

Vervolgdoctoraal Fysische Geografie,  
Tentamen Matrix algebra en Multivariate Statistiek (MUVS)  
28 nov 1997, Tr I, 102

Maak vraag 1, en kies er twee uit de vragen 2, 3 en 4. Antwoord kort.

VRAAG 1

a bereken het inproduct van de vectoren (1, 3, 1, 2) en (4, 5, -2, 8)

gegeven  $A = \begin{bmatrix} 1 & 2 & 1 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix}$ ,  $b = [7, 5, 9]$ ,  $c = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$ .

b bereken  $bc$ ,  $cb$  en  $|cb|$ .

c inverteer  $A$ , controleer je antwoord.

d los  $Ax = b'$  op voor  $x$  (hint: gebruik hiervoor het antwoord van vraag c, als je dit hebt).

e is  $c$  een eigenvector van  $A$ ?

VRAAG 2

a Leg de verschillen en overeenkomsten tussen discriminantanalyse en arbitrary origin (k-means) clusteranalyse kort uit.

b Beschouw de dataset met de furanen en dioxines uit de twee havens in Rotterdam en de Rijn, bekeken in het practicum. Welke indeling is, in termen van Wilk's criterium, de beste: de indeling in herkomst (Chemiehaven, Laurenhaven, Rijn) of die gevonden met behulp van k-means clusteranalyse?

c Had je het antwoord op vraag 2b ook kunnen beantwoorden zonder dat je de dataset had gekend?

d Welke van de twee indelingen bij vraag 2b is het meest waardevol voor de interpretatie van de gegevens, en waarom?

### VRAAG 3

Gegeven de dataset in onderstaande kruistabel

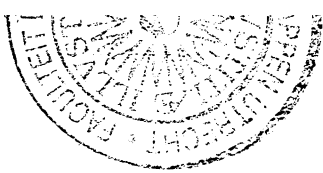
	locatie 1	locatie 2	locatie 3
soort a	15	9	0
soort b	6	2	4
soort c	1	1	4
soort d	8	8	2

- a bereken de tabel met verwachte frequenties, en de tabel met de chi-kwadraat bijdragen.
- b welke soort en welke locatie wijken het sterkste af van de marginale verdeling? Waaraan zie je dit?
- c Welke soort en welke locatie lijken het sterkste op de marginale verdeling?
- d Stel, je wilt bovenstaande tabel ordenen. Welke techniek zou je hiervoor gebruiken? Op welk criterium worden rijen en kolommen dan gesorteerd?

### VRAAG 4

Een onderzoeker (Steven de Jong) heeft een set van 70 satellietbeelden van een enkel gebied, elk voor een smal bandje uit het zichtbare of near infrared spectrum. Om een globale indruk te krijgen van het patroon, aanwezig in al deze beelden, gebruikt hij hoofdcomponentenanalyse.

- a Waarom?
- b Hoe ziet de data matrix er in dit geval uit?
- c De eerste drie hoofdcomponenten lijken de beelden vrij goed samen te vatten, maar de interpreteerbaarheid van deze hoofdcomponenten laat te wensen over: het is niet duidelijk hoe een hoofdcomponent samenhangt met een groep beelden. Is hier verbetering in aan te brengen? Zo ja, hoe?
- d wat is, bij hoofdcomponentenanalyse, een loading (of loding)? Wat is een score?



a.  $(1, 3, 1, 2) \cdot (4, 5, -2, 3) = 4 + 15 + (-2) + 16 = 33$

b.  $bc = [7 \ 5 \ 9] \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} = [14 + 5 + 9] = [28]$

$cb = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} [7 \ 5 \ 9] = \begin{bmatrix} 14 & 10 & 18 \\ 7 & 5 & 9 \\ 7 & 5 & 9 \end{bmatrix}$

$|cb| = \begin{vmatrix} 14 & 10 & 18 \\ 7 & 5 & 9 \\ 7 & 5 & 9 \end{vmatrix} = \begin{vmatrix} 14 & 10 & 18 \\ 7 & 5 & 9 \\ 0 & 0 & 0 \end{vmatrix} = \begin{vmatrix} 0 & 0 & 0 \\ 7 & 5 & 9 \\ 0 & 0 & 0 \end{vmatrix} = 0$

c.  $\begin{bmatrix} 1 & 2 & 1 & | & 1 & 0 & 0 \\ 2 & 1 & 0 & | & 0 & 1 & 0 \\ 3 & 2 & 1 & | & 0 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 1 & | & 1 & 0 & 0 \\ 2 & 1 & 0 & | & 0 & 1 & 0 \\ 0 & -4 & -2 & | & -3 & 0 & 1 \end{bmatrix} \rightarrow$

$\begin{bmatrix} 1 & 2 & 1 & | & 1 & 0 & 0 \\ 0 & -3 & -2 & | & -2 & 1 & 0 \\ 0 & -4 & -2 & | & -3 & 0 & 1 \end{bmatrix} \rightarrow$   ~~$\begin{bmatrix} 1 & 2 & 1 & | & 1 & 0 & 0 \\ 0 & -3 & -2 & | & -2 & 1 & 0 \\ 0 & -4 & -2 & | & -3 & 0 & 1 \end{bmatrix}$~~   $\rightarrow$

~~$\begin{bmatrix} 1 & 0 & 1 & | & -1 & -2 & 2 \\ 0 & 1 & 0 & | & 1 & 1 & -1 \\ 0 & -2 & 1 & | & 4 & -3 & -3 \end{bmatrix}$~~   $\rightarrow$   ~~$\begin{bmatrix} 1 & 0 & 0 & | & 1 & 0 & 1/2 \\ 0 & 1 & 0 & | & 1 & 1 & -1 \\ 0 & 0 & -2 & | & -1/3 & -4/3 & 1 \end{bmatrix}$~~

$\begin{bmatrix} 1 & 2 & 1 & | & 1 & 0 & 0 \\ 0 & -3 & -2 & | & -2 & 1 & 0 \\ 0 & 0 & 2/3 & | & -1/3 & -4/3 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 0 & | & 1/2 & 2 & -3/2 \\ 0 & -3 & 0 & | & -3 & -3 & 3 \\ 0 & 0 & 2/3 & | & -1/3 & -4/3 & 1 \end{bmatrix} \rightarrow$

$\begin{bmatrix} 1 & 0 & 0 & | & -1/2 & 0 & 1/2 \\ 0 & -3 & 0 & | & -3 & -3 & 3 \\ 0 & 0 & 2/3 & | & -1/3 & -4/3 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 & | & -1/2 & 0 & 1/2 \\ 0 & 1 & 0 & | & 1 & 1 & -1 \\ 0 & 0 & 1 & | & -1/2 & -2 & 3/2 \end{bmatrix} =$

$A^{-1} = \begin{bmatrix} -1/2 & 0 & 1/2 \\ 1 & 1 & -1 \\ -1/2 & -2 & 3/2 \end{bmatrix}$

controle:  $A^{-1}A = I$

$$\frac{1}{2} \begin{bmatrix} -\frac{1}{2} & 0 & \frac{1}{2} \\ 1 & 1 & -1 \\ -\frac{1}{2} & -2 & \frac{3}{2} \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \text{lept!}$$

1.  $b = [7 \ 5 \ 9]$      $b' = \begin{bmatrix} 7 \\ 5 \\ 9 \end{bmatrix}$

$$Ax = b' \rightarrow A^{-1}Ax = A^{-1}b' \rightarrow Ix = A^{-1}b'$$

$$A^{-1}b' = \begin{bmatrix} -\frac{1}{2} & 0 & \frac{1}{2} \\ 1 & 1 & -1 \\ -\frac{1}{2} & -2 & \frac{3}{2} \end{bmatrix} \begin{bmatrix} 7 \\ 5 \\ 9 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \\ 0 \end{bmatrix} = x \quad \rightarrow x = A^{-1}b'$$

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

$$\begin{vmatrix} 1-\lambda & 2 & 1 \\ 2 & 1-\lambda & 0 \\ 3 & 2 & 1-\lambda \end{vmatrix} = (1-\lambda) \begin{vmatrix} 1-\lambda & 0 \\ 2 & 1-\lambda \end{vmatrix} - 2 \begin{vmatrix} 2 & 0 \\ 3 & 1-\lambda \end{vmatrix} + 1 \begin{vmatrix} 2 & 1-\lambda \\ 3 & 2 \end{vmatrix}$$

$$= (1-\lambda) \left( (1-\lambda)^2 - 0 \right) - 2 \left( (2-2\lambda) - 0 \right) + (4 - (3-3\lambda))$$

$$= (1-\lambda)^3 + 4\lambda - 4 + 4 - 3 + 3\lambda = (1-\lambda)^3 + 7\lambda - 3 = 0$$

als  $c$  is eigenvector dan:  $\begin{bmatrix} 1 & 2 & 1 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} = \lambda \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$  (immers:  $Ax = \lambda x$ )

$$\Rightarrow \left. \begin{array}{l} 5 = 2\lambda \rightarrow \lambda = \frac{5}{2} \\ 5 = \lambda \rightarrow \lambda = 5 \\ 9 = \lambda \rightarrow \lambda = 9 \end{array} \right\} \text{ Dit kan niet, dus } c \text{ is geen eigenvector van } A$$

③ a. Waargenomen:

	locatie 1	locatie 2	locatie 3	Totaal (Dr)
soort a	15	9	0	24
soort b	6	2	4	12
soort c	1	1	4	6
soort d	8	8	2	18
Totaal (De)	30	20	10	60

Verwacht:

	locatie 1	locatie 2	locatie 3	Totaal
soort a	12	8	4	24
soort b	6	4	2	12
soort c	3	2	1	6
soort d	9	6	3	18
Totaal	30	20	10	60

Waargenomen - Verwacht

	locatie 1	locatie 2	locatie 3	Totaal
soort a	3	1	-4	0
soort b	0	-2	2	0
soort c	-2	-1	3	0
soort d	-1	2	-1	0
Totaal	0	0	0	0

Chi-kwadraat bijdrages

	locatie 1	locatie 2	locatie 3	<del>Kwadraat</del> kwadraatsom
soort a	0,112	0,046	-0,258	0,081
soort b	0	-0,129	0,183	0,050
soort c	-0,149	-0,091	0,387	0,180
soort d	-0,043	0,105	-0,075	0,018
<del>Kwadraat</del> kwadraatsom	0,037	0,038	0,255	

(Chi-kwadraat bijdrage :  $D_r^{-95} (\text{waargenomen} - \text{verwacht}) D_k^{-95}$ )

- b. \* soort c wijkt het sterkst af van de vier soorten: de <sup>kwadraatsom van de</sup>  $\chi^2$ -bijdrages is het grootst tevens is al in de waargenomen waarden te zien dat de frequentie van soort c toeneemt van links naar rechts, terwijl de kolomtotaal juist afneemt.
- \* locatie 3 wijkt het sterkst af van de drie locaties: de kwadraatsom van de  $\chi^2$ -bijdrages is het grootst. In de waargenomen waarden is te zien dat de frequentie van ~~de~~ locatie 3 toe en afneemt, waar de rijtotaal juist af- en toeneemt.

c. \* locatie 1 lijkt het sterkst op de marginaal: laagste kwadraatssom en het verloop van de waargenomen frequenties is ~~g~~ gelijk aan dat van de rijtotalen

\* soort d lijkt het sterkst op de marginaal: laagste kwadraatssom en het verloop van de waargenomen frequenties lijkt op dat van de kolomtotalen.

d. Omdat er sprake is van kwalitatieve gegevens ~~en~~ ~~het~~ (de frequentie in voorkomen van een soort op een locatie is gegeven) is het het beste om een correspondentie-analyse (CA) uit te voeren.

~~De rijen en kolommen worden gesorteerd op de ladingen van de~~  
~~rijen en kolommen~~  
PCA

De rijen en kolommen worden gesorteerd op de ladingen van een bepaalde dimensie. Deze ladingen worden berekend op basis van de  $X^2$ -bijdragen.



7) a. Het doel van een hoofdcomponenten-analyse (PCA) is:

- \* data-reductie
- \* data-exploratie (in dit geval: is er een patroon zichtbaar?)
- \* het voorkomen van data-redundancy
- \* het samenvatten van de data

Omdat er per pixel van het gebied 70 ~~er~~ variabelen zijn (het de dataset is 70-dimensionaal) lijkt het handig om een data-reductie te bewerkstelligen. Het zou handig zijn als de dataset samengevat kan worden in minder dimensies. Het is dan makkelijker een patroon te ontdekken.

b. De data-matrix ziet er als volgt uit:

<del>pixelnr.</del> pixelnr.	spectrale band (genummerd)			
	1	2	...	70
1	$DN_{1,1}$	$DN_{1,2}$	...	$DN_{1,70}$
2	$DN_{2,1}$	$DN_{2,2}$	...	$DN_{2,70}$
...	...	...	...	...
n	$DN_{n,1}$	$DN_{n,2}$	...	$DN_{n,70}$

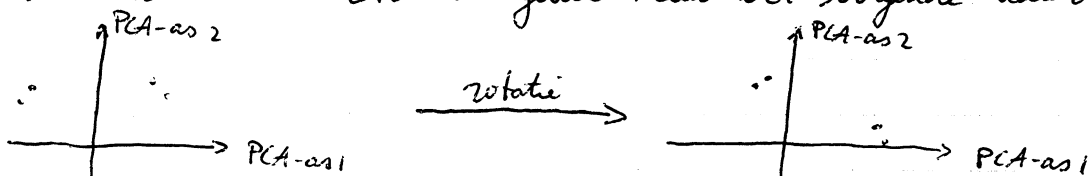
hierin geldt:

$n$  = aantal pixels per beeld  
 $DN_{i,j}$  = Digital Number voor pixel  $i$  en spectrale band  $j$

c. Er zou een verbetering aangebracht kunnen worden door rotatie van de drie PCA-assen.

Zonder rotatie is er weinig verschil in de loadings van de ~~er~~ variabelen op de PCA-assen. Het is dan niet duidelijk welke as door welke ~~er~~ variabelen gevormd wordt. Door rotatie komen bepaalde variabelen dicht bij een bepaalde as te liggen (er ontstaat dan meer verschil in loadings tussen de variabelen).

In een 2-dimensionaal geval kan het volgende aan de hand zijn:



Na rotatie hebben de ~~variabelen~~ <sup>variabelen</sup> een hoge ~~loading~~ <sup>loading</sup> voor de ene PCA-as en een lage loading voor de andere (voor een 2-dimensionaal voorbeeld is gekozen omdat dit makkelijker te tekenen is dan een meer-dimensionale situatie).

Een loading van een variabele ~~geeft de~~ op een PCA-as geeft de samenhang ~~van~~ <sup>van</sup> de variabele met die PCA-as. Een loading die sterk positief of negatief is, geeft een resp. sterk positief of sterk negatief verband tussen de variabele en de PCA-as.

Een PCA-as kan beschreven als:  $y = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n$   
waarbij  $\alpha_i$  is loading van variabele  $x_i$  op de PCA-as

in score <sup>voor een PCA-as</sup> is een nieuwe waarde <sup>van een</sup> die berekend wordt op basis van de ~~variabelen~~ <sup>loadings</sup>. De berekening is dan:  $y = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n$

waarbij  $x_i$  = ~~score~~ waarde voor variabele  $i$

$n$  = totaal aantal ~~van~~ variabelen

$\alpha_i$  = loading van variabele  $i$  op de PCA-as

$y$  = score <sup>van</sup> de PCA-as

~~Dit kan samengevat worden als  $y = X \cdot C$   
waarbij  $X$  = matrix met ~~de~~ waarden voor de variabelen  
 $C$  = transformatie matrix met~~

op die manier kan voor elke PCA-as een score uitgerekend worden.