

Review of data and statistics

1. Classification of data

→ ① cross sectional / time series / panel data

① Cross sectional data

→ collecting the data in single frame of time

→ conducting a survey only once.

Respondent individuals	Variable 1 age	Variable 2 salary	...	variable k gender
---------------------------	-------------------	----------------------	-----	----------------------

Purpose

→ identification of current status of statistical data.

→ comparison by analytical domain

by age, gender, province etc

→ identify relationship b/w variables

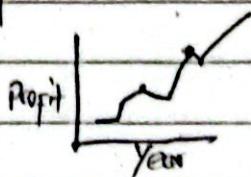
→ prediction / estimation

→ to make policy / strategy

② Time series

→ Collection of data over period of time

Year/GDP, Year/revenue, Year/profit



Purpose

→ To observe trend

→ Forecasting the future

→ construction of cycle

→ study of seasonal impact

④ Panel data

→ combined of cross-sectional & time series

2. Categorical vs Numerical

Categorical

→ any question having alphabetical or grouping response

→ Freq / %

→ any question having numeric response.

→ min/max/mean/std.dev.

3. According to no of variables

Univariate

→ analysis of single variables

Bivariate

→ ... two ...

multivariate

→ ... three or more

→ dependence technique

 ↳ one dependent variable

 ↳ two or more independent variables

\bar{x} - sample
 w - population

Date: _____
Page: _____

Population:

Totality of all items under the study. (Study area - coverage)

Census:- Complete enumeration

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

Parameter: Result/outcome of population data

Characteristics of population data

Population size $\rightarrow N$

" mean $\rightarrow \mu$

" proportion $\rightarrow \pi$ (or P)

" variance $\rightarrow \sigma^2$

" correlation $\rightarrow \rho$

" regression coefficient $\rightarrow \beta_i$

Sample:

Subset of the population.

Sampling: collecting the information/data from sample.

Estimator/~~Statistic~~^{statistic}: Result/outcome of sample data

Characteristic of sample data

Sample size $\rightarrow n$ (\hat{N})

" mean $\rightarrow \bar{x}$ ($\hat{\mu}$)

" proportion $\rightarrow p$ (\hat{p})

sample variance $\rightarrow s^2 (\hat{\sigma}^2)$

" correlation $\rightarrow r (\hat{\rho})$

" regression coefficient $\rightarrow b_i (\hat{\beta}_i)$

Sampling error $\rightarrow |\text{Estimator - Parameter}|$

\rightarrow used to determine sample size

* Sample design



II Research

Quantitative research



Survey
↓

Questionnaire (Instrument)

(Instrument should be valid)

(i.e. should measure the what
research wants to measures
→ objectives of the research)

→ should be correct & objective
reliable (validity)



answer/response → Correct

(data)

(reliability)



Data Management/Processing



Result → Data analysis



Knowledge



Policy/Plan/strategy/Program.

Null and alternative are

→ mutually exclusive

→ exclusive (belong to one)

Date:
Page:

II Statistics

Statistics

Descriptive

→ Describing the status
Table / chart

Central values / freq / ...
cross tabulation

Inferential

→ Technique of giving conclusion about
the population parameter on the basis
of sample data / estimator / statistic

→ Estimation

→ Hypothesis

* Hypothesis:

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

With the help of :

→ Past literature

→ Qualitative research / pilot

Null hypothesis (H_0): Statement of no difference or no relationship
 $'=', '≥', '≤'$ (equality)

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

\neq . \geq , \leq
two One

As $n \rightarrow N$

estimator \rightarrow parameter

Critical value \rightarrow point approach / value at the point
p-value \rightarrow area approach

Date: _____
Page: _____

* Level of significance (α):

Probability of rejecting the true statement

1% | 5% | 10%

* Confidence level ($1-\alpha$)

Probability of accepting the true statement

99% | 95% | 90%

* Decision Criteria:

\rightarrow Critical value approach

\rightarrow If $|t\text{est statistic}| \leq \text{critical value (tab value)}$
do not reject H_0

\rightarrow P-value approach

If $p\text{-value} < \alpha$

reject H_0

else

do not reject H_0

Example:

Q1 Gender [2]	Q2. Income <u>34,000</u>	Q3	Q4	Q5	Q6
Female -1		2	1	34000	15000
Male -2	Q2 Expenditure <u>15000</u>				

Q2. Post [1]

Assistant -1

Officer -2

Manager -3

$\theta_1 \rightarrow$ nominal

$\theta_2 \rightarrow$ ordinal

$\theta_3 \rightarrow$ ratio/nominal data

* Chivariate

Case a: Categorical data (nominal/ordinal)

↳ Frequency / % analysis
↳ Graph

↳ Bar
↳ pie

Distribution by gender

Gender	Freq	%
Male		
Female		
Total		

Case b: Numerical data

↳ min/max/mean/std dev

↳ one sample t-test

↳ Graph

↳ Histogram

↳ Box-plot

* Bivariate analysis

Two variables at a time

Case a: If both are categorical/grouping

↳ cross-tabulation betw two variables with freq and

appropriate v.

↳ Test

↳ Chi square test

Post Gender	Assistant	Officer	Manager	
Female	70 (0.20)	90 (0.1)	100 (0.1)	100
Male	100 (0.3)	50 (0.1)	50 (0.2)	200
	170	70.	60	300

Case b: If one is categorical and one is numerical

→ comparison of numerical by categorical

→ independent sample t-test (Two samples)

→ One way ANOVA (f-test) (More than two samples)

Gender	Average	Min	Max	St.dev
Female				
Male				
Total				

Comparison betwⁿ two means

Comparison of income by post

Post\Salary	Average	Min	Max	St.dev
Assistant				
Officer				
Manager				

Case c: If both are numerical

→ correlation (relationship betwⁿ two variables)

→ regression (we need dependent and independent variable)

48

Regression analysis

Interpretation of dependent and independent variables

- no. of dependent variable → one
- no. of independent variables →
 - one
 - two or more
(multiple)
 - (x_1, x_2, \dots, x_n)

Example :-

Price x	Sales volume y	Sales volume	Price	no. of outlets
100	100	100	100	100
110	90	90	110	100
120	80	80	120	100
130	70	70	130	100
140	60	60	140	100
150	50	50	150	100
160	40	40	160	100
170	30	30	170	100
180	20	20	180	100
190	10	10	190	100
200	0	0	200	100

→ Developing an equation

$$y = (\beta_0 + \beta_1 x) + \epsilon \quad (\text{Popn})$$

y = real value / actual / collected

$$\hat{y} = b_0 + b_1 x \quad (\text{sample})$$

\hat{y} = estimated value / predicted

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

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon, \quad (\text{Popn})$$

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p + \epsilon \quad (\text{sample})$$

$$\hat{y} = b_0 + b_1 x_1 + \dots + b_p x_p$$

$$y - \hat{y} = \text{error/residuals}$$

Objectives:

① Prediction of dependent variable by using the information of independent variable.

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 dependent variable.

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

$$b_1 = \frac{\Delta y}{\Delta x} = \frac{\text{amount change in } y}{\text{per unit change in } x}$$

$$\ln \hat{y} = b_0 + b_1 \ln x$$

$$b_1 = \frac{\Delta \ln y}{\Delta \ln x} = \frac{\% \text{ change in } y}{\% \text{ change in } x} \quad [\text{Elasticity}]$$

$$\ln \hat{y} = b_0 + b_1 x$$

$$b_1 = \frac{\Delta y}{\Delta x} = \frac{\% \text{ change in } y}{\text{unit change in } x}$$

④ Computation of growth rate/amount per time period.

⑤ Proportion of variation in y (dependent variable) that is explained by independent variable.

[coeff. of determination - R^2]

Price (x) | sales volume

$$b_1 = -0.75$$

$$R^2 = 0.82$$

82% of variations in sales is explained by Price.

⑥ Policy/program/pln/decision

Simple regression analysis

dependent variable \rightarrow one (y)

independent " \rightarrow one (x)

Equations:

$$y = \beta_0 + \beta_1 x + e \text{ (popn)}$$

$$y = b_0 + b_1 x + e \text{ (simple)}$$

By using least sq. method

$$b_1 = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2} = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

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

Example:

x (income) in 1000	y (expenditure) in 1000	xy	x^2	y^2
15	10	150	225	100
18	12	216	324	144
20	12	240	400	144
20	17	340	400	289
25	17	425	625	289
29	18	522	841	324
127	86	1893	2815	1250

$$\sum xy = 1893 \quad \sum x^2 = 2815 \quad \sum y^2 = 1250$$

$$\sum x = 127 \quad \bar{x} = 86$$

$$b_1 = \frac{n \sum xy - \bar{x}(\bar{y})}{n \sum x^2 - (\bar{x})^2}$$

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

= observed value - estimated value

$$> \frac{6 \times 1873 - 127 \times 86}{6 \times 2815 - 127^2}$$

$$= \frac{136}{761} = 0.573$$

standard error of estimate

$$SE(\hat{y}) = \sqrt{\frac{\sum e^2}{n-2}}$$

$$b_0 = \frac{\sum y}{n} - b_1 \frac{\sum x}{n} = 0$$

$$= \frac{86}{6} - \left(\frac{136}{761} \right) \times \frac{127}{6}$$

$$= 2.206$$

The estimation equation

$$\hat{y} = 2.206 + 0.572x$$

$$\hat{y} = 2.206 + 0.572x$$

$\hat{y} = 2.206 + 0.572x$	$(\hat{y} - y)$	$(\hat{y} - y)^2$
16.786	-0.786	0.617
12.502	-0.502	0.252
13.146	-1.696	2.709
13.146	-3.359	11.293
16.506	0.494	0.244
18.799	-0.799	0.630
	$0.32 \approx 0$	15.701
		15.701

→ error terms should be randomly distributed

→ if there is pattern in sign of error it is auto-correlation

d) Standard error of estimate

$$s_e(\hat{y}) = s_{yx} = s_e(\hat{y})$$

$$= \sqrt{\frac{\sum (y - \hat{y})^2}{n-p-1}}$$

$$= \sqrt{MSE}$$

$p \rightarrow$ no. of independent variable

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

Multiple regression analysis

Example:

Sales Volume	Price	no of stores	no of orders	P = 3
y	x ₁	x ₂	x ₃	k = 3
				} independent variable

Note: If the graph of error follows certain pattern there is auto-correlation which means there is problem in data.

* In matrix form

$$\begin{array}{|c|cccc|} \hline y_i & x_{1i} & x_{2i} & \dots & x_{ki} \\ \hline x_1 & x_{11} & x_{21} & \dots & x_{k1} \\ x_2 & x_{21} & x_{22} & \dots & x_{k2} \\ \vdots & & & & \\ \hline \end{array}$$

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \epsilon_i$$

$$\text{i.e. } Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + E_i$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{21} & \dots & x_{k1} \\ 1 & x_{22} & x_{22} & \dots & x_{k2} \\ \vdots & \vdots & \vdots & & \\ 1 & x_{1n} & x_{2n} & \dots & x_{kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} E_1 \\ E_2 \\ \vdots \\ E_n \end{bmatrix}$$

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

$$= E(x - \bar{x})^2$$

$$= E[(x - \bar{x})(\bar{x} - \bar{x})^T]$$

ANOVA table:

Software \rightarrow F-value & p-value

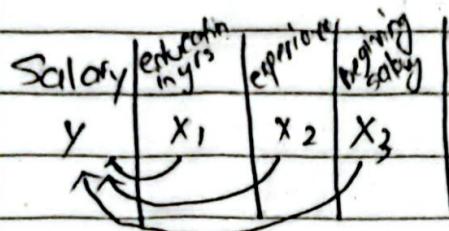
coefficient table

summary \rightarrow coefficient table

\hookrightarrow t-value + p-value

F-test \rightarrow it is used for overall test

t-test \rightarrow it is used for individual test



impact of $\underbrace{x_1 \text{ on } y}$, $\underbrace{x_2 \text{ on } y}$, $\underbrace{x_3 \text{ on } y} \rightarrow$ t-test (individual test)

$$\beta_1 = 0$$

$$\beta_2 = 4$$

$$\beta_3 = 0$$

impact of x_1, x_2 and x_3 on $y \rightarrow$ overall test

$$\beta_1 = \beta_2 = \beta_3 = 0$$

F-test $\Rightarrow H_0: \beta_1 = \beta_2 = \beta_3 = 0$

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

$H_1: \text{At least } \beta_i \neq 0$

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

if $p\text{-value} < \alpha (0.05) \rightarrow \text{reject } H_0$

$p\text{-value} \geq \alpha (0.05) \rightarrow \text{do not reject } H_0$

Format of ANOVA-table

Source	df	sum of squares	MSE	F-value	P-value
Regression	$k(0, p)$	SSR	$\frac{SSR + MSE}{k}$	$F = \frac{MSE}{MSE}$	
Error	$n-k-1$	SSE	$\frac{SSE}{n-k-1} = MSE$		
Total	$n-1$	TSS			

Format of coeff. table

Predictor	Regression coeff.	st. error	t-value ($t = \frac{b_i}{S_{b_i}}$)	P-value
Constant	b_0	S_{b_0}		
x_1	b_1	S_{b_1}		
x_2	b_2	S_{b_2}		
:	:	:		
x_k	b_k	S_{b_k}		

From ANOVA table and coeff. table

→ Developing and equation & prediction

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

i) Computing multiple coeff. 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

SSE = Sum of squares due to error / residuals

$$TSS = SSR + SSE$$

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

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

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

81% of variations in dependent variable is explained by 'k' independent variables.

iii) Adjusted R^2

→ It is 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

Increment in independent variable increases R^2 .

So we used adjusted R^2 .

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

iv) Standard error of estimate

$$SE(\hat{y}) = S_{yx} = \sqrt{MSE}$$

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

v) Overall test/goodness of fit

Hypothesis testing for all regression coefficient simultaneously.

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

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

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

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

Test statistics

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

$$\text{Cal } F = F$$

$\left\{ \begin{array}{l} \text{at } \alpha \% \\ F_{k, n-k-1} \end{array} \right.$

$K \rightarrow$ df. for numerator
 $n-k-1 \rightarrow$ " denominator

If $\text{cal } F \leq \text{Tab } F$
 we do not reject H_0

If $\text{cal } F > \text{Tab } F$
 we reject H_0

P-value approach

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

If $p\text{-value} \geq \alpha$
 we don't reject H_0

Individual test (t-test) [coeff. table]

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

There is no significant impact of x_i on y

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

There is significant impact of x_i on y

Test statistics

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

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

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

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

P-value $\geq \alpha \Rightarrow$ do not reject

Confidence interval estimate

$$b_i \pm t_{n-k-1}, \alpha/2 s_{b_i}$$

↓ ↓ ↓

coeff table coeff
table ↓
statistical table.

ii) Example:

By using following ANOVA table obtained from 30 obs.

Source	df	SS	MSS	F-value
Regression	4	300	?	?
Error	?	?	?	
Total	?	500		

- a) Compute the given ANOVA table
- b) Compute standard error of estimate & interpret its meaning
- c) Compute multiple coeff. of determination and interpret its meaning
- d) Compute Adj. R²
- e) Set up null and alternative hypothesis and carry out F-test.

g) $n=30$

Source	df	SS	MSS	F-value
Regression	4	300 (SSR)	$300/4 = 75$	$75/25 = 3.37$
Error	$25(n-k-1)$	$200(\text{TSS}-\text{SSR})$	$200/25 = 8$	
Total	$29(n-1)$	$500 (\text{TSS})$		

g) Standard error of estimate (S_{yx}) = \sqrt{MSE}
 $= \sqrt{8}$
 $= 2.828$

The avg. variation of observed values of dependent variable around fitted eqⁿ is 2.828

c) Multiple coeff of determination

$$R^2 = \frac{\text{SSR}}{\text{SST}}$$

$$= \frac{300}{500} = \frac{3}{5} = 0.60$$

60% of variation in dependent variable is explained by 9 independent variables.

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

$$= 1 - \left\{ (1 - 0.6)^2 \cdot \frac{20}{25} \right\}$$

$$= 0.536$$

e) Null hypothesis (H_0): $\beta_1 = \dots = \beta_4$

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

Alternative hypothesis (H_1): $\beta_i \neq 0 \quad i = 1 \dots 4$

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

Test-statistics

$$F = \frac{\text{MSR}}{\text{MSE}} = 9.37$$

$$\text{Col } F = 2.37$$

$$\text{Tab } F = 2.759$$

col F > tab F

so, we reject H_0