

K O N I N K L I J K N E D E R L A N D S  
M E T E O R O L O G I S C H I N S T I T U U T

D e B i l t

WETENSCHAPPELIJK RAPPORT

W.R. 75-1

P.J. Rijkoort

Een algemeen lineair statistisch model  
voor de bepaling van relaties  
tussen een aantal grootheden.

De Bilt, 1975

Publikationsnummer: W.R. 75-1 (Stat. Bur.)

U.D.C.: 519.2

EEN ALGEMEEN LINEAIR STATISTISCH MODEL VOOR DE  
BEPALING VAN RELATIES TUSSEN EEN AANTAL GROOTHEDEN

---

---

	pag.
Algemene inleiding	1
<u>1</u> Theoretische behandeling van een algemeen statistisch "multivariate" model	4
1.1. Inleiding	4
1.2. De beste schattingen van de parameters	4
1.3. Toetsing van hypothesen t.a.v. de parameters	7
1.4. Matrixvoorstelling van het model	10
1.5. De frequentieverdeling van de toetsingsgrootheid	14
1.6. Een aequivalente formulering van de F-toets	17
<u>2</u> Enkele voorbeelden	18
2.1. Enkelvoudige variantie-analyse	18
2.1.1. Inleiding	18
2.1.2. Kleinste kwadratische oplossing	19
2.1.3. Formele matrixoplossing	21
2.1.4. Numeriek voorbeeld	27
2.2. Regressie-analyse	29
2.2.1. Toetsing van kwadratische regressie	29
2.2.2. Numeriek voorbeeld	33
<u>3</u> Een bewerkingstechniek voor "multivariate" analyse	35
3.1. Inleiding	35
3.2. Orthogonalisatie	37
3.3. De "Sweep" operator	41
<u>4</u> Toepassing van de "Sweep" operator techniek bij "optimalisering" van multiple lineaire regressie	46
4.1. Inleiding	46
4.2. Algemene beschrijving van de "optimaliserings" methode	47
4.3. Toetsing tijdens het "optimaliserings" procedé	48
4.4. Enkele bijzonderheden in het artikel "Multiple regression analysis" bij M.A. Efroymsen	51

-----  
-----

## Samenvatting.

Wetenschappelijk onderzoek is voor een groot deel het zoeken naar relaties tussen één bepaald verschijnsel en een aantal andere verschijnselen.

Mathematisch geformuleerd, hoe ziet de functie eruit in:  $y_j = f(x_{ij}; \alpha_r)$  ( $i = 1, \dots, k; r = 1, \dots, t; j = 1, \dots, N$ ) waarin  $y_i$  resp.  $x_{ij}$  kwantificaties van de verschijnselen voorstellen en  $\alpha_r$  parameterwaarden zijn. Het probleem is in zijn algemeenheid slechts oplosbaar als men ten aanzien van de vorm van de functie een bepaalde onderstelling invoert, waarna het verder gaat om de schatting van de waarden van de parameters  $\alpha_r$  en het geven van een statistische uitspraak betreffende de realiteit van de veronderstelde relatie.

In die gevallen, waarin  $r = i$  en  $t = k$  kan  $f$  vaak als een lineaire functie in de  $\alpha_i$  worden beschouwd:

$$y_j = \sum_{i=1}^k \alpha_i x_{ij}.$$

In principe zijn in dit model zeer veel relaties samengevat. Ten eerste omdat de  $x_{ij}$  onderling afhankelijk mogen zijn, bijvoorbeeld  $x_{3j} = x_{1j} x_{2j}^2$  enz. Ten tweede omdat de  $x_{ij}$  zowel continue variabelen kunnen voorstellen als ook z.g. indicatorvariabelen. Dit laatste wil zeggen, dat slechts de waarde 1 of 0 (voor iedere  $j$  is één  $x_{ij} = 1$  en de overigen zijn 0) worden aangenomen, wat er op neer komt, dat de  $y_j$  in klassen  $i$  worden ingedeeld. In geval van continue variabelen heet de statistische behandeling regressie-analyse; in geval van indicatorvariabelen variantie-analyse. De gemengde vorm heet covariantie-analyse.

Voor de oplossing van het hier aangeduide probleem zijn mathematische-statistische methoden ontwikkeld die o.a. te vinden zijn in de boeken van Scheffé en van Dempster (zie literatuur).

In dit rapport wordt de algemene theorie volgens bovengenoemde boeken zo kort mogelijk behandeld in hoofdstuk 1, terwijl daarna in 2 numerieke toepassingen worden gegeven. De hoofdstukken 3 en 4 behandelen een bijzonder geval van regressie-analyse namelijk een methode om "met een computer" tot de optimale keuze van  $r$  variabelen  $x_i$  uit de totaal  $k$  variabelen te komen ter verklaring van de variatie van  $y$ . Deze methode is te vinden in een artikel van Efroymson.

De bedoeling van dit rapport is de behandelde materie toegankelijk te maken voor diegenen, die niet de tijd hebben de uitvoerige literatuur door te werken.

### Summary.

Scientific research is mainly concerned with attempts to find relations between one phenomenon and different other phenomena.

Mathematically formulated, what can be said about:  
 $y_j = f(x_{ij}; \alpha_r)$  ( $i = 1, \dots, k; r = 1, \dots, t; j = 1, \dots, N$ ) in which  $y_j$  resp.  $x_{ij}$  represent the phenomena and the  $\alpha_r$  are parameters to be chosen. In general a solution is only possible if some form of the function is postulated. The problem is reduced then to finding estimates of the  $\alpha_r$  and to giving a statistical judgement with regard to the reality of the relation so formed.

If  $r = i$  and  $t = k$  it is often possible to use a linear relation (linear with regard to the  $\alpha_i$ ):

$$y_j = \sum_{i=1}^k \alpha_i x_{ij}.$$

In principle many sorts of relations are contained in this model. Firstly because the  $x_{ij}$  need not to be independent of each other, for instance  $x_{3j} = x_{1j}x_{2j}^2$  ect. may be possible and secondly because the  $x_{ij}$  can be continuous variables as well as so called, indicator variables. The latter means that for each  $j$  only  $x_{ij} = 1$  and all other  $x_{ij} = 0$ ; so that the  $y_j$  are arranged in different classes. The statistical treatment is called regression analysis in the case of continuous variables and analysis of variance in the case of indicator variables. Analysis of covariance is a mixed form.

The mathematical statistical solution of the problem can be found in different textbooks for instance in those written by Scheffé and by Dempster (see literature).

In the present report the general theory according to Scheffé and Dempster is treated shortly in chapter 1, while chapter 2 gives some numerical examples. The chapter 3 and 4 contain a treatment of a special regression problem viz. the solution (which a computer) of the problem of finding an optimum choice of only  $r$  from the  $k$  variables  $x_j$  to explain the variation of  $y$ . This method has been developed by Efroymsen, using results of Orden (see literature).

This report is meant to give an introduction to the method to be used by investigators who do not have time to study the books of Scheffé, Dempster and other themselves.

Een algemeen lineair statistisch model voor de bepaling van relaties tussen een aantal grootheden.

Dr. P.J. Rijkoort

Algemene inleiding.

Een groot gedeelte van het onderzoek dat in wetenschappelijke instellingen wordt uitgevoerd, bestaat in het zoeken naar het verband tussen één bepaalde grootheid ( $y$ ) en één of meer grootheden ( $x_i$ ).

De statistische behandeling van dergelijke problemen, voor zover die in dit verslag zal worden uiteengezet, zal zich beperken tot die gevallen waarbij het verband de volgende lineaire vorm heeft:

$$y = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \quad (0.1)$$

De lineariteit van het model is beperkt tot de parameters  $\beta_i$  en niet wezenlijk tot de  $x_i$ , daardoor kunnen met dit model meer relaties beschreven worden dan men in eerste instantie zou denken. Naast die gevallen waarin inderdaad een rechtstreekse lineaire relatie bestaat, kan men voor gevallen waarbij dit niet zo is, door transformaties toch tot bovenstaande vorm komen; voor een exponentieel verband bijvoorbeeld kan dit met een logaritmische transformatie. Bovendien kan men een gecompliceerde vorm van  $x_i$ 's opvatten als een nieuwe  $x_i$ ; het kwadratische verband  $y = ax^2 + bx + c$  kan bijvoorbeeld ook in de vorm (0.1) gezien worden als men  $x^2$  door  $x_1$  en  $x$  zelf door  $x_2$  voorstelt; verder resulteert het opvatten van  $x_3$  als de constante 1 in het opnemen van een constante in de regressievergelijking.



De behandeling van de in het voorgaande ter sprake komende problemen, waarbij de  $x_i$ 's evenals  $y$  continue variabele grootheden zijn, wordt enkelvoudige regressie genoemd als er één  $x$  is, of meervoudige (multiple) regressie als er meerdere  $x_i$ 's zijn.

Het is ook mogelijk een ander soort problemen in de vorm (0.1) te brengen n.l. die waarbij  $y$  een grootheid is die onder verschillende omstandigheden wordt waargenomen en waarbij het er dan om gaat vast te stellen of  $y$  zich verschillend gedraagt onder deze verschillende omstandigheden. Als de verschillende omstandigheden met klassen A, B, C enz. worden aangeduid dan kan gesteld worden :

$$x_{(A)} = 1 ; x_{(B)} = 0 \text{ enz. als } y \in A \text{ (} y \text{ behoort tot A)}$$

$$x_{(A)} = 0 ; x_{(B)} = 1 ; x_{(C)} = 0 \text{ enz. als } y \in B$$

enz.

De grootheid  $x$  wordt nu een "indicator" variabele genoemd. Onze aanpak van dit soort problemen zal blijken tot de bekende "variantie-analyse" te voeren.

In geval er zowel continue als indicator variabelen in een probleem voorkomen, wordt de statistische behandeling in de literatuur met covariantie-analyse aangeduidt.

De oplossing van de gestelde problemen kan men in diverse statistische handboeken vinden. Een vrij algemene behandeling, later speciaal op de variantie-analyse gericht, is te vinden in H. Scheffé: "Analysis of variance" Wiley - New York 1959 - 1967, terwijl een algemene behandeling met matrices en verder toepassing op de regressie-analyse uitvoerig wordt uiteengezet in A.P. Dempster: "Elements of continuous multivariate analysis" - Addison - Wesley - 1969.

Hoewel beide boeken niet bepaald bijzonder moeilijk zijn, is het toch voor onderzoekers die niet speciaal statistisch zijn onderlegd, bezwaarlijk eerst deze boeken te moeten doorwerken, alvorens de resultaten te kunnen toepassen. Dit is vooral zo, omdat de methode van behandeling in deze boeken, zoals vaak in de mathematische literatuur, zo is, dat eerst allerlei stellingen en hulpstellingen worden ontwikkeld en bewezen, waarvan de betekenis pas veel later duidelijk wordt.

Men kan natuurlijk wel volstaan met korte recepten van statistische methoden die in artikelen in de statistische literatuur wel te vinden zijn, maar dan kan de onderzoeker niet beoordelen of er correct gewerkt wordt en bovendien is het bij de interpretatie van de resultaten veel beter als men inzicht heeft in de achtergrond van de gebruikte methode, in het hoe en waarom.

Om deze redenen zal in dit verslag getracht worden een zoveel mogelijk recht op het doel afgaande uiteenzetting van de algemene theorie te geven. Ter toelichting zullen één geval van variantie-analyse en één geval van regressie-analyse worden behandeld. Bovendien zal een artikel waarin multiple regressie als een soort recept wordt behandeld nader worden beschouwd.

Uiteraard moet voor het volgen van de behandeling van de theorie in dit rapport wel enige statistische basiskennis aanwezig zijn en tevens bekendheid met het werken met vectoren in een  $n$ -dimensionale ruimte en met matrices.

In eerste instantie zal het boek van Scheffé worden gevolgd en later Dempster. De notatie zal uniform worden gemaakt.

Dit rapport is bedoeld als basisrapport, waarin alleen de algemene zaken ter sprake komen. Later zal getracht worden in vervolgrapporten op bijzondere gevallen in te gaan.

1. Theoretische behandeling van een algemeen "multivariabele" model

1.1. Inleiding.

In de algemene inleiding is de lineaire vorm (0.1) gegeven voor het onderstelde verband tussen een variabele  $y$ , die afhankelijk variabele zal worden genoemd en variabelen  $x_i$ , die onafhankelijke variabelen heten. (N.B. de uitdrukking onafhankelijk met betrekking tot de  $x_i$ 's is niet helemaal correct, er kan onderlinge afhankelijkheid tussen de  $x_i$  voorkomen; onafhankelijk betekent hier alleen dat het geheel van de  $x_i$ 's primair varieert en dat de  $y$  als gevolg daarvan varieert).

Waarnemingen van  $y$  zullen meestal niet exact aan (0.1) voldoen maar met een zekere fout  $\varepsilon$  behept zijn. Daarom zal als algemeen model, waarvoor de notatie  $\Omega$  zal worden gebruikt, worden ingevoerd:

$$\Omega: y_k = \sum_{i=1}^p \beta_i x_{ik} + \varepsilon_k \quad (1.1.1)$$

We onderstellen dat er  $n$  waarnemingen zijn van  $p + 1$  grootheden  $y$  en  $x_i$ . Verder nemen we aan dat  $\varepsilon_k$  onafhankelijk is van  $\varepsilon_m$  voor iedere  $k$  en  $m$  en dat  $\varepsilon_k$  voor iedere  $k$  een normale verdeling met verwachtingswaarde nul ( $E(\varepsilon_k) = 0$ ) en standaarddeviatie  $\sigma$  volgt. Dit wordt als:  $\varepsilon_k$  is  $N(0, \sigma)$ , genoteerd. De statistische verwachtingswaarde van  $y_k$  is derhalve:

$$E(y_k) = \sum_{i=1}^p \beta_i x_{ik} \quad (1.1.2)$$

1.2. De beste schattingen van de parameters.

Het eerste probleem is nu: Wat zijn voor gegeven waarden van  $y_k$  en  $x_{ik}$  voor  $i = 1, \dots, p$ ;  $k = 1, \dots, n$  de beste schattingen

voor de parameters  $\beta_i$ ? Er zijn voor de oplossing verschillende methoden beschikbaar, die echter alle equivalent zijn. De bekendste is de kleinste kwadratenmethode. Een andere methode berust op het principe van de z.g. waarschijnlijkheidsfunctie voor een gegeven steekproef; maximaliseren van deze functie geeft z.g. "Maximum Likelihood" schatters. In dit rapport wordt er de voorkeur aangegeven de afleiding van de formule voor de schattingen van de  $\beta_i$  te baseren op een geometrische voorstelling. De waarden  $y_k$  ( $k = 1, \dots, n$ ) zullen daartoe als componenten van een vector  $\vec{y}$  in een  $n$ -dim. ruimte  $R_n$  worden opgevat. Evenzo zijn er  $p$ -vectoren  $\vec{x}_i$  ( $i = 1, \dots, p$ ) in  $R_n$  die onafhankelijk variabelen voorstellen. Deze  $p$ -vectoren  $\vec{x}_i$  vormen een deelruimte van  $R_n$  die we met  $V_p$  zullen aangeven en die  $p$  of minder dimensies kan hebben. Het is namelijk mogelijk dat er tussen enkele der  $\vec{x}_i$ 's lineaire afhankelijkheid bestaat en dan zijn er dus minder dan  $p$  onafhankelijke vectoren en heeft de deelruimte  $V$  minder dan  $p$  dimensies. In principe is de algemene theorie zo te ontwikkelen, dat met deze mogelijkheid wordt rekening gehouden. In dit rapport wordt echter gemakshalve ondersteld, dat  $V$  de dimensie  $p$  heeft.

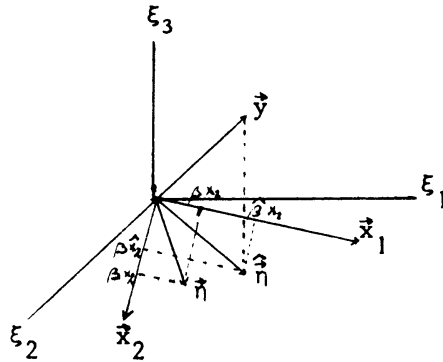
Zoals in (1.2) is gesteld vormen de sommen  $\sum_{i=1}^p \beta_i x_{ik}$  de verwachtingswaarden van  $y_k$ ; deze kunnen ook als componenten van een vector in  $R_n$  beschouwd worden die met  $\vec{\eta}(=E(\vec{y}))$  zal worden aangegeven. Aangezien de componenten van deze vector  $\vec{\eta}$  geheel zijn opgebouwd uit lineaire combinaties van de componenten van de vectoren  $\vec{x}_i$  zal de vector  $\vec{\eta}$  dus ook in  $V_p$  liggen. De vector  $\vec{y}$  stelt een bepaalde steekproef  $(y_1, \dots, y_n)$  voor, een andere steekproef van  $n$ -waarden  $(y_1, \dots, y_n)$  zal een andere vector  $\vec{y}$  leveren met dezelfde verwachtingswaarde  $\vec{\eta}$ .

De vector  $\vec{y}$  is een stochastische vector. De verschillende vectoren  $\vec{y}$  kunnen in principe overal in  $R_n$  liggen, maar zullen in het algemeen geconcentreerd zijn rond  $\vec{\eta}$ .

Aangezien er maar één steekproef  $(y_1, \dots, y_n)$  en dus één vector  $\vec{y}$  beschikbaar is kan  $\vec{\eta}$  niet precies bepaald worden. Het enige wat gedaan kan worden is een vector  $\vec{\eta}_0$  bepalen die zo dicht mogelijk bij  $\vec{y}$  ligt maar wel tot  $V_p$  behoort.  $\vec{y}$  zelf ligt vrijwel zeker niet in  $V_p$ , hoewel dit niet absoluut onmogelijk is; hiervoor zou moeten gelden:  $\epsilon_k = 0$  voor alle  $k$ , en de kans hierop is wel te verwaarlozen; en bovendien is er dan statistisch gezien geen probleem meer.

Nu is de vector in  $V_p$  die zo dicht mogelijk bij  $\vec{y}$  ligt de projectie van  $\vec{y}$  op  $V_p$ , deze zal als schatting van  $\vec{\eta}$  worden genomen en met  $\hat{\vec{\eta}}$  worden aangeduid.

fig. 1.



Ter illustratie fig. 1 waarbij om een tekening mogelijk te maken  $i$  tot 2 is beperkt en  $n$  tot 3. Het assenkruis van  $R_n$  is zo gekozen, dat  $\xi_1$  en  $\xi_2$  het vlak van  $\vec{x}_1$  en  $\vec{x}_2$  aangeven. Het assenkruis (de basis in  $n$ -dim. geometrie), waarop de componenten van  $\vec{y}$  en  $\vec{x}_i$  zijn aangegeven, kan een ander zijn. De vector  $\vec{\eta}$  is in werkelijkheid onbekend.

$$\text{Er geldt nu } \vec{\eta} = \sum_{i=1}^p \beta_i \vec{x}_i \quad \text{en} \quad \hat{\vec{\eta}} = \sum_{i=1}^p \hat{\beta}_i \vec{x}_i, \text{ waarin de } \hat{\beta}_i$$

de gezochte schattingen van de regressiecoëfficiënten zijn. Verder is  $\vec{y} - \hat{\vec{\eta}} \perp V_p$  en dus  $\perp \vec{x}_i$  voor  $i = 1, \dots, p$ . De verschillen  $\vec{y} - \vec{\eta}$  en  $\vec{y} - \hat{\vec{\eta}}$  stellen vectoren voor waarvan de componenten resp.  $\epsilon_k$  en  $\hat{\epsilon}_k$  kunnen worden genoemd. Het kwadraat van de lengten van deze

$$\text{vectoren is dus resp. } \sum_{k=1}^n \epsilon_k^2 \quad \text{en} \quad \sum_{k=1}^n \hat{\epsilon}_k^2.$$

Wegens  $(\vec{y} - \hat{\vec{\eta}}) \perp V_p$  is  $\sum_{k=1}^n \hat{\xi}_k^2$  de kleinste waarde die  $\sum_{k=1}^p \epsilon_k^2$  kan aannemen als  $\vec{\eta}$  over  $V_p$  varieert. Nu is

$$\sum_{k=1}^n \epsilon_k^2 = \sum_{k=1}^n (y_k - \sum_{i=1}^p \beta_i x_{ik})^2; \text{ als dit } \mathcal{S}(\vec{y}, \vec{\beta}) \text{ genoemd wordt dan}$$

zijn de  $\hat{\beta}_i$  die  $\beta_i$  waarden die  $\mathcal{S}(\vec{y}, \vec{\beta})$  minimaal maken, met andere woorden, de geometrische oplossing van het probleem is identiek aan de kleinste kwadratenmethode.  $\mathcal{S}(\vec{y}, \hat{\vec{\beta}})$  zal als  $\mathcal{S}_\Omega$  worden genoteerd.

Het minimaliseren van de kwadraatsom geeft de zogenaamde normaalvergelijkingen. Dus:

$$\frac{\partial \mathcal{S}(\vec{y}, \hat{\vec{\beta}})}{\partial \beta_i} = 0 \text{ geeft}$$

$$\sum_{k=1}^n y_k x_{jk} = \sum_{k=1}^n (x_{jk} \sum_{i=1}^p \hat{\beta}_i x_{ik}) = \sum_{i=1}^p (\hat{\beta}_i \sum_{k=1}^n x_{jk} x_{ik}) \quad (1.2.1)$$

voor  $j = 1, \dots, p$ . Hieruit kunnen de  $\hat{\beta}_i$  worden berekend.

Men kan de geometrische oplossing ook als volgt formuleren: De vector  $\vec{y} \in R_n$  die de waarnemingen voorstelt wordt ontbonden in een vector  $\vec{y}_\parallel$  die  $\parallel V_p$  is en die de beste aanpassing van de waarnemingen aan het model weergeeft ( $\vec{y}_\parallel \equiv \hat{\vec{\eta}}$ ) en een vector  $\vec{y}_\perp$  die  $\perp V_p$  is en waarvan de lengte een maat is voor de afwijking tussen waarnemingen en model  $\vec{y}_\perp \equiv \vec{y} - \hat{\vec{\eta}}$ .

### 1.3. Toetsing van hypothesen t.a.v. de parameters.

In het voorgaande zijn de beste schattingen  $\hat{\beta}_i$  van de  $\beta_i$  gevonden in de onderstelling dat het algemene model volledig geldt. De  $\beta_i$  zullen in het algemeen onderling ongelijk zijn. Het is echter mogelijk dat van de  $\beta_i$  een aantal nul is of dat er enige aan elkaar gelijk zijn.

De vraag is dan of deze onderstellingen ten aanzien van de  $\beta_i$  getoetst kunnen, met andere woorden kunnen de gevonden  $\hat{\beta}_i$ , de onnauwkeurigheid van deze schattingen in aanmerking genomen, in overeenstemming zijn met bepaald veronderstellingen ten aanzien van de  $\beta_i$ . Dit kan wat algemener gesteld worden, namelijk dat er  $q$  onderling onafhankelijke lineaire functies  $\Psi_j$  van de  $\beta_i$ 's zijn. Als nulhypothese  $H_0$  wordt dan ingevoerd:

$$H_0: \Psi_j(\beta_1, \dots, \beta_p) = 0 \quad j = 1, \dots, q \quad (1.3.1)$$

De onderstelling  $\beta_3 = \beta_4$  kan bijvoorbeeld ingevoerd worden met  $\Psi_1 = \beta_3 - \beta_4$  enz.

Onder deze nulhypothese zijn er aan de  $p$   $\beta_i$ 's  $q$  beperkingen opgelegd en dat wil zeggen dat de schattingen  $\hat{\beta}_i$  die aan  $H_0$  voldoen, gevonden worden door de deelruimte  $V_p$  te vervangen door een deelruimte  $V_{p-q}$  waarin  $q$  beperkingen tot uiting komen. Immers als bijvoorbeeld  $\beta_1 = 0$  moet gelden dan moet  $x_1$  buiten beschouwing blijven,  $V_p$  wordt vervangen door een  $V_{p-1}$  opgespannen door  $x_2, \dots, x_p$ ; als  $\beta_2 = \beta_3$  moet worden getoetst, dan moeten  $x_1$  en  $x_2$  vervangen worden door hun som enz.

De beste schattingen van  $\beta_i$  vindt men nu door  $\vec{y}$  te projecteren op  $V_{p-q}$ . Deze projectie zal door  $\hat{\eta}_\omega$  worden voorgesteld, waarbij  $\omega = \Omega \cap H_0$ , dat wil zeggen het algemene model  $\Omega$  in combinatie met de nulhypothese  $H_0$ . Wat de "ware" vector  $\vec{\eta}$  betreft, kan dus nu gesteld worden:

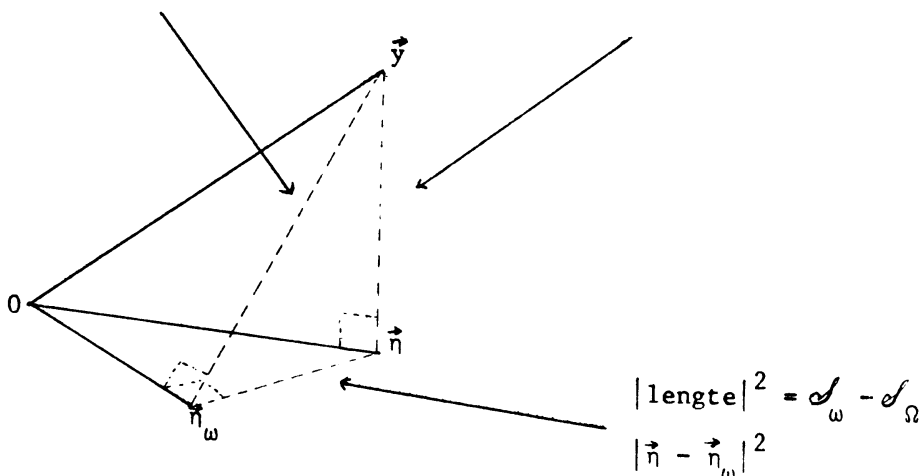
Onder  $\Omega$  geldt:  $\eta \in V_p$

Onder  $\omega$  geldt:  $\eta \in V_{p-q}$ .

Het kwadraat van de lengte van de verschilvector  $\vec{y} - \hat{\vec{n}}_\omega$  stelt de kwadraatsom van de "afwijkingen"  $\hat{\epsilon}$  voor onder  $\omega$ . Deze zal door  $\mathcal{J}_\omega$  worden aangeduid. In fig. 2 zijn de betrokken vectoren voorgesteld,

fig. 2

$$|\text{lengte}|^2 = |\vec{y} - \hat{\vec{n}}_\omega|^2 = \mathcal{J}_\omega \qquad |\text{lengte}|^2 = \mathcal{J}_\Omega = |\vec{y} - \hat{\vec{n}}|^2$$



waarbij dus bedacht moet worden dat  $\vec{y} \in R_n$ ,  $\hat{\vec{n}} \in V_p \in R_n$  en  $\hat{\vec{n}}_\omega \in V_{p-q} \in V_p \in R_n$ . Er geldt dus  $\vec{y} - \hat{\vec{n}} \perp \hat{\vec{n}}$  en ook  $\perp \hat{\vec{n}}_\omega$  omdat  $\vec{y} - \hat{\vec{n}} \perp V_p$  en  $\hat{\vec{n}}_\omega$  ook in  $V$ . Derhalve  $(\vec{y} - \hat{\vec{n}}) \perp (\hat{\vec{n}} - \hat{\vec{n}}_\omega)$  waaruit volgt:  $|\vec{y} - \hat{\vec{n}}_\omega|^2 = |\vec{y} - \hat{\vec{n}}|^2 + |\hat{\vec{n}} - \hat{\vec{n}}_\omega|^2$  of  $|\hat{\vec{n}} - \hat{\vec{n}}_\omega|^2 = \mathcal{J}_\omega - \mathcal{J}_\Omega$ . Er geldt bovendien nog  $\hat{\vec{n}}_\omega \perp (\vec{y} - \hat{\vec{n}}_\omega)$  en  $\perp (\vec{y} - \hat{\vec{n}})$  dus  $\perp (\hat{\vec{n}} - \hat{\vec{n}}_\omega)$ , met andere woorden  $\hat{\vec{n}}_\omega$  is ook de projectie van  $\hat{\vec{n}}$  op  $V_{p-q}$ .

Om nu een uitspraak omtrent het al of niet gelden van de nulhypothese  $H_0$  te doen, suggereert fig. 2 dat  $H_0$  geaccepteerd kan worden als  $\hat{\vec{n}}_\omega$  dicht genoeg bij  $\hat{\vec{n}}$  ligt, met andere woorden, naarmate  $\hat{\vec{n}}_\omega$  verder van  $\hat{\vec{n}}$  verwijderd ligt is het onwaarschijnlijker dat  $H_0$  geldt.



Hierbij moet natuurlijk ook de positie van  $\vec{y}$  ten opzichte van  $\hat{\eta}$  beschouwd worden, immers naarmate  $\vec{y}$  dichter bij  $\hat{\eta}$  ligt zal ook  $\hat{\eta}_\omega$  dichter bij  $\hat{\eta}$  moeten liggen om met gelijkblijvende waarschijnlijkheid een uitspraak te kunnen doen, dit is althans plausibel.

Derhalve zal het quotiënt 
$$\frac{J_\omega - J_\Omega}{J_\Omega} \tag{1.3.2}$$

bepalend moeten zijn voor het al of niet verwerpen van  $H_0$ . Hiertoe is nodig de kans te kunnen aangeven dat dit quotiënt bepaalde numerieke waarden bereikt; de frekwentieverdeling van  $\frac{J_\omega - J_\Omega}{J_\Omega}$  onder de nulhypothese moet worden bepaald. Alvorens hiertoe over te gaan zal eerst nog het probleem en de oplossing, met behulp van matrices worden geformuleerd.

1.4. Matrixvoorstelling van het model.

Als voor de grootheden die in het probleem een rol spelen de volgende notatie wordt ingevoerd

$$\vec{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}; \quad \vec{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}; \quad \vec{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \tag{1.4.1}$$

en het geheel van de  $p \times n$  grootheden  $x_{ik}$  als een matrix

$$X = \begin{pmatrix} x_{11} & \dots & x_{1n} \\ x_{21} & \dots & x_{2n} \\ \vdots & & \vdots \\ x_{p1} & \dots & x_{pn} \end{pmatrix} \tag{1.4.2}$$

wordt voorgesteld, is met behulp van de bekende rekenregels voor vectoren en matrices voor (1.1.1) te schrijven:

$$\vec{y} = X' \vec{\beta} + \vec{\epsilon} \tag{1.4.3}$$

$$\text{en } \vec{\hat{\eta}} = E(\vec{y}) = X'\hat{\beta} \quad (1.4.4)$$

waarin  $X'$  de gespiegelde matrix  $\begin{pmatrix} x_{11} & \dots & x_{p1} \\ \vdots & \ddots & \vdots \\ x_{1n} & \dots & x_{pn} \end{pmatrix}$  is.

Er moet nu om  $\hat{\beta}$  te vinden gelden:  $(\vec{y} - \vec{\hat{\eta}}) \perp \vec{\hat{\eta}}$  of met de notatie voor het inwendig vectorproduct:  $(\vec{y} - \vec{\hat{\eta}}, \vec{\hat{\eta}}) = 0$ , dus:  $(\vec{y} - X'\hat{\beta}, X'\hat{\beta}) = (\vec{y}, X'\hat{\beta}) - (X'\hat{\beta}, X'\hat{\beta}) = 0$ , waaruit:  $\vec{y}'X'\hat{\beta} = \hat{\beta}'XX'\hat{\beta} \equiv \hat{\beta}'S\hat{\beta}$  als:

$$XX' \equiv S \quad (1.4.5)$$

wordt gesteld. Er ontstaat dus:

$$\vec{y}'X'\hat{\beta} = \hat{\beta}'S \text{ of } S\hat{\beta} = X\vec{y} \quad (1.4.6)$$

Nu is  $S$  een  $p \times p$  symmetrische matrix waarvan de elementen bestaan uit de inwendige productsommen

$S_{ij} = \sum_{k=1}^n x_{ik}x_{kj}$ . Onder de voorwaarde dat  $S$  nonsingulier is, dat wil zeggen, dat de determinantwaarde  $S \neq 0$  is en er dus een inverse  $S^{-1}$  bestaat, kan (1.4.6) geschreven worden als:  $\hat{\beta}' = \vec{y}'X'S^{-1}$  of  $\vec{\hat{\beta}} = S^{-1}X\vec{y}$  (1.4.7)

$$\text{en } \vec{\hat{\eta}} = X'\hat{\beta} = X'S^{-1}X\vec{y} \quad (1.4.8)$$

In 1.2 was reeds ingevoerd  $\mathcal{J}(\vec{y}, \vec{\hat{\beta}}) = \mathcal{J}_\Omega = \|\vec{y} - \vec{\hat{\eta}}\|^2 =$

$$\|\vec{\epsilon}\|^2 = \sum_{k=1}^n \epsilon_k^2 \text{ dus:}$$

$$\mathcal{J}_\Omega = \|\vec{y} - X'\hat{\beta}\|^2 = \|\vec{y} - X'S^{-1}X\vec{y}\|^2 = \vec{y}'\vec{y} - \vec{y}'X'S^{-1}X\vec{y} \quad (1.4.9)$$

$$\text{met } D = X'S^{-1}X \text{ is ook te schrijven } \mathcal{J}_\Omega = \vec{y}'(I-D)\vec{y} \quad (1.4.10)$$

$$\frac{\partial \mathcal{J}(\vec{y}, \vec{\hat{\beta}})}{\partial (\vec{\hat{\beta}})} = X(\vec{y} - X'\hat{\beta}) = 0 \text{ geeft } X\vec{y} = XX'\hat{\beta} = S\hat{\beta} \text{ of } \vec{\hat{\beta}} = S^{-1}X\vec{y}.$$

Uitwerking geeft de normaalvergelijkingen zoals deze reeds in (1.2.1) gevonden waren.

Onder  $\Omega$  was ondersteld, dat variantie  $y_k$   
 $= E(y_k - Ey_k)^2 = \sigma^2$ . Als voor de covariantie matrix van  
 $\epsilon$  de notatie  $M_\epsilon$  wordt gebruikt, dan is deze te schrijven als:  
 $M_\epsilon = \sigma^2 I$  (1.4.11)

Nu geldt in de matrixtheorie nog de volgende stelling:

Als  $\vec{x} = M\vec{y}$  dan  $M_{\vec{x}} = MM_y M'$  (1.4.12)

Hiermede:  $M_{\vec{\beta}} = S^{-1} X M_y X' S^{-1} = \sigma^2 S^{-1} X X' S^{-1} = \sigma^2 S^{-1}$  (1.4.12)

Uit deze formule is de standaarddeviatie van de schattingen  
 $\hat{\beta}_i$  af te leiden.

In 1.2 waren de  $q$  lineaire functies  $\psi_j(\vec{\beta}_i)$  ingevoerd.

Deze zijn te schrijven:

$$\psi_j = \sum_{i=1}^p c_{ji} \beta_i$$

in de matrix notatie  $\vec{\psi} = C\vec{\beta}$  (1.4.13)

met:

$$C = \begin{pmatrix} c_{11} & \dots & c_{1p} \\ \vdots & & \vdots \\ c_{j1} & \dots & c_{jp} \end{pmatrix}$$

Voor de uit de waarnemingen verkregen schattingen van de para-  
 meters geldt:  $\hat{\psi} = C\hat{\beta} = CS^{-1}X\vec{y} = A\vec{y}$  (1.4.14)

met  $A = CS^{-1}X$ .

Nu invoeren  $B = AA' = CS^{-1}XX'S^{-1} = CS^{-1}C'$  (1.4.15)

dan wordt de covariantie matrix van  $\hat{\psi}$ :  $M_{\hat{\psi}} = \sigma^2 B$  (1.4.16)

Onder de voorwaarde dat  $S$  nonsingulier is, laat  $\mathcal{J}_\Omega$  zich  
 schrijven volgens (1.4.9) of (1.4.10).

Een analoge uitdrukking voor  $\mathcal{J}_\omega$  is als volgt te verkrijgen:  
 Nodig is die waarden van  $\vec{\beta}$  te vinden waarvoor  $\mathcal{J}(\vec{y}, \vec{\beta})$  een minimum  
 heeft onder de nevenvoorwaarden dat  $\vec{\psi} = 0$ . Dit is op te lossen  
 met behulp van Langrange multiplicatoren.

Hiertoe wordt ingevoerd:  $\mathcal{J}(\vec{y}, \vec{\beta}) = \mathcal{J}(\vec{y}, \vec{\beta}) + \vec{\lambda}' C \vec{\beta}$  (1.4.17)

waarin  $\vec{\lambda}'$  een  $1 \times q$  vector is, die  $q$  constanten (multiplicatoren)  $q_j$  als componenten heeft. Als (1.4.17) naar  $\vec{\beta}$  wordt

gedifferentieerd ontstaat:  $\frac{\delta \mathcal{J}(\vec{y}, \vec{\beta})}{\delta \vec{\beta}} = -2X(\vec{y} - X'\vec{\beta}) + C'\vec{\lambda}$ ,

$\frac{\delta \mathcal{J}}{\delta \vec{\beta}} = 0 \rightarrow 2S\vec{\beta} = 2X\vec{y} - C'\vec{\lambda}$  of  $\vec{\beta} = S^{-1}X\vec{y} - \frac{1}{2}S^{-1}C'\vec{\lambda}$  waarin  $\vec{\beta}$  de notatie is voor de schattingen van  $\vec{\beta}$  als  $\omega$  geldt. Nu moet aan  $\vec{\psi} = 0$  voldaan zijn. Dus  $C\vec{\beta} = CS^{-1}X\vec{y} - \frac{1}{2}CS^{-1}C'\vec{\lambda} = 0$ . De matrix  $CS^{-1}C'$  is in (1.4.15) reeds als  $B$  ingevoerd. Bewezen kan worden, dat  $B$  nonsingulier is (zie Scheffé, blz. 30), zodat geschreven kan worden:  $\vec{\lambda} = 2B^{-1}CS^{-1}X\vec{y}$  derhalve  $\vec{\beta} = S^{-1}X\vec{y} - S^{-1}C'B^{-1}CS^{-1}X\vec{y}$  waarna:

$$\begin{aligned} \mathcal{J}_\omega &= \mathcal{J}(\vec{y}, \vec{\beta}) = \vec{y}'\vec{y} - 2\vec{y}'X'(S^{-1}X\vec{y} - S^{-1}C'B^{-1}CS^{-1}X\vec{y}) + \\ & (\vec{y}'X'S^{-1} - \vec{y}'X'S^{-1}C'B^{-1}CS^{-1}) S (S^{-1}X\vec{y} - S^{-1}C'B^{-1}CS^{-1}X\vec{y}) = \\ & \vec{y}'\vec{y} - \vec{y}'X' \{2S^{-1} - 2S^{-1}C'B^{-1}CS^{-1} - S^{-1} + S^{-1}C'B^{-1}CS^{-1} + \\ & S^{-1}C'B^{-1}CS^{-1} - S^{-1}C'B^{-1}CS^{-1}C'B^{-1}CS^{-1}\} X\vec{y}. \text{ Tenslotte:} \\ \mathcal{J}_\omega &= \vec{y}'\vec{y} - \vec{y}'X'(S^{-1} - S^{-1}C'B^{-1}CS^{-1})X\vec{y} \end{aligned} \quad (1.4.18)$$

In 1.3 was gevonden dat voor de beslissing tot al of niet verwerpen van  $H_0$  een toetsingsgrootte nodig is, die niet  $\mathcal{J}_\omega$  maar het verschil  $\mathcal{J}_\omega - \mathcal{J}_\Omega$  bevat. Voor het verschil is een eenvoudige uitdrukking te vinden.

Uit (1.4.10) en (1.4.18) volgt:  $\mathcal{J}_\omega - \mathcal{J}_\Omega = \vec{y}'X'S^{-1}CB^{-1}CS^{-1}X\vec{y}$ . Nu is volgens (1.4.14)  $\hat{\vec{\psi}} = CS^{-1}X\vec{y}$  zodat tenslotte gevonden wordt:

$$\mathcal{J}_\omega - \mathcal{J}_\Omega = \hat{\vec{\psi}}'B^{-1}\hat{\vec{\psi}} \quad (1.4.19)$$

1.5. De frequentieverdeling van de toetsingsgrootheid.

De vectoren  $\vec{y}, \vec{\eta}$  en  $\vec{\eta}_\omega$  zijn gegeven in coördinaten ten aanzien van een of andere basis in  $R_n$  die met  $\{v_i\}$  zal worden aangegeven.

De verschilvectoren  $\vec{y}-\vec{\eta}$  en  $\vec{\eta}-\vec{\eta}_\omega$ , die loodrecht op elkaar staan zoals in 1.3 is aangetoond, zijn beide in principe dus uitdrukkingen in alle  $v_i$ 's. Het is nu handig een andere basis te kiezen en wel een orthonormale basis  $\alpha_1 \dots \alpha_n$  zodanig, dat 1<sup>o</sup>)  $p-q$  basisvectoren  $\vec{\alpha}_{q+1} \dots \vec{\alpha}_p$  de ruimte  $V_{p-q}$  opspannen, 2<sup>o</sup>)  $q$  andere basisvectoren  $\vec{\alpha}_1 \dots \vec{\alpha}_q$  met die van  $V_{p-q}$  samen,  $V_p$  opspannen en 3<sup>o</sup>) de resterende  $n-p$  vectoren  $\vec{\alpha}_{p+1} \dots \vec{\alpha}_n$  er voor zorgen, dat  $R_n$  gevormd wordt.

Als nu  $z_1, \dots, z_n$  de coördinaten van  $\vec{y}$  ten opzichte van de basis  $\{\vec{\alpha}_i\}$  zijn, dan is de coördinaat  $z_i$  gelijk aan de projectie van  $\vec{y}$  op  $\vec{\alpha}_i$  en daar  $\vec{\alpha}_i$  een eenheidsvector is, is deze projectie gelijk aan het inproduct  $(\vec{\alpha}_i, \vec{y})$ . De vector  $\vec{z}$ , die identiek is aan de vector  $\vec{y}$ , kan dus in de orthonormale basis  $\{\vec{\alpha}_j\}$  uitgedrukt, geschreven worden als  $P\vec{y}$ , waarin  $P$  een  $n \times n$  matrix is, waarvan de elementen de componenten zijn van de  $\vec{\alpha}_i$  ten opzichte van de oude basis  $\{v_i\}$ .

$$P = \begin{pmatrix} \alpha_{11} & \dots & \alpha_{1n} \\ \vdots & & \vdots \\ \alpha_{n1} & \dots & \alpha_{nn} \end{pmatrix} \quad (1.5.1.)$$

De rijen van  $P$  vormen dus de vectoren  $\vec{\alpha}_i$ . Volgens de regels van de matrixtheorie is:

$$P\vec{y} = \begin{pmatrix} \alpha_{11} & \dots & \alpha_{1n} \\ \vdots & & \vdots \\ \alpha_{n1} & \dots & \alpha_{nn} \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \sum_{j=1}^n \alpha_{1j} y_j \\ \vdots \\ \sum_{j=1}^n \alpha_{nj} y_j \end{pmatrix} = \begin{pmatrix} (\alpha_1, y) \\ \vdots \\ (\alpha_n, y) \end{pmatrix} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix}$$

De  $z_i$  zijn dus lineaire combinaties van de  $y_j$ 's. Nu zijn de  $y_j$ 's stochastische variabelen en dus zijn de  $z_i$ 's ook stochastisch. De verwachtingswaarde van  $z_i$  zal door  $\zeta_i$  worden voorgesteld:  $Ez_i = \zeta_i$ .

In vectornotatie verder:  $\vec{\zeta} = E(\vec{z}) = E(P\vec{y}) = PE(\vec{y}) = P\vec{\eta}$  (1.5.2)  
 Als combinatie van normaal verdeelde grootheden  $y_j$  is  $z_j$  zelf ook normaal verdeeld en wel met gemiddelde  $\zeta_i$  en standaarddeviatie  $\sigma$ . Dit laatste volgt zonder meer uit de eigenschap, dat geometrische grootheden invariant zijn t.o.v. de overgang op een andere basis. Het volgt echter ook uit de matrixnotatie met (1.4.9) namelijk:  $M_{\vec{z}} = PM_{\vec{y}}P' = M_{\vec{y}}PP' = M_{\vec{y}} = \sigma^2 I$  ( $PP' = I$ , omdat de  $\vec{\alpha}_i$  onderling loodrechte eenheidsvectoren zijn). De covariant is matrix van  $\vec{z}$  is dus  $\sigma^2 I$ .

Nu stellen dus  $\vec{\eta} = E(\vec{y})$  en  $\vec{\zeta} = E(\vec{z})$  één en dezelfde vector voor, maar genoteerd respectievelijk t.o.v. de basis  $\{v_i\}$  en  $\{\alpha_i\}$ ; daar  $\vec{\eta} \in V_p$  ook  $\vec{z} \in V_p$  en dat wil zeggen  $\zeta_{p+1} = \dots = \zeta_n = 0$ .

Als  $\omega$  geldt dan is  $E(\vec{y}) = \vec{\eta}_\omega$  en derhalve  $\vec{\zeta}_\omega = E_\omega(\vec{z}) = P\vec{\eta}_\omega$ . Daar  $\vec{\eta}_\omega \in V_{p-q}$  moet ook gelden  $\vec{\zeta}_\omega \in V_{p-q}$  en dus  $\zeta_1 = \dots = \zeta_q = 0$ . Voor de grootheden  $\mathcal{J}$  wordt nu respectievelijk gevonden:

$$\mathcal{J}_\Omega = \|\vec{y} - \hat{\vec{y}}\|^2 = \|\vec{z} - \hat{\vec{z}}\|^2 = \left\| \sum_{i=p+1}^n z_i \alpha_i \right\|^2 = \sum_{i=p+1}^n z_i^2 \quad (1.5.3)$$

$$\mathcal{J}_\omega = \|\vec{y} - \hat{\vec{y}}_\omega\|^2 = \|\vec{z} - \hat{\vec{z}}_\omega\|^2 = \left\| \sum_{i=1}^q z_i \alpha_i + \sum_{i=p+1}^n z_i \alpha_i \right\|^2 = \sum_{i=1}^q z_i^2 + \sum_{i=p+1}^n z_i^2 \quad (1.5.4)$$

$$\text{en dus: } \mathcal{J}_\omega - \mathcal{J}_\Omega = \|\hat{\vec{y}} - \hat{\vec{y}}_\omega\|^2 = \|\hat{\vec{z}} - \hat{\vec{z}}_\omega\|^2 = \sum_{i=1}^q z_i^2 \quad (1.5.5)$$

Volgens de mathematische statistiek volgt de kwadraatsom van  $v$  onderling onafhankelijk normaal verdeelde grootheden  $x_i$  met  $\sigma = 1$  een zogenaamde  $\chi^2$  (Chi kwadraat) verdeling met  $v$  vrijheidsgraden. Is  $E(x_i) = \xi_i \neq 0$  dan is er sprake van een niet centrale  $\chi^2$  met centraliteitsparameter

$\delta = (\sum_1^v \xi_i^2)^{\frac{1}{2}}$ , is  $\xi_i = 0$  dan wordt de verdeling centraal genoemd,

respectievelijk aangeduid met  $\chi_{v,\delta}^2$  of  $\chi_v^2$ . Er geldt:  $E\chi_{v,\delta}^2 = v + \delta^2$ .

Uit (1.5.5) volgt:

$$E(\mathcal{J}_\omega - \mathcal{J}_\Omega) = E \sum_1^q z_i^2 = \sum_1^q E z_i^2 = \sum_1^q (E(z_i - \zeta_i)^2 + \zeta_i^2) = q\sigma^2 + \sum_1^q \zeta_i^2 = q\sigma^2 + (\sigma\delta)^2 \text{ met } \delta = (\sum_1^q \zeta_i^2 / \sigma^2)^{\frac{1}{2}} \quad (1.5.6)$$

Dus  $E(\frac{\mathcal{J}_\omega - \mathcal{J}_\Omega}{\sigma^2}) = q + \delta^2$ . Verder  $E\mathcal{J}_\Omega = E \sum_{p+1}^n z_i^2 = (n-p)\sigma^2$  wegens

$$\zeta_{p+1} = \dots = \zeta_n = 0.$$

$$\text{Derhalve } E \frac{\mathcal{J}_\Omega}{\sigma^2} = (n-p). \quad (1.5.7)$$

$\frac{\mathcal{J}_\omega - \mathcal{J}_\Omega}{\sigma^2}$  volgt dus een  $\chi_{q,\delta}^2$  verdeling onder  $\Omega$  en een  $\chi_q^2$  verdeling onder  $\omega$ .  $\frac{\mathcal{J}_\Omega}{\sigma^2}$  volgt een  $\chi_{n-p}^2$  verdeling; immers als  $\Omega$  geldt, dan ligt  $\vec{h}$  in  $V_p$  en geldt  $\zeta_1, \dots, \zeta_q \neq 0$ ; geldt  $\omega$  dan ligt  $\vec{h}$  in  $V_{p-q}$  en geldt  $\zeta_1 = \dots = \zeta_q = 0$ .

Het quotiënt van twee  $\chi^2$  verdelingen is een F verdeling, die ook weer centraal of niet centraal kan zijn. Derhalve is

$$\frac{(n-p)(\mathcal{J}_\omega - \mathcal{J}_\Omega)}{q\mathcal{J}_\Omega} \quad (1.5.8)$$

verdeeld volgens  $F_{q,n-p,\sigma}$  onder  $\Omega$  en volgens  $F_{q,n-p}$  onder  $\omega$ .

Volgens (1.5.6) en (1.5.2) is  $\sigma^2 \delta^2 = \sum_{i=1}^q z_i^2 = \sum_{i=1}^q \left( \sum_{j=1}^n p_{ij} \eta_j \right)^2$

Verder is  $\mathcal{J}_\omega - \mathcal{J}_\Omega = \sum_{i=1}^q z_i^2 = \sum_{i=1}^q \left( \sum_{j=1}^n p_{ij} y_j \right)^2$  met andere

woorden de non-centraliteitsparameter  $\sigma$  is te vinden door in de uitdrukking  $\mathcal{J}_\omega - \mathcal{J}_\Omega$  de  $y$ 's te vervangen door  $Ey_i = \eta_j$

1.6. Een aequivalente formulering van de F-toets.

In 't voorgaande is aangetoond, dat  $\frac{\mathcal{J}_\omega - \mathcal{J}_\Omega}{\sigma^2}$  een  $\chi_{q,\sigma}^2$  verdeling volgt, terwijl ook geldt dat  $\mathcal{J}_\omega - \mathcal{J}_\Omega = \hat{\psi}' B^{-1} \hat{\psi}$ . Verder was  $\sigma$  te vinden door in de uitdrukking voor  $\mathcal{J}_\omega - \mathcal{J}_\Omega$  als functie van de  $y_j$ , deze  $y_j$  door  $\eta_j = Ey_j$  te vervangen. Dit is dus equivalent met  $E(\mathcal{J}_\omega - \mathcal{J}_\Omega) = E\hat{\psi}' B^{-1} \hat{\psi} = \hat{\psi}' B^{-1} \hat{\psi} = \sigma^2 \delta^2$ . Derhalve volgt  $\frac{(\hat{\psi} - \bar{\psi})' B^{-1} (\hat{\psi} - \bar{\psi})}{\sigma^2}$  een  $\chi_q^2$  terwijl  $\frac{\mathcal{J}_\Omega}{\sigma^2}$  een  $\chi_{n-p}^2$  verdeling volgt. Als  $S^2 = \frac{\mathcal{J}_\Omega}{n-p}$  wordt gesteld, dan volgt  $\frac{(\hat{\psi} - \bar{\psi})' B^{-1} (\hat{\psi} - \bar{\psi})}{qS^2}$  een  $F_{q,n-p}$  verdeling. Dit is ook zo op te vatten, dat  $(\hat{\psi} - \bar{\psi})' B^{-1} (\hat{\psi} - \bar{\psi}) \leq qS^2 F_{q,n-p}$  een vertrouwensgebied vormt, een ellipsoïde voor  $\hat{\psi}$  in de  $q$ -dimensionale  $\psi$ -ruimte met centrum  $\hat{\psi}_1, \dots, \hat{\psi}_q$ . In het geval  $q = 1$  gaat dit over in een vertrouwensinterval:  $b^{-1} (\hat{\psi} - \bar{\psi})^2 \leq S^2 F_{1,n-p}$  (1.6.1)

Nu is een F-verdeling met  $v_1 = 1$  equivalent met een t-verdeling dat wil zeggen, als  $\underline{t}$  een t-verdeling met  $v$  vrijheidsgraden volgt, dan volgt  $\underline{t}^2$  een  $F_{1,v}$  verdeling; derhalve is het vertrouwensinterval voor  $\hat{\psi}$  ook te schrijven als:

$$\hat{\psi} - t_{n-p} S \hat{\psi} \leq \hat{\psi} \leq \hat{\psi} + t_{n-p} S \hat{\psi} \tag{1.6.2}$$

waarbij  $bS^2$  door  $S^2 \hat{\psi}$  is vervangen omdat  $bS^2$  gelijk is aan de variantie van  $\hat{\psi}$ .



2. Voorbeelden.

2.1. Enkelvoudige variantie analyse.

Als een grootte  $y$  onder bijvoorbeeld  $k$  verschillende omstandigheden wordt waargenomen en men wil weten of deze verschillende omstandigheden invloed hebben op de waarde van  $y$  dan kan men dit onderzoeken met behulp van enkelvoudige variantie analyse. (Het is dan zo, dat een eventuele invloed van de omstandigheden zich uit in een niveauverandering van  $y$ ). De verschillende omstandigheden betekenen voor de waarnemingen en voor de theoretische grootheden van het betrokken probleem een indeling in klassen, die met een index  $i$  zullen worden aangegeven. In de termen van 1 is dit probleem als volgt te formuleren:

$$\Omega: y_{t,r} = \sum_{i=1}^k \beta_i x_{i,t} + \epsilon_{t,r} \quad \text{waarin } t = 1, \dots, k \quad r = 1, \dots, n_k$$

$$\sum_{t=1}^k n_t = n \text{ en } \delta_{i,t} = (\text{Kronecker-symbool; dat wil zeggen}$$

$$\delta_{i,t} = 1 \text{ als } t = i; \delta_{i,t} = 0 \text{ als } t \neq i) x_{i,t} \text{ is de in (0.3)}$$

genoemde indicator variabele. Korter kan ook  $y_{t,r} = \beta_t + \epsilon_{t,r}$  (2.1.1) geschreven worden. Voor de  $\epsilon_{t,r}$  geldt: onderling onafhankelijk normaal verdeeld met verwachtingswaarde nul en standaarddeviatie  $\sigma$ . Derhalve  $E y_{t,r} = \beta_t$ .

Er zijn dus  $k$  steekproeven van de grootte  $n_i$  ( $i=1, \dots, k$ ). De eerste vraag die men zich kan stellen is of er inderdaad een reëel verschil tussen de diverse  $\beta_i$  bestaat. In statistische zin ligt het dan voor de hand te toetsen of de  $\beta$ 's mogelijk alle aan elkaar gelijk zijn. Dat wil zeggen, onderzocht wordt de nulhypothese:  $H_0 : \beta_1 = \beta_2 = \dots = \beta_k$  (2.1.2) Nu is dit een vrij eenvoudig probleem, dat zonder meer rechtstreeks met de methode der kleinste kwadraten is op te lossen; daarom echter juist geschikt als illustratie.

Naast de rechtstreekse oplossing zal ook de formele oplossing worden gegeven.

2.1.2 Kleinste kwadratische oplossing.

De kwadraatsom  $S(\vec{y}, \vec{\beta})$  is onder  $\Omega$ :  $S(\vec{y}, \vec{\beta}) = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \beta_i)^2$

Deze is minimaal als  $\frac{\delta S(\vec{y}, \vec{\beta})}{\delta \beta_v} = -2 \sum_{j=1}^{n_v} (y_{vj} - \beta_v) = -2 \sum_{j=1}^{n_v} y_{vj} + 2n_v \beta_v = 0$ ,

waaruit:  $\hat{\beta}_v = \sum_{j=1}^{n_v} y_{vj} / n_v = \bar{y}_v$ . Met andere woorden, de beste

schattingen van de  $\beta_v$  is het gemiddelde van de  $y$  waarden onder de met  $v$  aangegeven omstandigheden. Als ook  $H_0$  geldt (notatie van 1.2:  $\omega = H_0 \cap \Omega$ ) dan zijn alle  $\beta_i$  gelijk, bijvoorbeeld  $= \beta$  en dan is  $S(\vec{y}, \vec{\beta}) = \sum_{ij} (y_{ij} - \beta)^2$  en dus

$$\frac{\delta S(\vec{y}, \vec{\beta})}{\delta \beta} = -2 \sum_i \sum_j (y_{ij} - \beta) = -2 \sum_{ij} y_{ij} + 2\beta \sum_j n_j = 0$$

waaruit:  $\hat{\beta} = \bar{y} (= \frac{1}{n} \sum_{ij} y_{ij})$ .

In de notatie van 1.2 is nu:

$\vec{y} = \{y_{11}, \dots, y_{1n_1}, y_{2,1}, \dots, y_{2n_2}, \dots, y_{k,1}, \dots, y_{k,n_k}\}$ . Verder is wegens

$\eta_{ij} = E y_{ij} = \beta_i$ :  $\hat{\eta} = \{\hat{\beta}_1, \dots, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_2, \dots, \hat{\beta}_k, \dots, \hat{\beta}_k\}$  en dus

$\hat{\eta}_\omega = \{\hat{\beta}, \dots, \hat{\beta}\}$  waarmee:

$$\mathcal{J}_\Omega = ||\vec{y} - \hat{\eta}||^2 = \sum_{ij} (y_{ij} - \hat{\beta}_i)^2 = \sum_{ij} (y_{ij} - \bar{y}_i)^2 = \sum_{ij} y_{ij}^2 - \sum_i n_i \bar{y}_i^2$$

$$\text{en } \mathcal{J}_\omega = ||\vec{y} - \hat{\eta}_\omega||^2 = \sum_{ij} (y_{ij} - \hat{\beta})^2 = \sum_{ij} (y_{ij} - \bar{y})^2 = \sum_{ij} y_{ij}^2 - n \bar{y}^2$$

Dus  $\mathcal{J}_\omega - \mathcal{J}_\Omega = \sum_i n_i \bar{y}_i^2 - n \bar{y}^2$  wat ook te vinden is als:

$$||\hat{\eta} - \hat{\eta}_\omega||^2 = \sum_{ij} (\bar{y}_i - \bar{y})^2 = \sum_i n_i (\bar{y}_i - \bar{y})^2.$$

$S_{\omega} - S_{\Omega} = \sum_i n_i (\bar{y}_i - \bar{y})^2$  is een gewogen maat voor de spreiding van de steekproefgemiddelden van de k steekproeven rond het algemeen gemiddelde. Dit wordt ook wel als  $S_{\text{tussen}}$  genoteerd en heet dan spreidingskwadraatsom tussen de steekproeven.  $\mathcal{E}_{\Omega} = \sum_i \sum_j (y_{ij} - \bar{y}_i)^2$  is een gecombineerde maat voor de spreiding binnen de steekproeven, genoteerd als  $S_{\text{binnen}}$ .

Volgens 1.5 volgt  $\mathcal{E}_{\Omega}$  een  $\chi^2$  verdeling met n-p vrijheidsgraden waarin p het aantal te schatten  $\beta$ 's is, dus k in bovenstaand geval.  $S_{\text{binnen}}$  heeft dus n-k vrijheidsgraden.  $\mathcal{E}_{\omega} - \mathcal{E}_{\Omega}$  heeft 9 vrijheidsgraden, waarin 9 het aantal functies  $\psi_j$  is dat  $H_0$  bepaalt. Met invoering van:  $\psi_1 = \beta_1 - \beta_k, \psi_2 = \beta_2 - \beta_k, \dots, \psi_{k-1} = \beta_{k-1} - \beta_k$  is  $H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$  te schrijven als  $\psi_i = 0$  voor  $i = 1, \dots, k-1$ . Derhalve is het aantal vrijheidsgraden dat voor  $\mathcal{E}_{\omega} - \mathcal{E}_{\Omega}$  in rekening moet worden genomen k-1.  $\mathcal{E}_{\omega}$  heeft tenslotte n-1 vrijheidsgraden.

Volgens (1.5.7) geldt ook:  $E \mathcal{E}_{\Omega}/n-p = \sigma^2$  ( dus  $\mathcal{E}_{\Omega}/n-p$  is een "zuivere" schatter van  $\sigma^2$  ),  $E \frac{\mathcal{E}_{\omega} - \mathcal{E}_{\Omega}}{q} = \sigma^2 + \frac{1}{q} \sigma^2 \delta^2 = \sigma^2 + \frac{1}{q} \sum_i c_i^2$

In bovenstaand geval wordt dit volgens de slotregel in (1.5):

$$E \frac{\mathcal{E}_{\omega} - \mathcal{E}_{\Omega}}{q} = \sigma^2 + \frac{1}{q} \sum_{i=1}^{k-1} n_i (\beta_i - \beta)^2.$$

De resultaten zijn in de volgende tabel samen te vatten:

$$y \quad S_t \equiv S_{\text{tussen}}$$

Bron	kwadraat-sommen	aantal vrijheidsgraden	variantie	E. (variantie)
tussen de groepen	$S_t = \sum_i n_i (\bar{y}_i - \bar{y})^2$	k-1	$S_t / k-1$	$\sigma^2 + \frac{\sum_i n_i (\beta_i - \beta)^2}{k-1}$
binnen de groepen	$S_b = \sum_i \sum (y_{ij} - \bar{y})^2$	n-k	$S_b / n-k$	$\sigma^2$
totaal	$S_{\text{"tot"}} = \sum_{ij} (y_{ij} - \bar{y})^2$	n-1	-	-

De uitspraak of  $H_0$  verworpen moet worden volgt uit de bepaling van de overschrijdingskans van

$$\frac{n-k}{k-1} \frac{S_{\text{tussen}}}{S_{\text{binnen}}} = \frac{n-k}{k-1} \frac{\sum_{i=1}^k n_i (y_{i\cdot} - \bar{y})^2}{\sum_{ij} (y_{ij} - y_{i\cdot})^2} \quad \text{op grond van}$$

de  $F_{k-1, n-k}$  verdeling.

### 2.1.3 Formele matrix oplossing.

Uitgangspunt:

$$\Omega : y_{ij} = \beta_i + \epsilon_{ij} \quad (i=1, \dots, k ; j=1, \dots, n_i ; \sum n_i = n ; \epsilon : N(0, \sigma^2))$$

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k$$

$$\psi_j = \beta_k - \beta_j \quad (j=1, \dots, k-1)$$

dus ook:

$$H_0 : \psi_j = 0 \quad (j=1, \dots, k-1)$$

$$\vec{y} : \{y_{11} \dots y_{1n_1} y_{21} \dots y_{2n_2} \dots y_{kn_k}\}$$

De in (0.1) ingevoerde grootheden  $\vec{x}_i$  zijn in dit geval

$$\vec{x}_1 = \{1, \dots, 1, 0, \dots, 0\} \text{ enz.}$$

$n_1$

$$\vec{x}_2 = \{ \underbrace{0, \dots, 0}_{n_1}, \underbrace{1, \dots, 1}_{n_2}, 0, \dots, 0 \} \text{ enz.}$$

dus:

$$X = \begin{pmatrix} 1 \dots \dots 1 & 0 \dots \dots 0 & & & & & & & & 0 \\ 0 \dots \dots 0 & 1 \dots \dots 1 & 0 & 0 & & & & & & 0 \\ \vdots & & & & & & & & & \vdots \\ 0 & & & & & & & & & 0 \\ & & & & & & 0 & 0 & 1 \dots \dots 1 & \end{pmatrix}$$

waarmee  $S = XX' = \begin{pmatrix} n_1 & & & & & & & & & 0 \\ & n_2 & & & & & & & & 0 \\ & & \ddots & & & & & & & \\ & & & n_k & & & & & & \\ 0 & & & & & & & & & \end{pmatrix}$ ;  $S^{-1} = \begin{pmatrix} \frac{1}{n_1} & & & & & & & & & 0 \\ & \frac{1}{n_2} & & & & & & & & \\ & & \ddots & & & & & & & \\ & & & \frac{1}{n_k} & & & & & & \\ 0 & & & & & & & & & \frac{1}{n_k} \end{pmatrix}$

volgens (1.4.7) is:

$$\hat{\vec{\beta}} = S^{-1}X\vec{y} = \begin{pmatrix} \frac{1}{n_1} & & & & & & & & & 0 \\ & \frac{1}{n_2} & & & & & & & & \\ & & \ddots & & & & & & & \\ & & & \frac{1}{n_k} & & & & & & \\ 0 & & & & & & & & & \end{pmatrix} \begin{pmatrix} 1 \dots \dots 1 & & & & & & & & & 0 \\ & \ddots & & & & & & & & \\ & & \ddots & & & & & & & \\ & & & \ddots & & & & & & \\ & & & & 1 \dots \dots 1 & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \end{pmatrix} \begin{pmatrix} y_{11} \\ \vdots \\ y_{kn_k} \end{pmatrix} =$$

$$\begin{pmatrix} \frac{1}{n_1} & & & & & & & & & 0 \\ & \frac{1}{n_2} & & & & & & & & \\ & & \ddots & & & & & & & \\ & & & \frac{1}{n_k} & & & & & & \\ 0 & & & & & & & & & \end{pmatrix} \begin{pmatrix} \sum_j y_{j1} \\ \sum_j y_{j2} \\ \vdots \\ \sum_j y_{jk} \end{pmatrix} = \begin{pmatrix} \bar{y}_1 \\ y_2 \\ \vdots \\ \bar{y}_k \end{pmatrix} \quad \text{waaruit } \hat{\beta}_j = \bar{y}_j$$

verder kan men  $\hat{\beta}$  bepalen als  $\|\vec{y} - \hat{\vec{\eta}}\|^2$  of als  $\vec{y}'(I-D)\vec{y}$  volgens (1.4.10).

Nu is  $\hat{\eta}' = (\bar{y}_1, \dots, \bar{y}_1, \bar{y}_2, \dots, \bar{y}_2, \dots, \bar{y}_k, \dots, \bar{y}_k)$

Dus  $\|\hat{y} - \hat{\eta}\|^2 = \sum_{ij} (y_{ij} - \bar{y}_i)^2 = \sum_{ij} y_{ij}^2 - \sum_i n_i \bar{y}_i^2$ .

ook  
 $D = X'S^{-1}X = \begin{pmatrix} 1 & & & 0 \\ \vdots & \ddots & & \\ & & 1 & \\ 0 & & & \ddots \\ & & & & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{n_1} & & & 0 \\ & & \frac{1}{n_2} & \\ & & & \ddots \\ 0 & & & & \frac{1}{n_k} \end{pmatrix} \begin{pmatrix} 11\dots 1 & & & 0 \\ & & & \\ & & & \\ 0 & & & & 11\dots 1 \end{pmatrix} =$

$$\left[ \begin{array}{cc|c} \begin{pmatrix} \frac{1}{n_1} & & \\ \vdots & \ddots & \\ \frac{1}{n_1} & & \end{pmatrix} & & \text{O} \\ & \begin{pmatrix} \frac{1}{n_2} & & \\ \vdots & \ddots & \\ \frac{1}{n_2} & & \end{pmatrix} & & \\ \text{O} & & \begin{pmatrix} \frac{1}{n_k} & & \\ \vdots & \ddots & \\ \frac{1}{n_k} & & \end{pmatrix} \end{array} \right]$$

Dus:  
 $I-D = \left[ \begin{array}{ccc|c} \begin{pmatrix} (1 - \frac{1}{n_1}) & & \\ -\frac{1}{n_1} & \ddots & \\ -\frac{1}{n_1} & & (1 - \frac{1}{n_1}) \end{pmatrix} & & \text{O} \\ & \begin{pmatrix} (1 - \frac{1}{n_2}) & & \\ -\frac{1}{n_2} & \ddots & \\ -\frac{1}{n_2} & & (1 - \frac{1}{n_2}) \end{pmatrix} & & \\ \text{O} & & \begin{pmatrix} (1 - \frac{1}{n_k}) & & \\ \vdots & \ddots & \\ (1 - \frac{1}{n_k}) & & \end{pmatrix} \end{array} \right]$

(I - D)y wordt nu

$$\begin{pmatrix} y_{11} - \bar{y}_1 \\ y_{12} - \bar{y}_1 \\ \vdots \\ y_{21} - \bar{y}_2 \\ \vdots \end{pmatrix}$$

zodat  $\mathcal{J}_\Omega = \sum_{ij} (y_{ij} - \bar{y}_i) y_{ij} = \sum_{ij} y_{ij}^2 - \sum_i n_i \bar{y}_i^2$ .

wegens  $H_\Omega: \psi_j = \beta_k - \beta_j$  wordt de matrix C:

$$C = \begin{pmatrix} -1 & 0 & \dots & 1 \\ 0 & -1 & & \\ \vdots & & \ddots & \\ 0 & & & -1 & 1 \end{pmatrix}$$

Dan is  $B = CS^{-1}C' =$

$$\begin{pmatrix} -1 & 0 & \dots & 0 & 1 \\ 0 & -1 & & & \\ \vdots & & \ddots & & \\ 0 & & & -1 & 1 \end{pmatrix} \cdot \begin{pmatrix} \frac{1}{n_1} & & & 0 \\ & \frac{1}{n_2} & & \\ & & \ddots & \\ 0 & & & \frac{1}{n_k} \end{pmatrix} \begin{pmatrix} -1 & 0 \\ 0 & -1 \\ \vdots & \vdots \\ 0 & 1 \end{pmatrix}$$

Hieruit:

$$B^{-1} = \begin{pmatrix} n_1 - \frac{n_1^2}{n} & -\frac{n_1 n_2}{n} & \dots & -\frac{n_1 n_{k-1}}{n} \\ -\frac{n_2 n_1}{n} & n_2 - \frac{n_2^2}{n} & & \\ \vdots & & \ddots & \\ \vdots & & & \vdots \end{pmatrix}$$

Verder is:

$$\hat{\psi} = \begin{pmatrix} v_k - \bar{y}_1 \\ y_k - \bar{y}_2 \\ \vdots \\ y_k - \bar{y}_{k-1} \end{pmatrix}$$

Derhalve  $\mathcal{L}'_{\omega} - \mathcal{L}'_{\Omega} = \hat{\psi}' B^{-1} \hat{\psi} =$

$$(\bar{y}_k - \bar{y}_1, \dots) \begin{pmatrix} n_1 - \frac{n_1^2}{n} & -\frac{n_1 n_2}{n} & \dots \\ \frac{n_1 n_2}{n} & n_2 - \frac{n_1^2}{n} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} \bar{y}_k - \bar{y}_1 \\ \vdots \\ \vdots \end{pmatrix} =$$

$$(\bar{y}_k - \bar{y}_1, \dots) \begin{pmatrix} n_1 (\bar{y}_k - \bar{y}_1) - \frac{n_1^k}{n} \sum_{j=1}^k n_j (\bar{y}_k - \bar{y}_j) \\ n_2 (\bar{y}_k - \bar{y}_2) - \frac{n_2^k}{n} \sum_{j=1}^k n_j (\bar{y}_k - \bar{y}_j) \\ \vdots \\ \vdots \end{pmatrix} =$$

$$(\bar{y}_k - \bar{y}_1, \dots) \cdot \begin{pmatrix} n_1 (\bar{y} - \bar{y}_1) \\ n_2 (\bar{y} - \bar{y}_2) \\ \vdots \\ \vdots \end{pmatrix} = \sum_{i=1}^k n_i (y - y_i) (y_k - y_i) =$$

$$\sum_{i=1}^k n_i (\bar{y} - \bar{y}_i) (\bar{y} - \bar{y}_i) + \bar{y}_k - \bar{y} = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2 = \sum_{i=1}^k n_i \bar{y}_i^2 - n \bar{y}^2.$$

Berekening van  $\zeta_{\omega}$  volgens (1.4.18) levert resp.:

$$S^{-1} C' = \begin{pmatrix} -\frac{1}{n_1} & & & & \\ & -\frac{1}{n_2} & & & \\ & & \ddots & & \\ & & & \frac{1}{n_{k-1}} & \\ \frac{1}{n_k} & & & & \frac{1}{n_k} \end{pmatrix}$$



$$S^{-1}C'B^{-1} = \begin{pmatrix} \frac{n_1}{n} - 1 & \frac{n_2}{n} & \dots & \frac{n_{k-1}}{n} \\ \frac{n_1}{n} & \frac{n_2}{n} - 1 & \dots & \frac{n_{k-1}}{n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{n_1}{n} & \dots & \dots & \frac{n_{k-1}}{n} - 1 \\ \frac{n_1}{n} & \dots & \dots & \frac{n_{k-1}}{n} \end{pmatrix}$$

$$S^{-1}C'B^{-1}CS^{-1} = \begin{pmatrix} \frac{1}{n_1} - \frac{1}{n} & -\frac{1}{n} & \dots & -\frac{1}{n} & -\frac{1}{n} \\ -\frac{1}{n} & \frac{1}{n_2} - \frac{1}{n} & \dots & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \vdots & \vdots & \dots & \frac{1}{n_{k-1}} - \frac{1}{n} & -\frac{1}{n} \\ -\frac{1}{n} & \dots & \dots & -\frac{1}{n} & \frac{1}{n_k} - \frac{1}{n} \end{pmatrix}$$

dus

$$S^{-1}S^{-1}C'B^{-1}CS^{-1} = \begin{pmatrix} \frac{1}{n} & \dots & \frac{1}{n} \\ \vdots & \ddots & \vdots \\ \frac{1}{n} & \dots & \frac{1}{n} \end{pmatrix}$$

waarmee

$$\vec{y}'X'(S^{-1} - S^{-1}C'B^{-1}CS^{-1})X\vec{y} = (n_1\bar{y}_1 \quad n_2\bar{y}_2 \quad \dots \quad n_k\bar{y}_k) \cdot$$

$$\begin{pmatrix} \frac{1}{n} & & & \\ & \ddots & & \\ & & \frac{1}{n} & \\ & & & \ddots \\ \frac{1}{n} & & & & \frac{1}{n} \end{pmatrix} \cdot \begin{pmatrix} n_1\bar{y}_1 \\ n_2\bar{y}_2 \\ \vdots \\ n_k\bar{y}_k \end{pmatrix} = \frac{1}{n}(\sum n_i y_i)^2 = n\bar{y}^2. \quad \text{Dus } \mathcal{J}_\omega = \sum y_i^2 - n\bar{y}^2.$$

Zodat inderdaad  $\mathcal{J}_\omega - \mathcal{J}_\Omega = \sum n_i y_i^2 - n\bar{y}^2$ .

Hiermede zijn dus dezelfde resultaten verkregen als met de rechtstreekse methode.

#### 2.1.4 Numeriek voorbeeld.

Een adequaat voorbeeld van enkelvoudige variantie-analyse is voor meteorologische waarnemingen niet makkelijk te vinden. Meestal zijn de verschijnselen zo complex, dat meerdere effecten door elkaar een rol spelen, of is er sprake van tijdreeksen, waarvoor andere methoden ten onderzoek beter geschikt zijn. Toch is een voorbeeld uit een tijdreeks gekozen namelijk de vorstexcentriciteiten die door F. Ynsen zijn samengesteld ter karakterisering van de winters en gepubliceerd in W.R.74-2. Deze gegevens hebben het voordeel, dat ze te beschouwen zijn als onafhankelijk verdeelde normale grootheden.

Gebruikt zullen worden de vorstexcentriciteiten van 210 jaar namelijk van 1760 tot en met 1969, waarvan 7 steekproeven van 30 worden samengesteld. De vraag is nu of de gemiddelde vorstexcentriciteiten van de 7 opeenvolgende 30-jaar-perioden gelijk zijn. De nulhypothese luidt in de notatie van 2.2.1 dus:

$\beta_1 = \beta_2 = \dots = \beta_7 = \beta$ . De schattingen van deze  $\beta$ 's blijken te zijn:

$$\hat{\beta}_1 = 0,0733$$

$$\hat{\beta}_2 = -0,0823$$

$$\hat{\beta}_3 = -0,1226$$

$$\hat{\beta}_4 = 0,1256$$

$$\hat{\beta}_5 = 0,1710$$

$$\hat{\beta}_6 = -0,2066$$

$$\hat{\beta}_7 = 0,2690$$

Het algemeen gemiddelde is  $\bar{\beta} = 0,0324$

De kwadraatsom tussen de groepen  $S_{\text{tussen}} = \mathcal{J}_{\omega} - \mathcal{J}_{\Omega} = 5,398$

De kwadraatsom binnen de groepen  $S_{\text{binnen}} = \mathcal{J}_{\omega} = 225,704$

De toetsingsgrootheid:  $\frac{n-k}{k-1} \cdot \frac{S_{\text{tussen}}}{S_{\text{binnen}}} = \frac{203}{6} \cdot \frac{5,398}{225,704} = 0,829$

De overschrijdingskans is  $> 20\%$ , zodat er geen enkele reden is de nulhypothese te verwerpen. Er is dus geen verschil in niveau van de winterkarakterisering in de zeven 30-jaar-tijdvakken te constateren. Natuurlijk in overeenstemming met de constatering in W.R.74-2 dat de reeks vorst-excentriciteiten als een volkomen toevalsreeks is te beschouwen.

Opmerking: Een uitkomst zoals hier, dat de nulhypothese niet verworpen hoeft te worden, geeft verder geen problemen. Een wel verwerpen van de nulhypothese geeft direct aanleiding tot de vraag : Welke waarden voor de  $\beta_i$  moeten dan als reëel worden aanvaard? Het is natuurlijk helemaal niet nodig de berekende schattingen  $\hat{\beta}_i$  zonder meer als de beste schattingen te gebruiken, er kan in beperktere mate onderlinge gelijkheid bestaan. Het is de bedoeling hierop in een vervolgrapport nader in te gaan.

2.2. Regressie-analyse.

2.2.1. Toetsing van een kwadratische regressie.

Als in geval van een relatie tussen  $y$  en  $x$  vermoed wordt, dat deze relatie kwadratisch is, namelijk  $y = a + bx + cx^2$  dan kan men toetsen of dit inderdaad het geval is of dat men toch met een lineair verband kan volstaan. Volgens de algemene theorie kan gesteld worden:

$$\underline{\Omega}: y_i = \alpha + \beta x_{1i} + \gamma x_{2i} + \epsilon_i \quad (2.2.1)$$

met  $x_1 = x$  en  $x_2 = x^2$  verder weer  $E\epsilon_i = 0$ ,  $E\epsilon_i\epsilon_j = \sigma^2\delta_{ij}$  waarin  $\delta_{ij} = 1$  als  $i = j$ ,  $\delta_{ij} = 0$  als  $i \neq j$ .

$$\underline{\omega}: y_i = \alpha + \beta x_i + \epsilon_i \quad (i = 1, \dots, n) \quad (2.2.2)$$

Nu heeft de matrix  $X$  de volgende vorm:

$$X = \begin{pmatrix} 1 & 1 & \dots & \dots & 1 \\ x_1 & x_2 & \dots & \dots & x_n \\ x_1^2 & x_2^2 & \dots & \dots & x_n^2 \end{pmatrix}$$

$$S = XX' = \begin{pmatrix} n & \Sigma x & \Sigma x^2 \\ \Sigma x & \Sigma x^2 & \Sigma x^3 \\ \Sigma x^2 & \Sigma x^3 & \Sigma x^4 \end{pmatrix}$$

Als de determinant van  $S$  door  $D_s$  wordt voorgesteld, dan geldt:

$$D_s = n\Sigma x^2 \Sigma x^4 - n(\Sigma x^3)^2 + 2\Sigma x \Sigma x^2 \Sigma x^3 - (\Sigma x)^2 \Sigma x^4 - (\Sigma x^2)^3$$

met  $t_i = \frac{\sum x^i}{\sqrt{D_s}}$  is  $S^{-1}$  te schrijven als :

$$S^{-1} = \begin{pmatrix} t_2 t_4 - t_3^2 & t_2 t_3 - t_1 t_4 & t_1 t_3 - t_2^2 \\ t_2 t_3 - t_1 t_4 & t_0 t_4 - t_2^2 & t_1 t_2 - t_0 t_3 \\ t_1 t_3 - t_2 & t_1 t_2 - t_0 t_3 & t_3 t_2 - t_1^2 \end{pmatrix} = \begin{pmatrix} T_{24} & T_{23} & T_{13} \\ T_{23} & T_{04} & T_{12} \\ T_{13} & T_{12} & T_{02} \end{pmatrix}$$

als  $T_{24} = t_2 t_4 - t_3^2$  enz. Hiermede worden de kolommen van  $S^{-1}X$ :

$$T_{24} + T_{23}x_i + T_{13}x_i^2$$

$$T_{23} + T_{04}x_i + T_{12}x_i^2$$

$$T_{13} + T_{12}x_i + T_{02}x_i^2$$

$$\text{en derhalve: } \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \\ \hat{\gamma} \end{pmatrix} = S^{-1}X\vec{y} = \begin{pmatrix} T_{24}\Sigma y + T_{23}\Sigma yx + T_{13}\Sigma yx^2 \\ T_{23}\Sigma y + T_{04}\Sigma yx + T_{12}\Sigma yx^2 \\ T_{13}\Sigma y + T_{12}\Sigma yx + T_{02}\Sigma yx^2 \end{pmatrix} \quad (2.2.4)$$

De elementen van de matrix  $D = X'S^{-1}X$  worden nu:

$$D_{ij} = T_{24} + T_{23}x_i + T_{13}x_i^2 + T_{23}x_j + T_{04}x_i x_j + T_{12}x_i x_j + T_{13}x_j^2 + T_{12}x_i x_j^2 + T_{02}x_i^2 x_j^2$$

$$\text{en } \mathcal{J}_\Omega = ||\vec{y} - \hat{n}||^2 = \vec{y}'\vec{y} - \vec{y}'D\vec{y} = \Sigma y^2 - \Sigma \Sigma y_i D_{ij} y_j = \Sigma y_{ij}^2 - T_{24}(\Sigma y)^2 - 2T_{23}\Sigma y \Sigma xy - 2T_{13}\Sigma y \Sigma x^2 y - T_{04}\Sigma x^2 y - 2T_{12}\Sigma x^2 y \Sigma xy - T_{02}(\Sigma x^2 y)^2 \quad (2.2.5)$$

Rechtstreeks berekening van  $\mathcal{J}_\omega$  gaat als volgt: De nulhypothese onder  $\omega$  luidt  $\psi = \gamma$ . Dus  $C = (001)$ , waarmee  $B = C'S^{-1}C = T_{02}$  derhalve  $B^{-1} = T_{02}^{-1}$ .

$$S^{-1}C^1 \text{ wordt nu } \begin{pmatrix} T_{13} \\ T_{12} \\ T_{02} \end{pmatrix}$$

$$\text{en } S^{-1}C'B^{-1}CS^{-1} = \begin{pmatrix} T_{13}^2/T_{02} & T_{12}T_{13}/T_{02} & T_{13} \\ T_{12}T_{13}/T_{02} & T_{12}^2/T_{02} & T_{12} \\ T_{13} & T_{12} & T_{02} \end{pmatrix}$$

$$\text{zodat } \mathcal{E}_\omega = \Sigma y_{ij}^2 \vec{y}' X' \begin{pmatrix} T_{24} - T_{13}^2/T_{02} & T_{23} - T_{12}T_{13}/T_{02} & 0 \\ T_{23} - T_{12}T_{13}/T_{02} & T_{04} - T_{12}^2/T_{02} & 0 \\ 0 & 0 & 0 \end{pmatrix} Xy$$

$$= \Sigma y_{ij}^2 - (T_{24} - \frac{T_{13}^2}{T_{02}})(\Sigma y)^2 - 2(T_{23} - \frac{T_{12}T_{13}}{T_{02}})\Sigma y \Sigma xy - (T_{04} - \frac{T_{12}^2}{T_{02}})(\Sigma xy)^2 \quad (2.2.7)$$

$$\text{en dus: } \mathcal{E}_\omega - \mathcal{E}_\Omega = \frac{T_{13}^2}{T_{02}} (\Sigma y)^2 + 2 \frac{T_{12}T_{13}}{T_{02}} \Sigma y \Sigma xy + \frac{T_{12}^2}{T_{02}} \Sigma x^2 y + 2T_{13} \Sigma y (\Sigma xy)^2 +$$

$$2T_{12} \Sigma xy \Sigma x^2 y + T_{02} (\Sigma x^2 y)^2 = \frac{1}{T_{02}} (T_{13} \Sigma y + T_{12} \Sigma xy + T_{02} \Sigma x^2 y)^2$$

Dit resultaat is sneller te vinden met (1.4.19):  $\mathcal{E}_\omega - \mathcal{E}_\Omega = \hat{\psi}' B^{-1} \hat{\psi}$

waarin  $\hat{\psi} = \hat{r} = T_{13} \Sigma y + T_{12} \Sigma xy + T_{02} \Sigma x^2 y$  volgens (2.2.4) en

$B^{-1} = T_{02}^{-1}$  volgens (2.2.6)

Er valt op te merken, dat  $\mathcal{S}_{\Omega}$  ook te berekenen is op analoge wijze als  $\mathcal{S}_{\omega}$ , omdat in dit geval de X-matrix die geldt als  $\omega$  in plaats van  $\Omega$  wordt genomen direct is op te schrijven, namelijk:

$$\hat{X} = \begin{pmatrix} 1 & 1 & & 1 \\ x_1 & x_2 & & x_n \end{pmatrix}$$

$$\text{en dus } \hat{S} = \hat{X}\hat{X}' = \begin{pmatrix} n & \Sigma x \\ \Sigma x & \Sigma x^2 \end{pmatrix}$$

Als de determinant van  $\hat{S}$  met  $\hat{D}_s$  wordt aangegeven, dan:

$$\hat{D}_s = n\Sigma x^2 - (\Sigma x)^2 \text{ waarna met } r_i = \frac{\Sigma x^i}{\hat{D}_s}$$

$$S^{-1} = \begin{pmatrix} r_2 & -r_1 \\ -r_1 & r_0 \end{pmatrix} \text{ en } \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = S^{-1}Xy = \begin{pmatrix} r_2\Sigma y - r_1\Sigma xy \\ -r_1\Sigma y - r_0\Sigma xy \end{pmatrix}$$

Opmerking:

$$\text{Dus wordt } \hat{\beta} = \frac{n\Sigma xy - \Sigma x\Sigma y}{n\Sigma x^2 - \Sigma x^2} = \frac{\Sigma xy - nx\bar{y}}{\Sigma x^2 - nx^2} = \frac{\Sigma(x - \bar{x})(y - \bar{y})}{\Sigma(x - \bar{x})^2}$$

waarmee de bekende formule voor de lineaire regressie-coëfficiënt van y op x is gevonden. Verder is

$$\hat{\alpha} = \frac{\Sigma x^2 \Sigma y - \Sigma x \Sigma xy}{n\Sigma x^2 - (\Sigma x)^2} = \frac{v\Sigma x^2 - x\Sigma xy}{\Sigma x^2 - nx^2} = y - \hat{\beta}x$$

de bekende formule voor de constante term in de lineaire regressie.

De elementen van de matrix  $D = XS^{-1}X$  worden nu

$$\hat{D}_{ij} = r_2 - r_1x_i - r_1x_j + r_0x_ix_j \text{ en dus:}$$

$$\mathcal{S}_{\omega}^f = \Sigma y^2 - \Sigma_i \Sigma_j v_i \hat{D}_{ij} y_i = \Sigma v^2 - r_2 \Sigma v^2 + 2r_1 \Sigma y \Sigma xy - r_0 \Sigma x^2 y \quad (2.2.8)$$

Dat (2.2.7) identiek is aan (2.2.6) is tenslotte na enig rekenwerk aan te tonen.

De toetsing van  $H_0$  kan nu weer op de bekende manier worden uitgevoerd door  $(n - 1) \frac{s_{\omega}^2 - s_{\Omega}^2}{s_{\Omega}^2}$  te vergelijken met het gekozen significantie-niveau volgens de  $F_{1, n-1}$  verdeling.

Zoals in 1.6 is aangegeven kan ook gebruik worden gemaakt van een vertrouwensinterval:  $\hat{\psi} - t_{n-p} S_{\psi}^{\wedge} < \psi < \hat{\psi} + t_{n-p} S_{\psi}^{\wedge}$

wat nu wordt:  $\hat{\gamma} - t_{n-p} S_{\gamma}^{\wedge} < \gamma < \hat{\gamma} + t_{n-p} S_{\gamma}^{\wedge}$  (2.2.8)

waarin  $S_{\gamma}^{\wedge 2} = \frac{s_{\Omega}^2}{n - p} \cdot T_{02}$  volgens (1.4.15); immers  $\frac{s_{\Omega}^2}{n - p}$  is een schatting van  $\sigma^2$  en  $T_{02}$  is B.

### 2.2.2 Numeriek voorbeeld.

Voor een numeriek voorbeeld zullen gegevens uit een artikel van W. Gerstmann "Regressionsbeziehungen für Turbulenz parameter" worden gebruikt. In Abb. 2 van dit artikel geeft Gerstmann meetresultaten van de turbulente warmtestroom H voor onstabiele situaties in vergelijking met een soort Richardsongetal  $\frac{U_4 - U_{0,25}}{t_{0,25} - t_4} = r$ , waarin  $U_i$  de windsnelheid is op i meter hoogte en  $t_i$  de temperatuur op dezelfde hoogte. Er wordt lineaire regressie-analyse toegepast en getoetst of regressiecoëfficiënt en constantsterm significant van nul afwijken; wat echter niet getoetst wordt is of er misschien een niet lineaire relatie zou kunnen zijn. Dit zal als voorbeeld volgens de in 2.2.1 geschetste methode worden uitgevoerd.

Aangezien de basisgegevens waaruit Gerstmann Abb. 2 samenstelde, niet in het artikel voorkomen zijn de waarden van de grootheden H en r in centimeters uit de figuur afgelezen en y en x genoemd, waarbij  $x \sim 2,5 \frac{U_4 - U_{0,25}}{t_{0,25} - t_4}$  en  $y \sim 25H$  is. In bijgaande figuur 1 zijn x en y tegen elkaar uitgezet.



Het regressiemodel  $\Omega$ :  $y = \alpha + \beta x + \gamma x^2 + \epsilon$  wordt nu aangepast. De parameters  $\alpha$ ,  $\beta$  en  $\gamma$  worden volgens (2.2.3)

berekend, met als resultaat:

$$\hat{\alpha} = 0,206 \quad \hat{\beta} = 1,608 \quad \hat{\gamma} = -0,093$$

terwijl volgens (1.4.12) de standaard deviaties zijn:

$$S_{\hat{\alpha}} = 0,507 \quad S_{\hat{\beta}} = 0,347 \quad S_{\hat{\gamma}} = 0,042 .$$

Aangezien het de bedoeling is te toetsen of het lineaire model  $\omega$ :  $y = \alpha + \beta x + \epsilon$  voldoende is, wordt de nulhypothese  $H_0: \psi = \gamma = 0$  ingevoerd.

Voor de F-toets zijn de numerieke resultaten nu:

$$\begin{aligned} \mathcal{J}_{\Omega} &= 101,109 \quad \text{berekend volgens (2.2.4)} \\ \mathcal{J}_{\omega} - \mathcal{J}_{\Omega} &= 12,230 \quad \text{berekend volgens (2.2.7)} \\ \text{zodat } \frac{n-p}{q} \cdot \frac{\mathcal{J}_{\omega} - \mathcal{J}_{\Omega}}{\mathcal{J}_{\Omega}} &= \frac{43,3}{1} \cdot \frac{12,230}{101,109} = 4,84. \end{aligned}$$

De overschrijdingskans van deze waarde is volgens de  $F_{1,40}$  verdeling ongeveer 4%. Bij een 5% significantiedrempel moet de nulhypothese  $\gamma = 0$  dus verworpen worden, met andere woorden de onderstelling dat men met een lineaire relatie kan volstaan moeten worden verworpen ten gunste van een kwadratische relatie. Men verkrijgt hetzelfde resultaat als gebruik gemaakt wordt van het vertrouwensinterval volgens (1.6.2). De 5%  $t_{40}$  waarde is namelijk 2,02 zodat  $\hat{\psi} \pm t_{n-p} S_{\hat{\psi}}$  wordt  $-0,093 \pm 2,02 \times 0,042$ . Dus  $-0,178 < \psi = \gamma < -0,008$ . De waarde  $\gamma = 0$  ligt buiten het vertrouwensinterval; de nulhypothese moet verworpen worden.

$\mathcal{J}_{\omega}$  is ook rechtstreeks te berekenen als  $y = \tilde{\alpha} + \tilde{\beta}x$  wordt aangepast. Er wordt dan gevonden  $\tilde{\alpha} = 0,973$  en  $\tilde{\beta} = 0,878$  met  $S_{\tilde{\alpha}} = 0,386$  en  $S_{\tilde{\beta}} = 0,105$  terwijl dan  $\mathcal{J}_{\omega}$  berekend volgens  $\mathcal{J}_{\omega} = \|\tilde{y} - \tilde{h}\|^2$  gelijk aan 133,335 is.

In het voorgaande werd berekend  $\mathcal{J}_\Omega = 101,109$  en  $\mathcal{J}_\omega - \mathcal{J}_\Omega = 12,230$ , zodat  $\mathcal{J}_\omega = 101,109 + 12,230 = 113,339$ ; dit klopt dus zeer fraai met de rechtstreekse berekening.

In figuur 1 is zowel het lineaire verband  $y = 0,973 + 0,878x$ , als het kwadratische  $y = 0,206 + 1,608x - 0,093x^2$  getekend. Het lineaire verband moet natuurlijk hetzelfde zijn als Gerstmann in Abb. 2 geeft. Dit klopt op kleine verschillen na, die het gevolg zijn van afleesonauwkeurigheden.

Tenslotte kan men zich nog afvragen of het nodig is een constante term in de regressievergelijking mee te nemen. Het vertrouwensinterval (2.2.8) voor  $\alpha$  is  $0,206 - 2,02 \times 0,507 \leq \alpha < 0,206 + 2,02 \times 0,507$  of  $-0,892 < \alpha < 1,230$  zodat er geen reden is  $\alpha \neq 0$  te stellen. Wordt nu  $y = \beta^*x + \gamma^*x^2$  aangepast, dan blijkt  $\beta^* = 1,723$  en  $\gamma^* = -0,105$ . De kromme  $y = 1,723x - 0,105x^2$  is eveneens in figuur 1 getekend.

### 3. Een bewerkingstechniek voor multivariate analyse.

#### 3.1. Inleiding.

In hoofdstuk 1 is de oplossing gezocht voor de beste schatting van een lineaire relatie tussen een grootte  $y$  en  $p$  grootheden  $x_i$ . De geometrische voorstelling gaf aan, dat deze oplossing te vinden is door in een  $n$ -dim. ruimte een vector  $y$  te projecteren op een  $p$ -dim. deelruimte  $V_p$  gevormd door  $p$ -vectoren  $x_i$  ( $i = 1, \dots, p$ ). In feite is er daarbij, zoals reeds in 1.2. is gesteld, sprake van de ontbinding van een vector  $\vec{y}$  in een component evenwijdig aan  $V_p$  ( $\vec{y}_{//}$ ) en een component loodrecht op  $V_p$  ( $\vec{y}_\perp$ ). De eerste geeft de "aanpassing" aan het lineaire model en de tweede is een maat voor de afwijking tussen waarneming en model.

Nu is het meestal zo dat de  $p$ -grootheden  $x_i$  niet allen even belangrijk zijn, dat wil zeggen, dat de variatie van  $y$  grotendeels bepaald wordt door een beperkt aantal  $x_i$ 's terwijl de anderen nauwelijks een rol spelen. Voor het bepalen van de  $x_i$ 's die tezamen een optimaal resultaat geven ter "verklaring" van de variatie van  $y$  zijn diverse berekeningsmethoden ontwikkeld, onder anderen door Efroymson. Hierop zal in 4 nader worden ingegaan.

Voor de berekeningen van dergelijke problemen zijn in het algemeen zeer grote aantallen bewerkingen nodig, zodat een rekenmachine een noodzakelijk hulpmiddel is.

Het gaat daarbij dan om de oplossing van een stelsel lineaire vergelijkingen, waarvoor een vector, matrixvorm in (1.4.7) was gevonden, namelijk:  $\hat{\beta} = S^{-1}X\vec{y}$  (met voorwaarde, dat de matrix  $S$  nonsingulier is).

Mathematisch gezien is hiermede het probleem opgelost. Numeriek kunnen er echter moeilijkheden ontstaan, als de matrix  $S$  slechts geconditioneerd is, dat wil zeggen, dat kleine variaties in de componenten van de vector  $y$  tot zeer grote variaties in de  $\beta$  leiden; n.l. hetzij ten gevolge van de afrondingsfouten die bij de vele computerbewerkingen cumulatief een groot effect kunnen hebben, hetzij omdat ergens nagenoeg even grote getallen van elkaar worden afgetrokken, waardoor de relatieve fout zeer groot wordt. Er zijn methoden ontwikkeld waarmee deze moeilijkheden kunnen worden vermeden. Hierbij wordt de vectorvoorstelling van het probleem door ontbinding van de vectoren in een onderling orthogonaal stelsel getransformeerd.

Voor het regressieprobleem, dat in dit rapport aan de orde, is ook orthogonalisatie toegepast maar niet zozeer om het conditieprobleem te ontlopen, maar om over een rekentechniek te beschikken waarmee de optimale keuze van het beperkte aantal  $x_i$ 's kan worden bereikt.

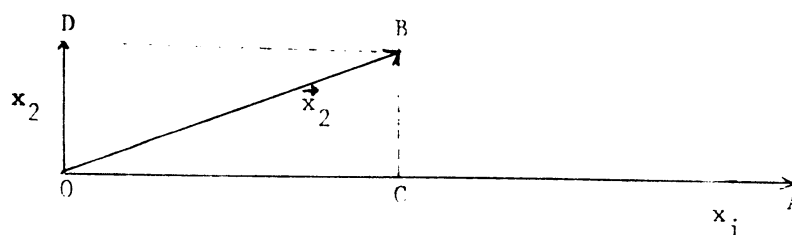
Het conditieprobleem is bij het relatief kleine aantal matrixelementen in dit geval doorgaans nauwelijks aanwezig (zie Orden p. 41).

In het volgende wordt nu de techniek, die de in het voorgaande genoemde orthogonalisatie realiseert, uiteengezet. Daarbij zullen ook de standaarddeviaties van de geschatte parameters van de lineaire relaties worden verkregen. Een uitvoerige beschrijving van de methode is te vinden in: A.P. Dempster: "Elements of continuous multivariate analysis".

### 3.2. Orthogonalisatie.

Uitgangspunt is een  $p$ -dim. ruimte  $V_p$  opgespannen door  $p$  vectoren  $\vec{x}_i$  ( $i = 1, \dots, p$ ). Deze  $p$  vectoren dienen te worden vervangen door  $p$  vectoren  $\vec{x}_i^*$  die alle onderling loodrecht op elkaar staan en dezelfde ruimte  $V_p$  bepalen. De  $\vec{x}_i$  vormen een willekeurig basis voor  $V_p$ ; de  $\vec{x}_i^*$  een orthogonale basis van  $V_p$ . Men kan de vectoren  $\vec{x}_i^*$  achtereenvolgens als volgt vinden:  $\vec{x}_1$  wordt  $\vec{x}_1^*$  genoemd, daarna wordt  $x_2$  ontbonden in een vector evenwijdig  $\vec{x}_1 = \vec{x}_1^*$  en één loodrecht op  $\vec{x}_1 = \vec{x}_1^*$ . Deze laatste wordt  $\vec{x}_2^*$  genoemd.  $\vec{x}_3^*$  wordt bepaald als de component van  $x_3$  die loodrecht is  $\vec{x}_2^*$  zowel als op  $\vec{x}_1^*$ . Op overeenkomstige wijze worden de overigen  $\vec{x}_i^*$  verkregen.

Een uitdrukking voor dit proces is te vinden met behulp van het inproduct van twee vectoren, zie bijgaande figuur:





Een matrix  $Q$  waarvan de elementen de inproducten  $(\vec{x}_i, \vec{x}_j)$  zijn, is dan te schrijven als:

$$Q = \{\vec{x}\vec{x}'\} = M\{x x'\} = MTM' \quad (3.2.6)$$

waarin  $T$  een diagonaalmatrix met elementen  $(\vec{x}_j, \vec{x}_j)$  in de diagonaal is; de elementen  $(\vec{x}_i, \vec{x}_j)$  zijn immers nul wegens de orthogonaliteit van de  $\vec{x}_i$ .

De matrix  $Q$  is het uitgangspunt voor een operatie waarbij stap voor stap de orthogonalisatie wordt doorgevoerd. Als met  $s$  het aantal vectoren  $\vec{x}_i$  wordt aangegeven dat op zeker moment georthogonaliseerd is, dan zal een operator worden gedefinieerd die uitgaande van de beginsituatie de orthogonalisatie tot een willekeurige  $s$  tot stand brengt. Daaruit is de operator af te leiden die het proces van  $s$  naar  $s + 1$  voert. Het ligt voor de hand voor vectoren en matrices een splitsing in te voeren op de volgende wijze:

Als  $\vec{v}$  en  $K$  respectievelijk een willekeurige  $(p \times 1)$  vector en  $(p \times p)$  matrix is:

$$\vec{v} = \begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_s \end{bmatrix} \\ \begin{bmatrix} v_{s+1} \\ \vdots \\ v_p \end{bmatrix} \end{bmatrix} \text{ en } K = \begin{bmatrix} K_{11}^{s \times s} & K_{12}^{s \times (p-s)} \\ K_{21}^{(p-s) \times s} & K_{22}^{(n-s) \times (p-s)} \end{bmatrix} \quad (3.2.7)$$

Als deze splitsing op (3.2.4) wordt toegepast ontstaat:

$$\begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \end{bmatrix} = \begin{bmatrix} M_{11} & 0 \\ M_{21} & M_{22} \end{bmatrix} \begin{bmatrix} \vec{x}_1^* \\ \vec{x}_2^* \end{bmatrix} = \begin{bmatrix} M_{11}\vec{x}_1^* \\ M_{21}\vec{x}_1^* + M_{22}\vec{x}_2^* \end{bmatrix}$$

$$\text{waarbij: } \vec{x}_1 = \begin{bmatrix} x_1 \\ \vdots \\ x_s \end{bmatrix} \text{ en } \vec{x}_2 = \begin{bmatrix} x_{s+1} \\ \vdots \\ x_p \end{bmatrix}.$$

$$\text{Dus: } \vec{x}_1 = M_{11} \vec{x}_1^* \quad (3.2.8)$$

$$\vec{x}_2 = M_{21} \vec{x}_1^* + M_{22} \vec{x}_2^* \quad (3.2.9)$$

Nu is iedere vector van  $\vec{x}_2^*$  loodrecht op iedere vector, die in  $\vec{x}_1^*$  is begrepen, dus loodrecht op de ruimte door deze vectoren  $\vec{x}_1^*$  gevormd en dit is dezelfde ruimte die door de  $\vec{x}_1$  vectoren wordt opgespannen. (3.2.9) betekent dus een splitsing van de  $\vec{x}_2^*$  vectoren in componenten  $M_{21} \vec{x}_1^*$  die evenwijdig zijn aan de ruimte van de  $\vec{x}_1$  vectoren en componenten  $M_{22} \vec{x}_2^*$  die loodrecht op die ruimte staan. Hiervoor wordt de notatie:

$$\vec{x}_{2.1} = M_{22} \vec{x}_2^* \quad (3.2.10)$$

ingevoerd.

Verder geldt:  $MN = I$ . Dus:

$$\begin{pmatrix} M_{11} & 0 \\ M_{21} & M_{22} \end{pmatrix} \cdot \begin{pmatrix} N_{11} & N_{12} \\ N_{21} & N_{22} \end{pmatrix} = \begin{pmatrix} N_{11} N_{11} & M_{11} N_{12} \\ M_{21} N_{11} + M_{22} N_{21} & M_{21} N_{12} + M_{22} N_{22} \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}$$

$$\text{Dus } M_{11} N_{11} = I, \text{ derhalve } N_{11} = M_{11}^{-1} \quad (3.2.11)$$

$$M_{11} N_{12} = 0, \text{ wegens } M_{11} \neq 0 \text{ dus } N_{12} = 0 \quad (3.2.12)$$

$$M_{21} N_{12} + M_{22} N_{22} = M_{22} N_{22} = I \text{ waaruit } N_{22} = -M_{22}^{-1} \quad (3.2.13)$$

$$M_{21} N_{11} + M_{22} N_{21} = 0 \text{ of } M_{21} N_{11} = -M_{22} N_{21}$$

Voor de laatste formule zal ingevoerd worden de notatie  $H_{21}$ :

$$H_{21} = M_{21} N_{11} = -M_{22} N_{21} \quad (3.2.14)$$

Nu zijn  $N_{11}$  en  $M_{22}$  beide driehoeksmatrices met diagonaal elementen die alle 1 zijn, derhalve is uit (3.2.11) af te lezen:

- (i) De eerste rij van  $H_{21}$  is gelijk aan het tegengestelde van de eerste rij van  $N_{21}$ .
- (ii) De laatste kolom van  $H_{21}$  is gelijk aan de laatste kolom van  $M_{21}$ .

Uit (3.2.9) en (3.2.10) volgt:  $\vec{x}_2 = M_{21}\vec{x}_1^* + \vec{x}_{2.1}$

Uit (3.2.8) en (3.2.11) :  $\vec{x}_1^* = M_{11}^{-1}\vec{x}_1 = N_{11}\vec{x}_1$

Dus:  $\vec{x}_{2.1} = \vec{x}_2 - M_{21}\vec{x}_1^* = \vec{x}_2 - M_{21}N_{11}\vec{x}_1 = \vec{x}_2 - H_{21}\vec{x}_1$  (3.2.15)

### 3.3. De "Sweep" Operator.

Nu kan de operator ingevoerd worden die de overgang van begin tot de toestand s moet bewerkstellen. Hiervoor zal de notatie SWP, die "sweep operator" 1) betekent, worden gebruikt:

$$\text{SWP}(0, 1, \dots, s) Q = \text{SWP}(0, 1, \dots, s) \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix} = \begin{pmatrix} -Q_{11}^{-1} & Q_{11}^{-1}Q_{12} \\ Q_{21}Q_{11}^{-1} & Q_{22} - Q_{21}Q_{11}^{-1}Q_{12} \end{pmatrix}$$

met  $\text{SWP}(0) Q = Q$  en  $\text{SWP}(0, 1, \dots, p) Q = -Q^{-1}$  (3.3.1)

Als het splitsingsprocedé op (3.2.6) wordt toegepast, ontstaat:

$$Q = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix} = \begin{pmatrix} M_{11} & 0 \\ M_{21} & M_{22} \end{pmatrix} \cdot \begin{pmatrix} T_{11} & 0 \\ 0 & T_{22} \end{pmatrix} \cdot \begin{pmatrix} M'_{11} & M'_{21} \\ 0 & M'_{22} \end{pmatrix} = \begin{pmatrix} M_{11}T_{11}M'_{11} & M_{11}T_{11}M'_{21} \\ M_{21}T_{11}M'_{11} & M_{21}T_{11}M'_{21} + M_{22}M'_{22} \end{pmatrix}$$

Hierop  $\text{SWP}(0, 1, \dots, s)$  toepassen levert:

1) Deze term is volgens Dempster afkomstig van A.E. Beaton (1964)



$$\text{SWP}(0, 1, \dots, s) Q = \begin{pmatrix} -M'_{11}{}^{-1} T'_{11} M'_{11}{}^{-1} & M'_{11}{}^{-1} M'_{21} & 1 \\ M_{21} M'_{11}{}^{-1} & M_{22} T'_{22} M'_{22} & \end{pmatrix}$$

Nu is  $M'_{11} T'_{11} M'_{11} = M'_{11} (\vec{x}'_1 \vec{x}'_1) M'_{11} = (\vec{x}'_1 \vec{x}'_1) = Q_{11}$ .

Verder is volgens (3.2.11) en (3.2.13):  $M_{21} M'_{11}{}^{-1} = H_{21}$  en

$M'_{11}{}^{-1} M'_{21} = H'_{21}$  en tenslotte met (3.2.10):

$M_{22} T'_{22} M'_{22} = M_{22} (\vec{x}'_2 \vec{x}'_2) M'_{22} = (\vec{x}'_{21} \cdot \vec{x}'_{21})$  hiervoor invoeren

$$M_{22.1}, \text{ dus: } M_{22.1} = (\vec{x}'_{2.1} \cdot \vec{x}'_{2.1}) = M_{22} T'_{22} M'_{22} \quad (3.3.2)$$

hieruit zien we dat het element op de eerste rij en in de eerste kolom van  $M_{22.1}$  gelijk is aan het eerste diagonaal element van  $T'_{22}$ . Nu kan geschreven worden:

$$\text{SWP}(0, 1, \dots, s) M = \begin{pmatrix} -M'_{11}{}^{-1} & H'_{21} \\ H_{21} & M_{22.1} \end{pmatrix} \quad (3.3.3)$$

Uit het voorgaande blijkt dus dat bij toepassing van  $\text{SWP}(0, 1, \dots, s)M$ , voor achtereenvolgens  $s = 1, \dots, p$ , uit de submatrices  $H_{21}$  de transformatie matrices  $M$  en  $N$  zijn af te leiden met de regels (i) en (ii), terwijl uit de submatrix  $M_{22.1}$  de matrix  $T$  volgt.

Als de formule (3.2.15) als  $\vec{x}'_2 = \vec{x}'_{2.1} + H_{21} \vec{x}'_1$  geschreven wordt, dan is in te zien dat de matrix  $H_{21}$  in zijn rijen de coëfficiënten levert van de projecties van de vectoren uit  $\vec{x}'_2$  op de ruimte van de geschatte regressiecoëfficiënten tussen één der  $x_j$  ( $j = s+1, \dots, p$ ) en alle  $x_k$  ( $k=1, \dots, s$ ).

Verder leveren de diagonaalelementen van  $M_{22.1}$  de kwadraten van de componenten van de  $\vec{x}'_2$  vectoren die loodrecht op  $\vec{x}'_1$  staan; dit is in de regressie-analyse de rest variantie nadat de regressie tussen de  $\vec{x}'_2$  en de  $\vec{x}'_1$  grootheden is tot stand gebracht. Tenslotte levert  $M'_{11}{}^{-1}$  informatie over de standaarddeviatie van de regressiecoëfficiënten (zie (1.4.12)).  $M$  is gelijk aan de matrix  $S$ . Dit zal nog nader worden bekeken.

Eerst moet nu nog worden aangegeven hoe de opeenvolgende operaties  $SWP(0, 1, \dots, s)$  voor  $s = 1, \dots, p$  uit elkaar worden afgeleid. Hiertoe de reciproke operator RSW invoeren volgens:

$$RSW(0, 1, \dots, s) \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} = \begin{pmatrix} -M_{11}^{-1} & -M_{11}^{-1}M_{12} \\ -M_{21}M_{11}^{-1} & M_{21}M_{11}^{-1}M_{12} \end{pmatrix} \quad (3.3.4)$$

Gemakkelijk is in te zien dat inderdaad:

$$RSW(0, 1, \dots, s) SWP(0, 1, \dots, s) M = M$$

evenals:

$$SWP(0, 1, \dots, s) RSW(0, 1, \dots, s) M = M$$

Er wordt nu een  $SWP(s)$  ingevoerd waarvoor geldt:

$$SWP(s) SWP(0, 1, \dots, s-1) = SWP(0, 1, \dots, s)$$

en dus:

$$SWP(s) M = SWP(0, 1, \dots, s) RSW(0, 1, \dots, s-1)M$$

M eerst als volgt schrijven:

$$M = \begin{pmatrix} \tilde{M}_{11} & m' & \tilde{M}_{12} \\ m & c & r \\ \tilde{M}_{21} & r' & \tilde{M}_{22} \end{pmatrix} \quad \text{Waarin } \tilde{M}_{11} \text{ een } (s-1) \times (s-1) \text{ matrix,}$$

m een  $1 \times (s-1)$  matrix,  
r een  $1 \times (p-s)$  matrix en  
c één enkel element is.

Dan is:

$$RSW(0, 1, \dots, s-1)M = \begin{pmatrix} -\tilde{M}_{11}^{-1} & -\tilde{M}_{11}^{-1}m' & -\tilde{M}_{11}^{-1}\tilde{M}_{12} \\ -m\tilde{M}_{11}^{-1} & c-m\tilde{M}_{11}^{-1}m' & r-m\tilde{M}_{11}^{-1}\tilde{M}_{12} \\ -\tilde{M}_{21}\tilde{M}_{11}^{-1} & r'-\tilde{M}_{21}\tilde{M}_{11}^{-1}m' & M_{22}-\tilde{M}_{21}\tilde{M}_{11}^{-1}\tilde{M}_{12} \end{pmatrix}$$

Nu hierop  $SWP(0, 1, \dots, s)$  toepassen; hiervoor is eerst de reciproke nodig van de matrix:

$$\begin{pmatrix} -M_{11}^{-1} & -M_{11}^{-1}m' \\ -mM_{11}^{-1} & c-mM_{11}^{-1}m' \end{pmatrix}$$

Stel dat deze te schrijven is als  $\begin{pmatrix} A & B' \\ B & D \end{pmatrix}$  dan moet gelden:

$$\begin{pmatrix} A & B' \\ B & D \end{pmatrix} \cdot \begin{pmatrix} \tilde{M}_{11}^{-1} & -\tilde{M}_{11}^{-1}m' \\ -m\tilde{M}_{11}^{-1} & c-m\tilde{M}_{11}^{-1}m' \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}$$

of:

$$\begin{pmatrix} -\tilde{A}\tilde{M}_{11}^{-1} - B'm\tilde{M}_{11}^{-1} & -\tilde{A}\tilde{M}_{11}^{-1}m' + B'(c-m\tilde{M}_{11}^{-1}m') \\ -B\tilde{M}_{11}^{-1} - Dm\tilde{M}_{11}^{-1} & -B\tilde{M}_{11}^{-1}m' + D(c-m\tilde{M}_{11}^{-1}m') \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}$$

A en B oplossen uit gelijkstelling van de elementen van de eerste rijen geeft  $B' = -mc^{-1}$  en  $A = -\tilde{M}_{11}^{-1} + m'c^{-1}m$ .

B en D uit elementen van de tweede rijen geeft  $D = c^{-1}$  en  $B = -c^{-1}m$ .

Derhalve:

$$\begin{pmatrix} \tilde{M}_{11}^{-1} & -\tilde{M}_{11}^{-1}m' \\ -m\tilde{M}_{11}^{-1} & c-m\tilde{M}_{11}^{-1}m' \end{pmatrix}^{-1} = \begin{pmatrix} -\tilde{M}_{11}^{-1} + m'c^{-1}m & -m'c^{-1} \\ -c^{-1}m & c^{-1} \end{pmatrix}$$

SWP(0, 1, ..., s) RSW(0, 1, ..., s-1)M wordt nu:

$$SWP(s)M = \begin{pmatrix} \tilde{M}_{11}^{-1} - m'c^{-1}m & m'c^{-1} & \tilde{M}_{12} - m'c^{-1}r \\ c^{-1}m & -c^{-1} & c^{-1}r \\ \tilde{M}_{21} - r'c^{-1} & r'c^{-1} & Q_{22} - r'c^{-1}r \end{pmatrix}$$

Met het bovenstaande is de operator SWP(s) afgeleid die de opeenvolgende operaties SWP(0, 1, ..., s) voor  $s = 1, \dots$  uit elkaar kan afleiden. Deze is ook als volgt te definiëren, als  $m_{ij}$  de elementen van M en  $\tilde{m}_{ij}$  die van SWP(s)M resp. RSW(s)M zijn:

$$\begin{aligned}
 \tilde{m}_{ss} &= -1/m_{ss} \\
 \tilde{m}_{is} &= m_{is}/m_{ss} \\
 \tilde{m}_{sj} &= m_{sj}/m_{ss} \\
 \tilde{m}_{ij} &= m_{ij} - m_{is}m_{sj}/m_{ss}
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \begin{array}{l} i \text{ en } j \neq s \\ \\ \\ \end{array}$$

SWP(s) =

en

$$\begin{aligned}
 \tilde{m}_{ss} &= -1/m_{ss} \\
 \tilde{m}_{is} &= -m_{is}/m_{ss} \\
 \tilde{m}_{sj} &= -m_{sj}/m_{ss} \\
 \tilde{m}_{ij} &= m_{ij} - m_{is}m_{jj}/m_{ss}
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} (3.3.5)$$

RSW(s)

N.B. Door Dempster wordt een geheel andere volgorde bij de afleiding van de methode aangehouden. Eerst wordt de overgang van  $s$  naar  $s + 1$  in het proces moeizaam omschreven, waarbij nagegaan wordt, wat er met de elementen van  $M_{11}^{-1}$ ,  $H_{21}$  en  $M_{22.1}$  gebeurt bij de orthogonalisatie. Daarna wordt de operator (3.3.5) ingevoerd en wordt aangetoond, dat deze aan de eisen voldoet. (3.3.3) komt ergens uit de lucht vallen, hiervan blijkt pas later de bedoeling, terwijl tenslotte (3.3.1) als eigenschap blijkt te gelden.

Alternatieve Sweepoperatie definitie.

In afwijking van de definitie, zoals die door Dempster is ingevoerd, kan men ook

$$\text{SWP}^*(0, 1, \dots, s) \begin{pmatrix} M_{21} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} = \begin{pmatrix} M_{11}^{-1} & M_{11}^{-1}M_{12} \\ -M_{21}^{-1}M_{11}^{-1} & M_{22} - M_{21}^{-1}M_{11}^{-1} & M_{12} \end{pmatrix}$$

definieren. En

$$\begin{aligned}
 \tilde{m}_{ss} &= 1/m_{ss} \\
 \tilde{m}_{is} &= -m_{is}/m_{ss} \\
 \tilde{m}_{sj} &= m_{sj}/m_{ss} \\
 \tilde{m}_{ij} &= m_{ij} - m_{is}m_{sj}/m_{ss}
 \end{aligned}
 \quad \begin{array}{l} \\ \\ \\ \end{array} \quad \begin{array}{l} i; j \neq s \\ \\ \\ \end{array} \quad (3.3.6)$$

SWP\*(s) :

Dit heeft het voordeel dat nu geen afzonderlijke inverse operator behoeft te worden ingevoerd.  $SWP^*(s) SWP^*(s) M = M$  zoals gemakkelijk is in te zien. Wel moet bedacht worden, dat de rijen van  $H_{21}$  nu de regressiecoëfficiënt met tegengesteld teken leveren.

4. Toepassing van de Sweepoperator techniek bij "optimalisering" van multiple lineaire regressie.

4.1 Inleiding.

Van de techniek die in 3 is ontwikkeld, wordt gebruik gemaakt voor de oplossing van het probleem van de multiple regressie, waarbij het er om gaat voor de regressie van  $y$  op  $x_1, \dots, x_p$  te bepalen, welke van de  $p$ -grootheden  $x_i$  men in de regressie moet betrekken of met andere woorden, welke  $x_i$ 's geven een significante bijdrage ter verklaring van de variatie van  $y$  en welke  $x_i$ 's kunnen buiten beschouwing blijven.

De optimale oplossing is natuurlijk die welke mogelijke combinaties van 1 tot en met  $p$  variabelen uit de  $p$   $x_i$ 's in beschouwing neemt en voor iedere combinatie toetst of de nulhypothese, dat de regressiecoëfficiënten voor de niet in de combinatie voorkomende  $x_i$ 's gelijk aan nul zijn, moet worden verworpen. Van de combinaties waarvoor deze nulhypothese niet verworpen behoeft te worden kiest men dan die welke de kleinste restvariantie geeft. Nu is dat aantal combinaties  $2^p - 1$  en dit wordt bijvoorbeeld bij  $p = 10$  al meer dan 1000. De hoeveelheid rekenwerk wordt dan zeer groot. Een methode die minder rekenwerk eist is bijvoorbeeld door Efroymsen beschreven in het artikel "Multiple Regression Analysis".

In de eerste plaats zal hier een korte algemene beschrijving van de methode worden gegeven, daarna zullen in 4.4 enkele details en onduidelijkheden in het artikel van Efroymsen nader worden besproken.

4.2. Algemene beschrijving van de "optimaliserings" methode.

De methode bestaat in het successievelijk toevoegen aan een verwijdering uit de regressie van  $x_i$ 's en wel als volgt:

In ieder stadium wordt nagegaan hoe groot de restvariantie-reductie is die de niet in de regressie meespelende  $x_i$ 's ieder afzonderlijk bij toevoegen zouden veroorzaken. Van de maximale waarde van deze restvariantie-reductie wordt getoetst of de bijdrage significant is. Is dit het geval dan wordt de betrokken  $x_i$  in de regressie opgenomen. Vervolgens wordt nagegaan of de reeds meespelende  $x_i$ 's moeten blijven meedoen. Het is namelijk zo, dat door toevoeging van een nieuwe  $x_i$  de bijdrage van de reeds aanwezig  $x_i$ 's verandert, in het algemeen kleiner wordt. Het is dus mogelijk dat een bepaald  $x_i$  nu niet langer een significante bijdrage levert. Van de meedoende  $x_i$ 's wordt dus nagegaan welke de kleinste restvariantie reductie geeft en vervolgens wordt getoetst of deze reductie significant is; zo ja, dan wordt de betrokken  $x_i$  uit de regressie verwijderd. Daarna wordt opnieuw onderzocht of nog een nieuwe  $x_i$  moet worden toegevoegd. Dit proces wordt zolang voortgezet totdat van de niet meespelende  $x_i$ 's er geen meer is die een significante bijdrage tot restvariantie reductie kan leveren.

De vraag die nu naar voren komt is natuurlijk: "Vindt men op bovenstaande wijze ook de optimale oplossing?" Helaas is op deze vraag geen antwoord te geven. Het is mogelijk dat bepaalde gunstige combinaties in bovenstaande methode niet aan bod komen. Te verwachten is wel dat deze heen en weer gaande methode een beter resultaat geeft dan 1<sup>o</sup> de simpele "voorwaarts"methode die alleen maar stopsgewijs nagaat welke  $x_i$ 's moeten worden toegevoegd, dan wel 2<sup>o</sup> de "achterwaarts"methode die, uitgaande van volledige regressie op alle  $p$   $x_i$ 's nagaat welke  $x_i$  verwijderd kunnen worden.

4.3. Toetsing tijdens het "optimaliserings"procedé.

Uitgaande van  $\Omega: y_k = \sum_{i=1}^p b_i x_{ik} + \epsilon_k ; \epsilon_k \Rightarrow N(0, \sigma^2); E\epsilon_k \epsilon_l = 0;$   
 $k = 1, \dots, n$

wordt de matrix

$$X = \begin{pmatrix} x_{11} & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{p1} & \dots & x_{pn} \\ y_1 & \dots & y_n \end{pmatrix}$$

gevormd, waaruit:

$$M = XX' = \begin{pmatrix} \sum_{j=1}^n x_{1j}^2 & \sum_{j=1}^n x_{1j}x_{2j} & \dots & \sum_{j=1}^n x_{1j}\sum x_{pj} & \sum_{j=1}^n x_{1j}y_j \\ \vdots & \vdots & & \vdots & \vdots \\ \sum_{j=1}^n x_{pj}x_{1j} & \dots & \dots & \sum_{j=1}^n y_j x_{pj}^2 & \sum_{j=1}^n x_{pj}y_j \\ \sum_{j=1}^n y_j x_{1j} & \dots & \dots & \sum_{j=1}^n y_j x_{pj} & \sum y_j^2 \end{pmatrix}$$

Stel nu dat in een zeker stadium van de bewerking de groot-  
 heden  $x_{r_1}, \dots, x_{r_t}$  in de regressie zijn betrokken, dat wil zeggen,  
 dat  $SWP^*(0, r_1, \dots, r_t)$  is uitgevoerd. M is veranderd in een matrix  
 voor te stellen is door:  $\tilde{M}^t = SWP(0, r_1, \dots, r_t)M$

$$\begin{pmatrix} (\tilde{M}^t)^{-1} & \tilde{H}_{12}^t \\ -\tilde{H}_{21}^t & \tilde{M}_{22.1}^t \end{pmatrix}$$

waarvoor (3.3.6) is gebruikt.

De elementen van  $M$  zullen door  $m_{ij}$  en die van  $\hat{M}^t$  door  $\hat{m}_{ij}^t$  worden voorgesteld, derhalve: als  $i$  en  $j$  beide tot  $r_1, \dots, r_t$  behoren dan  $\hat{m}_{ij}^t \in M_{11}^{-1}$ ; als  $i$  wel en  $j$  niet tot  $r_1, \dots, r_t$  behoren dan  $\hat{m}_{ij}^t \in H_{12}$ , enz.

Er wordt nu voor alle  $s$  de waarde berekend van

$$V_s = \frac{\hat{m}_{p+1,s}^t \hat{m}_{s,p+1}^t}{\hat{m}_s^t \hat{m}_s^t}. \text{ Dit is volgens (3.3.6) de verandering die}$$

de term  $\hat{m}_{p+1,p+1}^t$  ondergaat als  $SWP^*(s)$  op  $\hat{M}^t$  wordt toegepast en dit is dus volgens de regressie-analyse de verandering van de restvariantie. Behoort  $s$  niet tot  $r_1, \dots, r_t$  dan is  $x_s$  niet in de regressie opgenomen en zijn  $\hat{m}_{p+1,s}^t$  en  $\hat{m}_{s,p+1}^t$  elementen van  $\hat{M}_{22.1}^t$  en dus aan elkaar gelijk, derhalve is  $V_s > 0$ ; behoort

$s$  wel tot  $r_1, \dots, r_t$  dan geldt  $\hat{m}_{p+1,s}^t \in -\hat{H}_{21}^t$  en  $\hat{m}_{s,p+1}^t \in \hat{H}_{12}^t$  dus  $\hat{m}_{p+1,s}^t = -\hat{m}_{s,p+1}^t$  derhalve  $V_s < 0$ . Toepassing van  $SWP(s)$  betekent

in het geval  $s$  niet tot  $r_1, \dots, r_t$  behoort toevoeging van  $x_s$  aan de regressie en volgens (3.3.6) dus verkleining van  $\hat{m}_{p+1,p+1}^t$ , dus verkleining van de restvariantie; terwijl toepassing op  $s$  als deze wel tot  $r_1, \dots, r_t$  behoort betekent dat  $x_s$  uit de regressie wordt verwijderd, waarbij de restvariantie dus toeneemt.

Is de situatie  $r_1, \dots, r_t$  ontstaan na toevoeging van een nieuwe  $x_i$ , dan wordt voor alle negatieve  $V_i$ , het minimum  $|V_m|$  van  $|V_i|$  bepaald. Stel  $|V_m|$  behoort bij  $x_m$ . Getoetst wordt dan of de variantiereductie nog significant van nul afwijkt (in feite wordt getoetst of de bijbehorende regressiecoëfficiënt van nul afwijkt).

Is dit zo, dan volgt verwijdering van  $x_m$ . Voor de toetsing geldt dat  $\hat{m}_{p+1,p+1}^t$  overeenkomt met  $\mathcal{J}_\Omega$  (zie 1.2) en  $|V_m|$  met  $\mathcal{J}_\omega - \mathcal{J}_\Omega$ , verder is  $p = t$  en  $q = 1$ , dus met (1.5.8) wordt  $x_m$  verwijderd als



$$\frac{(n-t)|V_m|}{\tilde{m}_{p+1,p+1}^t} < F_{(1,n-t)\alpha} \quad (4.3.1)$$

waarbij  $F_{(1,n-t)\alpha}$  de  $\alpha\%$  waarde is van een  $F_{1,n-t}$  verdeling.

Is de  $r_1, \dots, r_t$  ontstaan na verwijdering van een  $x_i$ , dan wordt uit de positieve  $V_i$  de maximale waarde bepaald stel  $V_M$  vervolgens getoetst of deze restvariantiereductie significant groter dan nul is. Nu is  $\mathcal{J}_\Omega$  niet zonder meer de reeds berekende  $\tilde{m}_{p+1,p+1}^t$  maar de niet berekende  $\tilde{m}_{p+1,p+1}^{t+M}$  die echter gelijk is aan  $\tilde{m}_{p+1,p+1}^t - V_M$ . Geldt nu: 
$$\frac{(n-t-1) \cdot V_M}{\tilde{m}_{p+1,p+1}^{-V_M}} > F_{(1,n-t-1)\alpha} \quad (4.3.2)$$

dan volgt toevoeging van  $r_m$  aan de regressie. Is dit niet het geval dan stop het proces. (4.2.1) komt overeen met 17 in het artikel van Efroymsen en (4.2.2) met 19.

Verwijdering van  $x_m$  wordt uitgevoerd door toepassing van SWP ( $r_m$ ) en opname van  $X_M$  door SWP ( $r_M$ ).

Is het proces klaar, dan kunnen de gewenste regressie-coëfficiënten afgelezen worden als  $\hat{\beta}_i = -\tilde{m}_{p+1,r_i}^t$ , waarin de elementen zijn uit de rij  $p+1$  van de submatrix  $-H_{21}$ .

De standaarddeviatie  $S_y^v$  waarin  $\tilde{y} = \tilde{y} - \sum_{i=1}^t \beta_{r_i} x_{r_i}$  (4.3.3)  
 is  $\sqrt{\tilde{m}_{p+1,p+1}^t/t}$  dit is volgens (1.5.7) een zuivere schatting van  $\sigma$ .

De standaarddeviatie van  $\hat{\beta}_i$  is volgens (1.4.12) te schatten als

$$s_{\hat{\beta}_i} = S_y^v \sqrt{\tilde{m}_{r_i,r_i}^t} \quad (4.3.4)$$

waarin  $\hat{m}_{r_i r_i}^t$  een element is van submatrix  $M_{11}^{-1}$ , die zoals reeds in 3.2. opgemerkt is, identiek aan  $S^{-1}$  in (1.4.12).

4.4. Enkele bijzonderheden in het artikel "Multiple Regression Analysis" van Efroymson.

1. Bij het onderstellen van lineaire regressie wordt veelal aangenomen, dat ook een constante term aanwezig kan zijn. Nu is dit in feite in  $y = \sum_{r=1}^p \beta_r x_r$  opgesloten; men kan namelijk zonder bezwaar onderstellen dat bijvoorbeeld voor  $r=1$  alle  $x_1 \equiv 1$  zijn, dan zal  $\beta_1$  de constante term voorstellen, waarvan de schatting  $\hat{\beta}_0$  als  ${}_{p+1,1}$  is af te lezen, terwijl ook de standaarddeviatie van  $\hat{\beta}_0$  bepaald kan worden.

Door Efroymson wordt echter die constante term eerst verwijderd door  $y$  en  $x_i$  te vervangen door  $y - \bar{y}$  en  $x_i - \bar{x}_i$ . Hierbij is  $\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$ , enz. De constante term wordt dan aan het eind van het proces als  $\hat{\beta}_0 = \bar{y} - \sum_{i=1}^p \hat{\beta}_i \bar{x}_i$  bepaald. In dit geval krijgt men echter niet de standaarddeviatie van  $\hat{\beta}_0$  ter beschikking. Efroymson definieert zijn uitgangsmatrix met de elementen

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

Interessant is het op te merken, dat de matrix  $S$  voor de op de gemiddelde gereduceerde grootheden door toepassing van de SWP-operator uit de matrix van de ongereduceerde grootheden is af te leiden.

Stel deze laatste is

$$M = \begin{pmatrix} \Sigma x_{ok}^2 & \Sigma x_{ok} x_{1k} & \dots & \Sigma x_{ok} y_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma y_k x_{ok} & \dots & \dots & \Sigma y_k^2 \end{pmatrix} = \begin{pmatrix} n & \Sigma x_1 & \dots & \Sigma y \\ \Sigma x_1 & \Sigma x_1^2 & \dots & \Sigma x_1 y \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma y & \Sigma x_1 y & \dots & \Sigma y^2 \end{pmatrix}$$

Waarin de constante term als een  $o^{de}$  term is toegevoegd in plaats van als  $x_1$ . Nu kan de Sweepoperator op het constante element worden toegepast:

$$SWP(0)M = \begin{pmatrix} -\frac{1}{n} & \bar{x}_1 & \dots & \bar{x}_p & \bar{y} \\ \bar{x}_1 & \Sigma x_1^2 - \frac{(\Sigma x_1)^2}{n} & \dots & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \bar{y} & \dots & \dots & \dots & \Sigma y^2 - \frac{(\Sigma y)^2}{n} \end{pmatrix}$$

$$= \begin{pmatrix} -\frac{1}{n} & \bar{x}_1 & \dots & \bar{x}_p & \bar{y} \\ \bar{x}_1 & \Sigma (x_1 - \bar{x}_1)^2 & \dots & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \bar{x}_p & \dots & \dots & \dots & \dots \\ \bar{y} & \dots & \dots & \dots & \Sigma (y - \bar{y})^2 \end{pmatrix}$$

Behoudens de eerste rij en kolom is dit de matrix waar Efromson van uitgaat.

2. Dat Efroymson niet (3.3.5.) als Sweep operator gebruikt, maar de vorm met tekenwisselen in de rijen is reeds behandeld. Daarnaast voert Efroymson nog een uitbreiding van de matrix in, waarvan de betekenis bij de computerberekening niet duidelijk is.

Efroymson noteert namelijk de uitgangsmatrix als  $\begin{pmatrix} S & T' \\ T & Z \end{pmatrix}$  waarin S

zoals reeds onder 1<sup>o</sup> werd opgemerkt uit de elementen  $\sum_k (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$  bestaat, terwijl verder  $Z = \sum_k (y_k - \bar{y})^2$  en T uit  $\sum_k (x_{ik} - \bar{x}_i)(y_k - \bar{y})$  bestaat. Z is dus één enkel element. T één rij van elementen en T' één kolom van elementen.

$$\text{De matrix wordt uitgebreid tot } \begin{pmatrix} S & T' & I \\ T & Z & D \\ -I & B & C \end{pmatrix} \quad (4.4.1.)$$

waarin B, C en D aanvankelijk nul zijn. I is een identiteitsmatrix met evenveel elementen als S. Er wordt nu gesteld dat toepassing van het algoritme op de elementen van S opname van elementen in de regressie betekent en toepassing op elementen van C verwijdering van  $x_i$ 's uit de regressie. Na iedere stap bevat B de regressiecoëfficiënten en C de inverse elementen van de S matrix, dat wil zeggen de  $c_{ii}$  die nodig zijn om de standaarddeviatie van de regressiecoëfficiënten te bepalen. Wat B betreft, is het inderdaad zo, dat deze de regressiecoëfficiënten bevat, maar op de plaatsen die betrekking hebben op  $x_i$ 's die verwijderd zijn door toepassing van de Sweep operator op C elementen komen geen nullen; daar blijft iets staan. Dit is als volgt in te zien. Als bovenstaande matrix gesplitst wordt in :

$$\begin{pmatrix} S_{11} & S_{12} & T'_1 & I & 0 \\ S_{21} & S_{22} & T'_2 & 0 & I \\ T_1 & T_2 & Z & 0 & 0 \\ -I & 0 & 0 & 0 & 0 \\ 0 & -I & 0 & 0 & 0 \end{pmatrix}$$

en hierop wordt SWP\*(0, 1, ..., s) toegepast, ontstaat:  
rijen

$$\begin{array}{l}
 \left. \begin{array}{l} 1 \\ \vdots \\ s \end{array} \right\} \left( \begin{array}{cccccc}
 S_{11}^{-1} & S_{11}^{-1}S_{12} & S_{11}^{-1}T_1' & S_{11}^{-1} & 0 \\
 -S_{21}S_{11}^{-1} & S_{22}-S_{21}S_{11}^{-1}S_{12} & T_2-S_{21}S_{11}^{-1}T_1' & -S_{22}S_{11}^{-1} & 0 \\
 -T_1S_{11}^{-1} & T_2-T_1S_{11}^{-1}S_{12} & Z-T_1S_{11}^{-1}T_1' & -T_1S_{11}^{-1} & 0 \\
 S_{11}^{-1} & S_{11}^{-1}S_{12} & S_{11}^{-1}T_1' & S_{11}^{-1} & 0 \\
 0 & -I & 0 & 0 & 0
 \end{array} \right) \\
 \left. \begin{array}{l} s+1 \\ \vdots \\ p \end{array} \right\} \\
 p+1 \\
 \left. \begin{array}{l} p+2 \\ \vdots \\ p+1+s \end{array} \right\} \\
 \left. \begin{array}{l} p+2+s \\ \vdots \\ 2p+1 \end{array} \right\}
 \end{array}$$

De regressiecoëfficiënten zijn te vinden in de eerste groep elementen van rij p+1 namelijk  $-T_1S_{11}^{-1}$ ; maar met negatief teken. Ze zijn ook te vinden zowel in de rijen 1, ..., s als in p+r, ..., p+1+s, namelijk als  $S_{11}^{-1}T_1'$  inderdaad in de B submatrix van (4.4.1). Wordt de Sweepoperator nu op de elementen (p+1, p+2)...(p+1+s, p+1+s) toegepast, dan ontstaat:

$$\left( \begin{array}{ccccc}
 0 & 0 & 0 & I & 0 \\
 0 & S_{21} & T_2' & -S_{21} & I \\
 0 & T_2 & Z & -T_1 & 0 \\
 -I & -S_{12} & -T_1' & S_{11} & 0 \\
 0 & -I & 0 & 0 & 0
 \end{array} \right)$$

B is dus niet weer 0 geworden maar  $\begin{pmatrix} -T_1' \\ 0 \end{pmatrix}$  terwijl C nu de oorspronkelijke  $S_{11}$  bevat.

Men vraagt zich echter wel af waarom Efroymson deze uitbreiding invoert, want in het computerprogramma wordt dit niet toegepast; de matrix blijft beperkt tot  $\begin{pmatrix} S & T' \\ T & Z \end{pmatrix}$ .

3.

Tenslotte nog een laatste verandering ten opzichte van de algemene theorie. Er worden in het computerprogramma niet de oorspronkelijke gedefinieerde elementen  $q_{ij} = \Sigma(x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$  enz. ingevoerd maar de correlatiecoëfficiënten

$$r_{ij} = \frac{\Sigma(x_{ik} - \bar{x}_i)(x_{jt} - \bar{x}_j)}{\sqrt{\Sigma(x_{ik} - \bar{x}_i)^2 \cdot \Sigma(x_{jk} - \bar{x}_j)^2}} = \frac{q_{ij}}{S_i \cdot S_j}$$

hierin is ook  $y$  begrepen door  $y$  als  $x_{p+1}$  te schrijven en  $i$  van  $1$  t/m  $p+1$  te laten lopen. De matrix

$\begin{bmatrix} S & T' \\ T & Z \end{bmatrix}$  bestaat dus nu uit de elementenrij voor  $i, j = 1, \dots, p+1$ , en zal door  $\begin{bmatrix} \tilde{S} & \tilde{T}' \\ \tilde{T} & \tilde{Z} \end{bmatrix}$  worden aangegeven. Toepassing van SWP\*(0, 1, ..., s)

op beide heeft tot resultaat dat resp.:

$$\begin{bmatrix} S^{-1} & S^{-1}T' \\ -TS^{-1} & Z-TS^{-1}T' \end{bmatrix} \text{ en } \begin{bmatrix} \tilde{S}^{-1} & \tilde{S}^{-1}\tilde{T}' \\ -\tilde{T}\tilde{S}^{-1} & \tilde{Z}-\tilde{T}\tilde{S}^{-1}\tilde{T}' \end{bmatrix}$$

ontstaat. De elementen van  $S$  zijn  $q_{ij}$  voor  $i, j = 1, \dots, p$ , die van

$\tilde{S}$  zijn  $\frac{q_{ij}}{s_i s_j}$ ; van  $T$  en  $\tilde{T}$  resp.  $q_{p+1, j}$  en  $\frac{q_{p+1, j}}{s_{p+1} s_j}$  en van  $Z$  en  $\tilde{Z}$

resp.  $q_{p+1, p+1}$  en  $\frac{q_{p+1, p+1}}{s_{p+1}^2}$ . Als de elementen van  $S^{-1}$  met  $\tilde{q}_{ij}$

worden aangegeven, dan is gemakkelijk in te zien dat die van  $\tilde{S}^{-1}$  gelijk zijn aan  $q_{ij}^* = s_i s_j \tilde{q}_{ij}$ . Nu zijn de elementen van  $TS^{-1}$  de regressiecoëfficiënten voor het oorspronkelijke regressiemodel, dus  $\hat{\beta}_i$ ; de elementen van  $\tilde{T}\tilde{S}^{-1}$  die van het op standaarddeviaties genormeerde model, die met  $b_i$  zullen worden aangegeven. Dan is:

$$\hat{\beta}_i = \sum_{k=1}^p q_{p+1, k} \tilde{q}_{k, j} \text{ en derhalve } b_i = \sum_{k=1}^p \frac{q_{p+1, k}}{s_{p+1} s_k} \cdot s_k s_j \tilde{q}_{k, j} = \frac{s_j}{s_{p+1}} \cdot \hat{\beta}$$

$$\text{of } \hat{\beta}_i = \frac{s_{p+1}}{s_j} b_i.$$

Dit is in het artikel van Efroymsen, formule (15), genoteerd als  $\hat{\beta}_i = \frac{\sigma_n}{\sigma_i} b_{in}$ ;  $s_{p+1}$  is namelijk de standaarddeviatie van  $y$ , wat door Efroymsen met  $\sigma_n$  wordt aangegeven.

De standaarddeviatie van de volgens de regressie gereduceerde  $y$  (in feite van  $\hat{\epsilon}$ ) is volgens (4.3.3) te schrijven als

$$s_y^{\sim} = \sqrt{\frac{q_{p+1,p+1}^*}{b}} \quad \text{waarin } q_{p+1,p+1}^* = Z-TS^{-1}T'; \text{ als}$$

$$Z-TS^{-1}T' = \tilde{q}_{p+1,p+1} \text{ wordt genoemd, geldt } \tilde{q}_{p+1,p+1}^* = \frac{q_{p+1,p+1}^*}{s_{p+1}^2}$$

derhalve  $s_y^{\sim} = s_{p+1} \sqrt{\frac{\tilde{q}_{p+1,p+1}^*}{b}}$ .

Dit is formule 14 van Efroymsen:  $s_y = \sigma_n \sqrt{r_{nn}/\phi}$ .

Tenslotte is  $s_{\hat{\beta}_i} = s_y \sqrt{\tilde{q}_{ii}}$  volgens (4.3.4) waarin  $\tilde{q}_{ii}$  een element van  $S^{-1}$  is. Nu geldt voor de elementen  $q_{ii}^*$  en  $\tilde{S}^{-1}$  dat  $q_{ii}^* = S_i^{2\sim} \tilde{q}_{ii}$  derhalve  $S_{\hat{\beta}_i} = \frac{s_y}{s_i} \sqrt{q_{ii}^*}$ . Dit is de formule (16) van Efroymsen:  $\frac{s_y}{\sigma_i} \sqrt{c_{ii}}$ .

## LITERATUUR

- H. Scheffé : The Analysis of Variance  
Wiley - New York 1959
- A.P. Dempster : Elements of continuous multivariable analysis  
Addison - Wesley Publ. 1969
- M.A. Efroymsen: Multiple regression analysis  
Nr. 17 in Ralston and Wilf (e.d.):  
Mathematical Methods for Digital Computers N.Y. 1962 (1967)
- A. Orden : Matrix inversion and related topics by direct methods  
Nr. 2 in Ralston and Wilf (e.d.):  
Mathematical Methods for Digital Computers N.Y. 1962 (1969)



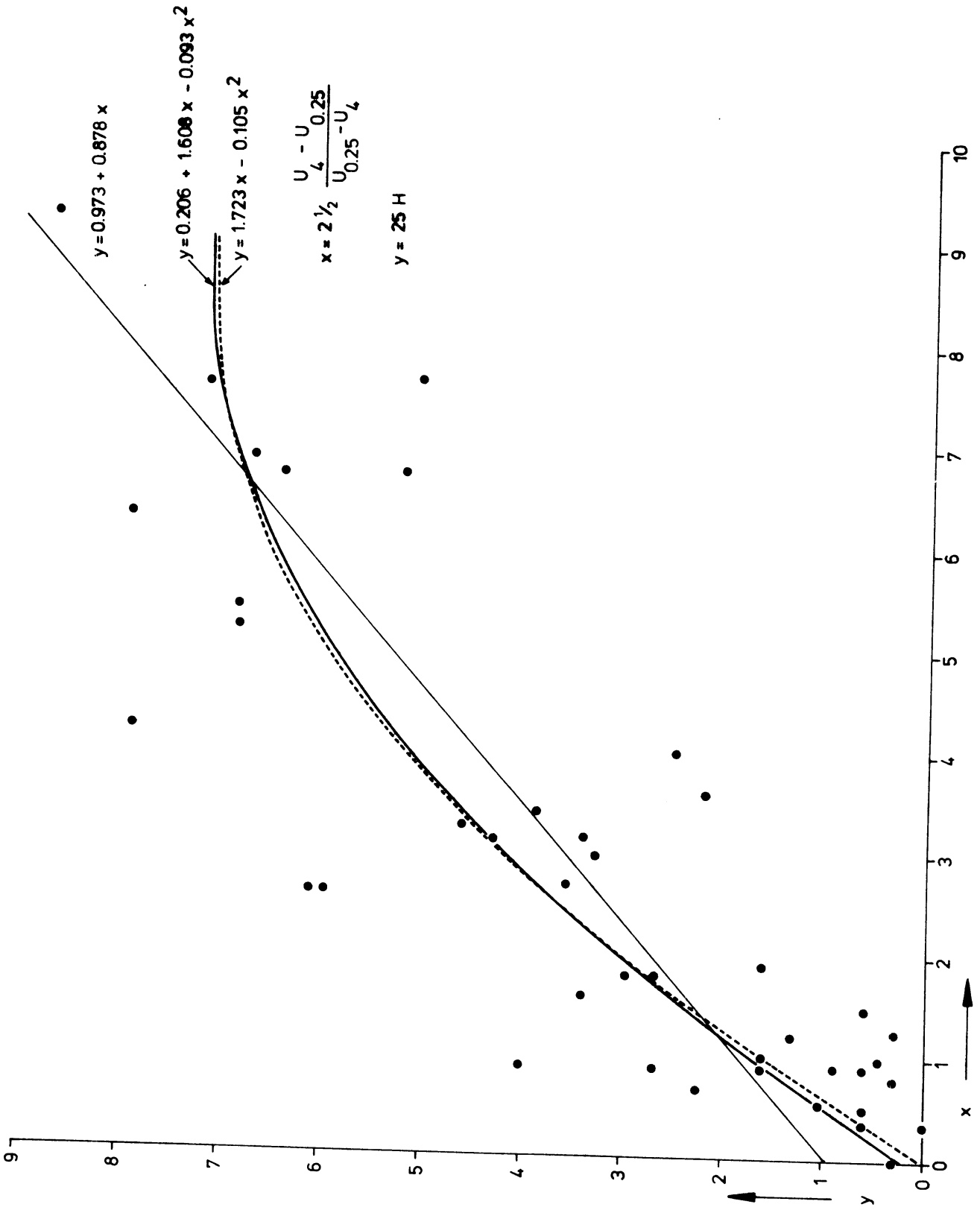


Fig.1 Relatie tussen warmtestroom en stabiliteitsparameter