

Enkele opmerkingen over het gebruik van
correlatie-coëfficiënten

Ph.Th.Stol

INSTITUUT VOOR CULTUURTECHNIEK EN WATERHUISHOUDING
Lindendreef 25
Postbus 10
6700 AB Wageningen

De mate van samenhang tussen twee variabelen kan op eenvoudige wijze worden weergegeven door de correlatie-coëfficiënt.

Het is bekend dat hoge correlatie-coëfficiënten niet zonder meer behoeven te betekenen dat er een causale samenhang tussen de basisgegevens bestaat. Steeds kan zich bijvoorbeeld al de mogelijkheid voordoen dat twee grootheden beide met een derde grootheid zijn gecorreleerd en zodoende zelf ook een hogere correlatie vertonen dan wanneer beide onafhankelijk van deze derde grootheid waren geweest of waren waargenomen.

Het behoort tot de specifieke taak van de onderzoeker op grond van inzicht in de eigen probleemstelling na te gaan of in zijn materiaal dergelijke correlaties kunnen optreden die moeten worden geëlimineerd.

Een ander voorbeeld waarbij hoge correlaties worden berekend die oorspronkelijk niet aanwezig waren is het geval waarbij een variabele per eenheid van een andere variabele wordt uitgedrukt (bijvoorbeeld per oppervlakte-eenheid, per tijdseenheid) of waarbij van verhoudingsgetallen gebruik wordt gemaakt. Zelfs behoeft het niet steeds bekend te zijn dat een resultaat op een dergelijke wijze is ontstaan.

Uit de gebruikelijke correlatie-coëfficiënten te weten de gewone (of totale), de partiële en de multipale zal de onderzoeker zelf een keus moeten doen. (FISHER, 1958, pagina 190 en 191). Zo kan het bijvoorbeeld voorkomen dat het constant veronderstellen van een derde opgenomen variabele voor een gesteld probleem irrelevant is. Hieruit volgt dat er geen algemene regel kan worden gegeven die voorschrijft welk van de correlatie-coëfficiënten moet worden gebruikt.

Soms kan het nuttig zijn alle correlatie-coëfficiënten te berekenen teneinde een volledig inzicht in alle samenhangen te verwerven. In terminologie van vectorrekening betekent dit dat van alle vectoren die een variabele voorstellen de onderlinge ligging wordt vastgesteld door het berekenen van de hoeken tussen de vectoren (figuur 4).

Ook kan het van belang zijn noemers van verhoudingsgetallen als zelfstandige variabele op te nemen teneinde in staat te zijn de invloed er van te elimineren.

De bedoeling van de in deze nota vermelde berekeningen is aan de hand van enkele numerieke voorbeelden het gebruik van de partiële correlatie-coëfficiënt toe te lichten en op het verschijnsel van hoge correlaties tussen afgeleide reeksen nog eens de aandacht te vestigen.

Verondersteld wordt thans dat in bijlage 1 de kolommen X_1 , X_2 en X_3 waarnemingsresultaten zijn. Hierin zijn X_1 en X_2 zodanig dat de correlatie tussen beide gelijk aan 0 is.

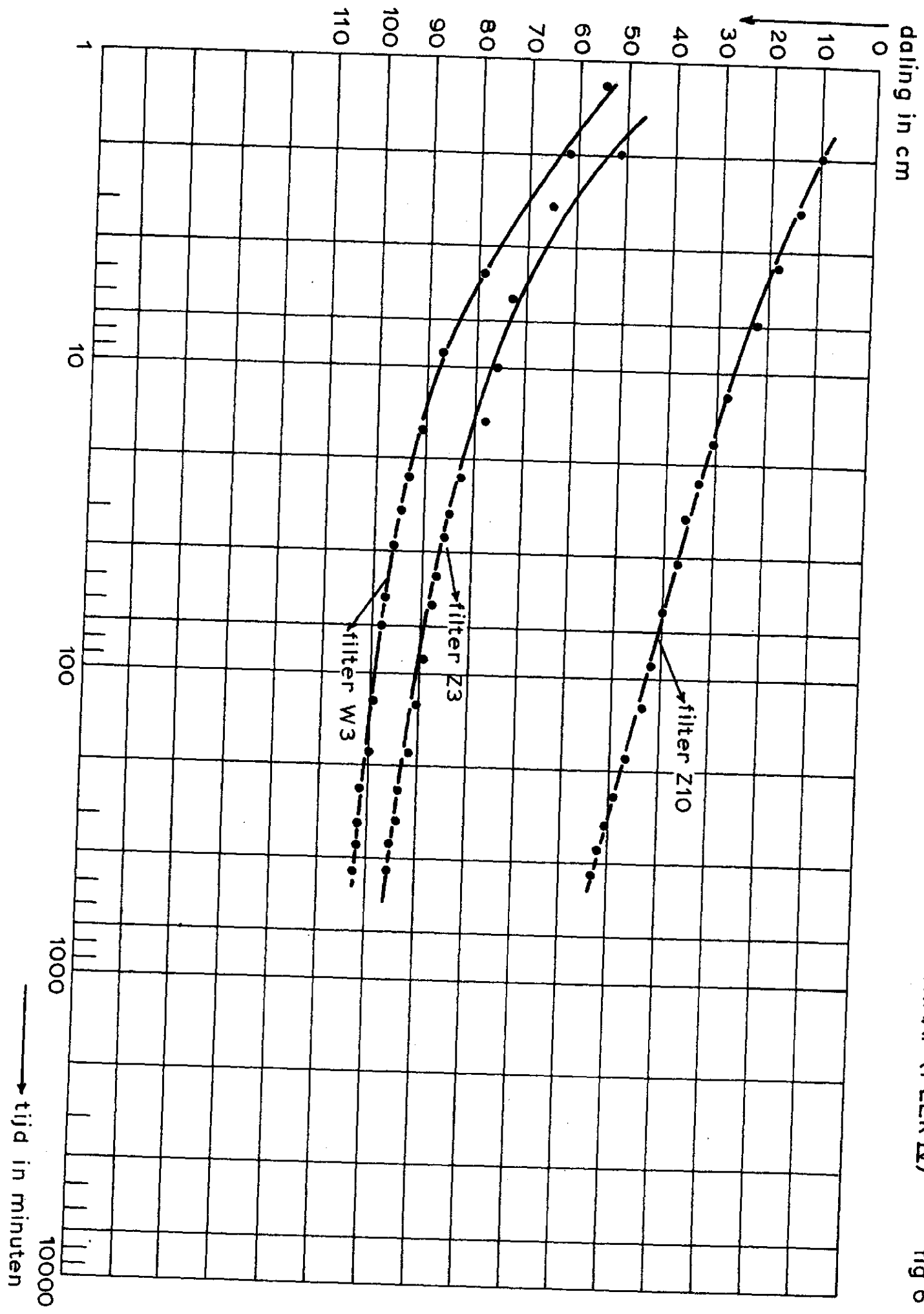
De kolom X_3 is verkregen door de rangnummers 1 tot en met 25 in willekeurige volgorde te plaatsen.

De matrix van correlatie-coëfficiënten ziet er nu als volgt uit:



TIJD-POTENTIAL KROMMEN VAN DE PEILPUTTEN IN DE ZUID- EN WESTRAAI (PLEK IV)

fig 6



$$\begin{array}{c}
 x_1 \\
 x_2 \\
 x_3
 \end{array}
 \begin{pmatrix}
 & x_1 & x_2 & x_3 \\
 & 1 & 0 & -.24 \\
 & 0 & 1 & -.21 \\
 & -.24 & -.21 & 1
 \end{pmatrix}$$

Als symbool voor de correlatie tussen X_1 en X_2 wordt gebruikt $r_{x_1 x_2}$ of korter r_{12} . Voor het gegeven geval zijn de significantie-niveaus voor 25 waarnemingen:

$$r = .396 \text{ met } \alpha = 5\%$$

$$r = .505 \text{ met } \alpha = 1\%$$

De partiële correlatie tussen X_1 en X_2 , de correlatie dus die bestaat wanneer de invloed van X_3 wordt geëlimineerd bedraagt $-.05$. Het symbool hiervoor is $r_{x_1 x_2 \cdot x_3}$ of $r_{12 \cdot 3}$.

Van de in de bijlage 1 gegeven waarden zijn nieuwe reeksen afgeleid waarbij eens is nagegaan wat het effect van het werken met verhoudingsgetallen is. De vraag doet zich voor of door correlatie-rekening toch weer een inzicht in de samenhang tussen de oorspronkelijke variabelen afzonderlijk kan worden verkregen.

Stel dat bijvoorbeeld wordt overgegaan op

$$\begin{aligned}
 Y_1 &= X_1/X_3 \\
 Y_2 &= X_2/X_3
 \end{aligned}$$

dan wordt de correlatie-matrix tussen Y_1 en Y_2

$$\begin{pmatrix}
 1 & .52 \\
 .52 & 1
 \end{pmatrix}$$

De 'oorspronkelijke' reeksen zijn dus thans hoog gecorreleerd en significant met $\alpha = 1\%$.

Geprobeerd kan nu worden de correlatie tussen Y_1 en Y_2 te berekenen bij constant houden van X_3 , teneinde te trachten de invloed van x_3 te elimineren.

Het volgende resultaat werd verkregen:

$$r_{x_1 x_2} = .0$$

$$r_{y_1 y_2 \cdot x_3} = .48$$

$$r_{x_1 x_2 \cdot x_3} = -.05$$

$$r_{y_1 y_2} = .52$$

$$r_{y_1 y_2 \cdot (1/x_3)} = .30$$

Voor de partiële correlaties in de laatste kolom geldt voor de significantie

$$.404 \text{ met } \alpha = 5\%$$

$$.515 \text{ met } \alpha = 1\%$$

De eliminatie van de variabele X_3 heeft niet veel effect gehad. De partiële correlatie is $.48$ tegen de gewone correlatie $.52$. Een duidelijker effect heeft eliminatie van de invloed van de reciproke waarde van X_3 . De partiële correlatie is nu $.30$ en niet meer significant. In een dergelijke uitkomst kan men de motivering zoeken over te gaan op

$$X_1 = Y_1 \cdot X_3$$

en

$$X_2 = Y_2 \cdot X_3$$

waarmee de oorspronkelijke variabelen zijn terugverkregen.

Overigens is hiermee niet aangetoond of bewezen dat op deze wijze steeds het juiste inzicht wordt verkregen.

DE JONGE (1960) vermeldt nog een voorbeeld van SNEDECOR waarin eveneens twee reeksen gegevens (X en Y) een correlatie-coëfficiënt gelijk aan 0 hebben (bijlage 2). De correlatie van X met Z = X + Y is echter hoog en bedraagt .94. Ook dit is een voorbeeld van een hoge correlatie die ontstaat door van de basisreeksen een nieuwe reeks af te leiden en deze weer met de oorspronkelijke gegevens te correleren. Indentiek met dit geval is het berekenen van een gemiddelde om dan deze gemiddelden als referentiereeks te gebruiken.

Van belang is te constateren dat het combineren van waarnemingsreeksen veelal leidt tot een sterke verhoging van de correlatie. Veelal zal dus de nodige voorzichtigheid in acht moeten worden genomen bij het interpreteren van correlaties berekend uit afgeleide reeksen.

Bij factor- respectievelijk aspectenanalyse tracht men eveneens een zo zuiver mogelijk beeld te verkrijgen door noemers van verhoudingsgetallen steeds als variabele op te nemen.

In de volgende paragrafen zullen nog een aantal formules die de betrekkingen tussen de correlatie-coëfficiënten weergeven worden besproken en toegelicht.

Correlatie-coëfficiënten, formulering en afleiding

Naast de gewone correlatie-coëfficiënt die de correlatie tussen twee variabelen X_1 en X_2 weergeeft (r_{12}) wordt toegepast de multipele correlatie-coëfficiënt die inlichtingen verstrekt over de samenhang tussen meer variabelen (bijvoorbeeld $r_{1,235}$) in welk laatste geval wordt bedoeld de correlatie tussen de variabele X_1 met de variabelen X_2, X_3 en X_5 gezamenlijk. Zie hiervoor bijvoorbeeld NOTA 134 pagina 7 en NOTA 147 pagina 12.

Voor de gevallen waarin wordt gevraagd de invloed van een derde variabele te elimineren wordt de partiële correlatie-coëfficiënt gebruikt. Deze geeft aan welke correlatie tussen de variabelen X_1 en X_2 bijvoorbeeld bestaat indien de invloed van X_3 wordt geëlimineerd ($r_{12,3}$).

Alvorens speciaal nader op de partiële correlatie in te gaan worden enkele formuleringen gegeven.

Stel gegeven een aantal metingen van de variabelen X_1, X_2 en X_3 . De correlatie tussen bijvoorbeeld X_1 en X_2 wordt dan gegeven door

$$r_{12} = \frac{\sum (X_1 - \bar{X}_1)(X_2 - \bar{X}_2)}{\sqrt{\sum (X_1 - \bar{X}_1)^2} \sqrt{\sum (X_2 - \bar{X}_2)^2}}$$

Wordt overgegaan op variabelen die zijn uitgedrukt ten opzichte van hun gemiddelde waarde door de transformatie

$$x_i = X_i - \bar{X}_i$$

dan komt er

$$r_{12} = \frac{\sum x_1 x_2}{\sqrt{\sum x_1^2} \sqrt{\sum x_2^2}}$$

In deze zin zullen variabelen steeds worden opgevat, ook in de figuren wordt aangenomen dat de herleiding op het gemiddelde heeft plaatsgevonden. In meetkundige termen betekent dit dat de voorgestelde ruimte die is welke loodrecht staat op het niveau, dat is de vector (1, 1, 1, ..., 1).

De correlatie-coëfficiënten worden samengevat in de correlatie matrix

$$\begin{pmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{pmatrix}$$

waarin $r_{11} = r_{22} = r_{33} = 1$ en $r_{ij} = r_{ji}$

De determinant van deze matrix is dan

$$R = \begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{vmatrix}$$

De cofactoren van r_{ij} aangegeven met R_{ij} zijn dan de van het juiste teken voorziene minoren. Bij ontwikkeling naar de eerste kolom ontstaat

$$R_{11} = + \begin{vmatrix} 1 & r_{23} \\ r_{23} & 1 \end{vmatrix}$$

$$R_{12} = - \begin{vmatrix} r_{12} & r_{13} \\ r_{23} & 1 \end{vmatrix}$$

$$R_{13} = + \begin{vmatrix} r_{12} & r_{13} \\ 1 & r_{23} \end{vmatrix}$$

zodat $R = R_{11} + r_{12} R_{12} + r_{13} R_{13}$

De meetkundige voorstelling van deze matrix is een wat merkwaardige tengevolge van het feit dat de hoofddiagonaal uit 1-en bestaat en de matrix symmetrisch is (figuur 1 en 2).

Het eindpunt van de eerste basis-vector bevindt zich dus steeds in de ruimte (lijn, vlak,)
loodrecht op de eerste as op afstand 1 van de oorsprong. Noem deze ruimte bijvoorbeeld vlak 1. Naarmate de correlatie van variabele 1 met variabele 2 groter is zal het eindpunt van de eerste basis-vector dichter bij het snijpunt, de snijlijn,, van vlak 1 met vlak 2 liggen.

In figuur 1 wordt voor een correlatie-matrix met twee variabelen een aantal situaties weergegeven. In figuur 2 eenzelfde voorstelling voor drie variabelen. De determinant heeft de waarde van het oppervlak, (respectievelijk de inhoud), van het parallellogram, (respectievelijk parallelpipidum), op de basisvectoren.

Een onderdeterminant kan nu als volgt worden weergegeven:

het schrappen van een kolom betekent het buiten beschouwing laten van de overeenkomstige vector, het schrappen van een rij betekent het verlagen van de dimensie door het buiten beschouwing laten van de betreffende kentallen. De basis-vectoren worden hierdoor op de overgebleven ruimte geprojecteerd (figuur 3).

De voorstellingswijze gebaseerd op de correlatie-matrix speelt een belangrijke rol bij de factor- of aspectenanalyse, welke analyse juist op deze matrix wordt uitgevoerd.

Voor het in deze nota aan de orde gestelde onderwerp kan de correlatie-matrix en -determinant dienen om een algemene formule van de multipele- en partiële correlatie-coëfficiënt te geven, ongeacht het aantal variabelen dat in de beschouwing is opgenomen. De gevraagde grootheden kunnen dan in de enkelvoudige correlatie-coëfficiënten worden uitgedrukt.

De bewijzen die op deze voorstellingswijze steunen worden gegeven in KENNY and KEEPING, deel II, (1959) op pagina 339 en volgende.

De partiële correlatie-coëfficiënt kan als volgt worden berekend:

Zij gegeven de variabelen x, y en z. Gevraagd wordt de correlatie tussen x en y, onder eliminatie van de invloed van z. Worden de meetuitkomsten van x, y en z op de gewone wijze als vectorvoorstelling

weergegeven (figuur 4) dan betekent de vraag dat de correlatie-coëfficiënt moet worden berekend tussen die vectoren x_z en y_z die ontstaan nadat x en y op de ruimte $\perp z$ zijn geprojecteerd. In dat geval geldt namelijk

$$\left. \begin{array}{l} x_z \perp z \\ y_z \perp z \end{array} \right\} x_z \text{ en } y_z \text{ onafhankelijk van } z$$

Vooropgesteld wordt dat de lengten van x , y en z gelijk aan 1 zijn gemaakt zodat $xx = yy = zz = 1$. De projecties worden nu als volgt berekend (figuur 4)

$$x_z = x - \lambda z$$

met de eis

$$(x - \lambda z) \perp z$$

dus

$$xz - \lambda zz = 0$$

waaruit volgt

$$\lambda = \frac{xz}{zz} = xz \quad (zz = 1)$$

zodat

$$x_z = x - (xz)z$$

Analoog

$$y_z = y - (yz)z$$

De partiële correlatie tussen x en y is nu dus de gewone correlatie tussen x_z en y_z zodat, in vectoren:

$$\begin{aligned} r_{xy.z} &= \frac{\{x - (xz)z\} \cdot \{y - (yz)z\}}{\sqrt{\{x - (xz)z\}^2 \{y - (yz)z\}^2}} \\ &= \frac{xy - (yz)(xz)}{\sqrt{\{1 - (xz)^2\} \{1 - (yz)^2\}}} \end{aligned}$$

Daar $xx = yy = zz = 1$ en de variabele reeds zijn gereduceerd ten opzichte van hun gemiddelden kan dit worden geschreven als

$$r_{xy.z} = \frac{r_{xy} - r_{yz} r_{xz}}{\sqrt{(1 - r_{xz}^2)(1 - r_{yz}^2)}}$$

of met indices

$$r_{12.3} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}} \quad (1)$$

Zijn er meer variabelen in het geding dan kan de bewerking verder worden doorgevoerd en ontstaat er bijvoorbeeld:

$$r_{12.34} = \frac{r_{12.4} - r_{13.4} r_{23.4}}{[(1 - r_{13.4}^2)(1 - r_{23.4}^2)]^{1/2}}$$

De partiële correlatie $r_{12.3}$ kan ook worden opgevat als de correlatie tussen x_1 en x_2 indien x_3 constant wordt gehouden. Op deze wijze gedefinieerd hangt de partiële correlatie-coëfficiënt echter af van de gekozen constante waarde van x_3 .

Zal de partiële correlatie-coëfficiënt onafhankelijk van het niveau van x_3 zijn dan moet aan een aantal voorwaarden omtrent lineariteit en het constant zijn van de standaard afwijkingen zijn voldaan (KENNY and KEEPING, deel II, pagina 352). In de praktijk zal aan deze voorwaarden veelal slechts bij benadering zijn voldaan wat inhoudt dat $r_{12.3}$ een soort gemiddelde waarde zal zijn voor de correlaties bij alle x_3 niveaus (idem), (FISHER, 1958, pagina 188) en (SNEDECOR, 1957, pagina 430).

Streng genomen geldt nog de eis dat x_1 en x_2 normaal zijn verdeeld, voor x_3 is deze eis niet noodzakelijk (FISHER, 1958, pagina 188). Per definitie laat men de gegeven formule ook wel gelden voor de partiële correlatie-coëfficiënt voor anders verdeelde grootheden (KENDALL and STUART, deel 2, pagina 318).

In de gegeven numerieke voorbeelden is er niet voor gezorgd dat aan bovengenoemde eisen is voldaan. De vermelde resultaten moeten dan ook worden beschouwd als een kwantitatieve weergave van het meetkundig model waarin een correlatie-coëfficiënt de betekenis van de cos van een hoek krijgt.

Eenvoudiger worden de gegeven formules nog met behulp van de determinant R geschreven. Een samenvatting volgt hieronder.

MULTIPELE CORRELATIE

2 variabelen

$$R = \begin{vmatrix} 1 & r_{12} \\ r_{12} & 1 \end{vmatrix}$$

$$R = 1 - r_{12}^2$$

dus

$$r_{12}^2 = 1 - R = 1 - \frac{R}{R_{11}}$$

3 variabelen

$$R = \begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{vmatrix}$$

$$r_{1.23}^2 = 1 - \frac{R}{R_{11}} = \frac{r_{12}^2 - r_{13}^2 - 2r_{12}r_{13}r_{23}}{1 - r_{23}^2}$$

$$r_{2.31}^2 = 1 - \frac{R}{R_{22}}$$

$$r_{3.12}^2 = 1 - \frac{R}{R_{33}}$$

algemeen voor k variabelen

$$r_{1.23 \dots k}^2 = 1 - \frac{R}{R_{11}}$$

algemeen voor < k variabelen

zou men kunnen definiëren

$$r_{1.234}^2 = 1 - \frac{R_{55.66 \dots kk}}{R_{11.55.66 \dots kk}}$$

waarin dan $R_{55.66}$ de cofactor is die uit R ontstaat, door daarin achtereenvolgens de 5e rij en kolom en 6e rij en kolom te schrappen. Daar dit steeds cofactoren zijn van diagonaalelementen blijft het teken ongewijzigd.

Daar correlatie-coëfficiënten worden berekend door de variabelen op niveau te herleiden (te projecteren op de ruimte loodrecht op het gemiddelde) wordt voor elke variabele die in de berekening is opgenomen bij de toetsing 1 vrijheidsgraad in mindering gebracht (vermindering van 1 dimensie) (FISHER, 1958, pagina 258).

PARTIELE CORRELATIE

2 variabelen

$$r_{12} = - \frac{-r_{12}}{1.1} = - \frac{R_{12}}{[R_{11} R_{22}]^{1/2}}$$

3 variabelen

$$r_{12.3} = - \frac{R_{12}}{[R_{11} R_{22}]^{1/2}}$$

algemeen voor k variabelen

$$r_{12.34 \dots k} = - \frac{R_{12}}{[R_{11} R_{22}]^{1/2}}$$

algemeen voor k variabelen kan men weer definiëren

$$r_{12.34} = - \frac{R_{12.55.66 \dots kk}}{\left[R_{11.55.66 \dots kk} R_{22.55.66 \dots kk} \right]^{1/2}}$$

Met deze symbolen geldt nog in het bijzonder voor vier variabelen als

$$R = \begin{vmatrix} 1 & r_{12} & r_{13} & r_{14} \\ r_{12} & 1 & r_{23} & r_{24} \\ r_{13} & r_{23} & 1 & r_{34} \\ r_{14} & r_{24} & r_{34} & 1 \end{vmatrix}$$

$$r_{12.34} = - \frac{R_{12}}{\left[R_{11} R_{22} \right]^{1/2}}$$

$$r_{12.3} = - \frac{R_{12.44}}{\left[R_{11.44} R_{22.44} \right]^{1/2}}$$

$$r_{12} = - \frac{R_{12.33.44}}{\left[R_{11.33.44} R_{22.33.44} \right]^{1/2}}$$

Bij de berekening van partiële correlaties wordt bovendien nog geprojecteerd op de ruimte loodrecht op de te elimineren variabelen. Voor elk wordt hiervoor bij het toetsen 1 vrijheidsgraad (dimensie) in mindering gebracht (FISHER, 1958, pagina 196).

Voor twee variabelen hebben de gegeven algemene formules minder betekenis daar dan het onderscheid tussen gewone, partiële en multiële correlatie irrelevant is. De formules blijven echter ook voor deze gevallen geldig.

Bijlage 1.

X_1	X_2	X_3	$Y_1 = \frac{X_1}{X_3}$	$Y_2 = \frac{X_2}{X_3}$	$\frac{1}{X_3}$
-12	92	5	-2.400	18.400	.200
-11	69	20	-.550	3.450	.050
-10	48	4	-2.500	12.000	.250
-9	29	23	-.391	1.261	.043
-8	12	22	-.364	.545	.045
-7	-3	10	-.700	-.300	.100
-6	-16	25	-.240	-.640	.040
-5	-27	14	-.357	-1.929	.071
-4	-36	2	-2.000	-18.000	.500
-3	-43	18	-.167	-2.389	.056
-2	-48	13	-.154	-3.692	.077
-1	-51	19	-.053	-2.684	.053
0	-52	11	.0	-4.727	.091
1	-51	24	.042	-2.125	.042
2	-48	15	.133	-3.200	.067
3	-43	17	.176	-2.529	.059
4	-36	7	.571	-5.143	.143
5	-27	9	.556	-3.000	.111
6	-16	8	.750	-2.000	.125
7	-3	3	2.333	-1.000	.333
8	12	21	.381	.571	.048
9	29	1	9.000	29.000	1.000
10	48	16	.625	3.000	.062
11	69	12	.917	5.750	.083
12	92	6	2.000	15.333	.167

Bijlage 2.

X	Y	X+Y=Z
32	18	50
31	13	44
23	22	50
24	17	41
44	11	55
53	19	72
9	16	25
35	23	58
33	23	56
31	18	49

$$r_{xy} = 0$$

$$r_{xz} = .94$$

(DE JONGE, 1960, deel II, pagina 534)

LITERATUUR

- FISHER, R.A., 1958 - Statistical Methods for Research Workers, London (I.C.W. 11/103)
- JONGE, H. DE, 1960 - Inleiding tot de medische statistiek, deel II, Leiden (I.C.W. 11/102 (2))
- KAMMIL, L.P., 1962 - Lineaire regressie. I.C.W.-nota 134
- KENDALL, M.G. and A.STUART, 1961 - The advanced theory of statistics. Deel II. Inference and relationship
Londen (I.C.W. 11/114)
- KENNY, J.F. and E.S.KEEPING, 1959 - Mathematics of Statistics, deel II, New York (I.C.W. 11/35 (2))
- SNEDDECOR, G.W., 1957. Statistical Methods. Iowa (I.C.W. 11/44)
- STOL, Ph.Th., 1962 - Een meetkundige toelichting op het oplossen van normaalvergelijkingen. I.C.W.-nota
147

fig.1. DE BASISVECTOREN VAN DE CORRELATIE-MATRIX (2 variabelen)

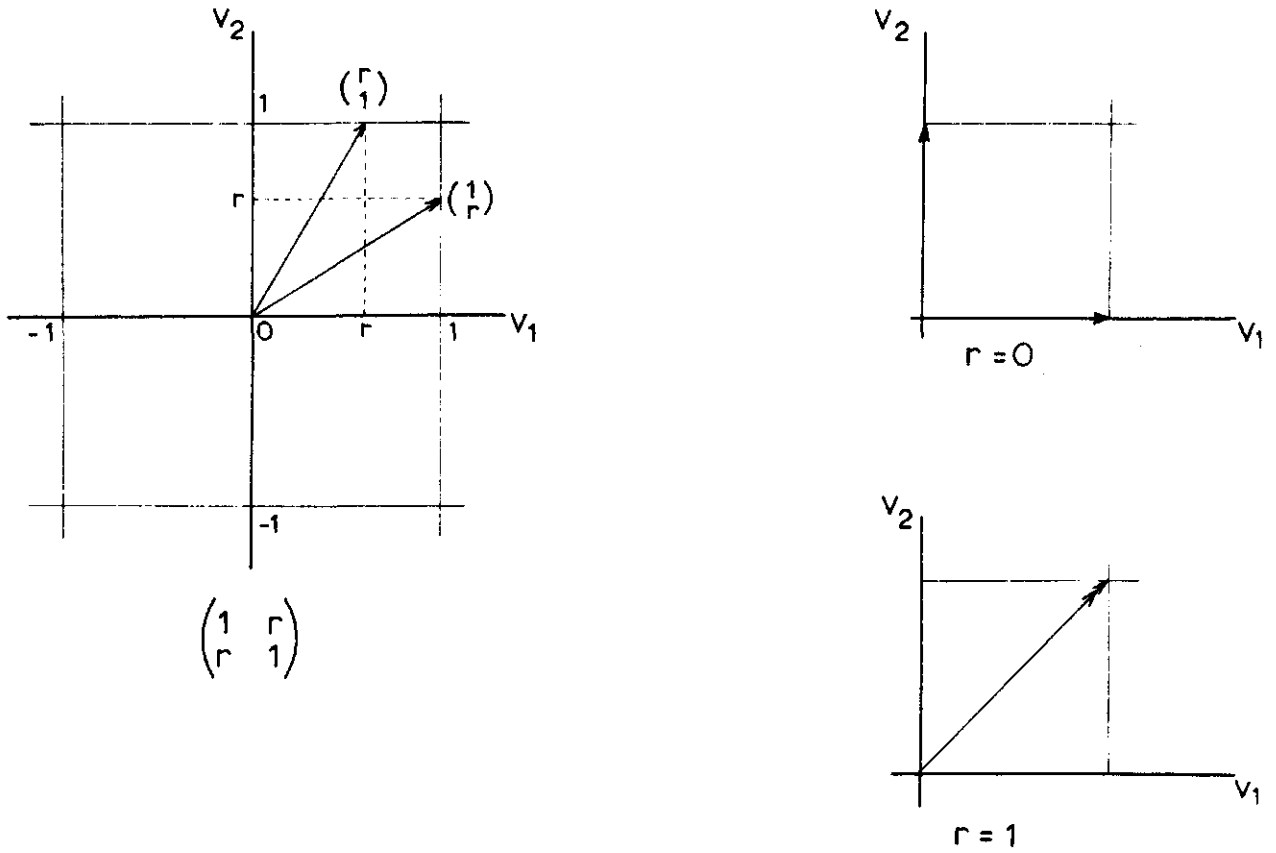


fig.2 DE BASISVECTOREN VAN DE CORRELATIE -MATRIX (3 variabelen)

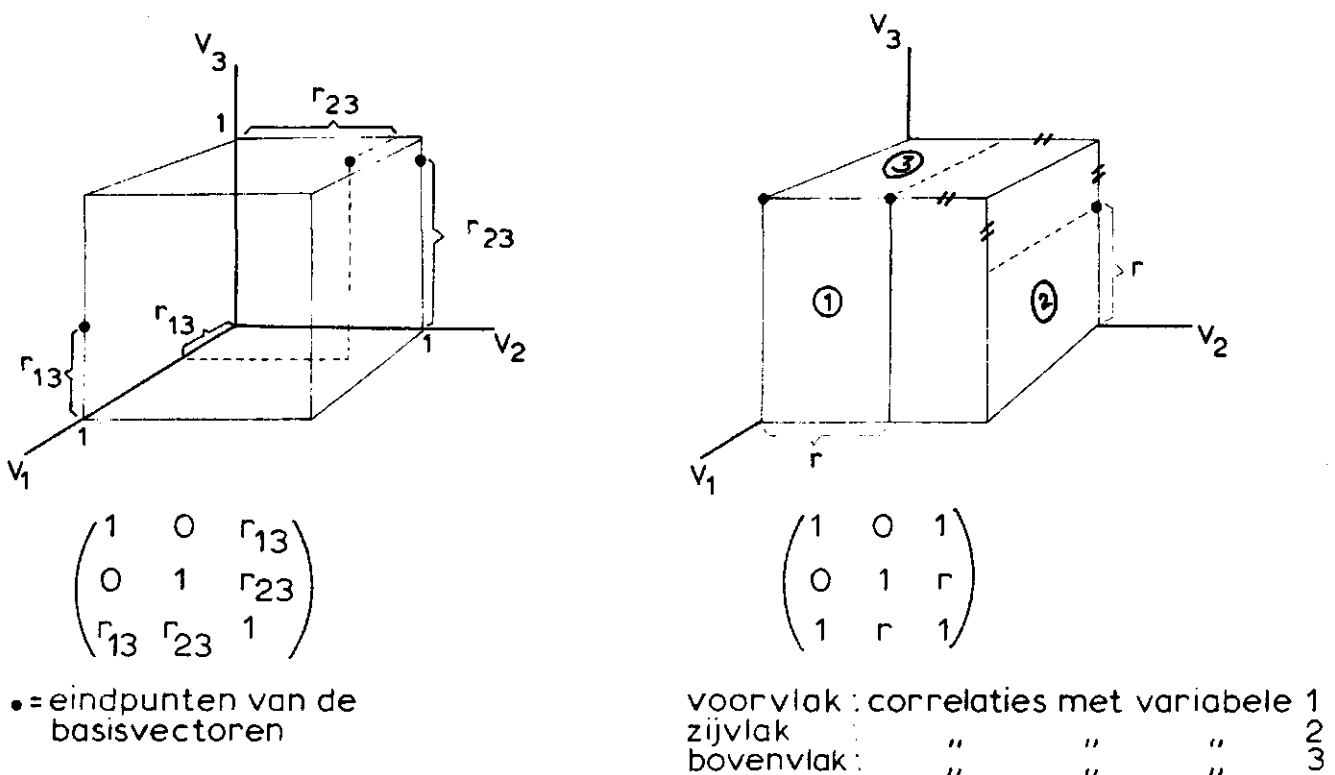


fig.3. HET SCHRAPPEN VAN EEN RIJ EN EEN KOLOM UIT DE CORRELATIE-MATRIX

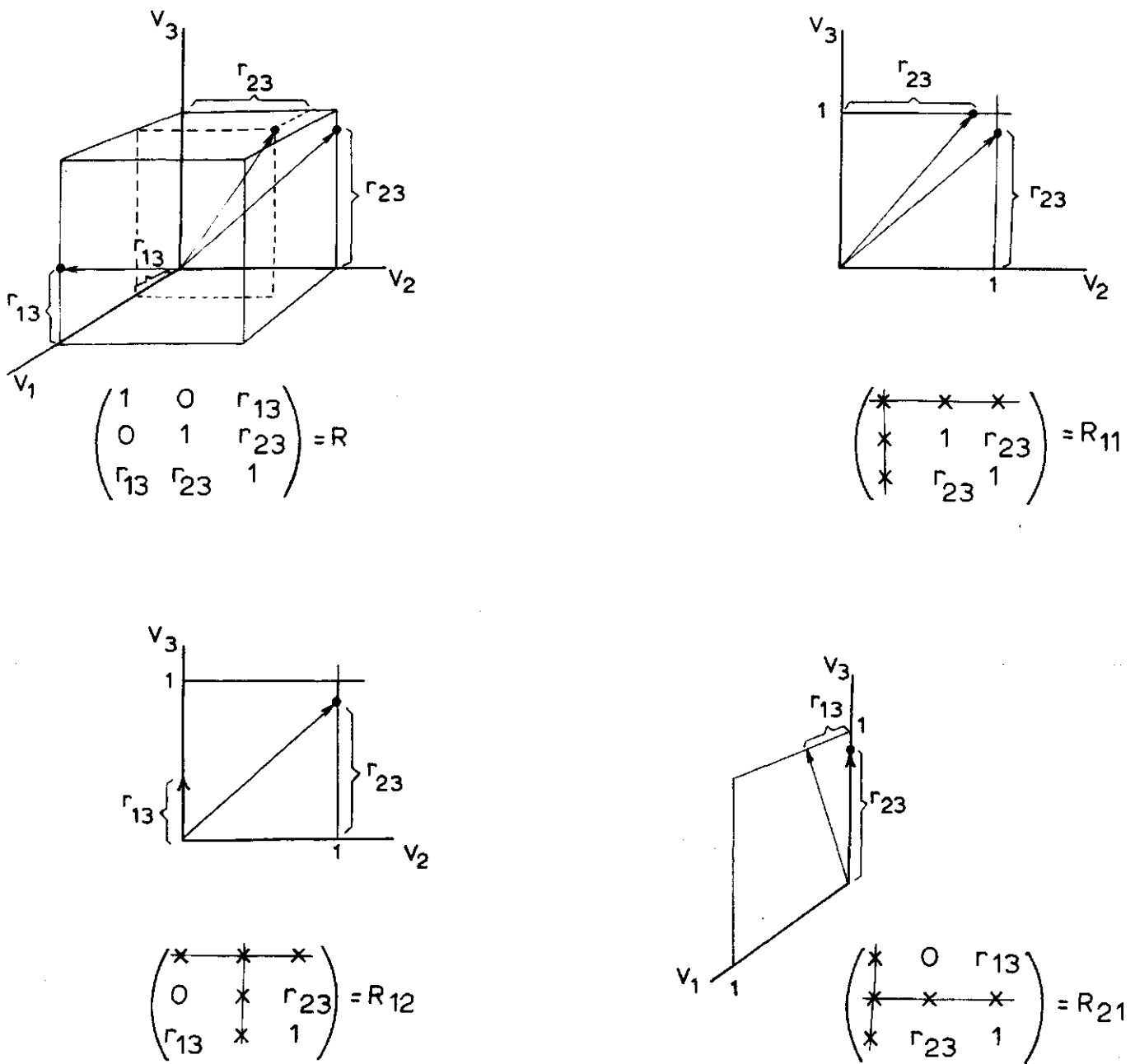


fig.4. PARTIËLE CORRELATIE

