



**Written exam in Multivariate Methods, 7.5 ECTS credits**

Thursday, 23<sup>th</sup> October 2014, 12:00 – 17:00

Time allowed: FIVE hours

Examination Hall: Ugglevikssalen

You are required to answer all **6 (six)** questions as well as motivate your solutions. The total amount of points is 80. In order to pass this part, you need to get at least 40 points. Points from this exam will be added to your results from the home assignment. The final grades are assigned as follows: **A** (91+), **B** (75-90), **C** (66-74), **D** (58-65), **E** (50-57), **Fx** (30-49), and **F** (0-29).

You are **allowed** to use a pocket calculator, a language dictionary, and a list of formulas (attached). In addition, you are **allowed** to use a one-sided A4 containing your own formulae, but excluding proofs and solutions. The A4 must be approved (signed) by the teacher; and, it must be submitted along with your solutions. If the A4 is not signed by the teacher and discovered, student might be accused in cheating on exam.

The teacher reserves the right to examine the students **orally** on the questions in this examination.

1. (10 points) The coordinates of three points with respect to the orthogonal basis vectors  $e_1$  and  $e_2$  are as follows:

$$A = (2, -1) \quad B = (1, 0.7) \quad C = (-3, 7.5)$$

Show that A, B and C lie on a straight line.

The set of oblique basis vectors  $f_1$  and  $f_2$  are related to vectors  $e_1$  and  $e_2$  as follows:

$$e_1 = .804f_1 + .629f_2$$

$$e_2 = .273f_1 - .588f_2$$

Compute the coordinates of A, B and C with respect to  $f_1$  and  $f_2$ . Do they still lie on a straight line? Justify your answer.

2. (18 points) Use the data in the contingency Table below to answer the following questions:

Blood Cholesterol (mg/100 cc)	Heart Disease	
	Present	Absent
<200	6	5
200-219	10	6
220-259	30	5
>259	45	7

- a. What is the probability that
- Heart disease is present?
  - Blood cholesterol is less than 219 mg/100 cc?

- b. What are the odds that
- (i) Heart disease is present given that blood cholesterol is less than 200 mg/100 cc?
  - (ii) Blood cholesterol is greater than 219 mg/100 cc given that heart disease is present?
- c. Compute the log of odds that
- (i) Heart disease is present given that blood cholesterol is greater than 259 mg/100 cc.
  - (ii) Heart disease is present given that blood cholesterol is between 200 and 219 mg/100 cc.

3. (15 points) Consider the two-indicator two-factor model represented by the following equations:

$$X_1 = 0.104F_1 + 0.824F_2 + U_1$$

$$X_2 = 0.065F_1 + 0.959F_2 + U_2$$

$$X_3 = 0.065F_1 + 0.725F_2 + U_3$$

$$X_4 = 0.906F_1 + 0.134F_2 + U_4$$

$$X_5 = 0.977F_1 + 0.116F_2 + U_5$$

$$X_6 = 0.827F_1 + 0.016F_2 + U_6$$

The usual assumptions hold for the above model. Answer the following questions assuming that the correlation between the common factors  $F_1$  and  $F_2$  is given by

$$\text{Corr}(F_1, F_2) = \phi_{12} = 0.3$$

- (a) What are the pattern loadings of indicators  $X_1, X_4$  and  $X_6$  on the factors  $F_1$  and  $F_2$ ? (3 points)
- (b) Compute the correlation between the indicators  $X_1$  and  $X_2$ . (4 points)
- (c) What percentage of the variance of indicators  $X_1$  and  $X_2$  is not accounted for by the common factors  $F_1$  and  $F_2$ ?

4. (12 points) Perform principal component analysis on the following data by hand. In other words, determine the angle (with best precision you can: 1%, 5%, 10%: your calculators should be of help) between the new axis and the old axis that would give a new variable, which accounts for the maximum variance in the data. What conclusions can you draw?

$X_1$	$X_2$
1	4
1	1
2	2
2	3
3	2
3	2
4	1
4	4

5. (15 points) Consider the following single-factor model

$$x_1 = \lambda_1 \xi + \delta_1$$

$$x_2 = \lambda_2 \xi + \delta_2$$

$$x_3 = \lambda_3 \xi + \delta_3$$

If the sample covariance matrix of the indicators is given by:

$$S = \begin{pmatrix} 1.20 & 0.93 & 0.45 \\ 0.93 & 1.56 & 0.27 \\ 0.45 & 0.27 & 2.15 \end{pmatrix}$$

Compute the estimates of the model parameters ( $\lambda_1, \lambda_2, \lambda_3, \text{Var}(\delta_1), \text{Var}(\delta_2), \text{Var}(\delta_3)$ ) using hand calculations. Are the parameter estimates unique?

6. (10 points)

(i) What problems cluster analysis is intended to solve? What types of cluster analysis you know: name them. What assumptions on data are imposed in order to apply "cluster analysis"? How strict you should be with those assumptions: speculate and exemplify "why?" to the best of your knowledge.

(ii) Cluster the following hypothetical data set into two groups using similarity matrixes. Intermediate similarity matrixes should be stated clearly. It is your choice to choose the type of distance used in calculations.

Hypothetical Data		
Subject ID	Income in EUR	Education (in years)
S1	5	5
S2	6	6
S3	15	14
S4	16	15
S5	25	20
S6	30	19

**Two-Factor Model:**  $x_1 = \lambda_{11}\xi_1 + \lambda_{12}\xi_2 + \varepsilon_1$

$$x_2 = \lambda_{21}\xi_1 + \lambda_{22}\xi_2 + \varepsilon_2$$

⋮

$$x_p = \lambda_{p1}\xi_1 + \lambda_{p2}\xi_2 + \varepsilon_p$$

The variance of  $x$ :  $E(x^2) = E(\lambda_1\xi_1 + \lambda_2\xi_2 + \varepsilon_1)^2$ ;  $Var(x) = \lambda_1^2 + \lambda_2^2 + Var(\varepsilon) + 2\lambda_1\lambda_2\phi$

The correlation between any indicator and any factor (the structure loading):

$$E(x\xi_1) = E[(\lambda_1\xi_1 + \lambda_2\xi_2 + \varepsilon_1)\xi_1]; \text{Corr}(x\xi_1) = \lambda_1 + \lambda_2\phi$$

The shared variance between the factor and an indicator: *Shared variance* =  $(\lambda_1 + \lambda_2\phi)^2$

The correlation between two indicators:

$$E(x_j x_k) = E[(\lambda_{j1}\xi_1 + \lambda_{j2}\xi_2 + \varepsilon_j)(\lambda_{k1}\xi_1 + \lambda_{k2}\xi_2 + \varepsilon_k)]$$

$$\text{Corr}(x_j x_k) = \lambda_{j1}\lambda_{k1} + \lambda_{j2}\lambda_{k2} + (\lambda_{j1}\lambda_{k2} + \lambda_{j2}\lambda_{k1})\phi$$

### Confirmatory Factor Analysis

The covariance matrix (one-factor model, two indicators):  $\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}$

Evaluating model fit:  $\chi^2$ -test  $H_0: \Sigma = \Sigma(\theta)$   $H_a: \Sigma \neq \Sigma(\theta)$  (test whether the difference between the sample and the estimated covariance matrix is a zero matrix)

$$\chi^2 = \sum_{i=1}^k \frac{|n_i - E(n_i)|^2}{E(n_i)}$$

### Cluster Analysis

Measure of similarity – squared Euclidean distance between two points

Hierarchical clustering:

**Centroid** method – each group is replaced by centroid

**Nearest-neighbor** or single-linkage method – the distance between two clusters is represented by the minimum of the distance between all possible pair of subjects in the two clusters

**Farthest-neighbor** or complete-linkage method - ... the maximum of the distances...

**Average-linkage** method - ... the average distance...

**Ward's** method – does not compute distances between clusters. Method tries to minimize the total within-group sums of squares.

### Discriminant Analysis

Assumptions: multivariate normality, equality of covariance matrices

Discriminant function:  $Z = w_1x_1 + w_2x_2$

$$\lambda = \frac{\text{between-group sum of squares}}{\text{within-group sum of squares}}$$

$\Sigma$ -variance-covariance matrix,  $T$ -total SSCP matrix.  $\gamma$ -vector of weights.

Discriminant function  $\xi = X' \gamma$ .  $B$  and  $W$  are between-groups and within-group SSCP matrices.

$$\text{Maximize } \lambda = \frac{\gamma' B \gamma}{\gamma' W \gamma}$$

$|W^{-1}B - \lambda I| = 0$ ;  $\gamma = \Sigma^{-1}(\mu_1 - \mu_2)$  - Fisher's discriminant function

### Logistic regression

$$\text{odds} = \frac{p}{1-p}$$

$$\ln \text{odds} = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}}$$

Maximum likelihood estimation:  $P(Y = 1) = p = \frac{e^{\beta X}}{1 + e^{\beta X}}$

$$L = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

**Quadratic equations:**  $ax^2 + bx + c = 0$ ;  $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

**Cubic equations:**

$$y^3 + ay^2 + by + c = 0; y = x - \frac{a}{3}; x^3 + px + q = 0; x_1 = \sqrt[3]{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} + \sqrt[3]{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}$$

## Formula Sheet, Multivariate Methods

### Matrices

Transpose – exchange rows and columns

Identity (I) – diag (1,1,...) of order n\*n

Inverse of A ( $A^{-1}$ ):  $AA^{-1} = A^{-1}A = I$

$A + B = B + A$ ;  $x(A + B) = xA + xB$ ;  $AB \neq BA$  (in general);

If order (A)=m\*n, order (B)=n\*p, then C=AB is of order m\*p

$$D = \det A = \begin{vmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{vmatrix}$$

$\det A = a_{i1}A_{i1} + a_{i2}A_{i2} + \dots + a_{in}A_{in}$  where cofactor  $A_{ij} = (-1)^{i+j} D_{ij}$  (i-row, j-column of D)

Cramer's rule:  $x_j = D_j / D$  where  $D = \det A$  and  $D_j$  is the determinant that arises when the j column of D is replaced by the column elements  $b_1, \dots, b_n$ . ( $Ax=b$ )

### Vectors

$$\mathbf{a} = (a_1 a_2 \dots a_p)$$

A right-angle triangle:  $\alpha$  - angle between a and c; c – hypotenuse;  $\cos \alpha = \frac{a}{c}$ ,  $\sin \alpha = \frac{b}{c}$

Length of vector  $\mathbf{a} = \|\mathbf{a}\| = \sqrt{a_1^2 + a_2^2}$

Basis vectors  $\mathbf{e}_1 = (1 \ 0)$ ,  $\mathbf{e}_2 = (0 \ 1)$

$$\mathbf{a} = a_1 \mathbf{e}_1 + a_2 \mathbf{e}_2$$

Scalar product  $\mathbf{ab} = a_1 b_1 + a_2 b_2 + \dots + a_p b_p$ ;  $\mathbf{ab} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \alpha$

Length of the projection:  $\|\mathbf{a}_p\| = \|\mathbf{a}\| \cos \alpha$

Variance of  $x_i$ :  $s_1^2 = \frac{\|x_i\|^2}{n-1}$ ; Generalized variance:  $GV = \left( \frac{\|x_1\| \|x_2\|}{n-1} \cdot \sin \alpha \right)^2$

### Distances

Euclidean:  $D_{AB} = \sqrt{\sum_{j=1}^p (a_j - b_j)^2}$

Statistical:  $SD_{ij}^2 = \left( \frac{x_i - x_j}{s} \right)^2$ , s-standard deviation

Manalanobis:  $MD_{ik}^2 = \frac{1}{1-r^2} \left[ \frac{(x_{i1} - x_{k1})^2}{s_1^2} + \frac{(x_{i2} - x_{k2})^2}{s_2^2} - \frac{2r(x_{i1} - x_{k1})(x_{i2} - x_{k2})}{s_1 s_2} \right]$

### Variance, Sum of Squares, and Cross Products

Variance:  $s_j^2 = \frac{\sum_{i=1}^n x_{ij}^2}{n-1} = \frac{SS}{df}$  (sum of squares/degrees of freedom)

Covariance:  $s_{jk} = \frac{\sum_{i=1}^n x_{ij} x_{ik}}{n-1} = \frac{SCP}{df}$  (sum of the cross products/degrees of freedom)

SSCP – sum of squares and cross products matrix  $\begin{pmatrix} SSX_1 & SCP \\ SCP & SSX_2 \end{pmatrix}$

S – covariance matrix  $S_t = \frac{SSCP_t}{df}$

**Within-Group Analysis:**  $SSCP_w = SSX_1 + SSX_2$  (pooled SSCP matrix)  $S_w = \frac{SSCP_w}{n_1 + n_2 - 2}$  (pooled cov m)

**Between-Group Analysis:**  $SS_j = \sum_{g=1}^G n_g (\bar{x}_{jg} - \bar{x}_j)^2$ ;  $SCP_{jk} = \sum_{g=1}^G n_g (\bar{x}_{jg} - \bar{x}_j)(\bar{x}_{kg} - \bar{x}_k)$

$SSCP_t = SSX_1 + SSX_2 + SS_j + SS_k$

### Principal Components Analysis

$x_1^* = \cos \theta * x_1 + \sin \theta * x_2$ ;  $x_2^* = -\sin \theta * x_1 + \cos \theta * x_2$

$\Sigma$  covariance matrix;  $\lambda$ -eigenvalues;  $|\Sigma - \lambda I| = 0$ ;  $\gamma$ -eigenvector;  $(\Sigma - \lambda I)\gamma = 0$ ;  $\gamma' \gamma = 1$ ;

### Factor Analysis

**Assumptions:** 1. Means of indicators, common factor, unique factors are zero.

2. Variances of indicators and common factors are one. 3.  $E(\xi_i \varepsilon_i) = 0$  and  $E(\varepsilon_i \varepsilon_j) = 0$

# MULTIVARIATE METHODS

## PCA

$$\begin{aligned} \xi_1 &= W_{11}X_1 + W_{12}X_2 + \dots + W_{1p}X_p \\ \xi_2 &= W_{21}X_1 + W_{22}X_2 + \dots + W_{2p}X_p \\ &\vdots \\ \xi_p &= W_{p1}X_1 + W_{p2}X_2 + \dots + W_{pp}X_p \end{aligned}$$

The weights:  $W_{i1}^2 + W_{i2}^2 + \dots + W_{ip}^2 = 1, i=1, \dots, p$   
 $W_{i1}W_{j1} + W_{i2}W_{j2} + \dots + W_{ip}W_{jp} = 0, \forall i \neq j$

Loadings = correlation between the original and the new variables

$$\lambda_{ij} = \frac{W_{ij}}{\sqrt{\lambda_i}}$$

Characteristic equation:  $\det(\lambda I - A) = 0$

## FA

$$\begin{aligned} X_1 &= \lambda_{11}\xi_1 + \lambda_{12}\xi_2 + \dots + \lambda_{1m}\xi_m + \epsilon_1 \\ X_2 &= \lambda_{21}\xi_1 + \lambda_{22}\xi_2 + \dots + \lambda_{2m}\xi_m + \epsilon_2 \\ &\vdots \\ X_p &= \lambda_{p1}\xi_1 + \lambda_{p2}\xi_2 + \dots + \lambda_{pm}\xi_m + \epsilon_p \end{aligned}$$

Assumptions:

- Means of indicators, common factors and unique factors are zero.
- Variations of indicators and common factors are one.
- The unique factors are not correlated among themselves or with the common factors.  $\Rightarrow E(\xi_i \epsilon_j) = E(\epsilon_i \epsilon_j) = 0, \forall i \neq j$

Model	Equations	The variance of any indicator $X_j$	Structure loading	Shared variance = (structure loading) <sup>2</sup>	Correlation between indicators
One-factor model	$X_1 = \lambda_1 \xi + \epsilon_1$ $\vdots$ $X_p = \lambda_p \xi + \epsilon_p$	$V(X_j) = \lambda_j^2 + V(\epsilon_j)$	$\text{Cor}(X_j, \xi) = \lambda_j$	$\lambda_j^2$	$\text{Cor}(X_j, X_k) = \lambda_j \lambda_k$
Two-factor model	$X_1 = \lambda_{11}\xi_1 + \lambda_{12}\xi_2 + \epsilon_1$ $\vdots$ $X_p = \lambda_{p1}\xi_1 + \lambda_{p2}\xi_2 + \epsilon_p$	$V(X_j) = \lambda_{j1}^2 + \lambda_{j2}^2 + 2\lambda_{j1}\lambda_{j2} + V(\epsilon_j)$	$\text{Cor}(X_1, X_2) = \lambda_{12} + \lambda_{21}\phi$ $\text{Cor}(X_j, \xi_i) = \lambda_{ji} + \lambda_{j2}\phi$	$(\lambda_{j2} + \lambda_{j1}\phi)^2$ $(\lambda_{j1} + \lambda_{j2}\phi)^2$	$\text{Cor}(X_j, X_k) = \lambda_{j1}\lambda_{k1} + \lambda_{j2}\lambda_{k2} + [\lambda_{j1}\lambda_{k2} + \lambda_{j2}\lambda_{k1}]\phi$

$$\text{Cor}(\xi_1, \xi_2) = \phi$$

## GFA

Underidentified model: # equations < # variables  
 Just-identified model: # equations = # variables  
 Overidentified model: # equations > # variables

$$P(X_1, X_2, X_3) = P(X_1) \cdot P(X_2|X_1) \cdot P(X_3|X_1, X_2)$$

## CA

$$D_{AB}^2 = \sum_{j=1}^p (a_j - b_j)^2 \quad SD_{X_k}^2 = \sum_{j=1}^p \left( \frac{X_{ij} - X_{kj}}{S_j} \right)^2 \quad MD_{ik}^2 = \frac{1}{1-r^2} \left[ \frac{(X_{i1} - X_{k1})^2}{S_1^2} + \frac{(X_{i2} - X_{k2})^2}{S_2^2} - \frac{2r(X_{i1} - X_{k1})(X_{i2} - X_{k2})}{S_1 S_2} \right]$$

(two groups)

$$\lambda = \frac{SS_B}{SS_W} \rightarrow \max$$

$$\text{If } \xi = \bar{X}^T \bar{Y} \Rightarrow \text{The estimation of } \bar{Y}: \bar{Y}^T = (\bar{\mu}_1 - \bar{\mu}_2)^T \Sigma^{-1}$$

$$SSCP_W = SSCP_1 + SSCP_2$$

$$SSCP_T = SSCP_W + SSCP_B$$

## LOG-REG

$$\text{odds} = \frac{p}{1-p}$$

$$\ln \frac{p}{1-p} = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

$$p = \frac{\text{odds}}{1 + \text{odds}}$$

$$P(Y=1) = p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}}$$

$$\text{If } Y = \beta_0 + \beta_1 X:$$

$$\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Rotation  $\Rightarrow$   
 $a_1^* = \cos \theta a_1 + \sin \theta a_2$   
 $a_2^* = -\sin \theta a_1 + \cos \theta a_2$

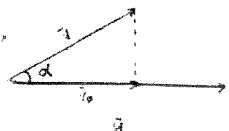
$$S^2 = \frac{SS}{df} = \frac{\sum (X_i - \bar{X})^2}{n-1} = \frac{\sum X_i^2 - n\bar{X}^2}{n-1}$$

$$S_{xy} = \frac{SCP}{df} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n-1} = \frac{\sum X_i Y_i - n\bar{X}\bar{Y}}{n-1}$$

Projection: of  $\vec{v}$  onto  $\vec{u}$

The projection vector:  $\vec{v}_p = \frac{\|\vec{v}_p\|}{\|\vec{u}\|} \vec{u}$

$$\|\vec{v}_p\| = \|\vec{v}\| \cos \alpha = \frac{\vec{v} \cdot \vec{u}}{\|\vec{u}\|}$$



Direction cosines = The cosines of the angle between a vector and the axes



Stockholms  
universitet

Statistiska institutionen

## Rättningsblad

**Datum:** 23/10 -2014

**Sal:** Ugglevikssalen

**Tenta:** Multivariata metoder

**Kurs:** Multivariata metoder

**ANONYMKOD:**

MM-0022

Jag godkänner att min tenta får läggas ut anonymt på hemsidan som studentsvar.

**OBS! SKRIV ÄVEN PÅ BAKSIDAN AV SKRIVBLADEN**

Markera besvarade uppgifter med kryss

1	2	3	4	5	6	7	8	9	Antal inl. blad
X	X	X	X	X	X				6 R
Lär.ant.									
10	18	15	12	14	10				

POÄNG

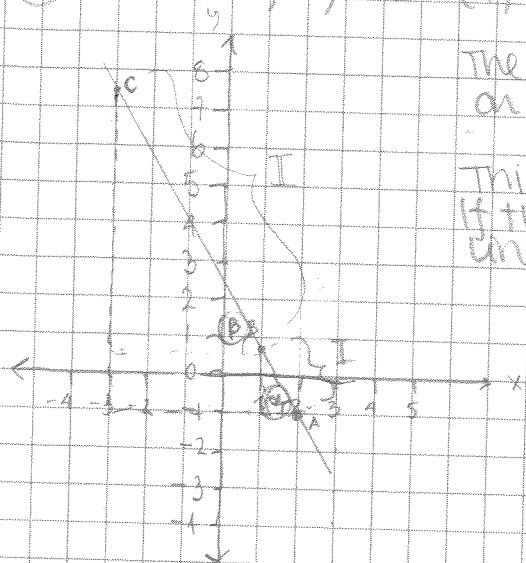
79  
AA+

BETYG

Lärarens sign.

①

$A = (2, -1)$     $B = (1, 0.7)$     $C = (-3, 7.5)$



The graphic solution shows that they are on a straight line.

This can also be proved analytically: If the points are in a straight line, then  $\cos(\alpha) = \cos(\beta)$ .

$$\cos(\beta) = \frac{3+1}{II}$$

$$II = D_{BC} = \sqrt{4^2 + 6.8^2} = 7.8892$$

$$\cos(\beta) = \frac{4}{7.8892} = 0.5070$$

$$\cos(\alpha) = \frac{5}{I+II}$$

$$I+II = \sqrt{5^2 + 8.5^2} = 9.8615$$

$$\cos \alpha = \frac{5}{9.8615} = 0.5070$$

$$\cos(\alpha) = \cos(\beta)!$$

We can also show this by calculating the linear equation:

$$-1 = 2x + m \quad 0.7 = x + m \rightarrow x = 0.7 - m$$

$$-1 = 2(0.7 - m) + m = 1.4 - 2m + m \rightarrow m = 2.4 \quad x = -1.7$$

BY CALCULATING Y FROM THE LINEAR EQUATION:

$y = -1.7x + m$	$\rightarrow$	$x_A = 2$	$\rightarrow$	$y_A = -1$	$\checkmark$	THEY ARE ON A STRAIGHT LINE!
		$x_B = 1$	$\rightarrow$	$y_B = 0.7$	$\checkmark$	
		$x_C = -3$	$\rightarrow$	$y_C = 7.5$	$\checkmark$	



① cont.

$$e_1 = 0,804f_1 + 0,629f_2$$

$$e_2 = 0,273f_1 - 0,588f_2$$

$$\begin{cases} A = 2e_1 - e_2 \\ B = e_1 + 0,7e_2 \\ C = -3e_1 + 7,5e_2 \end{cases}$$

$A = f(f_1, f_2)$ ?

$$A = 2e_1 - e_2 = 2(0,804f_1 + 0,629f_2) - (0,273f_1 - 0,588f_2)$$

$$A = 1,608f_1 + 1,258f_2 - 0,273f_1 + 0,588f_2$$

$$A = 1,335f_1 + 1,846f_2 \quad A = (1,335, 1,846)$$

$B = f(f_1, f_2)$ ?

$$B = e_1 + 0,7e_2 = 0,804f_1 + 0,629f_2 + 0,7(0,273f_1 - 0,588f_2)$$

$$B = 0,804f_1 + 0,629f_2 + 0,1911f_1 - 0,4116f_2$$

$$B = 0,9951f_1 + 0,2174f_2 \quad B = (0,9951, 0,2174)$$

$C = f(f_1, f_2)$ ?

$$C = -3e_1 + 7,5e_2 = -3(0,804f_1 + 0,629f_2) + 7,5(0,273f_1 - 0,588f_2)$$

$$C = -2,412f_1 - 1,887f_2 + 2,0475f_1 - 4,41f_2$$

$$C = -0,3645f_1 - 6,297f_2 \quad C = (-0,3645, -6,297)$$

To test if these are on a linear line I've chosen the method of finding our linear regression to see if it is suitable for A, B and C.

$$1,846 = 1,335x + m \quad 0,2174 = 0,9951x + m$$

$$m = 0,2174 - 0,9951x \rightarrow 1,846 = 1,335x + 0,2174 - 0,9951x$$

$$1,6286 = 0,3399x$$

$$x = 4,7914$$

$$m = -4,5505$$

$$y = 4,7914x - 4,5505$$

BY CALCULATING Y FROM EQUATION  $\rightarrow$

$$A: x_A = 1,335 \rightarrow y_A = 1,846 \quad \checkmark$$

$$B: x_B = 0,9951 \rightarrow y_B = 0,2174 \quad \checkmark$$

$$C: x_C = -0,3645 \rightarrow y_C = -6,297 \quad \checkmark$$

(ON A)  
THEY ARE LINEAR LINE!

We know that this is true since the axis have changed from  $e_1$  and  $e_2$  to  $f_1$  and  $f_2$ , it does not mean that the relationship between A, B and C has changed. They are still on a straight line.

OK

(2)

Blood cholesterol (mg/100cc)	Heart Disease		
	Present	Absent	
<200	6	5	= 11
200-219	10	6	= 16
220-259	30	5	= 35
>259	45	7	= 52
	91	23	114

a) i)  $P(\text{present}) = \frac{91}{114} = 0.7982$

ii)  $P(BC < 219) = \frac{16 + 11}{114} = 0.2368$

b) i)  $P(\text{Present} | BC < 200) = \frac{6}{11}$

odds(Present |  $BC < 200$ ) =  $\frac{\frac{6}{11}}{1 - \frac{6}{11}} = 1.20$

ii)  $P(BC > 219 | \text{Present}) = \frac{75}{91}$

odds( $BC > 219$  | Present) =  $\frac{\frac{75}{91}}{1 - \frac{75}{91}} = 4.6875$

c) i)  $P(\text{Present} | BC > 259) = \frac{45}{52}$

odds(Present |  $BC > 259$ ) =  $\frac{\frac{45}{52}}{1 - \frac{45}{52}} = 6.429$

$\ln(\text{odds}) = \ln(6.429) \approx 1.861$

ii)  $P(\text{present} | 200 < BC < 219) = \frac{10}{16}$

odds(present |  $200 < BC < 219$ ) =  $\frac{\frac{10}{16}}{1 - \frac{10}{16}} = 1.667$

$\ln(\text{odds}) = \ln(1.667) = 0.511$

3

$$\begin{aligned} X_1 &= 0,104F_1 + 0,824F_2 + U_1 \\ X_2 &= 0,065F_1 + 0,959F_2 + U_2 \\ X_3 &= 0,065F_1 + 0,725F_2 + U_3 \\ X_4 &= 0,906F_1 + 0,134F_2 + U_4 \\ X_5 &= 0,977F_1 + 0,116F_2 + U_5 \\ X_6 &= 0,827F_1 + 0,016F_2 + U_6 \end{aligned}$$

USUAL ASSUMPTIONS:

- 1) Mean indicator, common factor, unique factor = 0
- 2) Variance indicator, common factor = 1
- 3)  $E(X, U) = 0$
- 4)  $E(U_i, U_j) = 0$

$$\text{corr}(F_1, F_2) = \rho_{12} = 0,3$$

a) Pattern loadings:

	$F_1$	$F_2$
$X_1$	$\lambda_{11} = 0,104$	$\lambda_{12} = 0,824$
$X_4$	$\lambda_{41} = 0,906$	$\lambda_{42} = 0,134$
$X_6$	$\lambda_{61} = 0,827$	$\lambda_{62} = 0,016$

b)  $\text{corr}(X_1, X_2) = \lambda_{11} \cdot \lambda_{21} + \lambda_{12} \cdot \lambda_{22} + (\lambda_{11} \cdot \lambda_{22} + \lambda_{12} \cdot \lambda_{21}) \rho_{12}$   
 $= 0,104 \cdot 0,065 + 0,824 \cdot 0,959 + (0,104 \cdot 0,959 + 0,065 \cdot 0,824) \cdot 0,3$   
 $\text{corr}(X_1, X_2) = 0,843$  or

The correlation between the 2 indicators  $X_1$  and  $X_2$  is 0,8430.

c)  $\text{Var}(X) = \lambda^2 + \lambda^2 + \text{Var}(U) + 2\lambda_1\lambda_2\rho$

The variance in  $X$  that is not accounted for by the common factors  $F_1$  and  $F_2$

$\cdot \text{Var}(X_1) = 0,104^2 + 0,824^2 + \text{Var}(U_1) + 2 \cdot 0,104 \cdot 0,824 \cdot 0,3$

$\text{Var}(X_1) = 0,7412 + \text{Var}(U_1)$  !  $\text{Var}(X) = 1$  from assumption 2)!

% of the variance not accounted for by  $F_1$  and  $F_2$  =

$\text{Var}(U_1) = 1 - 0,7412 = 0,2588 = \underline{25,9\%}$

$\cdot \text{Var}(X_2) = 0,065^2 + 0,959^2 + \text{Var}(U_2) + 2 \cdot 0,065 \cdot 0,959 \cdot 0,3$

$\text{Var}(X_2) = 0,9613 + \text{Var}(U_2)$

% of the variance not accounted for by  $F_1$  and  $F_2$

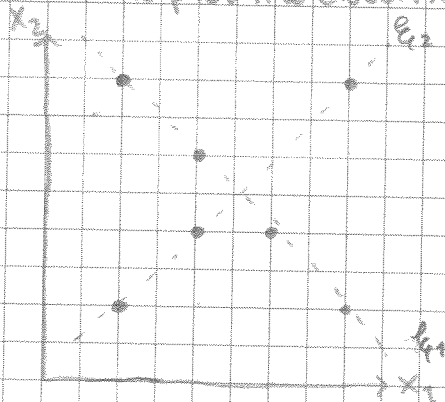
$\text{Var}(U_2) = 1 - 0,9613 = 0,0387 = \underline{3,87\%}$

OK

④

	$X_1$	$X_2$
1	1	4
2	1	1
3	2	2
4	2	3
5	3	2
6	3	2
7	4	1
8	4	4
MEAN:	2,5	2,375

First, I want to plot the observations:



just a first glimpse of what the answer should be.

I first want to create an SSCP-matrix:

$$SS_1^2 = \sum_{i=1}^8 (X_{i1} - \bar{X}_1)^2 = \sum_{i=1}^8 (X_{i1} - 2,5)^2 = 10$$

$$SS_2^2 = \sum_{i=1}^8 (X_{i2} - \bar{X}_2)^2 = \sum_{i=1}^8 (X_{i2} - 2,375)^2 = 9,875$$

$$SCP_{12} = \sum_{i=1}^8 (X_{i1} - \bar{X}_1)(X_{i2} - \bar{X}_2) = -0,5$$

$$SSCP_{12} = \begin{pmatrix} 10 & -0,5 \\ -0,5 & 9,875 \end{pmatrix}$$

COVARIANCE MATRIX

$$\Sigma = \frac{SSCP}{df} = \begin{pmatrix} 1,43 & -0,07 \\ -0,07 & 1,41 \end{pmatrix} \quad n-1=7$$

Want to calculate the  $\lambda_i$  (eigenvalue). But how?

$$\xi = \gamma' X \quad \Sigma = E(XX')$$

In principal component analysis we want to maximize the variance.

$$\text{var}(\xi) = E(\xi\xi') = E(\gamma\gamma' E(XX')) = \gamma' \Sigma \gamma$$

$$\text{MAX. } \gamma' \Sigma \gamma \quad \text{s.t. } \gamma' \gamma = 1 \rightarrow z = \gamma' \Sigma \gamma - \lambda(\gamma' \gamma - 1)$$

$$\frac{\partial z}{\partial \gamma} = 2 \Sigma \gamma - 2 \lambda \gamma = 0$$

$$\gamma(\Sigma - \lambda I) = 0$$

$$|\Sigma - \lambda I| = 0 \quad \text{The characteristic equation!}$$



④  $|\Sigma - \lambda I| = 0$

$$\left| \begin{pmatrix} 1,43 & -0,07 \\ -0,07 & 1,41 \end{pmatrix} - \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} \right| = 0$$

$$\begin{vmatrix} 1,43 - \lambda & -0,07 \\ -0,07 & 1,41 - \lambda \end{vmatrix} = (1,43 - \lambda)(1,41 - \lambda) + 0,0049 = 0$$

$$2,0163 - 1,43 \lambda - 1,41 \lambda + \lambda^2 + 0,0049 = 0$$

$$\lambda^2 - 2,84 \lambda + 2,0114 = 0$$

$$\lambda_{1,2} = \left( \frac{2,84}{2} \right) \pm \sqrt{\left( \frac{2,84}{2} \right)^2 - 2,0114} = \left( \frac{2,84}{2} \right) \pm 0,0707$$

$$\lambda_1 = 1,49 \quad \lambda_2 = 1,35 = \text{THE EIGENVALUES.}$$

We want to find the eigenvector for the PC that accounts for the maximum variance, i.e.  $\lambda_1 = 1,49$ .

$$(\Sigma - 1,49 I) \mathbf{x} = 0$$

$$\begin{pmatrix} 1,43 & -0,07 \\ -0,07 & 1,41 \end{pmatrix} - \begin{pmatrix} 1,49 & 0 \\ 0 & 1,49 \end{pmatrix} \mathbf{x} = 0$$

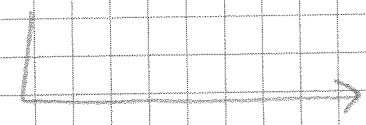
$$\begin{pmatrix} -0,06 & -0,07 \\ -0,07 & -0,08 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$-0,06 x_1 - 0,07 x_2 = 0 \rightarrow x_1 = \frac{-0,07}{0,06} x_2 \rightarrow \underline{x_1 = 1 \quad x_2 = -0,8571}$$

$$-0,07 x_1 - 0,08 x_2 = 0 \quad \mathbf{e}_1 = (1, -0,8571)$$

The eigenvector for the principal component that accounts for the maximum variance is  $\mathbf{e}_1 = (1; 0,8571)$

So, what is the angle between the new axis and the old axis?



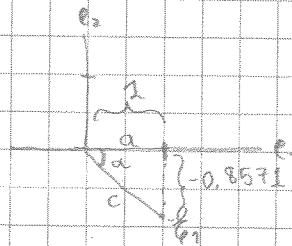


④.  $e_1 = (1, -0.8571)$

Angle  $\alpha$ ?

$$\cos \alpha = \frac{a}{c} = \frac{1}{\sqrt{1^2 + 0.8571^2}} = 0.75927$$

$$\alpha = 40.6^\circ$$



The angle between the new axis and the old axis is  $40.6^\circ$

This means that the new values will be:

$$\begin{aligned} X_1^* &= \cos \alpha \cdot X_1 + \sin \alpha \cdot X_2 \\ X_2^* &= -\sin \alpha \cdot X_1 + \cos \alpha \cdot X_2 \end{aligned}$$

	$X_1^*$	$X_2^*$
1	3.36	2.39
2	1.41	0.11
3	2.82	0.22
4	3.47	0.98
5	3.58	-0.43
6	3.58	-0.43
7	3.69	-1.84
8	5.64	0.43

see covariance matrix

Total variance in the dataset =  $1.43 + 1.41 = 2.84$

The new axis explains  $(1.49 / 2.84) \times 100 = 52.5\%$  of the total variance.

excelbut

5

$$\begin{aligned} X_1 &= \lambda_1 \xi + \delta_1 \\ X_2 &= \lambda_2 \xi + \delta_2 \\ X_3 &= \lambda_3 \xi + \delta_3 \end{aligned}$$

$$S = \begin{pmatrix} 1,2 & 0,93 & 0,45 \\ 0,93 & 1,56 & 0,27 \\ 0,45 & 0,27 & 2,15 \end{pmatrix}$$

$$\begin{aligned} \lambda_1^2 + \text{Var}(\delta_1) &= 1,2 & \lambda_2 \lambda_1 &= \lambda_1 \lambda_2 = 0,93 \\ \lambda_2^2 + \text{Var}(\delta_2) &= 1,56 & \lambda_3 \lambda_2 &= \lambda_2 \lambda_3 = 0,27 \\ \lambda_3^2 + \text{Var}(\delta_3) &= 2,15 & \lambda_3 \lambda_1 &= \lambda_1 \lambda_3 = 0,45 \end{aligned}$$

6 equations  
and  
6 unknowns  
↓  
JUST-IDENTIFIED! OK

We know the equations above since the covariance matrix is:

$$S = \begin{pmatrix} \frac{\lambda_1^2 + \text{Var}(\delta_1)}{\text{Var}(X_1)} & \frac{\lambda_1 \lambda_2}{\text{Cov}(X_1, X_2)} & \frac{\lambda_1 \lambda_3}{\text{Cov}(X_1, X_3)} \\ \frac{\lambda_2 \lambda_1}{\text{Cov}(X_2, X_1)} & \frac{\lambda_2^2 + \text{Var}(\delta_2)}{\text{Var}(X_2)} & \frac{\lambda_2 \lambda_3}{\text{Cov}(X_2, X_3)} \\ \frac{\lambda_3 \lambda_1}{\text{Cov}(X_3, X_1)} & \frac{\lambda_3 \lambda_2}{\text{Cov}(X_3, X_2)} & \frac{\lambda_3^2 + \text{Var}(\delta_3)}{\text{Var}(X_3)} \end{pmatrix}$$

where  $\lambda_i \lambda_j$  are covariances ( $\text{Cov}(i, j)$ ) and  $\lambda_i^2 + \text{Var}(\delta_i)$  is the variance of indicator  $j$  or ( $\text{Cov}(j, j)$ ).

$$\left. \begin{aligned} \lambda_1 &= \frac{0,93}{\lambda_2} \\ \lambda_1 &= \frac{0,45}{\lambda_3} \end{aligned} \right\} \frac{0,93}{\lambda_2} = \frac{0,45}{\lambda_3} \rightarrow 0,45 \lambda_2 = 0,93 \lambda_3 \rightarrow \lambda_2 = \frac{0,93 \lambda_3}{0,45}$$

$$\lambda_2 \lambda_3 = 0,27 \rightarrow \left( \frac{0,93 \lambda_3}{0,45} \right) \lambda_3 = 0,27 \rightarrow \frac{0,93}{0,45} \cdot \lambda_3^2 = 0,27 \rightarrow \lambda_3^2 = \frac{81}{620} \approx 0,1306$$

$$\lambda_3 \approx 0,3614 \quad \lambda_2 = \frac{0,93 \cdot 0,3614}{0,45} = 0,7470 \quad \lambda_1 = 1,2450 \quad \text{OK}$$

$$\begin{aligned} 1,2450^2 + \text{Var}(\delta_1) &= 1,2 & 0,7470^2 + \text{Var}(\delta_2) &= 1,56 & 0,3614^2 + \text{Var}(\delta_3) &= 2,15 \\ \text{Var}(\delta_1) &= -0,3500 & \text{Var}(\delta_2) &= 1,002 & \text{Var}(\delta_3) &= 2,019 \end{aligned}$$

→ how it can be negative?

$$\begin{aligned} \lambda_1 &= 1,245 & \text{Var}(\delta_1) &= -0,350 \\ \lambda_2 &= 0,7470 & \text{Var}(\delta_2) &= 1,002 \\ \lambda_3 &= 0,3614 & \text{Var}(\delta_3) &= 2,019 \end{aligned}$$

The solutions are unique since the equation is just-identified.  
(estimates of the parameters)

- 6) i) Cluster analysis can be used when we want to reduce data. It is done by clustering observations that are closest to each other. CA can also be used if we want to see which observations are the most similar to each other by grouping them.

There is a hierarchical analysis which is what will be used in the calculations below. The clusters are computed by either centroid method, nearest-neighbor method, furthest-neighbor method or Ward's method. There's also nonhierarchical analysis where you know the amounts of clusters a priori.

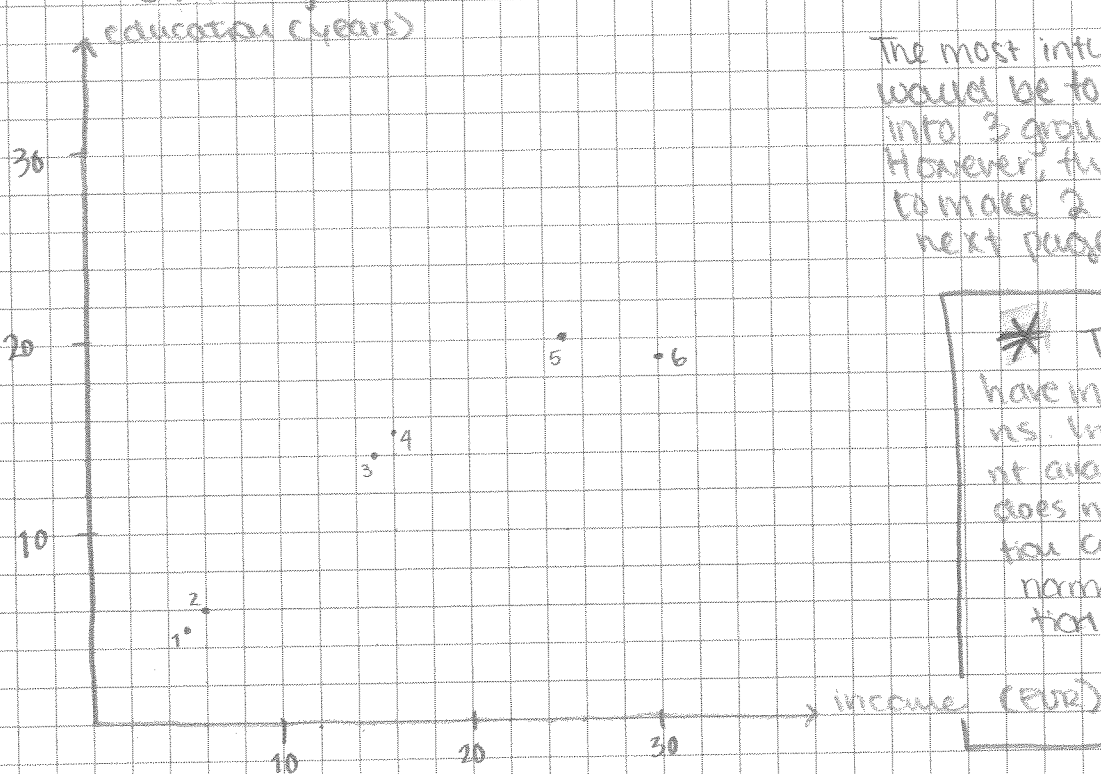
Assumptions for cluster analysis: \*

- 1) Each group is homogeneous with respect to a specific characteristic
- 2) Each group/cluster is different from the other groups with respect to the same characteristic.

These assumptions should be very strict in my opinion. When clustering observations there should not be other factors that combine them since that would mean that one observation might be misplaced in the wrong cluster due to this unobservable.

If we are not strict, then the analysis fails since we can't say that the clusters differ with respect to the observed variables. This would give weak conclusions and sample tests.

- (ii) I first want to plot the data to get a view of what answers I should get:



The most intuitive clustering would be to cluster them into 3 groups (1+2, 3+4, 5+6). However, the assignment is to make 2 groups. See next page →

\* The data must also have independent observations. In contrast to discriminant analysis, cluster analysis does not take misclassification costs, prior probability or normal distributional assumption into consideration.



(6.)

The distances are calculated by the squared euclidean distance:

$$D_{AB}^2 = \sum_{j=1}^p (a_j - b_j)^2$$

Similarity Matrix I

	S1	S2	S3	S4	S5	S6
S1	0	•	•	•	•	•
S2	2	0	•	•	•	•
S3	181	145	0	•	•	•
S4	221	181	2	0	•	•
S5	625	557	136	106	0	•
S6	821	745	250	212	26	0

The smallest distance is a tie between S1S2 and S3S4 where the distance is 2. These two observations will be clustered by the centroid method  $\Rightarrow$  the coordinates of the cluster = the average of the 2 clustered observations.

cluster S12: Income =  $\frac{(6+5)}{2} = 5,5$     Education =  $\frac{(16+5)}{2} = 5,5$

S34: Income =  $\frac{(16+15)}{2} = 15,5$     Education =  $\frac{(15+14)}{2} = 14,5$

New table

	Income	Educ
S12	5,5	5,5
S34	15,5	14,5
S5	25	20
S6	30	19

Similarity matrix II

	S12	S34	S5	S6
S12	0	•	•	•
S34	181	0	•	•
S5	590,5	120,5	0	•
S6	782,5	230,5	26	0

The smallest distance now is 26 so S5 and S6 will be clustered.

S56: Income:  $\frac{(30+25)}{2} = 27,5$     Education:  $\frac{(19+20)}{2} = 19,5$



(6)

New table:

	Income	Education
S12	5.5	5.5
S34	15.5	14.5
S56	29.5	19.5

Similarity Matrix

III

	S12	S34	S56
S12	0	.	.
S34	181	0	.
S56	688	169	0

The smallest distance now is between S34 and S56.

Now we have the 2 clusters:

cluster 1 = S1 and S2

cluster 2 = S3, S4, S5, S6.

OK



Stockholms  
universitet

Statistiska institutionen

## Rättningsblad

**Datum:** 23/10 -2014

**Sal:** Ugglevikssalen

**Tenta:** Multivariata metoder

**Kurs:** Multivariata metoder

**ANONYMKOD:**

MM-0016

Jag godkänner att min tenta får läggas ut anonymt på hemsidan som studentsvar.

**OBS! SKRIV ÄVEN PÅ BAKSIDAN AV SKRIVBLADEN**

Markera besvarade uppgifter med kryss

	1	2	3	4	5	6	7	8	9	Antal inl. blad
	X	X	X	X	X	X				15
Lär.ant.	7	18	15	12	14	9				

POÄNG

75  
AA+

BETYG

Lärarens sign.

Question 1:

$$a) A = (2; 1) \quad B = (1; 0,7) \quad C = (-3; 7.5)$$

$$\cos(\theta_{AB}) = \frac{a \cdot b}{|a| \cdot |b|} = \frac{(2 \cdot 1 + 0,7 \cdot 1)}{\sqrt{2^2 + 2^2} \sqrt{1^2 + 0,7^2}} = 1$$

$$\theta_{AB} = \cos^{-1}(1) = 0$$

Thus, A and B lie  
on a straight line,

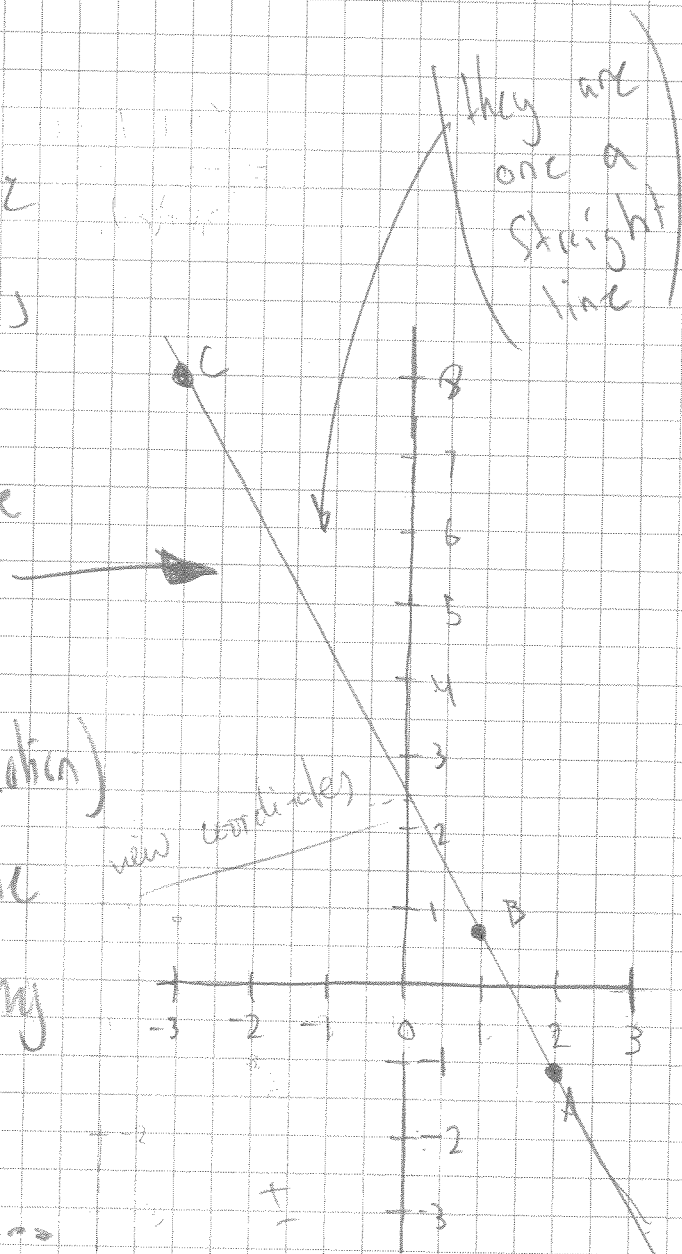
there is no angle.

From graph we see

that they do.

b) if the distance (Euclidean)  
do not change with the  
new coordinates then they  
lie on a straight line!

No more time to prove...



Question 2:

Blood (h)	Heart disease		$\Sigma$
	P	A	
< 200	6	5	11
200-219	10	6	16
220-259	30	5	35
> 259	45	7	52
$\Sigma$	91	23	114

a) What is the probability that:

$$(i) P(P) = \frac{91}{114} \approx 0,7982 = 79,82\%$$

$$(ii) P(< 219) = \frac{(11+16)}{114} \approx 0,2368 = 23,68\%$$

b) What are the odds that:

$$\boxed{\text{odds} = \frac{P}{1-P}} \quad \text{or}$$

$$(i) \text{odds}(P | < 200) = \frac{6}{5} = 1,2$$

$$(ii) \text{odds}(> 219 | P) = \frac{(30+45)}{(6+10)} = 4,6875 \quad \text{or}$$

c) Compute the log odds:

$$(i) \ln[\text{odds}(P | > 259)] = \ln\left(\frac{45}{7}\right) \approx 1,8668$$

$$(ii) \ln[\text{odds}(P | 200-219)] = \ln\left(\frac{10}{6}\right) \approx 0,5108 \quad \text{or}$$

Question 3:6-indikator 2-factor model!

$$X_1 = 0,104 \times F_1 + 0,824 \times F_2 + U_1$$

$$X_2 = 0,065 \times F_1 + 0,959 \times F_2 + U_2$$

$$X_3 = 0,065 \times F_1 + 0,725 \times F_2 + U_3$$

$$X_4 = 0,906 \times F_1 + 0,134 \times F_2 + U_4$$

$$X_5 = 0,977 \times F_1 + 0,116 \times F_2 + U_5$$

$$X_6 = 0,827 \times F_1 + 0,016 \times F_2 + U_6$$

$$\text{Corr}(F_1, F_2) = \phi_{12} = 0,3$$

g) What are the pattern loadings of indicators  $X_1, X_4$  and  $X_6$  on the factors  $F_1$  and  $F_2$ .

The pattern loading ( $\lambda_i$ ) for a common factor ( $F_i$ ) is simply the coefficient for that common factor.

Pattern loading of indicator ( $i$ ) and common factor ( $j$ ) =  $\lambda_{ij}$

$$X_1: \lambda_{1,1} = 0,104$$

$$\lambda_{1,2} = 0,824$$

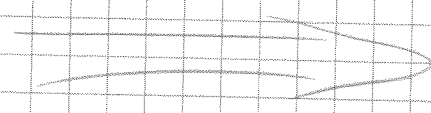
$$X_4: \lambda_{4,1} = 0,906$$

$$\lambda_{4,2} = 0,134$$

$$X_6: \lambda_{6,1} = 0,827$$

$$\lambda_{6,2} = 0,016$$

or

Next page 

b) Compute the correlation between  $X_1$  and  $X_2$ .

$$\begin{aligned} \text{Corr}(X_1, X_2) &= \lambda_{1j1} \lambda_{2j1} + \lambda_{1j2} \lambda_{2j2} + [\lambda_{1j1} \lambda_{2j2} + \lambda_{1j2} \lambda_{2j1}] \phi_{12} = \\ &= 0,104 \times 0,065 + 0,824 \times 0,959 + [0,104 \times 0,959 + 0,824 \times 0,065] \times 0,3 \\ &\approx \underline{\underline{0,8430}} \quad \text{or} \end{aligned}$$

c) What percentage of the variance of indicators  $X_1$  and  $X_2$  is not accounted for by the common factor  $F_1$  and  $F_2$ ?

$$\text{Var}(X_i) = \underbrace{\lambda_{ij1}^2 + \lambda_{ij2}^2 + 2(\lambda_{ij1})(\lambda_{ij2})\phi_{12}}_{\text{Variance from factors}} + \underbrace{\text{Var}(U_i)}_{\text{variance of unique factor}}$$

$\Leftrightarrow$

$$\text{Var}(U_i) = \text{Var}(X_i) - [\lambda_{ij1}^2 + \lambda_{ij2}^2 + 2(\lambda_{ij1})(\lambda_{ij2})\phi_{12}]$$

$\therefore$   $\uparrow$  is equal to 1 by assumption 2.

$$\Leftrightarrow \text{Var}(U_i) = 1 - [\lambda_{ij1}^2 + \lambda_{ij2}^2 + 2(\lambda_{ij1})(\lambda_{ij2})\phi_{12}]$$

$$X_1: \text{Var}(U_1) = 1 - [(0,104)^2 + (0,824)^2 + 2(0,104)(0,824)(0,3)] \approx 0,2588 = \underline{\underline{25,88\%}}$$

$$X_2: \text{Var}(U_2) = 1 - [(0,065)^2 + (0,959)^2 + 2(0,065)(0,959)(0,3)] \approx 0,0387 = \underline{\underline{3,87\%}}$$



Question 4

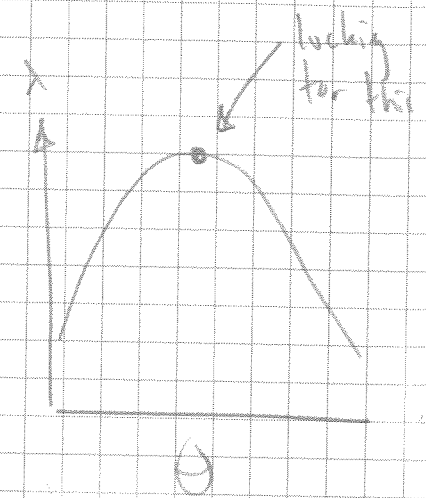
Finding the angle  $\theta$ , so that the first principal component (new variable) accounts for the maximum variance of the total variance of the data can be obtained by the equation;

$$\textcircled{1} \quad X_1^* = \cos(\theta) \times X_1 + \sin(\theta) \times X_2$$

Then one changes the angle  $\theta$  until a maximum of  $\lambda$  is reached,

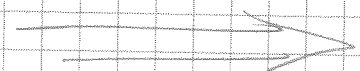
$$\text{Where } \lambda = \frac{\sigma_{X_1^*}^2}{\sigma_{X_1}^2 + \sigma_{X_2}^2} = \text{The problem}$$

here is to do all the computations.



I will instead compute the covariance matrix, and compute the eigenvalues and eigen vectors to get the  $\theta$  directly and confirm this point is the maximum by using equation  $\textcircled{1}$  two times for angles  $\theta$  close to  $\theta_{\max}$ .

a) Realize that covariance matrix is sufficient due to the fact that the data is in the same units.





obs	$x_1$	$x_2$
1	1	4
2	1	1
3	2	2
4	2	3
5	3	2
6	3	2
7	4	1
8	4	1

$$\bar{x}_i = \frac{1}{n} \sum_{k=1}^n x_{ik}$$

$$\sigma_i^2 = \frac{1}{(n-1)} \sum_{k=1}^n (x_{ik} - \bar{x}_i)^2$$

$$\sigma_{ij} = \frac{1}{(n-1)} \sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$$

$$\bar{x}_1 = \frac{1}{8} [2 \times 1 + 2 \times 2 + 2 \times 3 + 2 \times 4] = 2,5$$

$$\sigma_1^2 = \frac{1}{7} [2(1-2,5)^2 + 2(2-2,5)^2 + 2(3-2,5)^2 + 2(4-2,5)^2] \approx 1,43$$

$$\bar{x}_2 = \frac{1}{8} [2 \times 4 + 2 \times 1 + 3 \times 2 + 3] = 2,375$$

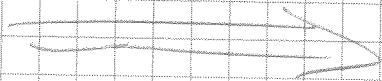
$$\sigma_2^2 = \frac{1}{7} [2(4-2,375)^2 + 2(1-2,375)^2 + 3(2-2,375)^2 + (3-2,375)^2] = 1,41$$

$$\sigma_{12} = \frac{1}{7} [(1-2,5)(4-2,375) + (1-2,5)(1-2,375) + (2-2,5)(2-2,375) + (2-2,5)(3-2,375) + (3-2,5)(2-2,375) + (3-2,5)(2-2,375) + (4-2,5)(1-2,375) + (4-2,5)(4-2,375)] = -0,07$$

$$\Sigma = \begin{pmatrix} 1,43 & -0,07 \\ -0,07 & 1,41 \end{pmatrix}$$

• Variances are similar which justify the use of the covariance matrix in PCA.

Really low covariance and thereby correlation between  $x_1$  and  $x_2$ .



Find eigenvalues  $\lambda_i$ :

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

$$\left| \begin{pmatrix} 1,43 & -0,07 \\ -0,07 & 1,41 \end{pmatrix} - \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} \right| = 0$$

$$\left| \begin{array}{cc} 1,43 - \lambda & -0,07 \\ -0,07 & 1,41 - \lambda \end{array} \right| = 0$$

$$(1,43 - \lambda)(1,41 - \lambda) - (-0,07)^2 = 0$$

$$\lambda^2 - 2,84\lambda + 2,0114 = 0$$

$$\left( \lambda - \frac{2,84}{2} \right)^2 = \left( \frac{2,84}{2} \right)^2 - 2,0114$$

$$(\lambda - 1,42)^2 = 0,005$$

$$\lambda = 1,42 \pm \sqrt{0,005}$$

$$\boxed{\begin{array}{l} \lambda_1 \approx 1,419 \\ \lambda_2 \approx 1,35 \end{array}}$$

Principal component 1 accounts for

$$\frac{1,419}{(1,419 + 1,35)} = 0,5246 \quad \underline{\underline{52,46\%}}$$

of the total variance in the data.  $\rightarrow$

⚠ This is really low, not good for PCA analysis!  $\rightarrow$

Find the eigen vector for  $\lambda_1$

$$(\Sigma - \lambda_1 I) \vec{w}_1 = \vec{0}$$

$$\begin{bmatrix} 1,43 & -0,07 \\ -0,07 & 1,41 \end{bmatrix} - \begin{bmatrix} 1,49 & 0 \\ 0 & 1,49 \end{bmatrix} \begin{bmatrix} w_{11} \\ w_{12} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} -0,06 & -0,07 & | & 0 \\ -0,07 & -0,08 & | & 0 \end{bmatrix} \times \left(-\frac{7}{6}\right) \downarrow \rightarrow \begin{bmatrix} -0,06 & -0,07 & | & 0 \\ 0 & 0 & | & 0 \end{bmatrix}$$

$$\Rightarrow \begin{cases} -0,06 w_{11} - 0,07 w_{12} = 0 \end{cases}$$

$$w_{12} = t_1, \text{ where } t_1 \in \mathbb{R}$$

$$\text{Set } t_1 = \frac{1}{\sqrt{21}}$$

(thus there are infinite number of solutions)

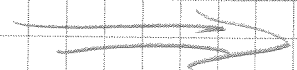
$$w_{11} = -\frac{7}{6} w_{12} = -\frac{7}{6} \times \frac{1}{\sqrt{21}} \approx -0,225$$

$$w_{12} = \frac{1}{\sqrt{21}} \approx 0,207$$

$$\cos(\theta) = -0,225$$

$$\theta = \frac{\cos^{-1}(-0,225) \times 180}{\pi} \approx \underline{\underline{34,4^\circ}}$$

Done



Conformation of  $34,4$  as  $\max \theta$ .

$$\boxed{\text{Case 1: } \theta = 30^\circ}$$

$$X_{41}^* = \cos\left(\frac{30^\circ \times \pi}{180}\right) \times X_1 + \sin\left(\frac{30^\circ \times \pi}{180}\right) \times X_2 =$$

$$= 0,966 \times X_1 + 0,15 \times X_2 \approx \begin{bmatrix} 2,87 \\ 1,37 \\ 2,73 \\ 3,23 \\ 3,60 \\ 3,60 \\ 3,96 \\ 5,46 \end{bmatrix}$$

$$\bar{X}_{41}^* = \frac{1}{8} [2,87 + 1,37 + 2,73 + 3,23 + 3,60 + 3,60 + 3,96 + 5,46] = 3,3525$$

$$\sigma_{41}^{*2} = \frac{1}{7} \left( (2,87 - \bar{X}_{41}^*)^2 + \dots + (5,46 - \bar{X}_{41}^*)^2 \right) = 1,36$$

$$\boxed{\text{Case 2 } \theta = 40}$$

$$X_{42}^* = \cos\left(\frac{40^\circ \times \pi}{180}\right) \times X_1 + \sin\left(\frac{40^\circ \times \pi}{180}\right) \times X_2 = 0,766 \times X_1 + 0,643 \times X_2 =$$

$$\bar{X}_{42}^* = \frac{1}{8} [3,34 + \dots + 5,62] = 3,48875$$

$$\sigma_{42}^2 = \frac{1}{7} \left[ (3,34 - \bar{X}_{42}^*)^2 + \dots + (5,62 - \bar{X}_{42}^*)^2 \right] = 1,34$$

$$\begin{bmatrix} 3,34 \\ 1,41 \\ 2,82 \\ 3,46 \\ 3,58 \\ 0,58 \\ 3,70 \\ 5,62 \end{bmatrix}$$

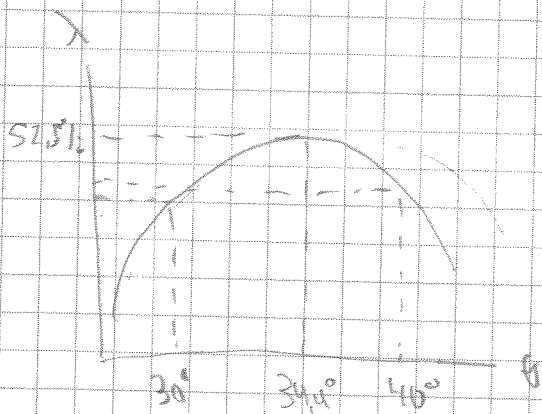
Vektork

Thus, we have obtained three points to locate the maximum angle  $\theta$ :

$$1) \underline{\theta = 30^\circ} : \frac{\sigma_{x_1}^*}{\sigma_1^2 + \sigma_2^2} = \frac{1,36}{(1,43 + 1,41)} \approx 0,479 = 47,9\%$$

$$2) \underline{\theta^{\max} = 34,4^\circ} : \frac{\sigma_{\max}}{\sigma_1^2 + \sigma_2^2} = \frac{\lambda_1}{\sigma_1^2 + \sigma_2^2} = \frac{1,49}{(1,43 + 1,41)} \approx 0,525 = 52,5\%$$

$$3) \underline{\theta = 40^\circ} : \frac{\sigma_{x_2}^*}{\sigma_1^2 + \sigma_2^2} = \frac{1,34}{1,43 + 1,41} = 0,472 = 47,2\%$$



However we can conclude that only accounting for 52,5% is not a good principal component. In this data PCA should not be used due to low correlation between  $X_1$  and  $X_2$ .

VG!

Question 3:

$$\begin{cases} x_1 = \lambda_1 \beta + \delta_1 \\ x_2 = \lambda_2 \beta + \delta_2 \\ x_3 = \lambda_3 \beta + \delta_3 \end{cases} \quad S = \begin{pmatrix} 1,20 & 0,43 & 0,45 \\ 0,43 & 1,56 & 0,27 \\ 0,45 & 0,27 & 2,15 \end{pmatrix}$$

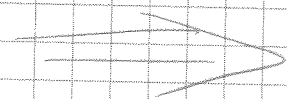
Are the parameter estimates unique?

$$S = \begin{bmatrix} \lambda_1^2 + \text{Var}(\delta_1) & \lambda_1 \lambda_2 & \lambda_1 \lambda_3 \\ \lambda_1 \lambda_2 & \lambda_2^2 + \text{Var}(\delta_2) & \lambda_2 \lambda_3 \\ \lambda_1 \lambda_3 & \lambda_2 \lambda_3 & \lambda_3^2 + \text{Var}(\delta_3) \end{bmatrix}$$

This matrix is symmetric and I only need to consider unique elements

We can conclude that there are 6 unique equations in the system and there is also 6 unique parameters to estimate. Because these are equal ( $6=6$ ) there will exist one solution for the system that is unique (the degrees of freedom are equal to zero ( $6-6=0$ )), the system is just-identified!

The unique parameter estimates are obtained on the next page!





$$\textcircled{1} \text{Var}(x_i) = \lambda_i^2 + \text{Var}(\delta_i)$$

$$\textcircled{2} \text{Cov}(x_i, x_j) = \lambda_i \lambda_j$$

The elements of these are in the matrices  $S$  and  $\hat{S}$ .

$$\begin{cases} 1,20 = \lambda_1^2 + \text{Var}(\delta_1) \\ 1,56 = \lambda_2^2 + \text{Var}(\delta_2) \\ 2,15 = \lambda_3^2 + \text{Var}(\delta_3) \\ 0,93 = \lambda_1 \lambda_2 \\ 0,27 = \lambda_2 \lambda_3 \\ 0,45 = \lambda_1 \lambda_3 \end{cases}$$

$$\Rightarrow \begin{cases} \text{Var}(\delta_1) = 1,20 - \lambda_1^2 & \textcircled{1} \\ \text{Var}(\delta_2) = 1,56 - \lambda_2^2 & \textcircled{2} \\ \text{Var}(\delta_3) = 2,15 - \lambda_3^2 & \textcircled{3} \\ \lambda_1 = \frac{0,93}{\lambda_2} & \textcircled{4} \\ \lambda_2 = \frac{0,27}{\lambda_3} & \textcircled{5} \\ \lambda_3 = \frac{0,45}{\lambda_1} & \textcircled{6} \end{cases}$$

Step 1: Sub.  $\textcircled{5}$  and  $\textcircled{6}$  into  $\textcircled{4}$ .

$$\lambda_1 = \frac{0,93}{\lambda_2} = \frac{0,93}{\left(\frac{0,27}{\lambda_3}\right)} = \frac{0,93}{\left(\frac{0,27}{\left(\frac{0,45}{\lambda_1}\right)}\right)} = \frac{1,55}{\lambda_1} \Leftrightarrow \lambda_1^2 = 1,55 \Rightarrow \lambda_1 = \sqrt{1,55} \approx 1,2450$$

Step 2: Solve  $\textcircled{6}$ .

$$\lambda_3 = \frac{0,45}{\lambda_1} = \frac{0,45}{\sqrt{1,55}} \approx 0,3614$$

Step 3: solve  $\textcircled{5}$ .

$$\lambda_2 = \frac{0,27}{\lambda_3} = \frac{0,27}{0,3614} \approx 0,7471$$

Step 4: solve  $\textcircled{1}$ ,  $\textcircled{2}$  and  $\textcircled{3}$

$$\text{Var}(\delta_1) = 1,20 - 1,55 = -0,35$$

$$\text{Var}(\delta_2) = 1,56 - (0,7471)^2 \approx 1,0018$$

$$\text{Var}(\delta_3) = 2,15 - (0,3614)^2 \approx 1,5918$$

Solutions:

$$\lambda_1 = 1,2450$$

$$\lambda_2 = 0,7471$$

$$\lambda_3 = 0,3614$$

negative why?

$$\text{Var}(\delta_1) = -0,3560$$

$$\text{Var}(\delta_2) = 1,0018$$

$$\text{Var}(\delta_3) = 1,5918$$

OK

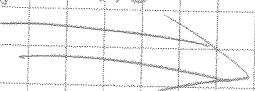
Question 6:

i) Cluster analysis intend accomplish data reduction by grouping observation into groups ("clusters") that share similar characteristics, and then, only analyze these groups. The similarity of characteristics is measured by distance between variables.

There are two types of clustering categories, hierarchical and non-hierarchical. Within the hierarchical there are the centroid, single linkage, complete linkage, average method and Ward's method to combine observations to clusters.

The assumptions are that, the data is multivariate normal, the groups have equal covariance matrices and that the observations are independent.

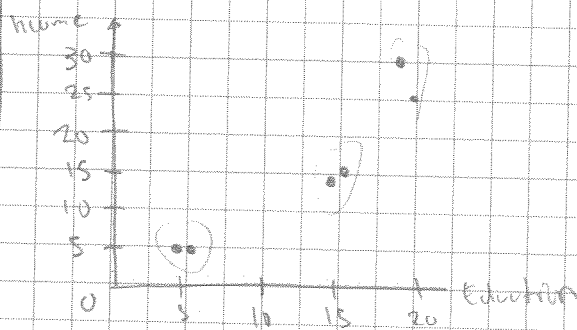
One have to have all these fulfilled to be able to trust and draw conclusions from the clusters, and if one wants to compute test statistics. It might be easier to cluster a variable on a categorical variable but the data should still be multivariate normal.





i) Squared Euclidean distance:  $D_{AB}^2 = \sum_{k=1}^p (a_k - b_k)^2$

The intermediate similarity matrix  $S$  consist of the squared Euclidean distances between observations, and thereby a zero diagonal.



• Visual inspection of data suggest three clusters.

$$D_{S1S2}^2 = (5-6)^2 + (5-6)^2 = 1$$

$$D_{S1S3}^2 = (5-15)^2 + (5-14)^2 = 181$$

$$D_{S1S4}^2 = (5-16)^2 + (5-15)^2 = 221$$

$$D_{S1S5}^2 = (5-25)^2 + (5-20)^2 = 625$$

$$D_{S1S6}^2 = (5-30)^2 + (5-14)^2 = 745$$

$$D_{S2S3}^2 = (6-15)^2 + (6-14)^2 = 145$$

$$D_{S2S4}^2 = (6-16)^2 + (6-15)^2 = 121$$

$$D_{S2S5}^2 = (6-25)^2 + (6-20)^2 = 557$$

$$D_{S2S6}^2 = (6-30)^2 + (6-14)^2 = 745$$

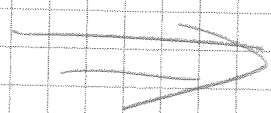
First column

Second column, and so on...

Until we get the intermediate similarity matrix

	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$
$s_1$	0					
$s_2$	1	0				
$s_3$	181	145	0			
$s_4$	221	121	1	0		
$s_5$	625	557	136	106	0	
$s_6$	745	745	250	212	26	0

Symmetry



Now we start the hierarchical clustering by the single linkage method, i.e. we choose the minimum distance of the two possibilities in the similarity matrix when forming a cluster. But we face a problem at the first step, should we start to cluster  $S_1, S_2$  or  $S_3, S_4$ ? They have the same squared Euclidean distance, I will start by  $S_1, S_2$  (it is convenient and they did it in the book).

$S_{(step 2)}$

$S_1, S_2$	0				
$S_3$	145	0			
$S_4$	181	1	0		
$S_5$	557	136	106	0	
$S_6$	745	250	212	26	0

Note that I have taken the minimum value of the "red" lines of  $S_i$  rowwise  
 ex  $145 = \text{Min}(181, 145)$

$S_{(step 2)}$

$S_1, S_2$	0			
$S_3, S_4$	145	0		
$S_5$	557	106	0	
$S_6$	745	212	26	0

$S_{(step 3)}$

$S_1, S_2$	0		
$S_3, S_4$	145	0	
$S_5, S_6$	557	106	0

I would stop here and use three clusters, as the plot and the increase in distance, from [1, 26] to 106, indicates but if it is required

$S_{(step 4)}$

$S_1, S_2$	0
$(S_3, S_4)(S_5, S_6)$	145

2 cluster solution:  $S_1, S_2$  and  $S_3, S_4, S_5, S_6$   
 OK, some calculations missing