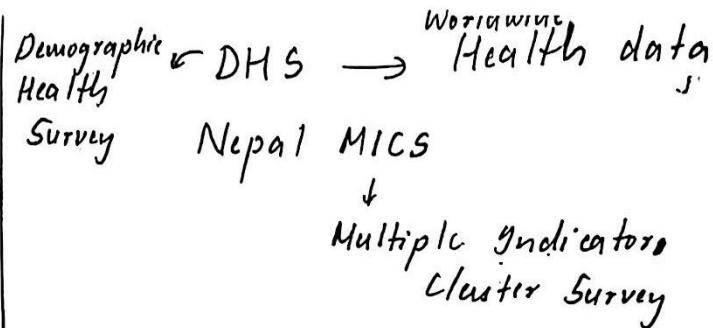
Share price (index)	Sun	Mon	Tues	Wed

Dummy Variable dependent Can't be estimated using OLS

↓ We use

logistic / Probit
Logit / Probit

Multinomial logistic Regression



April 5, 2025

Carter and Hue

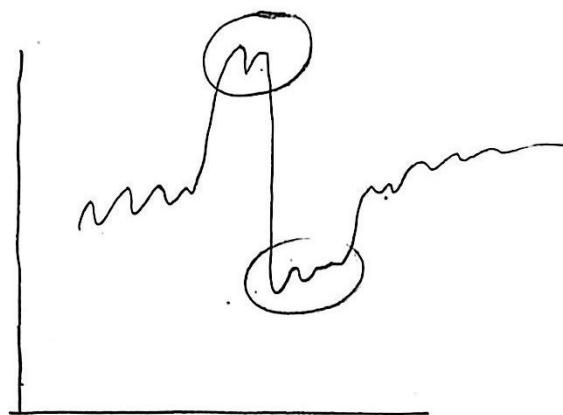
R Practicals

Categorical regression

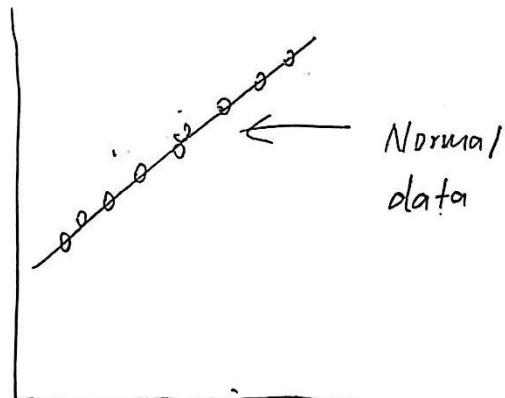
Autocorrelation

Heteroskedasticity

Dummy Variable example
(Wage)



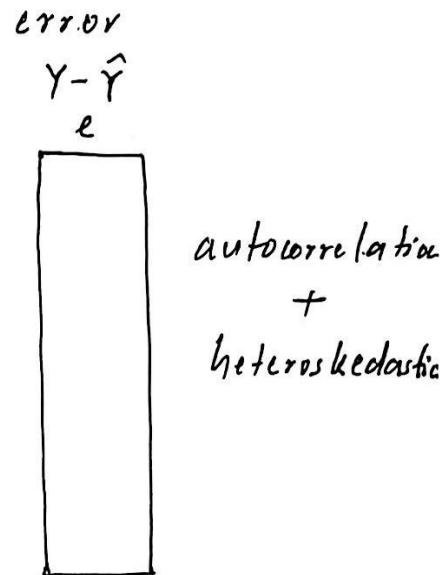
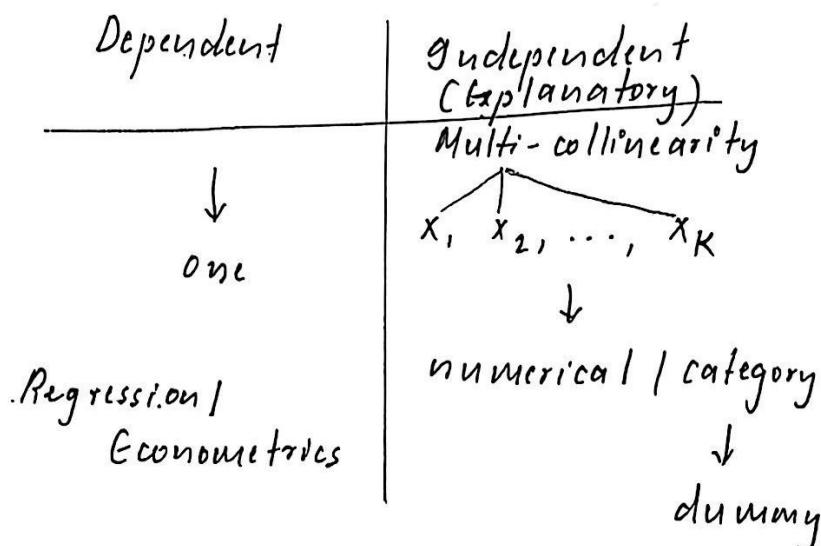
QQ-Plot



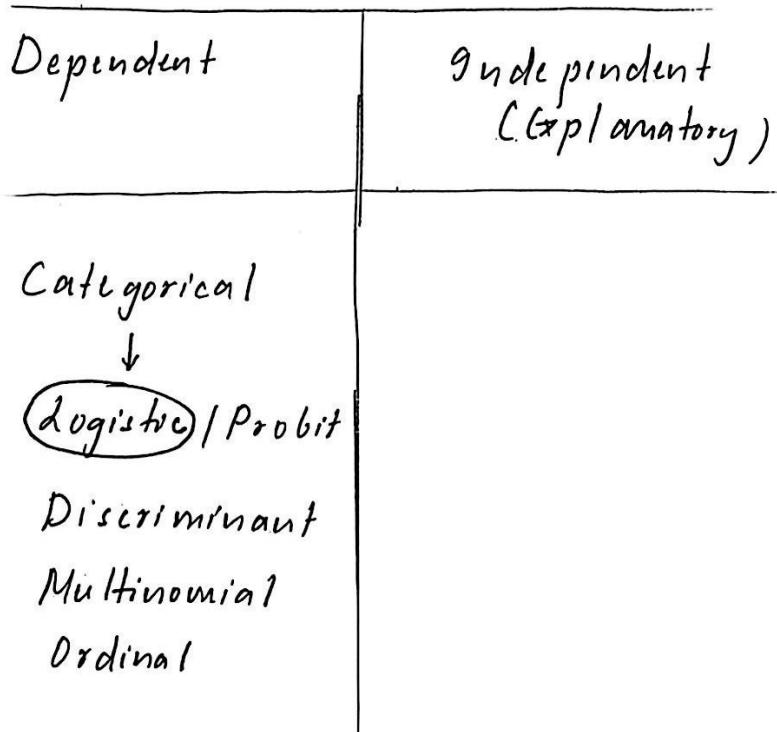
April 12, 2025 (Regression)

Econometrics Hilkaast
Book

Categorical Data Analysis



ARIMA: Auto-regression Integrated Moving Average



logistic / Probit



Dependent Variable



Categorical / Binary (0)

Dependent



Categorical

→ Discriminant (Multi-variate analysis)

Dependent Variable



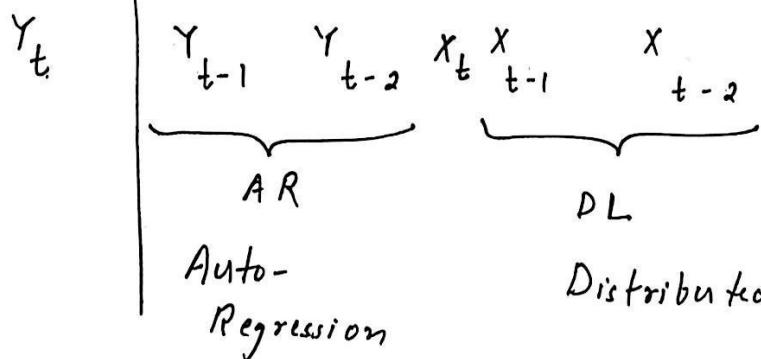
Ranking (ordinal)



Ordinal {for multinomial}

GDP	Remittance
y_t	x_t

Cointegration



Distributed lag Model

AR-DR model

Logistic Regression

Models with dummy dependent variable. Three important approaches:

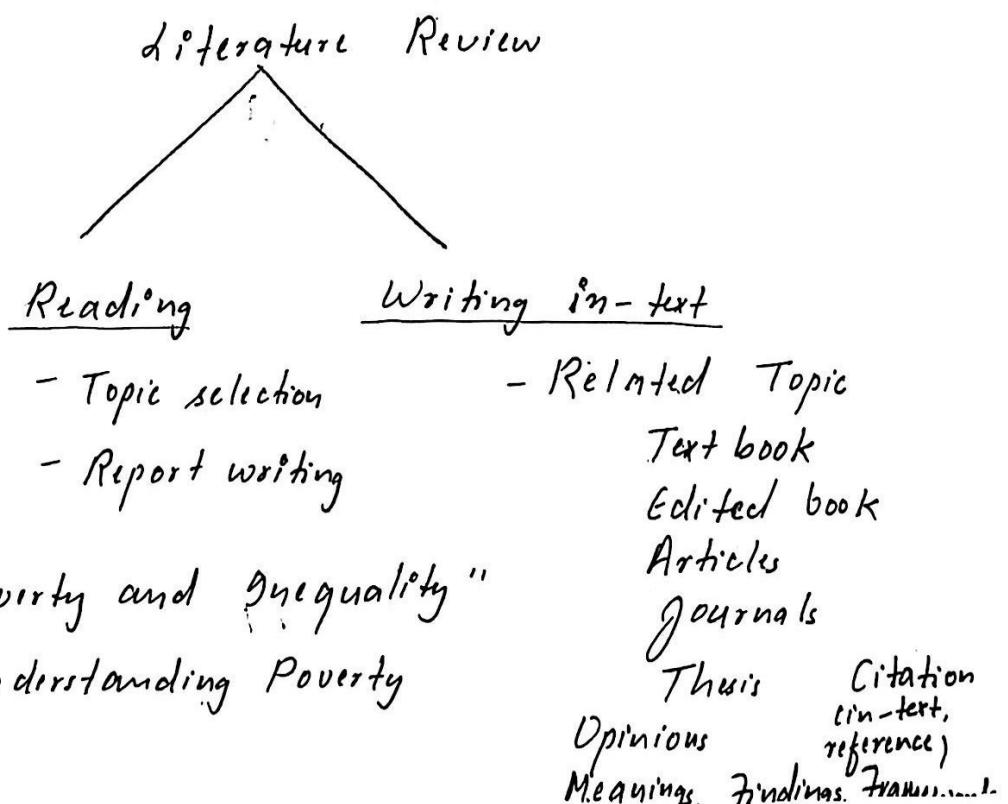
Linear Probability Model (LPM)

Logit Model

Probit Model

Example

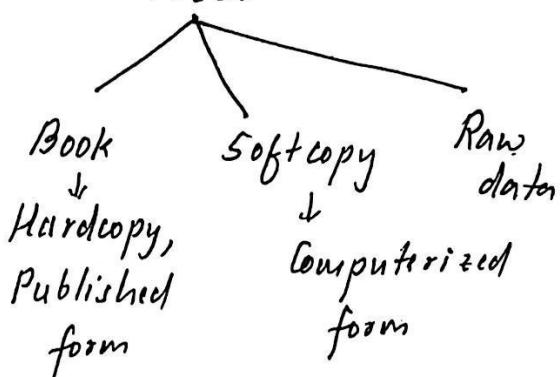
Factors affecting for poverty in Nepal



dependent variable	Independent variable
↓	
1	
0	

factors affecting employment status in Nepal: NLFS

factors affecting poverty in Nepal: NLSS



Computerized form:
 Internet - Free down loadable
 Online - Restricted
 Offline -
 Pendrive, CDs
 etc

Logistic Regression

→ Binary logistic regression is a form of regression which is used when the dependent variable is a true or forced dichotomy and the independent variables are of any type.

Employment Status : Not in labor force
 Employed
 Unemployed

- logistic regression can be used to predict a categorical dependent variable on the basis of continuous and/or categorical independent variables; to determine the effect size of the independent variables on dependent variable
- logistic regression applies maximum likelihood estimation after transforming the dependent into a logit variable.
- A logit is the natural log of the odds

(IMP)

Specification of the model

Y		No multicollinearity		
		x_1	x_2	\dots
1				
0				
0				

$$P_i = P(Y_i = 1) = F(z_i) = \frac{1}{1 + e^{-z_i}}$$

where, P_i is the probability

$$\begin{aligned}
 1 - P_i &= 1 - \frac{1}{1 + e^{-z_i}} = \frac{1 + e^{-z_i} - 1}{1 + e^{-z_i}} = \frac{e^{-z_i}}{1 + e^{-z_i}} \\
 &= \frac{1}{\frac{1 + e^{-z_i}}{e^{-z_i}}} = \frac{1}{1 + e^{z_i}}
 \end{aligned}$$

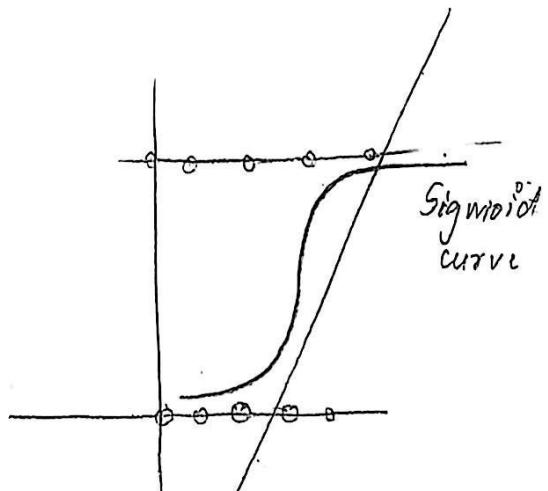
$$\begin{aligned}
 \frac{P_i}{1-P_i} &= \frac{\frac{1}{1+e^{-Z_i}}}{\frac{1}{1+e^{Z_i}}} = \frac{1+e^{Z_i}}{1+e^{-Z_i}} = \frac{1+e^{Z_i}}{1+\frac{1}{e^{Z_i}}} \\
 &= \frac{1+e^{Z_i}}{\frac{1+e^{Z_i}}{e^{Z_i}}} \\
 &= e^{Z_i}
 \end{aligned}$$

Now,

$$\ln\left(\frac{P_i}{1-P_i}\right) = Z_i = \alpha + \beta X_i \dots \dots \dots \quad (2)$$

Hence, $\frac{P_i}{1-P_i}$ is called odds ratio in favor of the event occurring

and $\ln\left(\frac{P_i}{1-P_i}\right)$ is the log odds ratio (also called logit of P).



Logistic Regression has many analogies to OLS regression:

Pseudo- R^2

Numerical
Y
T
 R^2
for categorical (high)
Pseudo- $R^2 = D_1/D$
(PL, High, significant)

Assumptions

- Models the natural log of odds (logits) of success probability.
- Associated dependent variable : Dichotomous (0/1)
- Probability distribution : Binomial
- link: logit
- Assumption
- Non autocorrelation
- linear relationship between predictors and log odds of dependent variable.
-

Measuring goodness-of-fit

* Effron's R^2

$$\text{Effron's } R^2 = 1 - \frac{n}{n_1 n_2} \sum (Y - \hat{Y})^2$$

$$\therefore \text{Wald's Test} = \left(\frac{b_i}{\hat{\sigma} b_i} \right)^2 = t^2$$

for individual test

* MacFadden's Pseudo R^2

$$\text{Pseudo } R^2 = 1 - \frac{d_n L}{d_n L_0}$$

Estimating overall significance of regression

$$LR = 2 \ln \frac{L}{L_0} = 2(\ln L - \ln L_0)$$

$LR \sim \chi^2$ with k degrees of freedom (explanatory variables)

$$\underline{\text{Odds Ratio}} : \frac{P(S)}{P(F)} = \frac{p_i}{1-p_i}$$

Odds Ratio: Ratio of two odds

log odds: Also called logit odds = $\log \left(\frac{p_i}{1-p_i} \right)$

e^{b_i} ^{slope} \Rightarrow Odds ratio for an independent variable

Smoking data (logistic Regression)

Smoker Y/N	age	income	Education	price of cigarettes (qt 10%)
***			***	
—	sign		— sign	— sign

Negative slope \Rightarrow odds ratio < 1

Positive slope \Rightarrow odds ratio > 1

April 25, 2025

Regression Analysis

Date _____
Page _____

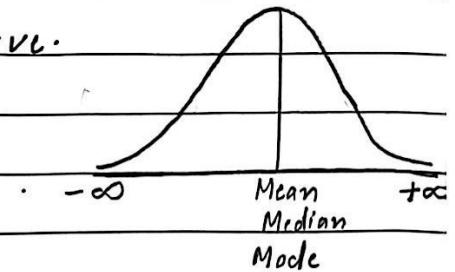
Checking Normality

(e)

Normal distribution

Normal
Distribution

It is a bell-shaped curve.



Properties

- * Symmetric
- * Mean = Median = Mode
- * Asymptotic curve

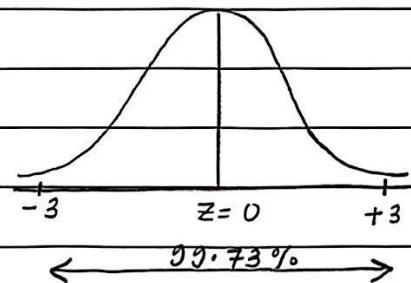
Probability function

Probability density function (pdf)

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

Standard normal variate

$$Z = \frac{x - \mu}{\sigma}$$



Data



Numerical

Σ

Discrete

\int

Continuous

Integer

(whole nos.)

Decimal

(fractional)

Counting processes (no. of colleges/books/students/visitors/defectives/calls)

measuring processes (price, wt, salary, eps)

$$\mu = 55$$

$$\sigma = 6$$

$$P(50 < x < 60) = ?$$

$$60$$

$$\int_{50}^{60} f(x) dx$$

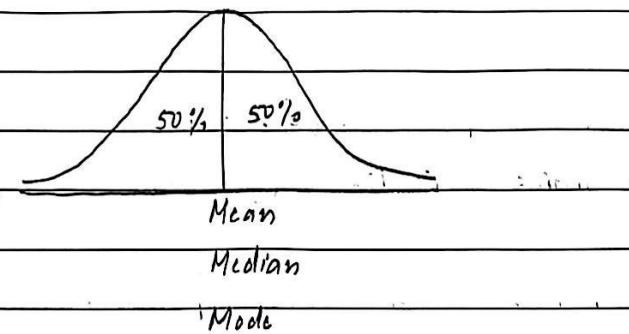
$$= \int_{50}^{60} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} dx$$

Detection

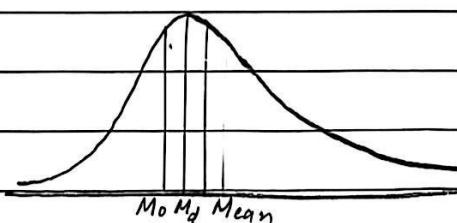
i) Mathematical method

$\text{Mean} \neq \text{Median} \neq \text{Mode}$

Normal \Rightarrow Symmetric



Skewed

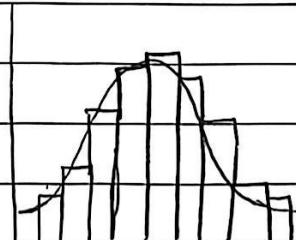


ii) Graphical method

a. Histograms

b. Frequency curve

(Density curve)

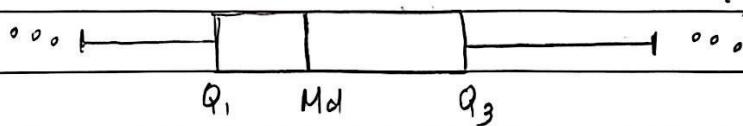
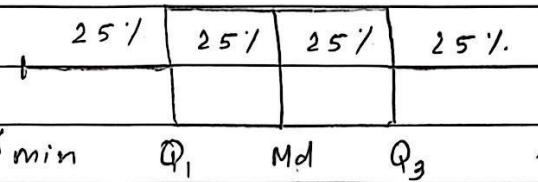


Histogram
Frequency curve

c. Box-and-Whisker plot

Five number summary

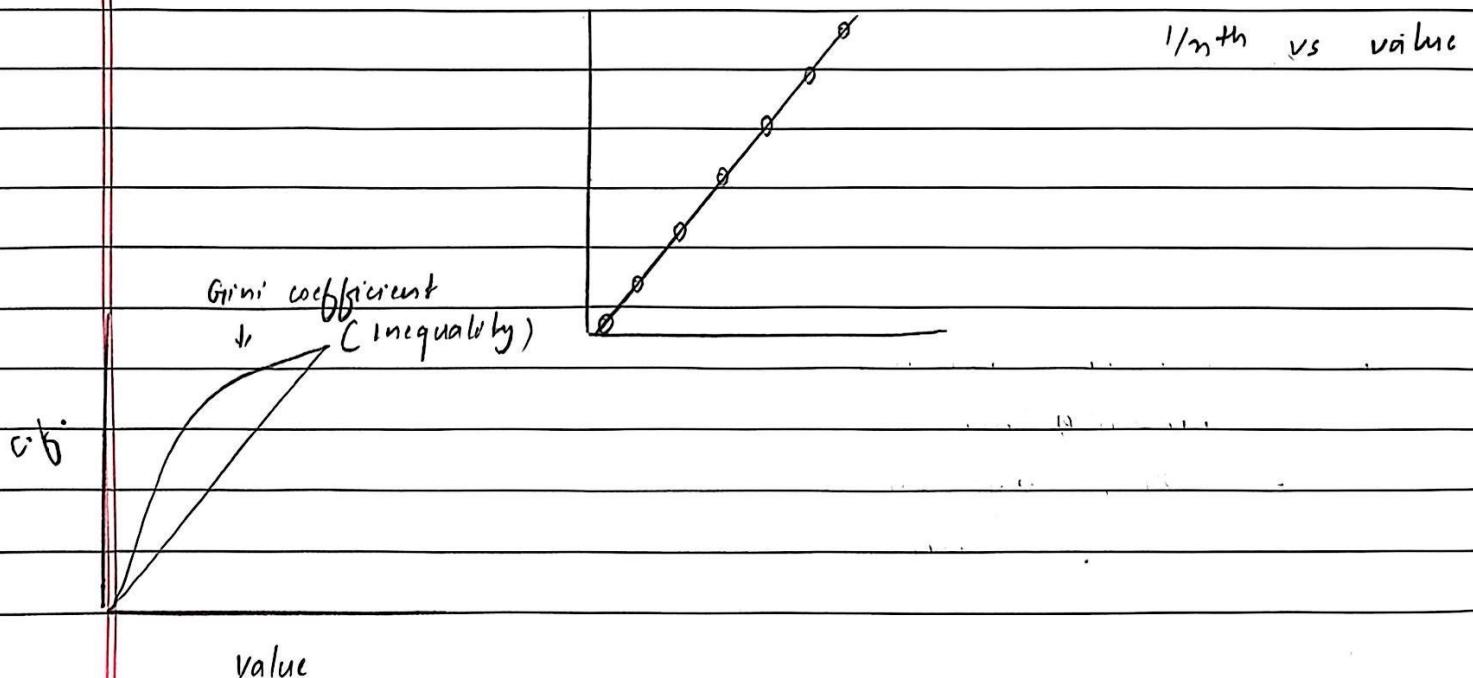
x_{\min}	Q_1	Md	Q_3	x_{\max}
------------	-------	------	-------	------------



$$\text{Lower limit} = Q_1 - 1.5(Q_3 - Q_1)$$

$$\text{Upper limit} = Q_3 + 1.5(Q_3 - Q_1)$$

d. Q-Q Plot (Quantile-Quantile Plot)



Testing

- Kolmogorov Smirnov Test
- Shapiro-Wilk's Test

} non-parametric

Null Hypothesis (H_0): Data are normally distributed

Alternative Hypothesis (H_1): Data are not normally distributed.

$p\text{-value} < \alpha (0.05) \Rightarrow$ we reject H_0 .

$p\text{-value} > \alpha (0.05) \Rightarrow$ we do not reject H_0 .

Impact of violation of normality

Almost all parametric tests become invalid.

t-test

z-test

F-test

→ no longer valid

Regression:

errors → normally distributed
 ↓

violated

tests are not valid

Data

Categorical



Non-parametric test

Numerical



Normality

Yes

No

Parametric

Non-parametric

Popular parametric test

Single mean \rightarrow t-test (z-test)

Two means \rightarrow t-test

Three or more means \rightarrow F-test

Correlation test \rightarrow t-test

Regression

Individual regression coefficient : t-test

$$H_0: \beta_i = 0$$

$$H_1: \beta_i \neq 0$$

Overall test : F-test

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1: \beta_i \neq 0$$

$$F = \frac{MSR}{MSt}$$

Cal F

Tab F

Decision.

Popular non-parametric test

Chi-Square test for independence

(Significant association between two categorical variables)

Profession

Brand of car

Goodness-of-fit tests.

Kolmogorov-Smirnov Test (Normality)

Mann-Whitney Test

To compare two groups (alternative technique of two independent sample t-test)
 $\mu_1 = \mu_2$

Kruskal-Wallis test.

To compare three or more groups (Alternative technique of F-Test or one way ANOVA)

$$[\mu_1 = \mu_2 = \mu_3 = \mu_4]$$

April 25, 2025

Date
Page

MARKOV CHAIN

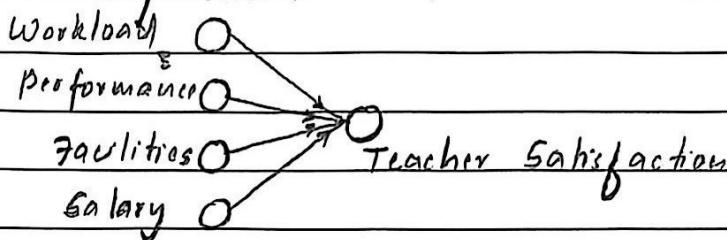
April 26,
2025 Stepwise Regression

Model Selection Criteria

Variable Selection (inclusion of ^{the} variable in final model)

↓
literature Review (Objectives / Hypothesis) → Topic
↓
Model-1 Variable Selection (Conceptual framework) Objectives
Hypothesis

Equation



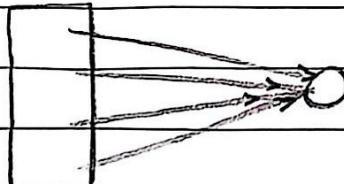
Model - 2

:

Model - 3



Own Conceptual framework



Literature Review



Variable Selection

→ Bivariate analysis (one independent and one dependent)
 significant
 insignificant

→ Multiple regression

Only significant in bivariate logit
 ordinal → VIF → Omitted
 multinomial

Stepwise Regression

	y	x_1	x_2	x_3	x_4	x_5 Rural/ x_6 Urban
Exp	income	family size	No. of child	no. of elder		$R=1$
						$U=0$

forward selection y | One variable



Result

R^2

y | Two variables *



Result

R^2

* Use significantly impacting variables

→ Standard Stepwise Regression selection
 Combines attributes of both forward, and backward elimination

Mallows' Cp

$$C_p = \frac{RSS_p / s^2}{n} - n + 2(p+1)$$

↓ Sample size

No. of predictor variables

lowest C_p is the best

Models

Define logistic model.

Which test is used for individual predictors in LR?
 Logistic Regression

Test for Individual Predictors

Individual Test

Gram

$$H_0: \beta_i = 0$$

$$H_1: \beta_i \neq 0$$

Test-statistic:

$$\text{Wald Statistic} \leftarrow W = \left(\frac{b_i}{s b_i} \right)^2$$

$$= \left(\frac{\hat{\beta}_i}{s \hat{\beta}_i} \right)^2$$

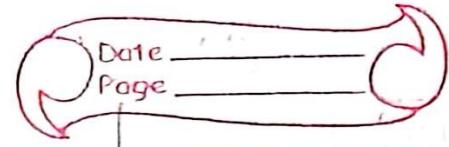


$p\text{-value} < \alpha \Rightarrow \text{reject } H_0$

$> \alpha \Rightarrow \text{do not reject } H_0$

TPR

FPR



Summarizing Predictive Power: Classification Model.

$$P(\hat{Y} = 1 | Y = 1)$$

Predicted

Y N

Actual Y O

N O

Sensitivity

1 - Specificity