# Unit IV

## Multivariate Analysis of Variance

### Introduction

Multivariate Analysis of variance is simply an ANOVA with several dependent variables. That is to say, ANOVA test for the difference in means between two or more groups, while MANOVA tests for difference in two or more vectors of means.

"MANOVA" stands for "Multivariate Analysis of variance". MANOVA in statistics contain multiple dependent variables. They help in determining the differences between either two or more than two independent variables. It assists in determining this difference simultaneously.

The MANOVA method determines if the dependent variables get significantly affected by

changes in the independent variables. It assists in determining. It also determines the interactions taking place amoungst dependent variables MANOVA finally determines the interaction taking place amoungst independent variables to.

Null hypothesis for MANOVA:

We would test $H_0$ to see if the vectors of means of the dependent variable is equal for multiple independent groups and our null hypothesis is

$$H_0 : \begin{bmatrix} \bar{x}_{11} \\ \bar{x}_{21} \\ \vdots \\ \bar{x}_{p1} \end{bmatrix} = \begin{bmatrix} \bar{x}_{12} \\ \bar{x}_{22} \\ \vdots \\ \bar{x}_{p2} \end{bmatrix} = \begin{bmatrix} \bar{x}_{13} \\ \bar{x}_{23} \\ \vdots \\ \bar{x}_{p3} \end{bmatrix} = \cdots = \begin{bmatrix} \bar{x}_{1k} \\ \bar{x}_{2k} \\ \vdots \\ \bar{x}_{pk} \end{bmatrix}$$

where $p$ represents the total no of dependent variables for $k$ levels.

Test statistics: $\Lambda = \dfrac{|W|}{|T|} = \dfrac{|W|}{|B+W|}$

where $W$ and $T$ are determinants of the within and Total sum of squares and cross product matrices.

This means that if between effect B, is very large then A approaches A, if B is small (or) even 0 then A approaches 1

Assumptions of MANOVA:

Normal data

The dependent variable should be normally distributed within groups. Overall, the F-test is robust to non-normality. If the non-normality is caused by skewness rather than by outliers. Tests for outliers should be run before performing a MANOVA and outliers should be transformed or removed.

Linearity

MANOVA assumes that there are linear relationship among all pairs of dependent variables, all pairs of covariates and all dependent variable covariate pairs in each all. Therefore when the relationship deviates from linearity the power of analysis will be compromised.

Homogenity of variances:

Homogenity of variances assumes that the dependent variable exhibit equal level of variance accross the range of predictor variable remember that the error variable is computed (ss error) by adding up the sum of squares within each group. If the variances in the two

groups are different from each other, then adding them together is not appropriate and will not yield an estimate of the common group variance. Homoscedastecity can be examined graphically or by means of a no. of statistical tools

## Homogeneity of variances and covariances

In multivariate designs, with multiple dependent measures, the homogeneity of variances assumption described earlier also applies. However since there are multiple dependent variables, it is also required that their inter correlation (covariances) are homogeneous across the cells of the design

## Special case of MANOVA:

Two special cases arise in in MANOVA the inclusion of within-subjects important variables and unequal sample size in cell

## Unequal sample size

As in ANOVA, when cells in a factorial MANOVA have different sample sizes, the sum of squares for effect plus error does not equal to total sum of squares. This causes tests of main effects and

interaction to be correlated. SPP effects the adjustment for unequal sample size MANOVA.

## Within subjects Design

Problem arises if the researcher measures several different dependent variables on different occasions. This situation can be viewed as a within subject independent variable with as many levels as occasions of it can be viewed as separate dependent variables for each occasion

## Additional limitations:

### Outliers

like ANOVA, MANOVA is extremely sensitive to outliers. Outliers may produce either a type I or type II error and give no indication as to which type of data error is occuring in the analysis. There are several programs available to test for univariate and multivariate outliers.

### Multicollinearity and singularity

When there is high correlation between dependent variables, one dependent variables becomes a near- linear combination of the other dependent variables. Under such circumstances it would become statistically reductant and sus- spect to includes both combinations.

## Multivariate Hotelling $T^2$ speciall issues concerning hotelling $T^2$ MANOVA

In statistical $T^2$ MANOVA is multivariate distn proportional to the F-distn and arises importantly as a distn of a set of statistics which are natural generalisation of the statistics.

Underlying student t-distn, Hotelling $T^2$ statistics is a generalisation of student t statistic that is used in multivariate hypothesis testing.

If a r.v 'x' has hotelling $T^2$ distn, $x \sim T^2_{p,m}$ then $\dfrac{m-p+1}{mp}$, $x \sim F_{p,m-p+1}$ where, $F_{p,m-p+1}$ is the F distn with parameter p and m-p+1.

Test versions: One sample: The multivariate vector means for a group equals a hypothetical vector of means.

Two sample: The multivariate vector means for a two group are equal.

Two sample $T^2$ Hotellings:-

If you know how to run a two sample t-test, then you know how to run a two sample.

Hotelling $T^2$:

The basic steps are the same although you will use a different formula to calculate the $T^2$ value and

you will use a different table (F table) to find the critical value

Hotelling $T^2$ has several advantages over the t-test:

The type I error rate is well controlled

The relationship b/w multivariate variables is taken into account.

It can generate an overall conclusion even if multiple (single) t-test are inconsistent while a t test will tell you which variable differ b/w groups. Hotelling $T^2$ summarizes the b/w group difference.

Test hypotheses: $H_0$: The two sample are from popln with the same multivariate means

$H_1$: The two samples are from popln with different multivariate means

Like the t-test you want to find a value for $T^2$ (in this case for $T^2$) and compare it with F table value

If the calculated values is greater than the table values you can reject the null hypotheses ($H_0$)

For ease of this calculation Hotelling $T^2$ is first transformed into an F-statistics.

$$F = \frac{n_1 + n_2 - p - 1}{p(n_1 + n_2 - 2)} T^2 \sim F_{p, n_1 + n_2 - p - 1}$$

where $n_1$ and $n_2$ sample sizes

$p$ = number of variable measured

$n_1 + n_2 - p - 1$ = degrees of freedom

Reject the null hypothesis (at a choosen significant level). If the calculated value is greater than the F table critical value

Rejecting the null hypothesis means that atleast one of the parameters is significantly different.

Note:

Testing the multiple dependent variables accomplish by creating new dependent variables that maximize group differences these artificial dependent variable are linear combination of the measured dependent variable.

MANOVA comparing two groups:-

A $T^2$ statistic for testing the equality of vector means from two multivariate populations can be developed by analogy with the univariate procedure. This $T^2$ statistic is appropriate for comparing responses from one set of

experimental select setting (population 1) with independent responses Vps from another set of experimental setting (population 2) The comparison can be without explicitly controlling for unit to unit variability as in the paired comparison case.

If possible, the experimental units should be randomly assigned to the sets of experimental conditions. Randomisation to some extent mitigate the effect at unit-to-unit variability in a sub-sequent comparison of treatments. Although some precision is lost relative to paired comparisons the inferences in the two population case are ordinarily, applicable to a more general collection of experimental units simply because unit homogenaity is not required.

Consider a random sample of size $n_1$ from population 1 and a sample of size $n_2$ from population 2

Principal component analysis.

A PCA is concerned with explaining the variance-covariance structure of a set of variables Its general objectives are i)Data reduction ii)Interpretation

Although p components are required to reproduce the total system variability often much of this variability

can be accounted for by a small number $k$ of the principal component. If so there is as information in the $k$ components as there is in the original p variables The $k$ principal component can then replace the initial P variables and original data set. Consisting of $n$ measurements on $p$ variables is reduced to a data set consisting of $n$ measurements on $k$ principal components.

An analysis of principal component often reveals relationship that were not previously suspected and there by allows interpretation that would not ordinarily result Good example of this is provided by stock market

Analysis of pc are more of a means to an end rather then themselves because they frequently serve as intermediate steps in much larger investigations pc may be input to a multiple regression (or) cluster analysis moreover pc are factoring of the covariance matrix for the factor analysis

Uses:-

PCA is a technique that is useful for the comparison and classification of data. The purpose is to reduce the dimensionality of a data set (sample) by finding a new set of variable smaller than the original

set of variables, that nonetheless returns most of the sample information.

## Computation of Principal Components in the population:-

Algebrically, principal components are particular linear combinations of the $p$ random variables $x_1, x_2 \ldots x_p$.

Geometrically, these linear combinations represent the selection of a new co-ordinate system obtained by relating the original system with $x_1, x_2, \ldots x_p$ as the co-ordinates axes. The new axes represent the directions with maximum variability and a simpler and more parsimonious description of the covariance structure.

Principal components depends solely on the covariance matrix $\Sigma$ ( or the correlation matrix $p$) $x_1, x_2 \ldots x_p$. Their development does not require a multivariate normal assumption.

Let the random vector $x' = [x_1, x_2 \ldots x_p]$ have the covariance matrix $\Sigma$ with eigen values $\lambda_1 > \lambda_2 > \lambda_3 \ldots > \lambda_p \geq 0$

Consider the linear combination

$$
\left.
\begin{aligned}
Y_1 &= a_1' x = a_{11} x_1 + a_{12} x_2 + \ldots + a_{1p} x_p \\
Y_2 &= a_2' x = a_{21} x_1 + a_{22} x_2 + \ldots + a_{2p} x_p \\
&\qquad\qquad\qquad \vdots \\
Y_p &= a_p' x = a_{p1} x_1 + a_{p2} x_2 + \ldots + a_{pp} x_p
\end{aligned}
\right\} \quad - \textcircled{1}
$$

Then we obtain $V(Y_i) = a_i' \Sigma a_i$, $i = 1, 2 \cdots p$ —— ②

$$Cov(Y_i, Y_k) = a_i' \Sigma a_{k_i}, \quad i = 1, 2 \cdots p \text{ —— ③}$$

These principal components are these uncorrelated linear combinations $Y_1, Y_2 \cdots Y_p$ whose variance in ② are as large as possible.

Thus 1st principal component is the linear combination with max variance (ie) it maximizes $Var(Y_i) = a_i' \Sigma a_i$

It is clear that $var(Y_i) = a_i' \Sigma a_i$ can be increased by multiplying any $a_i$ by some constant We therefore define

First principal component = linear combination $a_1' x$ that maxim. $V(a_1' x)$ subject to $a_1' a_1 = 1$ & $cov(a_1' x, a_1' x) = 1$

Second principal component = linear combination $a_2' x$ that maximizes $V(a_2' x)$ subject to $a_2' a_2 = 1$ & $Covar(a_1' x, a_2' x) = 0$

at the $i$th step $i$th principal component = linear combination $a_i' x$ that maximizes $V(a_i' x)$ subject to $a_i' a_i = 1$ and

$$cov(a_k' x_k, a_i' x_i) = 0 \text{ for } k < i$$

Extraction of principal components.

There is always the question of how many components to retain There is no definite answer to this ques

Things to consider include:
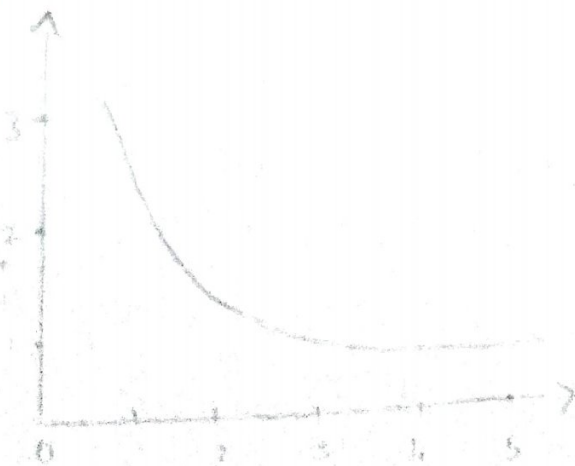
1) The amount of total sample variance explained

2) The relative sizes of the eigen values ( the variance of the sample components) and

3) The subject matter interpretations of the components

In addition, a component associated with an eigen value near zero, deemed to be unimportant, may indicate an unsuspected linear dependency in the data.

A useful visual aid to determine an appropriate no of principal components is a "scree plot" with the eigen values ordered from larger to smaller, a scree plot is a plot of $\lambda_i$, the magnitude of eigen values versus i, its number.

To determine the appropriate number of components, we look for an the scree plot. The no. of components is taken to be the point at which the remaining eigen values are relatively small and all about the same size. The following figure shows a scree plot for a situation with

An elbow occurs in the plot in the figure of at about i=3. That is the eigen values after $\hat{\lambda}_2$ are all relatively small and about the same size. In this case, it appears without any other evidence that two sample PC offers effectively summarise the total sample variance.

## Interpretation of principal components

The sample PC have several interpretations. First suppose the underlying distr of x is nearly : $N_p(\mu, \Sigma)$. Then the sample principal components $\hat{y}_i = \hat{e}_i'(x - \bar{x})$ are realizations of popln PC $y_i = e_i'(x - \mu)$, which have an $N_p(0, \Lambda)$ distr. The diagonal matrix $\Lambda$ has entries $\lambda_1, \lambda_2 \ldots \lambda_p$ and $(\lambda_i, e_i)$ are the eigen value - eigen vector pairs of $\Sigma$.

Also from the sample values $x_i$, we can approximate $\mu$ by $\bar{x}$ and $\Sigma$ by s. If s is +ve definite, the contour consisting of all px vectors x satisfying

$$(x - \bar{x})'s^{-1}(x - \bar{x}) = c^2 \quad \text{①}$$

estimates the constant density contours $(x - \mu)'\Sigma^{-1}(x - \mu) = c^2$ of the underlying normal density. The approximate contours can be drawn on the scatter plot to indicate the normal distr that generate the data. The normality assumption is useful for the inference procedures discussed but it is not required for the development of the properties of

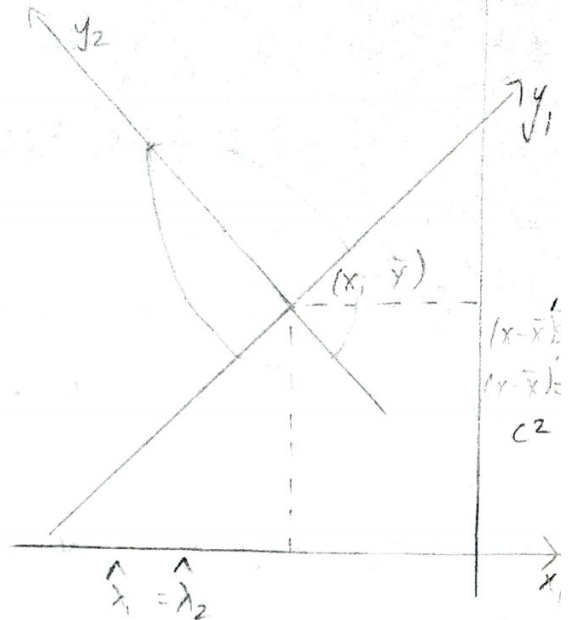estimates the constant density contour
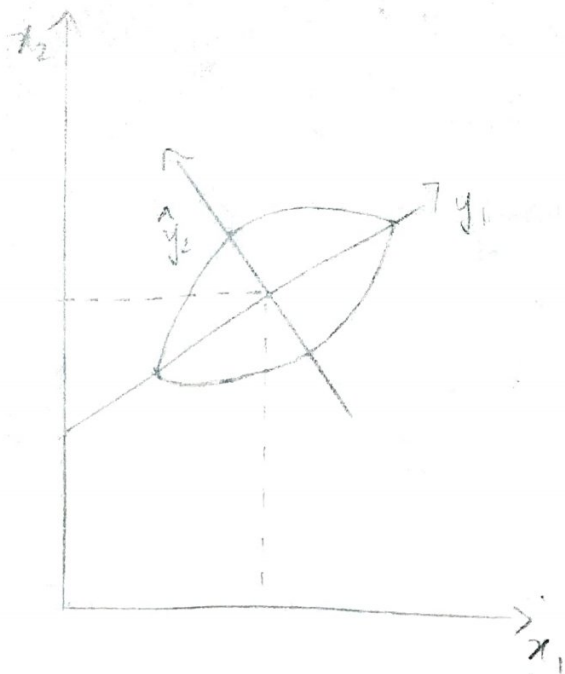
the sample pc summarized.

Even when the normal assumption is suspect and the scatter plot may depart somewhat from an elliptical pattern we can still extract the eigen values from s and obtain the sample pc Geometrically, the data may be plotted as n points in p space. The data can then be expressed in the new co-ordinates; which coincide with the axes of the contour ①. Now ① defines a hyperbolic elliptic said that is centered at $\bar{x}$ and whose axes are given by the eigen vectors of $s^{-1}$ or, equivalently of s. The length of these hyperbolic elliptic said axes are proportional to $\sqrt{\hat{\lambda}_i}$ $i = 1, 2 \ldots p$ where $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \ldots \geq \hat{\lambda}_p \geq 0$ are the given eigen values of s.

Because $\hat{e}_i$ has length 1, the absolute value of the $i$th principal component, $|\hat{y}_i| = |\hat{e}_i (x - \bar{x})|$, gives the length of the projection of the vector $(x - \bar{x})$ on the unit vector $\hat{e}_i$. Thus the sample pcs, $\hat{y}_i = \hat{e}_i (x - \bar{x})$ $i = 1, 2 \ldots p$ lie along the axes of the hyperbolic elliptic said and their absolute values are the length of the projection of $x - \bar{x}$ in the directions of the axes ①. Consequently, the sample pcs can be views as the results of translating the origin of the original coordinate

system to $\bar{x}$ and then rotating the coordinate axes untill they pass through the scatter in the directions of maximum variance.

The geometrical interpretation of the sample pcs for $p = e$ shows an ellipse of constant distance, centered at $\bar{x}$ with $\hat{\lambda}_1 > \hat{\lambda}_2$. The sample pcs are well determined. They lie along the axes of the ellipse in two perpendicular directions, including the direction of the original coordinate axes. similarly, the sample pcs can lie in any two perpendicular directions Including those of the original coordinate axes when the contours of constant distances are not unequally determined and can lie in any two perpendicular directions, including the direction of the original co-ordinate axes Similarly the sample pcs can lie in the any two perpendicular directions including those of the original co-ordinate axes. when the contours of constant distance are nearly circular or equivalently, when the eigen values of s are nearly equal. The sample variation is homogeneous in all directions. It is then not possible to represent data well in fewer than p dimensions.

Sample pcs and ellipse of constant distance.

If the last few eigen values $\hat{\lambda}_i$ are sufficiently small such that the variation in the corresponding $\hat{e}_i$ directions is negligible, the last few sample pcs can often be ignored and the data can be adequately approximated by their representations in the space of the retained components.

Finally. supplement

## Computation of the MLE of the principal components

There are several ways of computing the characteristic roots and characteristic vector [pc] of a matrix $\Sigma$ or $\hat{\Sigma}$. We shall indicate some of them.

One method involves expanding the dimens determinantal eqn, $0 = |\Sigma - \lambda I| \longrightarrow ①$ and solving the resulting pth degree eqn in $\lambda$.

Example, [By Newton's method or the second method] for the roots $\lambda_1 , \lambda_2 > \lambda_3 \cdots \lambda_p$. Then $\Sigma - \lambda I$ is of rank p-1 and a solution of $(\Sigma - \lambda_i I) \beta = 0$ can be obtained by taking

$a_j^{(i)}$ as the cofactor of the element in the 1st (or any other fixed) column and $j$ row of $\Sigma - \lambda_i I$.

The 2nd method iterates using the eqn for a characteristic vector

$$\Sigma \cdot x = \lambda x \qquad — ②$$

We have written the eqn for the popln. Let $x_{(0)}$ be any vector not orthogonal to the 1st characteristic vector, and define

$$x_{(i)} = \Sigma \, y_{i-1}, \quad y_{(i)} = \frac{x_{(i)}}{\sqrt{x_{(i)}' \, x_{(i)}}}, \quad i = 0, 1 \ldots \qquad — ③$$

It can be shown that

$$\lim_{i \to \infty} y_{(i)} = \pm \beta^{(i)} \quad \lim_{i \to \infty} x_{(i)}' = \lambda_i^2 \qquad — ④$$

The rate of convergence depends on the ratio $\lambda_2 / \lambda_1$. the closer this ratio is to 1, the slower the convergence.

To find the 2nd root and vector define

$$\Sigma_2 = \Sigma - \lambda_1 \beta^{(1)} \beta^{(1)'} \qquad — ⑤$$

Then $\Sigma_2 \beta^{(i)} = \Sigma \beta^{(i)} - \lambda_1 \beta^{(i)'} \beta^{(i)} = \Sigma \beta^{(i)} - \lambda_i \beta^{(i)} \qquad — ⑥$

if $i \neq 1$ and $\Sigma_2 \beta^{(1)} = 0. \qquad — ⑦$

Thus, $\lambda_2$ is the largest root of $\Sigma_2$ and $\beta^{(2)}$ is the corresponding vector. The iteration process is now applied to $\Sigma_2$ to find $\lambda_2$ and $\beta^{(2)}$. Defining $\Sigma_3 = \Sigma_2 - \lambda_2 \beta^{(2)} \beta^{(2)'}$ we can find $\lambda_3$ and $\beta^{(3)}$ and so forth.

There are several ways in which the labour of the iteration procedure may be reduced one is to raise $\Sigma$

a power before proceeding with the iteration. Thus one can use $\varepsilon^2$ defining,

$$x_{(i)} = \varepsilon^2 y_{(i-1)}, \quad y_{(i)} = \frac{x_{(i)}}{\sqrt{x'_{(i)} x_{(i)}}}, \quad i = 0, 1, \ldots \quad —(8)$$

This procedure will give twice as rapid convergence as the use of (3). Using $\varepsilon^4 = \varepsilon^2 \varepsilon^2$ will lead to convergence four times as rapid, and so on. It should be noted that since, $\varepsilon^2$ is symmetric, there are only $p(p+1)/2$ elements to be found.

Efficient computation, however, uses other methods. One method is the QR or QL algorithm. Let $\varepsilon_0 = \varepsilon$, define recuratively the orthogonal $Q_i$ and lower triangular $L_i$ by $\varepsilon_i = Q_i L_i$ and $\varepsilon_{i+1} = L_i Q_i \left(= Q_i' \varepsilon_i Q_i\right)$ $i = 1, 2 \ldots$

(The vam-schmidt orthogonalization is a way of finding $Q_i$ and $L_i$. the QR method replaces a lower triangular matrix.