

*Prepared By
Arya Kumar*

*Faculty of MBA, CIMA
aryatripathy@yahoo.com,
09853422575*

MULTIVARIATE ANALYSIS

Factor Analysis

- Centroid
- Principle Component Approach (PCA)

Prepared By ARYA KUMAR

Multivariate Analysis

- Many statistical techniques focus on just one or two variables. It simply means testing different independent variables that affect a dependent variable significantly.
- Multivariate analysis (MVA) techniques allow more than two variables to be analysed at once – Multiple regression is not typically included under this heading, but can be thought of as a multivariate analysis
- There are several techniques of conducting multivariate analysis:
 - i. **Multiple Regression**
 - ii. **Multiple Discriminant analysis**
 - iii. Multivariate analysis of variance (MANOVA)
 - iv. Canonical correlation analysis
 - v. **Factor Analysis**
 - a) **Centroid A and B**
 - b) **Principle Component Analysis(PCA)**
 - c) **Maximum Likelihood (ML)**
 - vi. Cluster Analysis
 - vii. Multidimensional scaling
 - viii. Latent structure Analysis

Note: As per the syllabus the contents mentioned in red colour are to be discussed

Factor Analysis

Prepared By ARYA KUMAR

- Factor analysis is a way to take a mass of data and shrinking it to a smaller data set that is more manageable and more understandable.
- It's a way to find hidden patterns, show how those patterns overlap and show what characteristics are seen in multiple patterns.

The two types: exploratory and confirmatory.

1. **Exploratory factor analysis** is if you don't have any idea about what structure your data is or how many dimensions are in a set of variables.
2. **Confirmatory Factor Analysis** is used for verification as long as you have a specific idea about what structure your data is or how many dimensions are in a set of variables.

some basic terms relating to factor analysis

Prepared By ARYA KUMAR

1. **Factor**: A factor is an underlying dimension that account for several observed variables. There can be one or more factors, depending upon the nature of the study and the number of variables involved in it.
2. **Factor-loadings**: Factor-loadings are those values which explain how closely the variables are related to each one of the factors discovered. They are also known as factor-variable correlations.
3. **Communality (h^2)**: Communality, symbolized as h^2 , shows how much of each variable is accounted for by the underlying factor taken together. A high value of communality means that not much of the variable is left over after whatever the factors represent is taken into consideration.
4. **Eigen value (or latent root)**: When we take the sum of squared values of factor loadings relating to a factor, then such sum is referred to as Eigen Value or latent root. Eigen value indicates the relative importance of each factor in accounting for the particular set of variables being analysed.

some basic terms relating to factor analysis

Prepared By ARYA KUMAR

5. Total sum of squares: When Eigen values of all factors are totalled, the resulting value is termed as the total sum of squares. This value, when divided by the number of variables (involved in a study), results in an index that shows how the particular solution accounts for what all the variables taken together represent.

6. Rotation: Rotation, in the context of factor analysis, is something like staining a microscope slide. From the standpoint of making sense of the results of factor analysis, one must select the right rotation. If the factors are independent orthogonal rotation is done and if the factors are correlated, an oblique rotation is made.

7. Factor scores: Factor score represents the degree to which each respondent gets high scores on the group of items that load high on each factor. Factor scores can help explain what the factors mean. With such scores, several other multivariate analyses can be performed.

Prepared By ARYA KUMAR

IMPORTANT METHODS OF FACTOR ANALYSIS

There are several methods of factor analysis, but they do not necessarily give same results. As such factor analysis is not a single unique method but a set of techniques. Important methods of factor analysis are:

- A. the centroid method;
- B. the principal components method;
- C. the maximum likelihood method.

Prepared By ARYA KUMAR

(A) Centroid Method of Factor Analysis

- This method of factor analysis, developed by L.L. Thurstone
- The centroid method tends to maximize the sum of loadings, disregarding signs, it is the method which extracts the largest sum of absolute loadings for each factor in turn.
- It is defined by linear combinations in which all weights are either + 1.0 or - 1.0.



Various steps involved in this method
are as follows

Prepared By ARYA KUMAR

1. **Computation of a matrix of correlations, R** , wherein unities are placed in the diagonal spaces
2. **The correlation matrix so obtained happens to be positive manifold**, in case the correlation matrix is not a positive manifold, then reflections must be made before the first centroid factor is obtained
3. The **first centroid** factor is determined as under
 - a) **The sum of the coefficients**
 - b) **Then the sum of these column sums (T) is obtained**
 - c) **The sum of each column obtained as per (a) above is divided by the square root of T obtained in (b) above, resulting in what are called centroid loadings.**

Prepared By ARYA KUMAR

4. To obtain second centroid factor (say B)- **obtain a matrix of residual coefficients**
 - a) the loadings for the **two variables on the first centroid factor are multiplied**
 - b) The resulting matrix of factor cross products may be **named as Q**
 - c) **Q, is subtracted** element by element **from the original matrix of correlation, R**
 - d) the result is the **first matrix of residual coefficients, R**
 - e) If the values are **found negative, then convert to positive and obtain a reflected matrix, R'**.
5. For subsequent factors (**C, D, etc.**) **the same process outlined above is repeated.** After the second centroid factor is obtained, cross products are computed forming, matrix, Q. This is then subtracted from R. (and not from R'.)

Prepared By ARYA KUMAR

Illustration 1

Given is the following correlation matrix, R, relating to eight variables with unities in the diagonal spaces:

		variables							
		1	2	3	4	5	6	7	8
variables	1	1.00	.709	.204	.081	.626	.113	.155	.774
	2	.709	1.00	.051	.089	.581	.098	.083	.652
	3	.204	.051	1.00	.671	.123	.689	.582	.072
	4	.081	.089	.671	1.00	.022	.798	.613	.111
	5	.626	.581	.123	.022	1.00	.047	.201	.724
	6	.113	.098	.689	.798	.047	1.00	.801	.120
	7	.155	.083	.582	.613	.201	.801	1.00	.152
	8	.774	.652	.072	.111	.724	.120	.152	1.00

Solution : Given correlation matrix, R is a positive manifold and such the weights for all variables be +1.0.

Calculating the first centroid factor(A)-

Solution

Prepared By **ARYA KUMAR**

Step 1. The first centroid factor is determined as under

a) *The sum of the coefficients*

		Variables							
	1	2	3	4	5	6	7	8	
1	1.00	.709	.204	.081	.626	.113	.155	.774	
2	.709	1.00	.051	.089	.581	.098	.083	.652	
3	.204	.051	1.00	.671	.123	.689	.582	.072	
4	.081	.089	.671	1.00	.022	.798	.613	.111	
5	.626	.581	.123	.022	1.00	.047	.201	.724	
6	.113	.098	.689	.798	.047	1.00	.801	.120	
7	.155	.083	.582	.613	.201	.801	1.00	.152	
8	.774	.652	.072	.111	.724	.120	.152	1.00	
Column sums	3.662	3.263	3.392	3.385	3.324	3.666	3.587	3.605	

b) Then the sum of these column sums (T) is obtained

$$T = 3.662 + 3.263 + 3.392 + 3.385 + 3.324 + 3.66 + 3.567 + 3.605 = 27.884$$

$$\sqrt{T} = 5.281$$

Prepared By ARYA KUMAR

c) The sum of each column obtained as per (a) above is divided by the square root of T obtained in (b) above, resulting in what are called centroid loadings.

First Centroid factor A=

$$\frac{3.662 \quad 3.662 \quad 3.263 \quad 3.392 \quad 3.385 \quad 3.324 \quad 3.66 \quad 3.567 \quad 3.605}{5.2281 \quad 5.2281 \quad 5.2281 \quad 5.2281 \quad 5.2281 \quad 5.2281 \quad 5.2281 \quad 5.2281 \quad 5.2281}$$

=.693, .618, .642, .641, .629, .694, .679, .683

First centroid factor A=

Columns	Factor A
1 -	0.693
2 -	0.618
3 -	0.642
4 -	0.641
5 -	0.629
6 -	0.694
7 -	0.679
8 -	0.683

Prepared By ARYA KUMAR

Step 2. The second centroid factor B

- a) the first matrix of factor cross product, **Q1- Multiply the two variables on the first centroid**

First matrix of factor cross product (Q1)

		First Centroid Factor A							
First Centroid Factor A		.693	.618	.642	.641	.629	.694	.679	.683
.693		.480	.428	.445	.444	.436	.481	.471	.473
.618		.428	.382	.397	.396	.389	.429	.420	.422
.642		.445	.397	.412	.412	.404	.446	.436	.438
.641		.444	.396	.412	.411	.403	.445	.435	.438
.629		.436	.389	.404	.403	.396	.437	.427	.430
.694		.481	.429	.446	.445	.437	.482	.471	.474
.679		.471	.420	.436	.435	.427	.471	.461	.464
.683		.473	.422	.438	.438	.430	.474	.464	.466

i.e. $.693 \times .693 = .480$, $.693 \times .618 = .428$, $.683 \times .683 = .466$

(cross multiple each numbers of centroid A which is presented in row and column at top and side wise)

Prepared By ARYA KUMAR

b) Now obtain residual coefficient (R1) by subtracting Q1 from R=

Row1 and column 1 of (R)- Row1 and column 1 of (Q1)
 $= 1.000 - 0.480 = 0.520$

Row1 and column 2 of (R)- Row1 and column 2 of (Q1)
 $= 0.709 - 0.428 = 0.281$ and so on for all the series

(i.e. the R value will be subtracted from respective Q1)

First matrix of residual coefficient (R_1)

		Variables							
		1	2	3	4	5	6	7	8
variables	1	0.520	0.281	-0.241	-0.363	0.190	-0.368	-0.316	0.301
	2	0.281	0.618	-0.346	-0.307	0.192	-0.331	-0.337	0.230
	3	-0.241	-0.346	0.588	0.259	-0.281	0.243	0.146	-0.366
	4	-0.363	-0.307	0.259	0.589	-0.381	0.353	0.178	-0.327
	5	0.190	0.192	0.281	0.381	0.604	-0.390	-0.217	0.294
	6	-0.368	-0.331	0.243	0.353	-0.390	0.518	0.330	-0.354
	7	-0.316	-0.337	0.146	0.178	-0.226	0.330	0.539	-0.312
	8	0.301	0.230	-0.366	-0.327	0.294	-0.354	-0.312	0.534

Prepared By **ARYA KUMAR**

c) Now obtain reflected matrix of residual coefficient ($R'1$) by converting – sign to + sign for the variables 3,4,6,7

Then add all the column variables
(repeat all the steps like Centroid-A)

Extraction of second centroid (B)

Variables		1	2	3*	4*	5	6*	7*	8
1	0.520	0.281	-0.241	-0.363	0.190	-0.368	-0.316	0.301	
2	0.281	0.618	-0.346	-0.307	0.192	-0.331	-0.337	0.230	
3*	-0.241	-0.346	0.588	0.259	-0.281	0.243	0.146	-0.366	
4*	-0.363	-0.307	0.259	0.589	-0.381	0.353	0.178	-0.327	
5	0.190	0.192	0.281	0.381	0.604	-0.390	-0.217	0.294	
6*	-0.368	-0.331	0.243	0.353	-0.390	0.518	0.330	-0.354	
7*	-0.316	-0.337	0.146	0.178	-0.226	0.330	0.539	-0.312	
8	0.301	0.230	-0.366	-0.327	0.294	-0.354	-0.312	0.534	
Column sums	2.580	2.642	2.470	2.757	2.558	2.887	2.375	2.718	

Then the sum of these column sums (T) is obtained

$$T = 2.580 + 2.642 + 2.470 + 2.757 + 2.558 + 2.887 + 2.375 + 2.718 = 20.987$$

$$\sqrt{T} = 4.581$$

Prepared By ARYA KUMAR

d) The sum of each column obtained as per (c) above is divided by the square root of T obtained in above, resulting in what are called centroid loadings.

First Centroid factor A=

$$\frac{2.580}{4.581}, \frac{2.642}{4.581}, \frac{2.470}{4.581}, \frac{2.757}{4.581}, \frac{2.558}{4.581}, \frac{2.887}{4.581}, \frac{2.375}{4.581}, \frac{2.718}{4.581}$$
$$=.563, .577, -.539, -.605, .558, -.630, -.518, .593$$

Step 3.

Representation of Factor loadings.

Factor loadings		
VARIABLES	Centroid Factor A	Centroid Factor B
1	.693	.563
2	.618	.577
3	.642	-.539
4	.641	-.602
5	.629	.558
6	.694	-.630
7	.679	.518
8	.683	.593

For subsequent factors(C,D, etc.)the same process outlined above is repeated.

- Cross products are computed forming, matrix Q2
- To obtain a third factor (C) one should operate R2 in the same way of R1.

Prepared By ARYA KUMAR

Step 4.

Find out the communality

Variables	Factor loadings		Communality (h^2)
	Centroid Factor A	Centroid Factor B	
1	.693	.563	$(.693)^2 + (.563)^2 = .797$
2	.618	.577	$(.618)^2 + (.577)^2 = .715$
3	.642	-.539	$(.642)^2 + (-.539)^2 = .703$
4	.641	-.602	$(.641)^2 + (-.602)^2 = .773$
5	.629	.558	$(.629)^2 + (.558)^2 = .707$
6	.694	-.630	$(.694)^2 + (-.630)^2 = .879$
7	.679	-.518	$(.679)^2 + (-.518)^2 = .729$
8	.683	.593	$(.683)^2 + (.593)^2 = .818$

Step 5.

Prepared By ARYA KUMAR

Find out the Eigen values from the final results and its interpretation

Variables	Factor Loadings		Communality
	Centroid A	Centroid B	
Eigen Value (Variance accounted for i.e., common variance)	3.490 (square all the values in centroid A and add)	2.631 (square all the values in centroid B and add)	6.121 (Total of communality values obtained or add Centroid A and Centroid B)
Proportion of total variance i.e. out of 8 As 8 is the no. of variables	.44 (44%)	.33 (33%)	.77 (77%)
Proportion of common variances i.e. 6.121	.57 (57%)	.43 (43%)	1.00 (100%)

Step 6 (Final).

Interpretation

Prepared By **ARYA KUMAR**

- Each communality in the above table represents the proportion of variance in the corresponding (row) variable and is accounted for by the two factors (A and B).
- For instance, 79.7% of the variance in variable one is accounted for by the centroid factor A and B
- The total value, 8.0, is partitioned into 3.490 as Eigen value for factor A and 2.631 as Eigen value for factor B and the total 6.121 as the sum of Eigen values for these two factors.
- 77% of the total variance is common variance whereas remaining 23% of it is made up of portions unique to individual variables and the techniques used to measure them.
- The last row shows that of the common variance approximately 57% are accounted for by factor A and the other 43% by factor B.
- Thus it can be concluded that the two factors together "explain" the common variance.
- **In one line the factor A and factor B both explains the dependent factor by 77%, while individually factor A explain 57% and B explains 43%.**

Prepared By ARYA KUMAR

B. Principal Component Analysis

Principal component analysis (PCA) is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components



Karl Pearson

- PCA was invented by him in 1901.



Harold Hotelling

- It was later independently developed and named by him in 1930s.

Prepared By ARYA KUMAR

Principal component analysis

- ▶ A backbone of modern data analysis.
- ▶ A black box that is widely used but poorly understood.
- ▶ It is a mathematical tool from applied linear algebra.
- ▶ It is a simple, non-parametric method of extracting relevant information from confusing data sets.
- ▶ It provides a roadmap for how to reduce a complex data set to a lower dimension

Prepared By ARYA KUMAR

Method

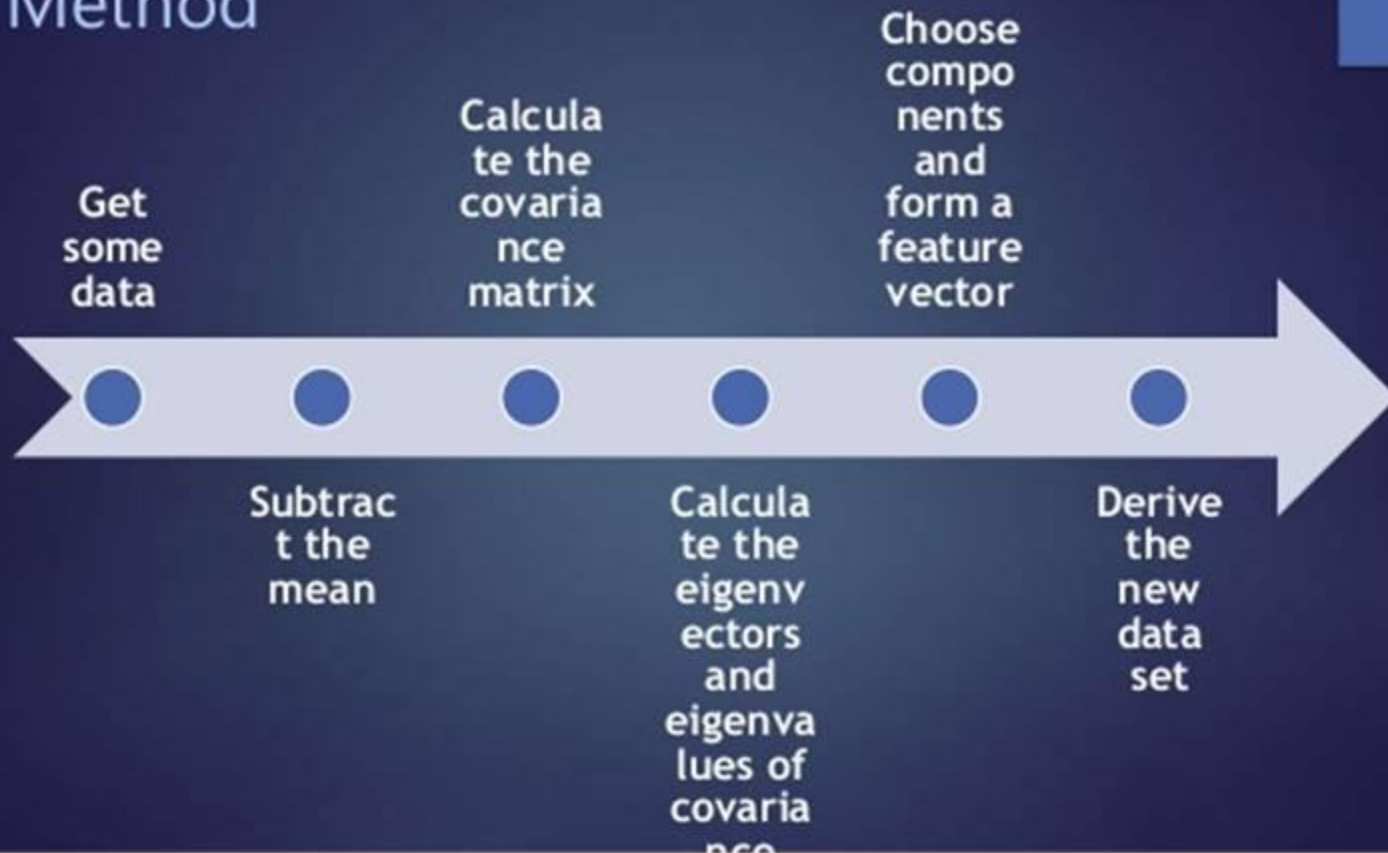


Illustration 2

Prepared By ARYA KUMAR

Take the correlation matrix, R , for eight variables of illustration 1 of this chapter and then compute:

- i. the first two principal component factors i.e. I and II
- ii. the communality for each variable on the basis of said two component factors;
- iii. the proportion of total variance as well as the proportion of common variance explained by each of the two component factors.

Note: all the three points are followed and explained in details

Prepared By ARYA KUMAR

Solution: Since the given correlation matrix is a positive manifold, we work out the first principal component factor (using trial vectors) as under:

Now calculate Principal Component I-

		variables							
		1	2	3	4	5	6	7	8
variables	1	1.00	.709	.204	.081	.626	.113	.155	.774
	2	.709	1.00	.051	.089	.581	.098	.083	.652
	3	.204	.051	1.00	.671	.123	.689	.582	.072
	4	.081	.089	.671	1.00	.022	.798	.613	.111
	5	.626	.581	.123	.022	1.00	.047	.201	.724
	6	.113	.098	.689	.798	.047	1.00	.801	.120
	7	.155	.083	.582	.613	.201	.801	1.00	.152
	8	.774	.652	.072	.111	.724	.120	.152	1.00

Step 4. To obtain Ua2 i.e.

Prepared By ARYA KUMAR

- a) the **first matrix** of factor cross product (follow the same steps like centroid)
- b) Represent in a table row wise and column wise then **multiply the row with column**

$$\text{i.e. } .371 \times .371 = .137, .371 \times .327 = .121 \dots \dots \dots, .365 \times .365 = .133$$

(cross multiple each numbers of centroid A which is presented in row and column at top and side wise)

- c) Add all the values in column Ua2=

$$1.296, 1.143, 1.201, 1.201, 1.165, 1.308, 1.280, 1.275$$

- d) Obtain normalizing factor Ua2 (similarly like the above slide)=

$$\text{i.e. square root of squaring the values of Ua2} = 3.493$$

- e) Find Va2, i.e. Ua2/NF=

$$.371, .327, .344, .344, .334, .374, .366, .365$$

Prepared By ARYA KUMAR

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Steps:

Step-1- Add all the values presented in column and name it as **Ua1**

$$=3.662,3.263,3.392,3.385,3.324,3.666,3.587,3.605$$

Step-2- find **normalization factor** i.e. square root of squaring the values of Ua1

$$\sqrt{(3.662)^2 + (3.263)^2 + (3.392)^2 + (3.385)^2 + (3.324)^2 + (3.666)^2 + (3.587)^2 + (3.605)^2} = 9.868$$

Step-3- **obtain Va1** i.e. Ua1/normalizing factor.

		variables							
		1	2	3	4	5	6	7	8
variables	1	1.00	.709	.204	.081	.626	.113	.155	.774
	2	.709	1.00	.051	.089	.581	.098	.083	.652
	3	.204	.051	1.00	.671	.123	.689	.582	.072
	4	.081	.089	.671	1.00	.022	.798	.613	.111
	5	.626	.581	.123	.022	1.00	.047	.201	.724
	6	.113	.098	.689	.798	.047	1.00	.801	.120
	7	.155	.083	.582	.613	.201	.801	1.00	.152
	8	.774	.652	.072	.111	.724	.120	.152	1.00
Ua1		3.662	3.263	3.392	3.385	3.324	3.666	3.587	3.605
Normalizing Factor		$\sqrt{(3.662)^2 + (3.263)^2 + (3.392)^2 + (3.385)^2 + (3.324)^2 + (3.666)^2 + (3.587)^2 + (3.605)^2}$ = 9.868							
Va1= Ua1/NF		.371	.33	.344	.343	.337	.372	.363	.365
			1						

Step 5. Compare V_{a1} and V_{a2} Prepared By ARYA KUMAR
Square root of normalizing factor of $U_{a2} = \sqrt{3.493}$

Comparing V_{a1} and V_{a2} , we find the two vectors are almost equal and this shows convergence has occurred. Hence V_{a1} is taken as the characteristic vector, V_a . Finally, we compute the loadings on the first principal component by multiplying V_a by the square root of the number that we obtain for normalizing U_{a2} . The result is as under:

Variables	(Characteristic vector V_a)	×	$\sqrt{\text{normalizing factor of } U_{a2}}$	=	Principal Component I
1	.371	×	1.868	=	.69
2	.331	×	1.868	=	.62
3	.344	×	1.868	=	.64
4	.343	×	1.868	=	.64
5	.337	×	1.868	=	.63
6	.372	×	1.868	=	.70
7	.363	×	1.868	=	.68
8	.365	×	1.868	=	.68

Prepared By ARYA KUMAR

Now calculate Principal Component II-

Repeat the similar steps (as stated in the context of obtaining centroid factor B earlier in this slides)

Variables	Principal Component II
1	.57
2	.59
3	-.52
4	-.59
5	.57
6	-.61
7	-.49
8	-.61

Step(Final)

Prepared By ARYA KUMAR

Find out the communality and the Eigen values from the final results and its interpretation

The other parts of the question can now be worked out (after first putting the above information matrix form) as given below:

Variables	Principal Components		Communality, h^2
	I	II	
1	.69	+.57	$(.69)^2 + (.57)^2 = .801$
2	.62	+.59	$(.62)^2 + (.59)^2 = .733$
3	.64	-.52	$(.64)^2 + (-.52)^2 = .680$
4	.64	-.59	$(.64)^2 + (-.59)^2 = .758$
5	.63	+.57	$(.63)^2 + (.57)^2 = .722$
6	.70	-.61	$(.70)^2 + (-.61)^2 = .862$
7	.68	-.49	$(.68)^2 + (-.49)^2 = .703$
8	.68	-.61	$(.68)^2 + (-.61)^2 = .835$
Eigen value i.e., common variance	3.4914	2.6007	6.0921
Proportion of total variance	.436 (43.6%)	.325 (32.5%)	.761 (76%)
Proportion of common variance	.573 (57%)	.427 (43%)	1.000 (100%)

Note: all this values can be interpreted as per the centroid A and B- (discussed in previous slide)

C. Maximum likelihood method

Prepared By ARYA KUMAR

1. The maximum likelihood (ML) method consists in obtaining sets of factor loading successively in such a way that each, in turn, explains as much as possible of the population correlation matrix as estimated from the sample correlation matrix.
2. If R_S stands for the correlation matrix actually obtained from the data in sample, R_p stands for the correlation matrix would be obtained if the entire population were tested, then the ML method seeks to extrapolate what is known from R_S is the best possible way to estimate R_p .
3. Thus the method is statistical approach in which one maximizes some relationship between the sample of data and the population from which the sample has drawn.
4. The loading obtained on the first factor are employed in the usual way to obtain a matrix of the residual coefficients.
5. A significance test is then applied to indicate whether it would be reasonable to extract a second factor.
6. This goes on repeatedly in search of one factor after another. One stops factoring after the significance fails to reject the null hypothesis for the residual matrix. The final product is a matrix of factor loading.

Note: this is similar to Centroid, but in Centroid the factor loading is mentioned immediately, while here we continue till reject null hyp.

Rotation in factor analysis

Simple structure is obtained by rotating the axes until:

Prepared By
ARYA KUMAR

- I. Each **row of the factor matrix has one zero.**
- II. Each **column of the factor matrix has p zero, where p is the number of factor.**
- III. For **each pair of factors, there are several variables for which the loading on one is virtually zero and the loading on the other is substantial.**
- IV. If **there are many factors, then for each pair of the factors there are many variables for which both loadings are zero.**
- V. For **every pair of factors, the number of variables with non- vanishing loading on both of them is small.** All these criteria imply that the factor analysis should reduce the complexity of all the variables.

R-type factor analysis

Prepared By ARYA KUMAR

In R-type factor analysis, high correlation occur when respondents who score high on variable 1 also score high on variable 2 and respondent who score low on variable 1 and also score low on variable 2.

Q-type factor analysis

In Q-type factor analysis, the correlation are computed between pairs of respondents instead of pairs of variables. High correlations occur when respondent 1's pattern of responses on all the variables is much like respondent 2's pattern of responses

MULTIVARIATE ANALYSIS

Prepared By
Arya Kumar

Faculty of MBA, CIME
aryantripathy@yahoo.com,
09853422575

Basic
Fundamentals

TECHNIQUES OF FACTOR ANALYSIS

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

1. Explanatory variable and criterion variable: If **X may be considered to be the cause of Y, then X is described as explanatory variable** (also termed as causal or independent variable) **and Y is described as criterion variable** (also termed as resultant or dependent variable). a set of many variables in which case set $(X_1, X_2, X_3, \dots, X_k)$ may be called a set of explanatory variables and the set $(Y_1, Y_2, Y_3, \dots, Y_k)$ may be called a set of criterion variables if the variation of the former may be supposed to cause the variation of the latter as a whole. *Prepared by Arya Kumar*
2. Observable variables and latent variables: **Explanatory variables described above are supposed to be observable directly in some situations**, and if this is so, the same are termed as observable variables. However, **there are some unobservable variables** which may influence the criterion variables. **We call such unobservable variables as latent variables.**

TECHNIQUES OF FACTOR ANALYSIS (CONTD.)

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

3) Discrete variable and continuous variable: Discrete variable is that variable which **when measured may take only the integer value whereas continuous** variable is one which, when measured, **can assume any real value (even in decimal points)**.

Prepared by Arya Kumar

4) Dummy variable (or Pseudo variable): This term is being used in a technical sense and is **useful in algebraic manipulations** in context of multivariate analysis. **We call X_i ($i = 1, \dots, m$) a dummy variable, if only one of X_i is 1 and the others are all zero.**

IMPORTANT MULTIVARIATE TECHNIQUES

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

1) Multiple regression: In multiple regression we form a linear composite of explanatory variables in such way that it has maximum correlation with a criterion variable.

- This is supposed to be a function of other explanatory variables.
- One can predict the level of the dependent phenomenon through multiple regression analysis model.
- Given a dependent variable, the linear-multiple regression problem is to estimate constants B_1, B_2, \dots, B_k and A such that the expression $Y = B_1X_1 + B_2X_2 + \dots + B_kX_k + A$ provides a good estimate of an individual's Y score based on his X scores.
- In practice, Y and the several X variables are converted to standard scores; $Z_y, Z_1, Z_2, \dots, Z_k$; each z has a mean of 0 and standard deviation of 1. Then the problem is to estimate constants, B_i , such that

$$z'_y = \beta_1 z_1 + \beta_2 z_2 + \dots + \beta_k z_k$$

- Where z' stands for the predicted value of the standardized Y score, z .

Prepared by Arya Kumar

IMPORTANT MULTIVARIATE TECHNIQUES (CONTD.)

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

2) **Multiple discriminant analysis:** Through discriminant analysis technique, **researcher may classify individuals or objects into one of two or more mutually exclusive and exhaustive groups** on the basis of a set of independent variables.

- Discriminant **analysis requires interval independent variables and a nominal dependent variable.** For example, suppose that brand preference (say brand x or y) is the dependent variable of interest and its relationship to an individual's income, age, education, etc. is being investigated, then we should use the technique of discriminant analysis.

Prepared by Arya Kumar

- **Regression analysis in such a situation is not suitable because the dependent variable is not interval scaled. Thus discriminant analysis is considered an appropriate technique when the single dependent variable happens to be non-metric** and is to be classified into two or more groups, depending upon its relationship with several independent variables which all happen to be metric

2) MULTIPLE DISCRIMINANT ANALYSIS

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

- For example, an individual is 20 years old, has an annual income of Rs 12,000, and has 10 years of formal education. Let b_1 , b_2 , and b_3 be the weights attached to the independent variables of age, income and education respectively. The **individual's score (z), assuming linear score**, would be:
$$z = b_1(20) + b_2(12000) + b_3(10)$$

Prepared by Arya Kumar

- The numerical values and signs of the **b's indicate the importance of the independent variables** in their ability to discriminate among the different classes of individuals. Thus through **the discriminant analysis, the researcher can as well determine which independent variables are most useful in predicting whether the respondent is to be put into one group or the other.**

2) MULTIPLE DISCRIMINANT ANALYSIS

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

- In case only two groups of the individuals are to be formed on the basis of several independent variables, we can then have a model like this

$$z_i = b_0 + b_1X_{1i} + b_2X_{2i} + \dots + b_nX_n$$

where X_{ji} = the i th individual's value of the j th independent variable;

b_j = the discriminant coefficient of the j th variable;

z_i = the i th individual's discriminant score;

Z_{crit} = the critical value for the discriminant score.

The classification procedure in such a case would be

If $z_i > Z_{crit}$, classify individual i as belonging to Group I

If $z_i < Z_{crit}$, classify individual i as belonging to Group II.

- When n (the number of independent variables) is equal to 2, we have a straight line classification boundary. Every individual on one side of the line is classified as Group I and on the other side; everyone is classified as belonging to Group II. When $n = 3$, the classification boundary is a two-dimensional plane in 3 space and in general the classification boundary is an $n - 1$ dimensional hyper-plane in n space.

Prepared by Arya Kumar

2) MULTIPLE DISCRIMINANT ANALYSIS

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

- In n-group discriminant analysis, a discriminant function is formed for each pair of groups. **If there are 6 groups to be formed, we would have $6(6 - 1)/2 = 15$ pairs of groups, and hence 15 discriminant functions**
- For **judging the statistical significance between two groups, we work out the Mahalanobis statistic, D^2** , which happens to be a generalized distance between two groups, where each group is characterized by the same set of n variables and where it is assumed that variance covariance structure is identical for both groups. It is worked out thus:

$$D^2 = (U_1 - U_2)v^{-1}(U_1 - U_2)$$

Where, U_1 = the mean vector for group I

U_2 = the mean vector for group II

v = the common variance matrix

- By transformation procedure, this D^2 statistic becomes an F statistic which can be used to see if the two groups are statistically different from each other

Prepared by Arya Kumar

2) MULTIPLE DISCRIMINANT ANALYSIS

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

- From all this, we can conclude **that the discriminant analysis provides a predictive equation, measures the relative importance of each variable and is also a measure of the ability of the equation to predict actual class-groups (two or more) concerning the dependent variable** *Prepared by Arya Kumar*

IMPORTANT MULTIVARIATE TECHNIQUES (CONTD.)

Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

3) Multivariate analysis of variance: Multivariate analysis of variance is **an extension of bivariate analysis of variance in which the ratio of among-groups variance to within-groups variance is calculated on a set of variables instead of a single variable.**

- This technique is considered appropriate **when several metric dependent variables are involved in a research study along with many non-metric explanatory variables.** *Prepared by Arya Kumar*
- multivariate analysis of variance is specially applied whenever the researcher wants **to test hypotheses concerning multivariate differences in group responses to experimental manipulations.**

For instance, the market researcher may be interested in using one test market and one control market to examine the effect of an advertising campaign on sales as well as awareness, knowledge and attitudes. In that case he should use the technique of multivariate analysis of variance for meeting his objective

MODULE-3 DATA ANALYSIS- II

Prepared by
Arya Kumar
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Regression
Analysis

Prepared By
ARYA KUMAR

MEANING OF REGRESSION:

The dictionary meaning of the word Regression is 'Stepping back' or 'Going back'. Regression is the measures of the average relationship between two or more variables in terms of the original units of the data. And it is also attempts to establish the nature of the relationship between variables that is to study the functional relationship between the variables and thereby provide a mechanism for prediction, or forecasting.

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

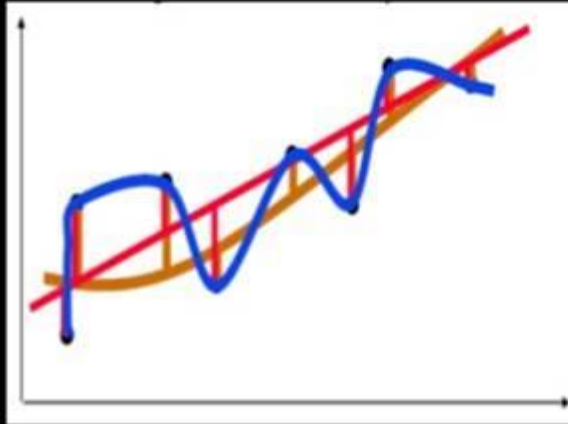
Prepared By ARYA KUMAR

- Father of Regression Analysis
Carl F. Gauss (1777-1855).
- contributions to physics, Mathematics & astronomy.
- The term “Regression” was first used in 1877 by Francis Galton.

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Prepared By ARYA KUMAR

Regression Analysis. . .



- It is the study of the relationship between variables.
- It is one of the most commonly used tools for business analysis.
- It is easy to use and applies to many situations.

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Importance of Regression Analysis KUMAR

Regression analysis helps in three important ways :-

- It provides estimate of values of dependent variables from values of independent variables.
- It can be extended to 2 or more variables, which is known as multiple regression.
- It shows the nature of relationship between two or more variable.

Faculty of MBA, CIME

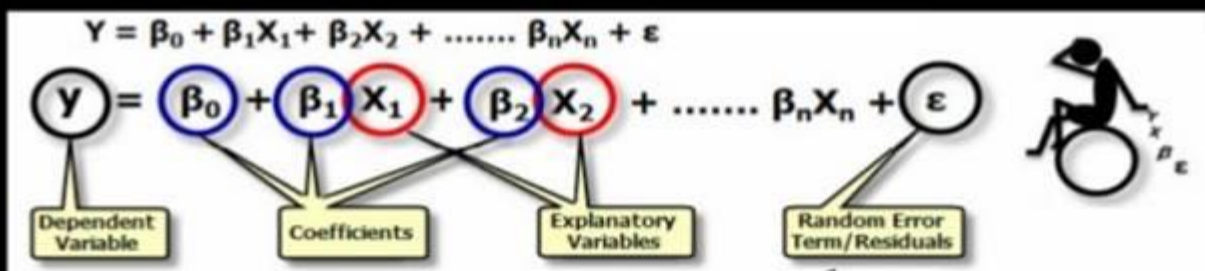
aryantripathy@yahoo.com, 09853422575

Regression types...

- Simple Regression: single explanatory variable
- Multiple Regression: includes any number of explanatory variables.

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Prepared By ARYA KUMAR



- Dependant variable: the single variable being explained/ predicted by the regression model
- Independent variable: The explanatory variable(s) used to predict the dependant variable.
- Coefficients (β): values, computed by the regression tool, reflecting explanatory to dependent variable relationships.
- Residuals (ϵ): the portion of the dependent variable that isn't explained by the model: the model under and over predictions.

Faculty of MBA, CIME

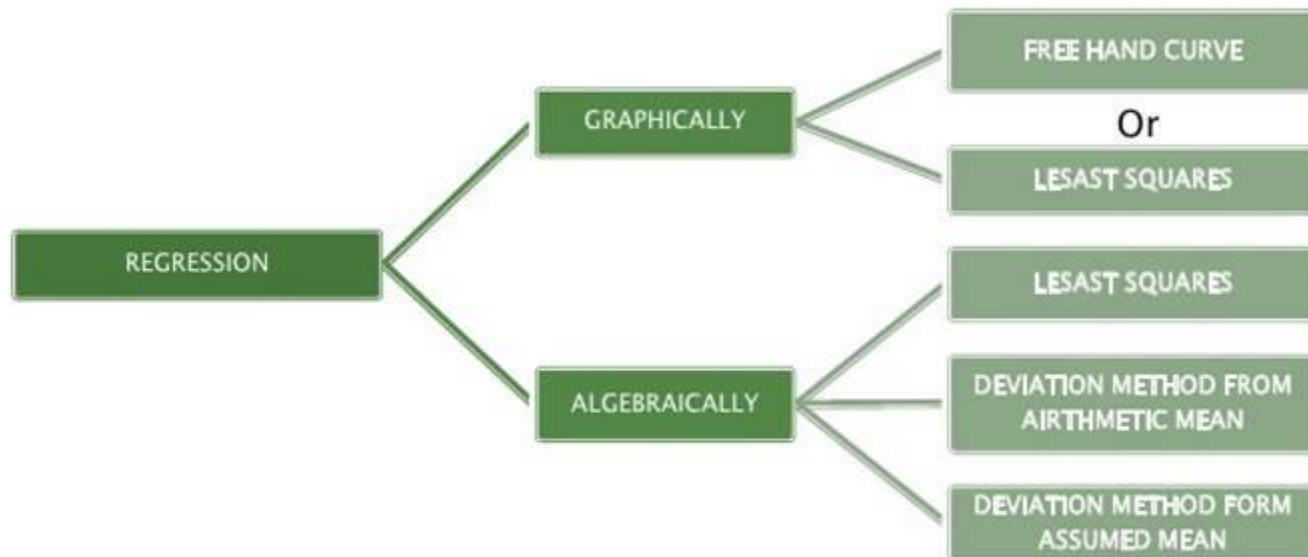
aryantripathy@yahoo.com, 09853422575

USE IN ORGANIZATION

In the field of business regression is widely used. Businessman are interested in predicting future production, consumption, investment, prices, profits, sales etc. So the success of a businessman depends on the correctness of the various estimates that he is required to make. It is also use in sociological study and economic planning to find the projections of population, birth rates. death rates etc.

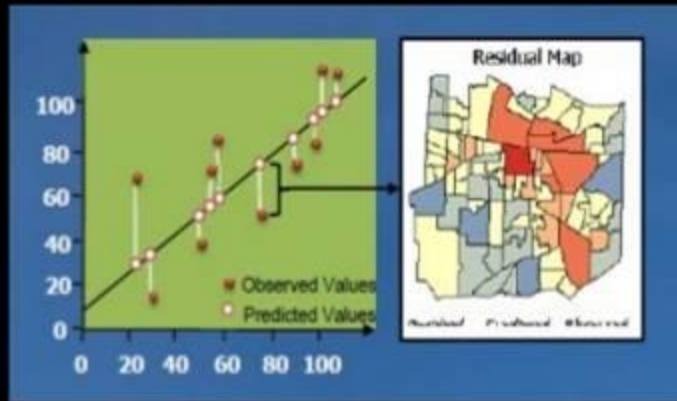
Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

METHODS OF STUDYING REGRESSION:



Prepared By ARYA KUMAR

Simple Linear Regression Model. . .



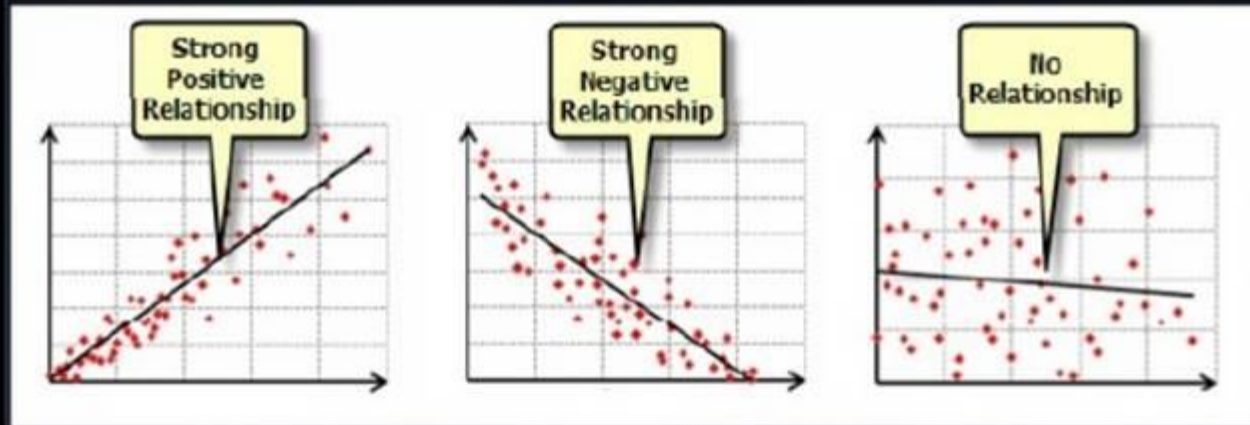
- Only **one** independent variable, x
- Relationship between x and y is described by a linear function
- Changes in y are assumed to be caused by changes in x

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Prepared By **ARYA KUMAR**

Types of Regression Models. . .



Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Estimated Regression Model. . .

The sample regression line provides an **estimate** of the population regression line

Estimated (or predicted) y value

Estimate of the regression intercept

Estimate of the regression slope

Independent variable

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

The individual random error terms e_i have a mean of zero

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Algebraically method-:

1. Least Square Method-:

Prepared By
ARYA KUMAR

The regression equation of X on Y is :

$$X = a + bY$$

Where,

X = Dependent variable

Y = Independent variable

The regression equation of Y on X is:

$$Y = a + bX$$

Where,

Y = Dependent variable

X = Independent variable

And the values of a and b in the above equations are found by the method of least of Squares-reference . The values of a and b are found with the help of normal equations given below:

$$\begin{array}{ll} \text{(I)} & \sum X = na + b \sum Y \\ & \sum XY = a \sum Y + b \sum Y^2 \end{array} \qquad \begin{array}{ll} \text{(II)} & \sum Y = na + b \sum X \\ & \sum XY = a \sum X + b \sum X^2 \end{array}$$

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Example 1 –: From the following data obtain the two regression equations using the method of Least Squares.

X	3	2	7	4	8
Y	6	1	8	5	9

Solution–:

X	Y	XY	X ²	Y ²
3	6	18	9	36
2	1	2	4	1
7	8	56	49	64
4	5	20	16	25
8	9	72	64	81
$\sum X = 24$	$\sum Y = 29$	$\sum XY = 168$	$\sum X^2 = 142$	$\sum Y^2 = 207$

Step-1: Formula: for solving Y on X

$$\sum Y = na + b \sum X$$

$$\sum XY = a \sum X + b \sum X^2$$

**Prepared By
ARYA KUMAR**

Substitution the values from the table we get

$$29 = 5a + 24b \dots\dots\dots(i)$$

168 = 24a + 142b As we get 2 is a common factor, so divide it with 2

$$84 = 12a + 71b \dots\dots\dots(ii)$$

Step-2: subtracting the equation by converting a/ b of similar number

Multiplying equation (i) by 12 and (ii) by 5

$$348 = 60a + 288b \dots\dots\dots(iii)$$

$$420 = 60a + 355b \dots\dots\dots(iv)$$

By solving equation(iii)and (iv) we get

$$a=0.66 \text{ and } b=1.07$$

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Step-3: By putting the value of a and b in the Regression equation Y on X we get

$$Y=0.66+1.07X$$

Step-4: Formula: for solving X on Y

Now to find the regression equation of X on Y ,
The two normal equation are

$$\begin{aligned}\sum X &= na + b \sum Y \\ \sum XY &= a \sum Y + b \sum Y^2\end{aligned}$$

Step-5: subtracting the equation by converting a/ b of similar number

Substituting the values in the equations we get

$$24=5a+29b\text{.....(i)}$$

$$168=29a+207b\text{.....(ii)}$$

Multiplying equation (i) by 29 and in (ii) by 5 we get

$$a=0.49 \text{ and } b=0.74$$

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

**Prepared By
ARYA KUMAR**

Step-6: Substituting the values of a and b in the Regression equation X and Y

$$X=0.49+0.74Y$$

Prepared By
ARYA KUMAR

2. Deviation from the Arithmetic mean method:

The calculation by the least squares method are quite cumbersome when the values of X and Y are large. So the work can be simplified by using this method.

The formula for the calculation of Regression Equations by this method:

Regression Equation of X on Y- $(X - \bar{X}) = b_{xy} (Y - \bar{Y})$

Regression Equation of Y on X-

$$(Y - \bar{Y}) = b_{yx} (X - \bar{X})$$

Where, b_{xy} and b_{yx} = Regression Coefficient

$$b_{xy} = \frac{\sum xy}{\sum y^2} \quad \text{and} \quad b_{yx} = \frac{\sum xy}{\sum x^2}$$

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Example2-: from the previous data obtain the regression equations by Taking deviations from the actual means of X and Y series.

X	3	2	7	4	8
Y	6	1	8	5	9

Solution-:

X	Y	$x = X - \bar{X}$	$y = Y - \bar{Y}$	x^2	y^2	xy
3	6	-1.8	0.2	3.24	0.04	-0.36
2	1	-2.8	-4.8	7.84	23.04	13.44
7	8	2.2	2.2	4.84	4.84	4.84
4	5	-0.8	-0.8	0.64	0.64	0.64
8	9	3.2	3.2	10.24	10.24	10.24
$\sum X = 24$	$\sum Y = 29$	$\sum x = 0$	$\sum y = 0$	$\sum x^2 = 26.8$	$\sum y^2 = 38.8$	$\sum xy = 28.8$

Prepared
By
ARYA
KUMAR

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Regression Equation of X on Y is

Formula $(X - \bar{X}) = b_{xy} (Y - \bar{Y})$

$$b_{xy} = \frac{\sum xy}{\sum y^2}$$

Prepared By
ARYA KUMAR

$$X - 4.8 = \frac{28.8}{38.8} Y - 5.8$$

$$X - 4.8 = 0.74 Y - 5.8$$

$$X = 0.74 Y + 0.49 \quad \dots\dots\dots(I)$$

Regression Equation of Y on X is

Formula $(Y - \bar{Y}) = b_{yx} (X - \bar{X})$

$$b_{yx} = \frac{\sum xy}{\sum x^2}$$

$$Y - 5.8 = \frac{28.8}{26.8} X - 4.8$$

$$Y - 5.8 = 1.07 (X - 4.8)$$

$$Y = 1.07 X + 0.66 \quad \dots\dots\dots(II)$$

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

It would be observed that these regression equations are same as those obtained by the direct method .

Prepared By
ARYA KUMAR

3.Deviation from Assumed mean method-:

When actual mean of X and Y variables are in fractions ,the calculations can be simplified by taking the deviations from the assumed mean.

The Regression Equation of X on Y-:

$$(X - \bar{X}) = b_{xy}(Y - \bar{Y})$$

The Regression Equation of Y on X-:

$$(Y - \bar{Y}) = b_{yx}(X - \bar{X})$$

But , here the values of b_{xy} and b_{yx} will be calculated by following formula:

$$b_{xy} = \frac{N \sum d_x d_y - \sum d_x \sum d_y}{N \sum d_x^2 - \sum d_x^2} \quad b_{yx} = \frac{N \sum d_x d_y - \sum d_x \sum d_y}{N \sum d_y^2 - \sum d_y^2}$$

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Example-: From the data given in previous example calculate regression equations by assuming 7 as the mean of X series and 6 as the mean of Y series.

Prepared By
ARYA KUMAR

Solution-:

X	Y	Dev. From assu. Mean 7 (d_x)= $X-7$	d_x^2	Dev. From assu. Mean 6 (d_y)= $Y-6$	d_y^2	$d_x d_y$
3	6	-4	16	0	0	0
2	1	-5	25	-5	25	+25
7	8	0	0	2	4	0
4	5	-3	9	-1	1	+3
8	9	1	1	3	9	+3
$\sum X = 24$	$\sum Y = 29$	$\sum d_x = -11$	$\sum d_x^2 = 51$	$\sum d_y = -1$	$\sum d_y^2 = 39$	$\sum d_x d_y = 31$

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

$$\bar{X} = \frac{\sum X}{N} \Rightarrow \bar{X} = \frac{24}{5} = 4.8$$

$$\bar{Y} = \frac{\sum Y}{N} \Rightarrow \bar{Y} = \frac{29}{5} = 5.8$$

The Regression Coefficient of X on Y-:

$$b_{xy} = \frac{N \sum d_x d_y - \sum d_x \sum d_y}{N \sum d_y^2 - \sum d_y^2}$$

$$b_{xy} = \frac{5(31) - (-11)(-1)}{5(39) - (-1)^2}$$

$$b_{xy} = \frac{155 - 11}{195 - 1}$$

$$b_{xy} = \frac{144}{194}$$

$$b_{xy} = 0.74$$

The Regression equation of X on Y- $(X - \bar{X}) = b_{xy}(Y - \bar{Y})$

$$(X - 4.8) = 0.74(Y - 5.8)$$

$$X = 0.74Y + 0.49$$

Prepared By
ARYA KUMAR

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

The Regression Coefficient of Y on X:-

$$b_{yx} = \frac{N \sum d_x d_y = \sum d_x \sum d_y}{N \sum d_x^2 - \sum d_x^2}$$

$$b_{yx} = \frac{5(31) - (-11)(-1)}{5(51) - (-11)^2}$$

$$b_{yx} = \frac{155 - 11}{255 - 121}$$

$$b_{yx} = \frac{144}{134}$$

$$b_{yx} = 1.07$$

The Regression Equation of Y on X:-

$$(Y - \bar{Y}) = b_{yx} (X - \bar{X})$$

$$(Y - 5.8) = 1.07 (X - 4.8)$$

$$Y = 1.07 X + 0.66$$

It would be observed that these regression equations are same as those obtained by the **least squares method** and deviation from **arithmetic mean**

Prepared By
ARYA KUMAR

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Meaning Of Report

- It is a statement that is prepared to present information in a structured manner about an investigation that has been undertaken.
- Report is a summary of findings and recommendations about a particular matter / problem.
- Report acts as guidance and higher authorities including company executives and directors.
- The main purpose of reports is to facilitate timely decisions and follow up measures.

Prepared By

ARYA KUMAR

Report writing: Formal

There are many different types of reports. This information is a basic outline only. Before you attempt to write a report, you should check the particular requirements for the subject.

A formal report should have the following arrangement.

1. TITLE PAGE — The Title Page must include the subject of the report, who the report is for, who the report is by and the date of submission.

2. ABSTRACT — An Abstract is usually 100 to 200 words and should include the following:

- why the report has been written (i.e. what question or problem is it addressing?)
- how the study was undertaken
- what the main findings were
- what the significance of the findings is.

Be specific and precise so that the reader can get a good understanding of the main points without having to read the whole report.

The abstract should be on a separate page with the centred heading ABSTRACT in capitals. It is usually written in a single paragraph with no indentation.

3. TABLE OF CONTENTS — The Table of Contents should be on a separate page. It helps the reader to find specific information and indicates how the information has been organised and what topics are covered. The table of contents should also include a **list of figures** and a **list of tables** if any are used in the report.

4. INTRODUCTION — The Introduction has three main components.

1. **The Background** which describes events leading up to the existing situation, what projects have been done previously, and why the project or study is necessary.
2. **The Purpose** which defines what the project or study is to achieve, who authorised it and the specific terms of reference.
3. **The Scope** which outlines any limitations imposed on the project such as cost, time etc.

5. BODY — The Body varies according to the type of report. Basically, it answers the questions — Who? Why? Where? When? What? How? In an investigative report, it would consist of all the information required to convince the reader that the conclusions and recommendations are valid/reliable. This information must be presented in a systematic way.

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Different Types of Report

Prepared By
ARYA KUMAR

According to function

1. Analytical Report
2. Informational Report
3. Research reports



According to formality

1. Statutory report
2. Non-statutory or voluntary report
3. Routine Report
4. Special Report
5. Technical Report
6. Popular Report

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Significance of Report

Prepared By
ARYA KUMAR

1. Provide sufficient information on various aspects of the business
2. Reports tells about the knowledge or skills required by the professionals to deal with it.
3. Decision making gets possible through report output
4. Solving a problem is possible through a well-structured report.
5. Reports communicate the planning, policies, and other matters regarding an organization to the masses.
6. News reports play the role of ombudsman and levy checks and balances on the establishment.

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

1. Letter or Memo of Transmittal (not actually part of the report but included with the report)

2. Title Page

**Prepared By
ARYA KUMAR**

- Title of the report
- Name of the person, department, or company commissioning the report
- Date submitted
- Authors and their corporate or departmental affiliation
- The title of the report should be long enough to describe the report's contents (two-line titles are acceptable) and
- Incorporate the key words from the report to allow indexing and retrieval.

3. Executive Summary

The executive summary is always one paragraph long and contains the following:

- problem identification addressed in the report
- significance of the problem to the company in brief
- brief description of the solution
- a very brief description of the time and funds required to implement the solution

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

- 4. Table of Contents**
- 5. Table of Figures**
- 6. Glossary of Terms**
- 7. Body**
- 8. Conclusion**
- 9. Recommendations**
- 10. Appendices**
- 11. Bibliography**
- 12. Page Numbering**

Prepared By
ARYA KUMAR

Faculty of MBA, CIME
aryantripathy@yahoo.com, 09853422575

Components of a project report

1. Cover page:

- Title of the project
- Degree for which submitted
- Name of Author
- Name of Supervisor
- Year of submission
- Name of the University Institution

Prepared By
ARYA KUMAR

2. Acknowledgement

3. Table of contents

4. Body of the Report

- Introduction
- Review of literature
- Objectives,
- Methodology - area of the study, sampling technique if used, type of the tools for data collection, method of analysis etc.
- Limitations of study and chapter planning.
- Conceptual framework in national / international scenario relating to the topic of the project.
- Processing of data analysis and findings.
- Conclusions and recommendations

5. Bibliography or References

6. Annexures: Questionnaires / Schedules if any, relevant / reports, etc.

Find the article: 6 ZIJMR VOL8 ISSUE6 JUNE 2018.pdf

Faculty of MBA, CIME

aryantripathy@yahoo.com, 09853422575

Oral Presentation

For successful presentation, the written report must be completed so as to be able to complete full presentation within the stipulated time span.

Basic strategies for effective oral presentation:

1. Preparation

- a) Understand the context of presentation.
- b) Analyse audience.
- c) Understand and articulate presentation's purpose.
- d) Choose and shape presentation's content.
- e) Organize presentation.

2. Delivery

- a) Choose an appropriate speaking style.
- b) Practise an effective delivery style
- c) Select and use visual aids effectively

Prepared By
ARYA KUMAR



Faculty of MBA, CIME
aryantripathy@yahoo.com,
09853422575