

LØSNING

Oppgave 1

a) Σ kjent:

$$U^2 = n(\bar{X} - \mu_0)' \Sigma^{-1} (\bar{X} - \mu_0) \sim \chi_p^2 \text{ under } H_0$$

Forkast hvis $U^2 > \chi_{p, \alpha}^2$

Σ ukjent:

$$T^2 = n(\bar{X} - \mu_0)' S^{-1} (\bar{X} - \mu_0) \sim \frac{(n-1)p}{n-p} F_{p, n-p}$$

Forkast hvis $T^2 > \frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)$

Konklusjon $p=5$, Σ ukjent:

Kritisk verdi er $\frac{49.5}{45} \cdot F_{5, 45, 0.05} = 13.17$
 ≈ 24.2

Observert verdi: $50 \cdot 0.167 = 8.35$ dos ikke forkast

Hvis Σ er kjent er kritisk verdi $\chi_{5, 0.05}^2 = 11.07$

dos forkast heller ikke da.

b. Bonferroni:

$$\text{Intervall for } \mu_i: \bar{X}_i \pm t_{49} \left(\frac{0.05}{4} \right) \sqrt{\frac{S_{ii}}{50}}$$

$$\text{dvs. } \bar{X}_i \pm 2.31 \sqrt{\frac{S_{ii}}{50}}$$

$$\text{For } \mu_1: 13.38 \pm 2.31 \sqrt{\frac{13.73}{50}}$$

$$13.38 \pm 1.21$$

$$\text{dvs. } (12.17, 14.59)$$

$$F^2\text{-metoden: } \bar{X}_i \pm \sqrt{\frac{p(n-t)}{n-p} F_{p,n-p}(0.05)} \sqrt{\frac{S_{ii}}{50}}$$

$$\bar{X}_i \pm \sqrt{\frac{5 \cdot 49}{45} \cdot 2.42} \cdot \sqrt{\frac{S_{ii}}{50}}$$

$$\text{dvs. } \bar{X}_i \pm 3.69 \sqrt{\frac{S_{ii}}{50}}$$

$$\text{dvs. store intervaller. (For } \mu_1: (13.38 \pm 1.90) \text{) } \\ \text{dvs. } (11.48, 15.28)$$

Konfidensellipse for (μ_1, μ_2) kan lages ved
 a) at sæt den øverste venstre 2×2 matricen
 i Σ^* , dvs.

$$\Sigma^* = \begin{pmatrix} 13.73 & 8.69 \\ 8.69 & 8.95 \end{pmatrix}$$

Konfidensellipsen vil da ha centrum i
 $\begin{pmatrix} \bar{X}_1 \\ \bar{X}_2 \end{pmatrix} = \begin{pmatrix} 13.38 \\ 14.09 \end{pmatrix}$ og akser langs egen-
 vektorene for Σ^* . (Se boka kap. 5).

c) Andel av total sampelvarians forklart ved \hat{y}_i er

$$\frac{180.70}{180.70 + 7.24 + 4.07 + 2.13 + 0.39} = \underline{0.929} \quad \underline{\underline{(92.9\%)}}$$

Korrelasjoner:

Formelsamlingen ("5.8 av 10") gir formlene for populasjons variansen av korrelasjonene. Her skal brukes sampel versjonene:

$$\frac{\sqrt{\hat{\lambda}_i} \hat{e}_{i1}}{\sqrt{s_{ii}}}$$

Telleren er her elementene i første oppgitte egenvektor, mens nevneren er diagonalelementene i $\hat{\Sigma}$

Demed: $i=1: \frac{3.30}{\sqrt{13.73}} = 0.891$

$i=2: 0.809$

$i=3: 0.885$

$i=4: 0.939$

$i=5: 0.994$

Alle korrelasjonene er høye. Korrelasjonen med X_5 er dominerende siden X_5 har den suverent største sampelvariansen. Siden X_5 har en mye høyere varians enn de andre X_i kunne korrelasjonsmatrisen R med X_5 i første søle og første rade...

$$\begin{aligned}
 d) \quad \tilde{\Sigma} &= E((\tilde{x} - \mu)(\tilde{x} - \mu)') \\
 &= E[(\tilde{l}F + \tilde{\varepsilon})(\tilde{l}F + \tilde{\varepsilon})'] \\
 &= \tilde{l}\tilde{l}'E(F^2) + E(\tilde{\varepsilon}\tilde{\varepsilon}') \\
 &= \tilde{l}\tilde{l}' + \tilde{\Psi}
 \end{aligned}$$

Prinzipal komponent - methoden:

$$\tilde{l} = \begin{pmatrix} 3.30 \\ 2.42 \\ 4.59 \\ 2.03 \\ 11.78 \end{pmatrix}, \quad \tilde{\Psi}_{ii} = -(\tilde{l}\tilde{l}')_{ii} + S_{ii}$$

des

$$\tilde{\Psi} = \begin{pmatrix} 2.84 & & & & \\ & 3.09 & & & \\ & & 5.80 & & \\ & & & 0.55 & \\ & & & & 1.54 \end{pmatrix}$$

siden

$$\tilde{l}\tilde{l}' = \begin{pmatrix} 10.89 & 7.99 & 15.15 & 6.20 & 38.87 \\ ? & 5.86 & ? & ? & ? \\ ? & ? & 21.07 & ? & ? \\ ? & ? & ? & 4.12 & ? \\ ? & ? & ? & ? & 138.77 \end{pmatrix}$$

Residualmatrisen:

$$S - \tilde{l}\tilde{l}' - \tilde{\Psi} = \begin{pmatrix} 0 & 0.70 & 0.50 & 0.33 & -1.22 \\ & 0 & & & \\ & & 0 & & \\ & & & 0 & \\ & & & & 0 \end{pmatrix}$$

Oppgave 2

$$\begin{aligned} a) \text{Cov}(\underline{X}, Y) &\stackrel{\text{def}}{=} E[(\underline{X} - \underline{\mu})(Y - \eta)] \\ &= E[(\underline{l}F + \underline{\varepsilon})(qF + \delta)] \\ &= q \underline{l}' E(F^2) = q \underline{l}' \end{aligned}$$

Dermed for formelsamling + innsettning:

$$\begin{aligned} \hat{Y}_{LS} &= \eta + \underline{\sigma}' \underline{\Sigma}^{-1} (\underline{X} - \underline{\mu}) \\ &= \eta + q \underline{l}' (\underline{l} \underline{l}' + \underline{\Psi})^{-1} (\underline{X} - \underline{\mu}) \end{aligned}$$

~~#~~ Generell definisjon av \hat{Y}_{LS} :

Den er den lineære funksjon av \underline{X} (dvs " $\hat{Y} = a + b \underline{X}$ ") som minimerer $E((\hat{Y} - Y)^2)$

$$\text{Svaret er } \hat{Y} = \eta + \underline{\sigma}' \underline{\Sigma}^{-1} (\underline{X} - \underline{\mu})$$

som ifølge formel på "s. 6 av 10" er den betingede forventning for Y gitt \underline{X} ved multivariatitet.

b) Ren insetting fra formelsamling
 (med $\underline{\omega}_1 = \underline{\sigma} = \gamma \underline{l}$) gir

$$\hat{Y}_{PLS} = \gamma + \gamma \underline{l}' \underline{l} (\underline{l}' (\underline{l} \underline{l}' + \underline{\Psi}) \underline{l})^{-1} \underline{l}' (\underline{X} - \underline{\gamma})$$

Hovedidé: Finne lineærkombinasjonen

$F_j = \underline{\omega}_j' X$ som har størst mulig
 korrelasjon med \underline{Y} (σ_j som er
 ukonstant, og med gitte lengder på $\underline{\omega}_j$.)

c) Vi har

$$\begin{aligned} & \underline{l}' (\underline{l} \underline{l}' + \underline{\Psi})^{-1} \\ &= \underline{l}' (\underline{\Psi}^{-1} - a \underline{\Psi}^{-1} \underline{l} \underline{l}' \underline{\Psi}^{-1}) \\ &= \underline{l}' (\underline{I} - a \underline{\Psi}^{-1} \underline{l} \underline{l}') \underline{\Psi}^{-1} \\ &= (\underline{l}' - a \underline{l}' \underline{\Psi}^{-1} \underline{l} \underline{l}') \underline{\Psi}^{-1} \\ &= (1 - a \underline{l}' \underline{\Psi}^{-1} \underline{l}) \underline{l}' \underline{\Psi}^{-1} \\ &= \frac{1}{1 + \underline{l}' \underline{\Psi}^{-1} \underline{l}} \underline{l}' \underline{\Psi}^{-1} \end{aligned}$$

$$\text{der } c = \frac{f}{1 + \underline{l}' \underline{\Psi}^{-1} \underline{l}}$$

Vil se

$$d = f \underline{l}' \underline{l} \left(\underbrace{\underline{l}' (\underline{l} \underline{l}' + \underline{\Psi}) \underline{l}}_{\text{et tall!}} \right)^{-1}$$

Vi ser dermed at \hat{Y}_{LS} skalerer observasjonen $X_i - \mu_i$ med ψ_i -ene (de spesifikke variansene). Legger altså størst vekt på de variable som har minst spesifikke varianser (og dermed er forklart mest av den felles faktoren F).

Bortsett fra denne skaleringen, er det den samme lineærkombinasjon med l_i -ene som går igjen.

$$\underline{\Psi} = \sigma^2 \underline{I}$$

Da er

$$c = \frac{f}{1 + \frac{1}{\sigma^2} \underline{l}' \underline{l}} = \frac{f \sigma^2}{\sigma^2 + \underline{l}' \underline{l}}$$

$$d = \underset{\sim}{g} \underset{\sim}{l}' \underset{\sim}{l} \left(\underset{\sim}{l}' \left(\underset{\sim}{l} \underset{\sim}{l}' + \sigma^2 \underset{\sim}{I} \right) \underset{\sim}{l} \right)^{-1}$$
$$= \underset{\sim}{g} \underset{\sim}{l}' \underset{\sim}{l} \left[\underset{\sim}{l}' \underset{\sim}{l} \left(\underset{\sim}{l}' \underset{\sim}{l} + \sigma^2 \right) \right]^{-1}$$
$$= \frac{\underset{\sim}{g}}{\sigma^2 + \underset{\sim}{l}' \underset{\sim}{l}}$$

Siden $\psi_i = \sigma^2$ for hver i , ser vi at

~~$\hat{\psi}_{LS} = \hat{\psi}_{PLS}$~~

$\hat{\psi}_{LS} = \hat{\psi}_{PLS}$

(d) $\begin{pmatrix} X \\ \underset{\sim}{y} \end{pmatrix}$ dannes via tilsammen den

3-dimensionale faktormodellen fra Opgave 1. Naturligt estimater er de

$$\hat{\underset{\sim}{l}} = \begin{pmatrix} 3.30 \\ 2.42 \\ 4.59 \\ 2.03 \end{pmatrix}$$

$$\hat{\underset{\sim}{\mu}} = \begin{pmatrix} 13.38 \\ 14.09 \\ 119.45 \\ 10.00 \end{pmatrix}$$

$$\hat{\underset{\sim}{\psi}} = \begin{pmatrix} 2.84 \\ 3.09 \\ 3.80 \end{pmatrix}$$

$$\hat{\underset{\sim}{g}} = 11.78$$

$$\hat{\underset{\sim}{\eta}} = 39.19$$

Prediksjoneene blir da

$$\begin{aligned} \hat{Y}_{LS} = & 39.19 + 0.6598 \cdot \left[\frac{3.30}{2.84} (10 - 13.38) \right. \\ & + \frac{2.42}{3.09} (10 - 14.09) \\ & + \frac{4.59}{5.80} (115 - 119.45) \\ & \left. + \frac{2.03}{0.55} (10 - 10) \right] = \underline{\underline{32.16}} \end{aligned}$$

$$\begin{aligned} \hat{Y}_{PLS} = & 39.19 + 0.2557 \left[3 \cdot (10 - 13.38) + 2.42(10 - 14.09) \right. \\ & \left. + 4.59(115 - 119.45) + 2.03 \cdot (10 - 10) \right] \\ = & \underline{\underline{28.58}} \end{aligned}$$

Oppgave 3

a) Avstander (brukes $d(\underline{x}, \underline{y}) = \sqrt{\sum (x_i - y_i)^2}$)

	a	b	c	d
a	0			
b	6.40	0		
c	6.71	2.83	0	
d	4.12	3.16	5.10	0

Single linkage:

1) bc danner cluster ved 2.83

Ny avst. matrise:

	bc	a	d
bc	0		
a	6.40	0	
d	3.16	4.12	0

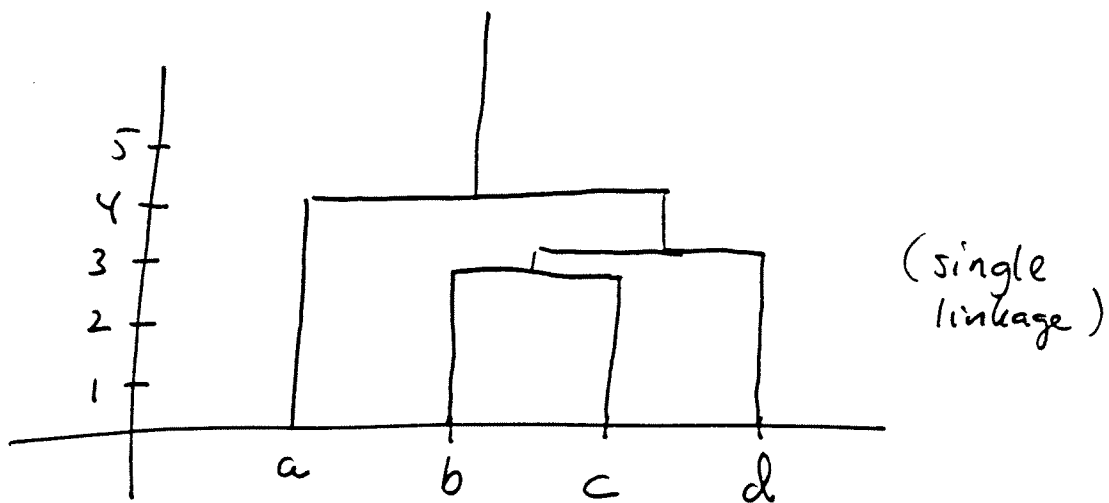
2) ^{0.5} bcd danner cluster ved 3.16

Ny avst. matrise

	bcd	a
bcd	0	
a	4.12	0

3) bcd og a danner cluster ved 4.12

Dendrogram single linkage:



Average linkage :

1) bc dannen cluster ved 2.83

My avst. matrise : ~~for single linkage~~

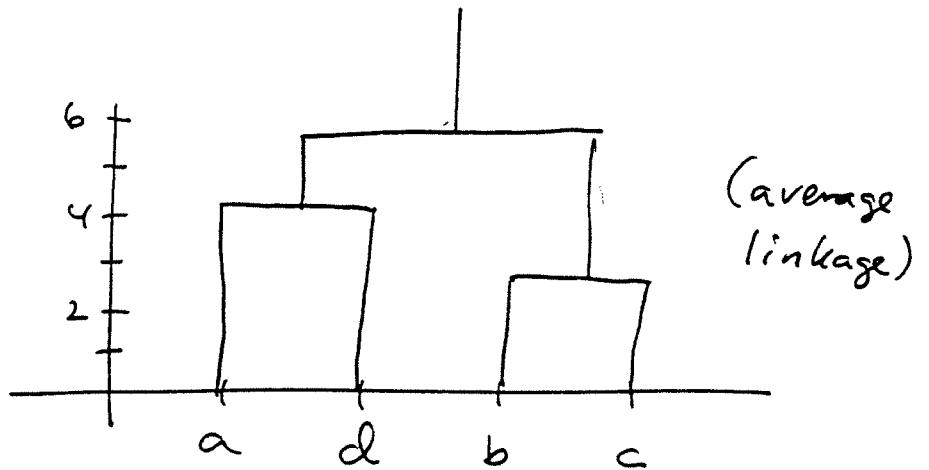
	bc	a	d
bc	0		
a	6.56	0	
d	4.13	4.12	0

2) a og d dannen nytt cluster ved 4.12

My avstandsmatrise :

	bc	ad
bc	0	
ad	5.35	0

Dendrogram



Dersom to cluster ønskes, ser vi av dendrogrammet for single at de to clusterne er {a} og {bcd} (mellom 3.16 og 4.12 består dendrogrammet av to cluster).

tilsvarende for average : Clusterne er {ad}, {bc} (avst. 4.12 - 5.35)

K-means: Anta vi starter med clustrene

{a}	Senter $\begin{pmatrix} 5 \\ 4 \end{pmatrix}$
{bcd}	Senter $\begin{pmatrix} 4/3 \\ 0 \end{pmatrix}$

Da vil a, b, c, d alle ligge nærmest ~~det~~ senteret i det cluster de selv ligger.

Altså vil K-means med de initiale cluster som ovenfor lede til den samme

opdelingen. (F.eks. Aust. for b til {a} er $\sqrt{(1-5)^2 + (-1-4)^2} = \sqrt{41}$, mens aust. for b til {bcd} er $\sqrt{(1-4/3)^2 + (-1)^2} = \sqrt{10/9} < \sqrt{41}$.)

Anta så vi starter med

{ad}	Senter $\begin{pmatrix} 9/2 \\ 2 \end{pmatrix}$
{bc}	Senter $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$

Også da vil a, b, c, d alle ligge nærmest ~~sitt~~ senteret i det cluster de selv ligger.