

1) Consider the following var-cov matrix and find the PCs. Also explain the proportion of variability explained by each PC and comment on the no of PC's to be retained if someone wishes to explain atleast 80% of the total system variability.

$$\Sigma = \begin{pmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix}$$

ans) By looking at  $\Sigma$  we can say that,

$$|\Sigma| = 2(5-4) = 2 \neq 0 \text{ i.e } \Sigma \text{ is non singular.}$$

Let  $\lambda_1, \lambda_2, \lambda_3$  be the eigenvalues of  $\Sigma$  i.e we need to solve  $|\Sigma - \lambda I| = 0$

$$\therefore |\Sigma - \lambda I| = 0$$

$$\Rightarrow \begin{vmatrix} 1-\lambda & -2 & 0 \\ -2 & 5-\lambda & 0 \\ 0 & 0 & 2-\lambda \end{vmatrix} = 0$$

$$\Rightarrow (2-\lambda)[(1-\lambda)(5-\lambda)-4] = 0$$

$$\Rightarrow (\lambda-2)[\lambda^2 - 6\lambda + 1] = 0$$

$$\Rightarrow \lambda = 2 \text{ or } \lambda = \frac{6 \pm \sqrt{36-4}}{2} = 3 \pm 2\sqrt{2}$$

$$\text{i.e } \lambda_1 = 3+2\sqrt{2}, \lambda_2 = 2, \lambda_3 = 3-2\sqrt{2}$$

$$\Rightarrow \lambda_1 = 5.828, \lambda_2 = 2, \lambda_3 = 0.172$$

$$\text{and } \lambda_1 > \lambda_2 > \lambda_3$$

Now we have to find orthonormal eigenvectors corresponding to  $\lambda_1, \lambda_2, \lambda_3$ .

Let  $\underline{g}_i$  be the eigenvector corresponding to  $\lambda_i$ ;  $i=1,2,3$

$$\text{i.e } \underline{g}_i = \begin{pmatrix} g_{i1} \\ g_{i2} \\ g_{i3} \end{pmatrix}; i=1,2,3$$

Now we have,

$$\sum \underline{a}_1 = \lambda_1 \underline{a}_1$$

$$\Rightarrow \begin{pmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \\ a_{13} \end{pmatrix} = \begin{pmatrix} 5.828 & a_{11} \\ 5.828 & a_{12} \\ 5.828 & a_{13} \end{pmatrix}$$

$$\Rightarrow -4.828 a_{11} - 2a_{12} = 0 \quad \text{and} \quad -2a_{11} - 0.828 a_{12} = 0 \quad \text{and} \quad a_{13} = 0$$

$$\Rightarrow a_{11} + 0.414 a_{12} = 0$$

Now if  $a_{11} = 1 \Rightarrow a_{12} = -2.415$

$$\therefore \underline{a}_1 = \begin{pmatrix} 1 \\ -2.415 \\ 0 \end{pmatrix} \quad \text{and} \quad \|\underline{a}_1\| = \sqrt{1^2 + (-2.415)^2} = 2.613$$

Normalized

$$\therefore \underline{e}_1 = \begin{pmatrix} 0.38 \\ -0.92 \\ 0 \end{pmatrix} \rightarrow \text{orthonormal eigenvector}$$

Hence the 1st PC is  $\gamma_1 = 0.38 x_1 - 0.92 x_2$

Similarly for  $\lambda_2 = 2$  we get,

$$\sum \underline{a}_2 = \lambda_2 \underline{a}_2$$

$$\Rightarrow \begin{pmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} a_{21} \\ a_{22} \\ a_{23} \end{pmatrix} = \begin{pmatrix} 2a_{21} \\ 2a_{22} \\ 2a_{23} \end{pmatrix}$$

$$\Rightarrow a_{21} + 2a_{22} = 0 \quad \text{and} \quad -2a_{21} + 3a_{22} = 0 \quad \text{and} \quad 2a_{23} = 2a_{23}$$

$$\Rightarrow 2a_{21} = 3a_{22} \quad \Rightarrow a_{23} = a_{23}$$

If  $a_{21} = 0$  then  $a_{22} = 0$  and  $a_{23} = 1$

$$\therefore \underline{a}_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \quad \text{and} \quad \|\underline{a}_2\| = \sqrt{0^2 + 0^2 + 1^2} = 1$$

$$\therefore \underline{e}_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

Hence the 2nd PC is  $\gamma_2 = x_3$

Now, for  $\lambda_3 = 0.172$  we have

$$\sum \underline{a}_3 = \lambda_3 \underline{a}_3$$

$$\Rightarrow \begin{pmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} a_{31} \\ a_{32} \\ a_{33} \end{pmatrix} = \begin{pmatrix} 0.172 a_{31} \\ 0.172 a_{32} \\ 0.172 a_{33} \end{pmatrix}$$

$$\Rightarrow 0.828 \alpha_{31} - 2\alpha_{32} = 0 \text{ and } -2\alpha_{31} + 4.823 \alpha_{32} = 0$$

and  $2\alpha_{33} = 0.172 \alpha_{33}$

Solving the above we get  $\underline{\alpha}_3 = \begin{pmatrix} 2.414 \\ 1 \\ 0 \end{pmatrix}$

$$\text{and } \| \underline{\alpha}_3 \| = \sqrt{2.414^2 + 1^2 + 0^2} = 2.613$$

Hence  $\underline{e}_3 = \begin{pmatrix} 0.923 \\ 0.383 \\ 0 \end{pmatrix}$

$\therefore$  The 3rd PC is  $Y_3 = 0.923 X_1 + 0.383 X_2$

Now variance explained by the 1st PC is  $\frac{5.828}{(5.828 + 2 + 0.172)} = 0.728 = 72.8\%$

Variance explained by 2nd PC is  $\frac{2}{(5.828 + 2 + 0.172)} = 0.25 = 25\%$

Variance explained by 3rd PC is  $\frac{0.172}{(5.828 + 2 + 0.172)} = 0.015 = 1.5\%$

and variance explained by

1st and 2nd PC is  $\frac{5.828 + 0.172}{(5.828 + 2 + 0.172)} = 0.75 = 75\%$

$\therefore$  If one wants to explain 80% of the total system variability one can retain the 1st and 3rd PC.

2) Solve the following var-cov matrix and find the PC's  
 Also convert the given matrix into correlation matrix,  
 find the PC's and interpret the results.

$$\begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix}$$

ans) By looking at  $\Sigma$  we can say that  $|\Sigma| = 84 \neq 0$  i.e  $\Sigma$  is non singular.

Let  $\lambda_1$  and  $\lambda_2$  be the eigenvalues of  $\Sigma$  i.e we need to solve  $|\Sigma - \lambda I| = 0$

$$\Rightarrow \begin{vmatrix} 1-\lambda & 4 \\ 4 & 100-\lambda \end{vmatrix} = 0$$

$$\Rightarrow (1-\lambda)(100-\lambda) - 16 = 0$$

$$\Rightarrow \lambda^2 - 101\lambda + 84 = 0$$

$$\therefore \lambda = \frac{101 \pm \sqrt{101^2 - (4 \times 84)}}{2} = 100.16 \text{ and } 0.84$$

$$\therefore \lambda_1 = 100.16 \text{ and } \lambda_2 = 0.84 \text{ and } \lambda_1 > \lambda_2$$

Now to have to find the orthonormal eigenvectors corresponding to  $\lambda_1$  and  $\lambda_2$ .

Let  $\underline{a}_i' = (a_{i1}, a_{i2})$ ;  $i=1,2$  be the eigenvector to  $\lambda_i$

Now we have, for  $\lambda_1 = 100.16$

$$\Sigma \underline{a}_1 = \lambda_1 \underline{a}_1$$

$$\Rightarrow \begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 100.16 a_{11} \\ 100.16 a_{12} \end{pmatrix}$$

$$\Rightarrow 4a_2 - 99.16a_1 = 0 \quad \text{and} \quad 4a_1 - 0.16a_2 = 0$$

$$\Rightarrow a_1 = 0.04a_2 = 0$$

$$\text{Let } a_2 = 1 \text{ then } a_1 = 0.04$$

$$\therefore \underline{a}_1 = \begin{pmatrix} 0.04 \\ 1 \end{pmatrix} \text{ and } \|\underline{a}_1\| = \sqrt{1^2 + 0.04^2} \approx 1$$

$$\therefore \underline{e}_1 = \begin{pmatrix} 0.04 \\ 1 \end{pmatrix}$$

Hence the 1st PC is  $y_1 = 0.04x_1 + x_2$

Now for  $\lambda_2 = 0.84$  we have

$$\Sigma \underline{a}_2 = \lambda_2 \underline{a}_2$$

$$\Rightarrow \begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 0.84 a_{11} \\ 0.84 a_{12} \end{pmatrix}$$

$$\Rightarrow 0.16a_{11} + 4a_{12} = 0 \quad \text{and} \quad 4a_{11} + 99.16a_{12} = 0$$

If  $a_{11} = 1$  then  $a_{12} = -0.04$

$$\therefore \underline{a}_1 = \begin{pmatrix} 1 \\ -0.04 \end{pmatrix} \quad \text{and} \quad \|a_1\| = \sqrt{1^2 + (-0.04)^2} \approx 1$$

$$\text{Hence } \underline{e}_1 = \begin{pmatrix} 1 \\ -0.04 \end{pmatrix}$$

$$\therefore \text{The 2nd PC is } Y_2 = X_1 + 0.04X_2$$

Now converting the given var-cov matrix into correlation matrix we get,

$$R = \begin{pmatrix} 1 & 0.4 \\ 0.4 & 1 \end{pmatrix}$$

Now  $\det R = 1 - 0.4^2 = 0.84 \neq 0$  ie R is non singular. Proceeding in a similar way we get,

$$|R - \lambda I| = 0 \Rightarrow \begin{vmatrix} 1 - \lambda & 0.4 \\ 0.4 & 1 - \lambda \end{vmatrix} = 0$$

$$\Rightarrow (1 - \lambda)^2 - 0.4^2 = 0$$

$$\Rightarrow (1 - \lambda + 0.4)(1 - \lambda - 0.4) = 0$$

$$\Rightarrow (1.4 - \lambda)(0.6 - \lambda) = 0$$

$$\therefore \lambda_1 = 1.4 \quad \lambda_2 = 0.6$$

$$\text{and } \lambda_1 > \lambda_2$$

Now for the first eigenvalue  $\lambda_1$  we get

$$R \underline{a}_1 = \lambda_1 \underline{a}_1$$

$$\Rightarrow \begin{pmatrix} 1 & 0.4 \\ 0.4 & 1 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 1.4 a_{11} \\ 1.4 a_{12} \end{pmatrix}$$

$$\Rightarrow 0.4 a_{12} - 0.4 a_{11} = 0 \quad \text{and} \quad 0.4 a_{22} - 0.4 a_{12} = 0$$

$$\Rightarrow a_{11} = a_{12} \quad \Rightarrow a_{11} = a_2$$

Hence  $\underline{a}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \| \underline{a}_1 \| = \sqrt{1+1} = \sqrt{2}$

$$\therefore \underline{e}_1 = \begin{pmatrix} 0.707 \\ 0.707 \end{pmatrix}$$

$\therefore$  The 1st PC is  $y_1 = 0.707x_1 + 0.707x_2$

Similarly for  $\lambda_2 = 0.6$  we get,

$$R \underline{a}_2 = \lambda_2 \underline{a}_2$$

$$\Rightarrow \begin{pmatrix} 1 & 0.4 \\ 0.4 & 1 \end{pmatrix} \begin{pmatrix} a_{21} \\ a_{22} \end{pmatrix} = \begin{pmatrix} 0.6 a_{21} \\ 0.6 a_{22} \end{pmatrix}$$

$$\Rightarrow a_{21} + 0.4 a_{22} = 0.6 a_{21} \quad \text{and} \quad 0.4 a_{21} + a_{22} = 0.6 a_{22}$$

$$\Rightarrow 0.4 a_{22} = -0.4 a_{21} \quad \text{and} \quad 0.4 a_{21} + 0.4 a_{22} = 0$$

$$\Rightarrow a_{22} = -a_{21} \quad \Rightarrow a_{21} = a_{22}$$

Hence  $\underline{a}_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \| \underline{a}_2 \| = \sqrt{1^2 + 1^2} = \sqrt{2}$

$$\therefore \underline{e}_2 = \begin{pmatrix} 0.707 \\ -0.707 \end{pmatrix}$$

$\therefore$  The 2nd PC is  $y_2 = 0.707x_1 - 0.707x_2$

$\therefore$  We can see that the PC's obtained from cov matrix and correlation matrix is not the same. It may also be noted that the corr-matrix is cov matrix of the standardized variables. PC's from corr-matrix is more preferable.

In cov matrix the 1st PC explains  $\frac{100.16}{101} = 0.992 = 99.2\%$

In corr matrix the 1st PC explains  $\frac{1.4}{2} = 0.707 \approx 70\%$

Moreover in corr matrix both the variables  $x_1$  and  $x_2$  contribute equally whereas in case of cov matrix it is not the case.

3) Consider the following correlation matrix and try to find a 1 factor model solution of it.

$$R = \begin{pmatrix} 1 & 0.705 & 0.731 \\ 0.705 & 1 & 0.716 \\ 0.731 & 0.716 & 1 \end{pmatrix}$$

ans) We know that for a 1 factor model solution, the correlation matrix can be written as,

$$R = \begin{pmatrix} l_{11}^2 & l_{11}l_{12} & l_{11}l_{31} \\ l_{21}l_{11} & l_{21}^2 & l_{21}l_{31} \\ l_{31}l_{11} & l_{31}l_{21} & l_{31}^2 \end{pmatrix} + \begin{pmatrix} \psi_1^2 & 0 & 0 \\ 0 & \psi_2^2 & 0 \\ 0 & 0 & \psi_3^2 \end{pmatrix} \quad \text{--- (1)}$$

$= LL' + \psi^2$  where  $l_{ij}$  = loading of  $i$ th variable on  $j$ th factor

Now comparing the given correlation matrix with the above form we get,

$$l_{11}l_{21} = 0.705 \quad l_{11}l_{31} = 0.731 \quad l_{21}l_{31} = 0.716$$

$$\text{Now, } l_{11} = \frac{0.705}{l_{21}} \quad \therefore l_{11}l_{31} = 0.731 \Rightarrow \frac{0.705}{l_{21}} l_{31} = 0.731$$

Now multiplying we get,

$$\frac{0.705}{l_{21}} l_{31} \times l_{21} l_{31} = 0.716 \times 0.731$$

$$\Rightarrow l_{31}^2 = \frac{0.716 \times 0.731}{0.705}$$

$$\Rightarrow l_{31} = \pm 0.861$$

$$\text{Taking } l_{31} = 0.861 \text{ we get, } l_{21} = \frac{0.705 \times 0.861}{0.731} = 0.831$$

$$\text{and } l_{11} = \frac{0.705}{0.831} = 0.848$$

If  $\ell_{31} = -0.861$  then the sign of the values of  $\ell_{11}$  and  $\ell_{21}$   
will be -ve

Now Putting the values of  $\ell_{11}, \ell_{21}, \ell_{31}$  in the model of (1)  
the values of  $\psi_1^2, \psi_2^2$  and  $\psi_3^2$  can be obtained.

$$\therefore \psi_1^2 = 1 - \ell_{11}^2 = 1 - (0.848)^2 = 0.280$$

$$\psi_2^2 = 1 - \ell_{21}^2 = 1 - (0.831)^2 = 0.309$$

$$\psi_3^2 = 1 - \ell_{31}^2 = 1 - (0.861)^2 = 0.258$$

$\therefore$  Given the correlation matrix a one factor model  
solution exists and can be written as,

$$x_1 = 0.848 F_1 + \epsilon_1$$

$$x_2 = 0.831 F_1 + \epsilon_2$$

$$x_3 = 0.861 F_1 + \epsilon_3$$

2) Consider the var-cov matrix and find one factor model  
solution if it exists using the principal component method

$$\Sigma = \begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix}$$

ans) Let us consider the one factor model as,

$$x_1 = \ell_{11} F_1 + \epsilon_1$$

$$x_2 = \ell_{21} F_1 + \epsilon_2$$

$$\text{Now we know, } \Sigma = \begin{pmatrix} \ell_{11}^2 & \ell_{11} \ell_{21} \\ \ell_{21} \ell_{11} & \ell_{22}^2 \end{pmatrix} + \begin{pmatrix} \psi_1^2 & 0 \\ 0 & \psi_2^2 \end{pmatrix}$$

$$= LL' + \Psi$$

Comparing the variance covariance with the above model  
we get,  $\ell_{11} \ell_{21} = 4 \Rightarrow \ell_{11} = \frac{4}{\ell_{21}}$

So here it may not be easy to find the one factor

model solution by this process. We consider the PCA method to solve this.

Using spectral decomposition we can say,

$$\Sigma = \lambda_1 e_1 e_1' + \lambda_2 e_2 e_2' + \dots + \lambda_p e_p e_p'$$

where  $\lambda_i$  = eigenvalues and  $e_i$  = corresponding eigenvectors  
 $i=1(1)p$

$$\text{or, } \Sigma = (\sqrt{\lambda_1} e_1, \sqrt{\lambda_2} e_2, \dots, \sqrt{\lambda_p} e_p) \begin{pmatrix} \sqrt{\lambda_1} e_1' \\ \sqrt{\lambda_2} e_2' \\ \vdots \\ \sqrt{\lambda_p} e_p' \end{pmatrix}$$

$\therefore$  The given matrix can be written as,

$$\begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix} = \left( \sqrt{100.16} \begin{pmatrix} 0.04 \\ 1 \end{pmatrix}, \sqrt{0.84} \begin{pmatrix} 1 \\ -0.04 \end{pmatrix} \right) \begin{pmatrix} \sqrt{100.16} (0.04 \ 1) \\ \sqrt{0.84} (1 \ -0.04) \end{pmatrix}$$

$$\text{where } \sqrt{\lambda_1} = \sqrt{100.16}, \quad \sqrt{\lambda_2} = \sqrt{0.84}$$

$$\text{and } e_1 = \begin{pmatrix} 0.04 \\ 1 \end{pmatrix}, \quad e_2 = \begin{pmatrix} 1 \\ -0.04 \end{pmatrix}$$

$$= \begin{pmatrix} 0.400 & 0.917 \\ 10.008 & -0.037 \end{pmatrix} \begin{pmatrix} 0.400 & 10.008 \\ 0.917 & -0.037 \end{pmatrix}$$

$$\simeq \begin{pmatrix} 0.400 \\ 10.008 \end{pmatrix} \begin{pmatrix} 0.400 & 10.008 \end{pmatrix} + \begin{pmatrix} \psi_1^2 & 0 \\ 0 & \psi_2^2 \end{pmatrix}$$

$$\simeq \begin{pmatrix} l_{11} \\ l_{12} \end{pmatrix} (l_{11} \ l_{12}) + \begin{pmatrix} \psi_1^2 & 0 \\ 0 & \psi_2^2 \end{pmatrix}$$

$\therefore$  It can be seen that  $l_{11} = 0.400$  and  $l_{12} = \frac{10.008}{\simeq 10}$

Now for obtaining  $\psi_1^2$  and  $\psi_2^2$  we can proceed as follows:

$$\Sigma = \begin{pmatrix} l_{11}^2 & l_{11}l_{12} \\ l_{21}l_{11} & l_{22}^2 \end{pmatrix} + \begin{pmatrix} \psi_1^2 & 0 \\ 0 & \psi_2^2 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} = \begin{pmatrix} l_{11}^2 & l_{11}l_{21} \\ l_{21}l_{11} & l_{22}^2 \end{pmatrix} + \begin{pmatrix} \psi_1^2 & 0 \\ 0 & \psi_2^2 \end{pmatrix}$$

$\therefore$  Comparing the above we can say that,

$$\sigma_{11} = l_{11}^2 + \psi_1^2 \Rightarrow \psi_1^2 = \sigma_{11} - l_{11}^2$$

$$\sigma_{22} = l_{21}^2 + \psi_2^2 \Rightarrow \psi_2^2 = \sigma_{22} - l_{21}^2$$

$$\text{Hence we get, } \psi_1^2 = 1 - 0.400^2 = 0.84$$

$$\text{and, } \psi_2^2 = 100 - (10)^2 \approx 0$$

$\therefore$  The one factor model solution is,

$$x_1 = 0.4 F_1 + \epsilon_1$$

$$x_2 = 10 F_1 + \epsilon_2$$

Q) 3) Consider the distance matrix given below. Perform single, complete and average linkage clustering on the distance matrix.

D =	A	B	C	D	E
A	1	0.8	0.4	0.2	0.1
B	0.8	1	0.6	0.3	0.2
C	0.4	0.6	1	0.7	0.5
D	0.2	0.3	0.7	1	0.9
E	0.1	0.2	0.5	0.9	1

ans) We follow the given algorithm to perform clustering.

- (1) Start with  $n$  clusters each containing 1 item
- (2) Club 2 most close clusters into one cluster. At this step there would be almost  $(n-1)$  clusters
- (3) Recompute or update the distances between the cluster
- (4) Repeat step 2 and 3 until one single cluster is obtained

■ Single clustering: we consider smallest distance i.e  $d(uv)_{\text{min}} = \min\{d_{uv}, d_{vw}\}$

There are 5 clusters A, B, C, D, E. Distance between A and C is minimum so they will be clubbed. The new cluster is,

	AE	B	C	D
AE	1			
B	0.2	1		
C	0.4	0.6	1	
D	0.2	0.3	0.7	1

$$\min\{d(A, B), d(E, B)\} = \min\{0.8, 0.2\} = 0.2$$

$$\min\{d(A, C), d(E, C)\} = \min\{0.4, 0.5\} = 0.4$$

$$\min\{d(A, D), d(E, D)\} = \min\{0.2, 0.9\} = 0.2$$

Now distance between AE and B is minimum so they are clubbed together.

	ABE	C	D
ABE	1		
C	0.4	1	
D	0.2	0.7	1

$$\min\{d(AC), d(B, C), d(E, C)\} = 0.4$$

$$\min\{d(A, D), d(B, D), d(E, D)\} = 0.2$$

Now distance between ABED and C is minimum

	A	B	E	D	C
ABED	1				
C	0.4	1			

$$\min\{d(Ac), d(AC), d(B,C), d(EC)\} = 0.4$$

Finally C will be clubbed with ABED and the process ends with one cluster.

- Complete linkage: We consider maximum distance here i.e  $d(uv)_w = \max\{d_{uw}, d_{vw}\}$

Here distance between AE is minimum so they will be merged.

	A	E	B	C	D
AE	1				
B	0.8	1			
C	0.5	0.6	1		
D	0.9	0.3	0.7	1	

$$\begin{aligned} \max\{d(A,B), d(E,B)\} &= \max\{0.8, 0.2\} = 0.8 \\ \max\{AC, EC\} &= 0.5 \\ \max\{AD, ED\} &= 0.9 \end{aligned}$$

Now distance between B and D is minimum so they are merged

	A	E	B	D	C
AE	1				
BD	0.9		1		
C	0.5		0.7	1	

$$\begin{aligned} \max\{B(AE), D(AE)\} &= \max\{0.9, 0.9\} = 0.9 \\ \max\{AC, EC\} &= \max\{0.4, 0.5\} = 0.5 \\ \max\{BC, DC\} &= \max\{0.6, 0.7\} = 0.7 \end{aligned}$$

Now distance between C and AE is minimum so they are merged.

	A	E	C	B	D
AEC	1				
BD	0.9			1	

$$\max\{B(AEC), D(AEC)\} = \max\{0.9, 0.7\} = 0.9$$

Finally BD will be merged with AEC and the process ends here with one cluster.

Average linkage: We take average distance between the cluster i.e

$$d(UV)_H = \frac{\sum_i \sum_k d_{ik}}{N_{(UV)} N_H}$$

Here distance between A and E is minimum so they are merged

	AE	B	C	D
AE	1			
B	0.5	1		
C	0.45	0.6	1	
D	0.55	0.3	0.7	1

$$\text{avg}\{AB, EB\} = 0.5$$

$$\text{avg}\{AC, EC\} = 0.45$$

$$\text{avg}\{AD, ED\} = 0.55$$

Now distance between B and D is minimum so they are merged together

	AE	B	D	C
AE	1			
BD	0.95	1		
C	0.45	0.65	1	

$$\text{avg}\{B(AE), D(AE)\} = 0.85$$

$$\text{avg}\{AC, ED\} = 0.45$$

$$\text{avg}\{BC, DC\} = 0.65$$

Now distance between AE and C is minimum so they are merged.

	AEC	BD
AEC	1	
BD	0.8	1

$$\text{avg}\{B(AEC), D(AEC)\} = 0.8$$

Finally BD will be merged with AEC and the process stops with one cluster

Q) Let there be 5 objects A, B, C, D, E where each of them has the coordinate as follows:

$$A(2,3), B(5,1), C(1,1), D(6,7), E(2,3)$$

Also let 2 clusters to be formed using K-means method.

Ans) Let us make 2 arbitrary clusters  $C_1$  and  $C_2$  with the following members,

$$C_1 = (A, B, C) \text{ and } C_2 = (D, E)$$

$$\text{Centroid of } C_1 = \left( \frac{2+5+1}{3}, \frac{3+1+1}{3} \right) = (3, 2)$$

$$\text{Centroid of } C_2 = \left( \frac{6+2}{2}, \frac{7+3}{2} \right) = (4, 5)$$

Now, distance of each obs from  $C_1$  and  $C_2$  are to be computed,

$$d(A, C_1) = [(3-2)^2 + (2-3)^2]^{1/2} = \sqrt{2}$$

$$d(A, C_2) = [(2-4)^2 + (3-5)^2]^{1/2} = \sqrt{8} \quad \therefore A \text{ stays in } C_1$$

$$d(B, C_1) = [(5-3)^2 + (1-2)^2]^{1/2} = \sqrt{5}$$

$$d(B, C_2) = [(5-4)^2 + (1-5)^2]^{1/2} = \sqrt{17} \quad \therefore B \text{ stays in } C_1$$

$$d(C, C_1) = [(3-2)^2 + (2-1)^2]^{1/2} = \sqrt{1}$$

$$d(C, C_2) = [(4-2)^2 + (5-2)^2]^{1/2} = \sqrt{13} \quad \therefore C \text{ stays in } C_1$$

$$d(D, C_1) = [(6-3)^2 + (7-2)^2]^{1/2} = \sqrt{50}$$

$$d(D, C_2) = [(6-4)^2 + (7-5)^2]^{1/2} = \sqrt{8} \quad \therefore D \text{ stays in } C_2$$

$$d(E, C_1) = [(2-3)^2 + (2-3)^2]^{1/2} = \sqrt{2}$$

$$d(E, C_2) = [(2-4)^2 + (3-5)^2]^{1/2} = \sqrt{8} \quad \therefore E \text{ goes to } C_1 \text{ as } d(E, C_1) \text{ is less than } d(E, C_2)$$

Now the clusters are  $C_1 = (A, B, C, E)$  and  $C_2 = D$

i.e. E is reassigned to  $C_1$  and others stay in their earlier clusters. Again the cluster centroid is computed.

$$\therefore \text{The new centroid of } C_1 = \left( \frac{1}{4}(2+5+2+3), \frac{1}{4}(3+1+2+3) \right) \\ = (2.75, 2.25)$$

and centroid of  $C_2 = (6, 7)$

Again the distance of each obs from  $C_1$  and  $C_2$  is computed

$$d(A, C_1) = [(2.75 - 2)^2 + (2.25 - 3)^2]^{1/2} = 1.06$$

$$d(A, C_2) = [(2 - 6)^2 + (3 - 7)^2]^{1/2} = 5.65$$

$$d(B, C_1) = [(5 - 2.75)^2 + (2.25 - 1)^2]^{1/2} = 2.52$$

$$d(B, C_2) = [(5 - 6)^2 + (1 - 7)^2]^{1/2} = 6.08$$

$$d(C, C_1) = [(2 - 2.75)^2 + (2 - 2.25)^2]^{1/2} = 0.79$$

$$d(C, C_2) = [(2 - 6)^2 + (2 - 7)^2]^{1/2} = 6.40$$

$$d(D, C_1) = [(6 - 2.75)^2 + (7 - 2.25)^2]^{1/2} = 5.75$$

$$d(D, C_2) = [(6 - 6)^2 + (7 - 7)^2]^{1/2} = 0$$

$$d(E, C_1) = [(2 - 2.75)^2 + (3 - 2.25)^2]^{1/2} = 1.06$$

$$d(E, C_2) = [(2 - 6)^2 + (3 - 7)^2]^{1/2} = 5.65$$

Now it is observed in the latest composition of  $C_1$  and  $C_2$  that each object is at the nearest centroid of the clusters. Hence the final 2 clusters made from the above 5 elements are as follows.

$$C_1 = (A, B, C, E), \quad C_2 = (D)$$

2) Consider the var-cov matrix given below and find out the canonical correlation.

$$R = \begin{pmatrix} 1 & 0.4 & 0.5 & 0.6 \\ 0.4 & 1 & 0.3 & 0.4 \\ 0.5 & 0.3 & 1 & 0.2 \\ 0.6 & 0.4 & 0.2 & 1 \end{pmatrix}_{4 \times 4}$$

(a) From the given matrix we can say that

$$\Sigma_{11} = \begin{pmatrix} 1 & 0.4 \\ 0.4 & 1 \end{pmatrix} \quad \Sigma_{22} = \begin{pmatrix} 1 & 0.2 \\ 0.2 & 1 \end{pmatrix}$$

$$\Sigma_{12} = \begin{pmatrix} 0.5 & 0.6 \\ 0.3 & 0.4 \end{pmatrix} \quad \Sigma_{21} = \begin{pmatrix} 0.5 & 0.3 \\ 0.5 & 0.4 \end{pmatrix}$$

Our objective is to find canonical correlation between the variables  $\underline{x}^{(1)}$  and  $\underline{x}^{(2)}$

where  $\underline{x}^{(1)} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ ,  $\underline{x}^{(2)} = \begin{pmatrix} x_3 \\ x_4 \end{pmatrix}$  and  $\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}$

We take linear combination such as,  $U = \underline{a}' \underline{x}^{(1)}$  and  $V = \underline{b}' \underline{x}^{(2)}$  such that  $\underline{a}$  and  $\underline{b}$  are so chosen that correlation coefficient between  $U$  and  $V$  is maximum

The canonical correlation can be found by calculating the eigenvalues of the matrix,

$$\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \quad \text{and} \quad \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Now,  $\Sigma_{22}^{-1} = \begin{pmatrix} 1 & 0.2 \\ 0.2 & 1 \end{pmatrix}^{-1}$ , Now  $\det \Sigma_{22} = 1 - 0.2^2 = 0.96$

$$\Sigma_{22}^{-1} = \frac{1}{\det \Sigma_{22}} \text{adj } \Sigma = \begin{pmatrix} 1.04 & -0.208 \\ -0.208 & 1.04 \end{pmatrix}$$

$$\text{similarly } \Sigma_{11}^{-1} = \frac{1}{|\Sigma_{11}|} \text{adj} |\Sigma_{11}) \quad \det |\Sigma_{11} = 1 - 0.4^2 \\ = 0.84$$

$$= \begin{pmatrix} 1.19 & -0.476 \\ -0.476 & 1.19 \end{pmatrix}$$

$$\therefore \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$$

$$= \begin{pmatrix} 1.04 & -0.208 \\ -0.208 & 1.04 \end{pmatrix} \begin{pmatrix} 0.5 & 0.3 \\ 0.6 & 0.4 \end{pmatrix} \begin{pmatrix} 1.19 & -0.476 \\ -0.476 & 1.19 \end{pmatrix} \begin{pmatrix} 0.5 & 0.6 \\ 0.3 & 0.4 \end{pmatrix}$$

$$= \begin{pmatrix} 247/625 & 143/625 \\ 13/25 & 221/625 \end{pmatrix} \begin{pmatrix} 1.19 & -0.476 \\ -0.476 & 1.19 \end{pmatrix} \begin{pmatrix} 0.5 & 0.6 \\ 0.3 & 0.4 \end{pmatrix}$$

$$= \begin{pmatrix} 0.361 & 0.084 \\ 0.450 & 0.173 \end{pmatrix} \begin{pmatrix} 0.5 & 0.6 \\ 0.3 & 0.4 \end{pmatrix}$$

$$= \begin{pmatrix} 0.2057 & 0.2502 \\ 0.2769 & 0.3392 \end{pmatrix}$$

Now the eigenvalues are given by,

$$|\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} - \lambda^2 I| = 0 \quad [\text{Consider } \lambda^2 \text{ as } \lambda]$$

$$\Rightarrow \begin{vmatrix} 0.2057 - \lambda & 0.2502 \\ 0.2769 & 0.3392 - \lambda \end{vmatrix} = 0 \Rightarrow (0.2057 - \lambda)(0.3392 - \lambda) - (0.2502 \times 0.2769) = 0$$

$$\Rightarrow 0.0850 - 0.2057\lambda - 0.3392\lambda + \lambda^2 - 0.069 = 0$$

$$\Rightarrow 0.0157 - 0.5899\lambda + \lambda^2 = 0$$

$$\therefore \lambda = \frac{0.5899 \pm \sqrt{0.5899^2 - (4 \times 0.0157)}}{2}$$

$$\therefore \lambda_1^2 = 0.0009, \lambda_2^2 = 0.545$$

$$\Rightarrow \lambda_1 = 0.03, \lambda_2 = 0.738$$

Hence the largest eigenvalue  $\lambda_2 = 0.738$  is the canonical correlation between the variables  $u$  and  $v$ .

2) Consider the data given below. Find equality of mean vectors of all the groups

(20, 185)	(15, 203)	(8, 215)
(12, 190)	(11, 210)	(7, 189)
(5, 187)	(9, 224)	(6, 201)
(30, 200)	(5, 211)	(5, 198)
(15, 203)	(10, 200)	(3, 207)
↓	↓	↓
G <sub>1</sub>	G <sub>2</sub>	G <sub>3</sub>

Ques) Our objective is to test  $H_0: \underline{\mu}_1 = \underline{\mu}_2 = \underline{\mu}_3$ , where

$$\underline{\mu}_i = \begin{bmatrix} \mu_{i1} \\ \mu_{i2} \end{bmatrix}; i=1, 2, 3$$

The model is,  $\tilde{x}_{ij} = \underline{\mu} + \underline{\alpha}_i + \underline{\epsilon}_{ij}; i=1, 2, \dots, j=1, 2, \dots, n$

where,  $\underline{\epsilon}_{ij} \sim N(0, \Sigma)$  and  $\sum_i \underline{\alpha}_i = 0$

$\underline{\mu}$  is the overall mean effect

$\underline{\alpha}_i$  is the  $i$ th treatment effect

The overall model can be written as,

$$\tilde{x}_{ij} = \bar{x} + (\bar{x}_i - \bar{x}) + (x_{ij} - \bar{x}_i)$$

↓      ↓      ↓      ↓  
 Obs    overall    sample    error  
 ↓      ↓      ↓      ↓  
 Obs    overall    estimated    error  
 ↓      ↓      ↓      ↓  
 mean    mean    treatment    effect

Now following the above idea we have,

$$\bar{x}_1 = \begin{bmatrix} 16 \\ 146.6 \end{bmatrix} \quad \bar{x}_2 = \begin{bmatrix} 10 \\ 212 \end{bmatrix} \quad \bar{x}_3 = \begin{bmatrix} 5.8 \\ 201.4 \end{bmatrix}$$

$$\text{and } \bar{x} = \begin{bmatrix} 10.6 \\ 186 \end{bmatrix}$$

Consider the 1st obs (1st variable) in the pair of variables. We can write them as,

$$\begin{pmatrix} 20 & 10 & 5 & 30 & 15 \\ 15 & 11 & 9 & 5 & 10 \\ 8 & 7 & 6 & 5 & 3 \end{pmatrix} = \begin{pmatrix} 10.6 & 10.6 & 10.6 & 10.6 & 10.6 \\ 10.6 & 10.6 & 10.6 & 10.6 & 10.6 \\ 10.6 & 10.6 & 10.6 & 10.6 & 10.6 \end{pmatrix} +$$

$$\begin{pmatrix} 4 & -6 & -11 & 14 & -1 \\ 5 & 1 & -1 & -5 & 0 \\ 2.2 & 1.2 & 0.2 & -0.8 & -2.8 \end{pmatrix} + \begin{pmatrix} 5.4 & 5.4 & 5.4 & 5.4 & 5.4 \\ -0.6 & -0.6 & -0.6 & -0.6 & -0.6 \\ -4.8 & -4.8 & -4.8 & -4.8 & -4.8 \end{pmatrix}$$

$\Rightarrow$  observation = mean + residual + treatment effect

i.e  $SS_{obs} = SS_{mean} + SS_{res} + SS_{int}$

Now  $SS_{obs} = 20^2 + 10^2 + 5^2 + 30^2 + \dots + 5^2 + 3^2 = 2385$

$$SS_{mean} = 15 \times 10.6^2 = 1685.4$$

$$SS_{int} = (5 \times 5.4^2) + (5 \times -0.6^2) + (5 \times -4.8^2) = 262.8$$

$$SS_{res} = 4^2 + 6^2 + 11^2 + \dots + 0.8^2 + 2.8^2 = 436.8$$

Similarly for the 2nd variable we can write,

$$\begin{pmatrix} 105 & 110 & 115 & 200 & 203 \\ 205 & 210 & 220 & 215 & 200 \\ 215 & 190 & 201 & 198 & 203 \end{pmatrix} = \begin{pmatrix} 186 & 186 & 186 & 186 & 186 \\ 186 & 186 & 186 & 186 & 186 \\ 186 & 186 & 186 & 186 & 186 \end{pmatrix} +$$

$$+ \begin{pmatrix} -41.6 & -36.6 & -31.6 & 53.4 & 56.4 \\ -5 & 0 & 10 & 5 & -10 \\ 13.6 & -11.4 & -0.4 & -3.4 & 1.6 \end{pmatrix}$$

$$+ \begin{pmatrix} -39.4 & -39.4 & -39.4 & -39.4 & -39.4 \\ 24 & 24 & 24 & 24 & 24 \\ 15.4 & 15.4 & 15.4 & 15.4 & 15.4 \end{pmatrix}$$

$\Rightarrow$  observation = mean + residual + treatment effect

i.e  $SS_{obs} = SS_{mean} + SS_{res} + SS_{int}$

$$\text{Now, } SS_{\text{obs}} = 105^2 + 110^2 + \dots + 198^2 + 203^2 = 54144.8$$

$$SS_{\text{mean}} = 15 \times 186^2 = 518940$$

$$SS_{\text{tnt}} = (5 \times -39.4^2) + (5 \times 24^2) + (5 \times 15.4^2) = 11827.6$$

$$SS_{\text{Res}} = 41.6^2 + 36.6^2 + 31.6^2 + \dots + 3.4^2 + 1.6^2 = 10680.4$$

Now proceeding row by row for the 2 variables the cross product are computed as below.

$$SSP_{\text{Mean}} = 15 \times (186 \times 10.6) = 29574$$

$$SSP_{\text{int}} = 5(-39.4 \times 5.4) + 5(0.6 \times 24) + 5(4.8 \times 15.4) = 1505.4$$

$$SSP_{\text{Res}} = (4 \times -39.4) + (-6 \times -36.6) + \dots + (-2.8 \times 1.6) = 1046.4$$

$$\text{Total SP} = (20 \times 105) + (10 \times 110) + \dots + (3 \times 203) = 29115$$

$$\begin{aligned}\text{Total (corrected) cross product} &= \text{Total SP} - SSP_{\text{mean}} \\ &= 29115 - 29574 = -459\end{aligned}$$

The Manova Table is given by,

### MANOVA TABLE

Source	Matrix of SS and CP	d.f
treatment	$\begin{bmatrix} 262.8 & 1505.4 \\ 1505.4 & 11827.6 \end{bmatrix}$	2
residual	$\begin{bmatrix} 436.8 & 1046.4 \\ 1046.4 & 10680.4 \end{bmatrix}$	6
Total	$\begin{bmatrix} 699.6 & 2551.8 \\ 2551.8 & 22508 \end{bmatrix}$	8

Therefore Wilk's Lambda is,

$$\Lambda = \frac{|WSS|}{|WSS + BSS|}$$

$$\Delta = \frac{\begin{vmatrix} 436.8 & 1046.4 \\ 1046.4 & 10680.4 \end{vmatrix}}{\begin{vmatrix} 699.6 & 2551.2 \\ 2551.2 & 22508 \end{vmatrix}} = \frac{3570245.76}{9234913.56} = 0.387$$

Now from the distribution of Wilk's Lambda we have,

No of variables = 2, No of groups = 3

$$\text{Hence, } \left( \frac{\sum_{i=1}^k n_i - k - 1}{k-1} \right) \left( \frac{1 - \sqrt{\Delta}}{\sqrt{\Delta}} \right)$$

$$= \left( \frac{9 - 3 - 1}{3 - 1} \right) \left( \frac{1 - \sqrt{0.387}}{\sqrt{0.387}} \right)$$

$$= 1.519 \sim F(2(k-1), 2(\sum_{i=1}^k n_i - k - 1)) = F(4, 10)$$

Now we can see that ,

$$F_{\text{obs}} = 1.519 < F_{0.05}(4, 10) = 3.48$$

Therefore, we accept the null hypothesis that there is no difference between the groups w.r.t to the mean vectors.