

October 14-18, 1985

Symposium on  
Nonlinear Multivariate Data Analysis

Jan de Leeuw  
Department of Data Theory  
University of Leiden  
Leiden, The Netherlands

Presented at  
The University of North Carolina at Chapel Hill

October 14-18, 1985

Symposium on  
Nonlinear Multivariate Data Analysis

Jan de Leeuw  
Department of Data Theory  
University of Leiden  
Leiden, The Netherlands

Presented at  
The University of North Carolina at Chapel Hill

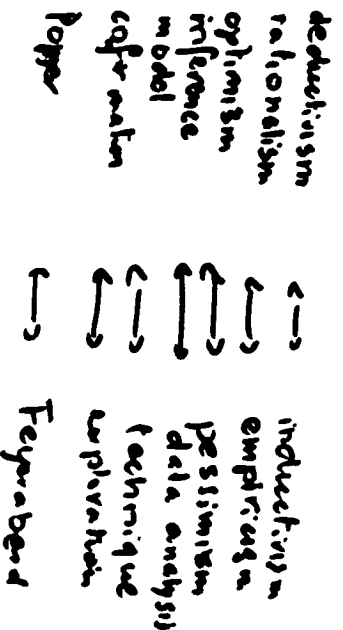
Monday, Oct. 14	
- Data analysis and statistics	1
- Approaches to multivariate analysis	7
Tuesday, Oct. 15	
- Coding and quantification	14
- Homogeneity analysis	20
Wednesday, Oct. 16	
- Correspondence analysis	53
- Nonlinear principal component analysis	83
Thursday, Oct. 17	
- Measurement levels and scaling criteria	100
- Sets of variables	126
- Nonlinear regression, canonical, and discriminant analysis	133
Friday, Oct. 17	
- General nonlinear multivariate analysis algorithms	163
- Statistical aspects of nonlinear multivariate analysis techniques	168
- Causal models, nonlinear path analysis	186
Annotated References	194
Illustrative Data Sets	197



②

① Gill's principles

There is no clear cut distinction between exploratory and confirmatory



2 In most of the situations we are familiar with the standard statistical prejudice (Kelman) does not make sense.

3 We must expect the many-many relationship between words and techniques.

(c) Robust and complex - inferential statistics

- replace models by super models
- replace assumptions by computations
- create your own strength - envision new
- ~> Robust statistics paper.
- ~> Faldoutie, Bootstrap, Subsampling.
- ~> Permutation and Randomization tests.

③ tools for data analysis

- 1 Gauging
- 2 Analysis of stability

ad1 A data set with known properties is a gauge. We want to find out how a particular technique represents these known properties.

- a Probabilistic gauges [multinomial]
- b Statistical gauges [samples from multinomial]
- c Monte Carlo gauges [idem, without formulas]
- d Algebraic gauges [Simplex, Guttman Scale, ...]
- e Empirical gauges [data with well-known structure]

④

ad2 A small and/or unimportant change in data, model, or technique should lead to a small and unimportant change in the results.

This is stability.

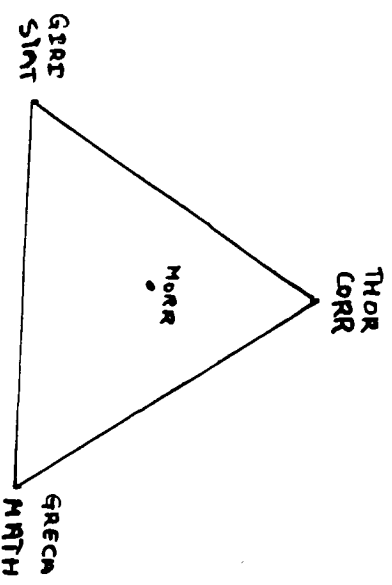
- a Replication stability
- b Statistical stability [replication without replication]
- c Stability under data selection [resampling]
- d Stability under model selection
- e Numerical stability
- f Analytical stability [differentiability]
- g Algebraic stability [... with algebraic tests]
- h Stability under selection of technique

Both tools are technique-oriented. But otherwise all of classical statistics fits in painlessly. Gauging is used to study the correspondence models  $\rightarrow$  techniques

## II) Multivariate analysis

- inventory
- problems
- prospects
- definitions

(a) Content analysis of 20 MVA Books  
[graphs & tables]



3

b) Current approaches

1 The multivariate normal [Muirhead, Eaton]

- usually a good description (?)
- The central limit effect (?)
- simple formulas (?)

- regressions are linear
  - independence  $\equiv$  orthogonality
  - equiprobability  $\equiv$  Euclidean distance
  - closed under linear combinations
- } (1)

- only one problem: usually a poor description [both the normality, the independence, and the identical distributions]

8

2 The log linear model

- more realistic

but

- poor at data description (interpretation)

- only usable if  $m \leq 4$  [empty cells]

- only usable if  $n$  very large [as least for the statistical part].

- model selection problem

3 Tabular analysis

- close to the data

but

- very heavy output

- problems with integration, reporting,

table selection, interpretation

4 The packages

RCA  
GUM  
etc ] by

SAS, SPSS, ....

convenient

but

limited to numerical variables.

So

We need techniques for data sets where

are

{ = mixed level  
= large  $m$   
= non-random  
= analyzed joint-bivariately



## Definition of MVA

Content analysis of books shows confusion

Is MDS a form of MVA?

Must we have a random sample in order to use MVA?

We answer both questions with a firm NO.

Multivariate analysis studies the structure of multivariables

$\phi_j : \Omega \rightarrow \mathcal{V}_j$  is a variable with domain  $\Omega$  and range  $\mathcal{V}_j$ .

We can have

$$\mathcal{V}_j = \mathbb{R} = (-\infty, +\infty)$$


$$\mathcal{V}_j = \{0, 1, 2, \dots\} = \mathbb{N}$$

$\mathcal{V}_j = \{ \text{protestant, catholic, buddhist} \}$   
 $\mathcal{V}_j = \{ \text{agree, don't know, disagree} \}$   
etc

We can have

$\Omega$  a finite set  $\{w_1, \dots, w_n\}$

$\Omega$  a probability space, and  $\phi_j$  measurable

[i.e.  by a random variable]

$(\phi_1, \dots, \phi_m)$  is a multivariable with domain  $\Omega$  and range  $\mathcal{V}_1 \times \dots \times \mathcal{V}_m$ .

Observe that rows and columns enter asymmetrically [conditionality]



TUESDAY OCT 15 13<sup>00</sup>

Yesterday we have seen that MVA techniques are needed which can be used

- on large data sets (large  $n$ )
- on mixed measurement level data
- on nonrandom data

Today we shall start building a system of

MVA techniques that

- has the existing linear techniques as special cases
- containing the necessary extensions

The system is built around the notion of homogeneity and around a particular least squares loss function measuring departure from homogeneity

The techniques will be discussed from the geometrical [MOS], algebraic [MVA], and statistical point of view.

### Quantification and coding

Suppose that  $\Omega = \{\omega_1, \dots, \omega_m\}$  is finite and suppose  $\mathcal{D}_j = \{1, \dots, k_j\}$  for all  $j=1, \dots, m$ .

### Quantification

A quantification of the objects (individuals, ...) is a mapping of  $\Omega$  into  $\mathbb{R}^p$ . The number  $p$  is the dimensionality of the quantification.

A quantification of (the categories of a) variable is a mapping of  $\mathcal{D}_j$  into  $\mathbb{R}^p$ . In some contexts we do not use quantification, but we use transformation.

Perfect quantification

A quantification of the objects  $X$ , with  $x_1, \dots, x_n \in \mathbb{R}^n$ , is perfect if, for all  $s, j, k$

$\phi_j(\omega_i) = \phi_j(\omega_k) \implies x_i = x_k$

i.e. if objects within are in the same category of a variable have the same score

Quantifications  $Y_1, \dots, Y_m$  of the variables, with  $Y_j$  matrices of order  $k \times p$ , are perfect if

$\phi_j(\omega_i) = a \wedge \phi_l(\omega_i) = b \implies (Y_j)_a = (Y_l)_b$

i.e. if the quantifications corresponding with the values of each object are the same.

1	a	p	u
2	b	q	v
3	a	r	v
4	a	p	u
5	b	p	v
6	c	p	v
7	a	p	u
8	a	p	v
9	c	p	v
10	a	p	v

perfect object series

from column 1

$\begin{cases} x_1 = x_8 = x_4 = x_7 = x_8 = x_9 \\ x_2 = x_5 \\ x_6 = x_3 \end{cases}$

from column 2

$x_1 = x_4 = x_5 = x_6 = x_7 = x_8 = x_9 = x_{10}$

etc

perfect category quantification

from row 1

$y_a = y_p = y_u$

from row 2

$y_b = y_q = y_v$

etc

Matrix notation

We need the concept of an indicator matrix

For variable  $j$  this is defined as the

$n \times k$  matrix  $G_j$  with

$$(G_j)_{ij} = \begin{cases} 1 & \text{if } \phi_j(w_i) = \ell \\ 0 & \text{otherwise} \end{cases}$$

	abc	pqr	uvw
1	100	100	100
2	010	010	010
3	100	001	010
4	100	100	100
5	010	100	010
6	001	100	000
7	100	100	000
8	100	100	000
9	001	100	010
10	100	100	010
	622	811	370

$D_j = G_j' G_j$   
is diagonal with  
univariate marginals

$C_j = G_j' G_j \ell$  contains  
bivariate marginals (i.e.  
cross tables).

This is a complete indicator matrix because

$$G_j u = u \quad \text{for all } j.$$

Table 2.1

a	p	u
b	q	v
a	r	v
a	p	v
b	p	v
c	p	v
a	p	v
c	p	v
a	p	v

Table 2.1.

Example of data matrix  $W$ .

Table 2.3

a	p	u	3
a	p	v	2
a	r	v	1
b	p	v	1
b	q	v	1
c	p	v	2

Table 2.3.

Reduced profile frequency matrix.

Table 2.2

a	p	u	3
a	p	v	2
a	p	w	0
a	q	w	0
a	q	v	0
a	r	v	1
a	r	w	0
a	r	w	0
b	p	v	1
b	p	w	0
b	q	v	1
b	q	w	0
b	r	w	0
b	r	v	0
b	r	w	0
c	p	v	2
c	p	w	0
c	q	w	0
c	q	v	0
c	r	w	0
c	r	v	0
c	r	w	0

Table 2.2.

Profile frequency matrix.

Table 2.4

	p	q	r		p	q	r		p	q	r
a	3	0	0		3	0	1		3	0	0
b	0	0	0		0	1	1		0	0	0
c	0	0	0		0	2	0		0	0	0
	0				v				v		

Table 2.4. Higher dimensional cross tabulation.

Table 2.5

	a b c	p q r	u v w
100	100	100	100
010	010	010	010
100	001	010	010
100	100	100	100
010	100	010	010
001	100	010	010
100	100	100	100
100	100	010	010
001	100	010	010
100	100	010	010

Table 2.5 Indicator matrix G for data matrix H of table 2.1

Table 2.6

	a b c	p q r	u v w
a	600	501	330
b	020	110	020
c	002	200	020
p	512	800	350
q	010	010	010
r	100	001	010
u	300	300	300
v	322	511	070
w	000	000	000

Table 2.6 Matrix C of bivariate marginals.

Table 2.7

	a b c	p q r	u v w
a	600	000	000
b	020	000	000
c	002	000	000
p	000	800	000
q	000	010	000
r	000	001	000
u	000	000	300
v	000	000	070
w	000	000	000

Table 2.7 Matrix D of univariate marginals.

Table 2.8

	a	b	c	d	e	f
old	1	0	0	0	0	0
	1	1	0	0	0	0
	0	1	1	1	0	0
	0	0	1	1	1	0
	0	0	0	1	1	1
new	0	0	0	1	1	1

Table 2.8A. Incomplete indicator matrix.

	a	b	c	d	e	f
old	10	01	01	01	01	01
	10	10	01	01	01	01
	01	10	10	10	01	01
	01	01	10	10	10	01
new	01	01	01	10	10	10

Table 2.8B. Completed indicator matrix

Perfect quantification in matrix notation

$(Y_1, \dots, Y_m)$  is perfect iff  $\exists X \Rightarrow G_1 Y_1 = \dots = G_m Y_m = X$

$(X)$  is perfect iff  $\exists Y_1, \dots, Y_m \ni G_1 Y_1 = \dots = G_m Y_m = X$

In summary  $(X; Y_1, \dots, Y_m)$  is perfect iff  $X = G_1 Y_1 = \dots = G_m Y_m$ .

Theorem

$(Y_1, \dots, Y_m)$  is perfect ff

$G_1 Y_1 = \dots = G_m Y_m = \frac{1}{m} \sum_{j=1}^m G_j Y_j$  iff

$\sum_{j=1}^m Y_j' G_j Y_j = m \sum_{j=1}^m Y_j' D_j Y_j$

between variables SSC

total SSC.

Theorem X is perfect iff

$$X = P_1 X = P_2 X = \dots = P_m X$$

iff

$$X'X = X'P_1 X \quad (\text{with } P_j = G_j(G_j G_j')^{-1} G_j')$$

total ssq  $\xrightarrow{P}$  between object ssq

Theorem X is perfect iff the p largest

eigenvalues of  $P_n$  are equal to +1.

$$Y = (Y_1, \dots, Y_m) \text{ is perfect iff } h \text{ is } p \text{ largest}$$

general eigenvalues of  $(C, M'D)$  are +1.

Perfection cannot be attained for real data sets. Thus we look for optimality. This defines homogeneity analysis.

less further measuring departure from perfect homogeneity (of given scores X and quantitative Y).

$$\mathcal{L}(X; Y) = \frac{1}{m} \sum_{j=1}^m t_j [X - G_j Y]' [X - G_j Y]$$

In homogeneity analysis (also known as

- multiple correspondence analysis [BENZÉCRIS]
- GUTMAN'S PROTRAM COMPONENT OF SCALE ANALYSIS [G 1941, L 1960]
- FACTORIAL ANALYSIS OF QUALITATIVE DATA [BURT, 1950]
- HARTIG'S n<sup>th</sup> METHOD OF QUANTIFICATION [H, 1956]

we want to find X and Y such that  $\mathcal{L}(X; Y)$  is minimized.

BUTwe do not want the trivial solution  $X=0 \wedge Y=0$ .NORdo we want the trivial dimension  $X = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad Y = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ 

Thus we must impose some normalization restriction.

$$N_1 : \begin{cases} W^T Y = 0 \\ X^T X = I. \end{cases}$$

$$N_2 : \begin{cases} W^T D Y = 0 \\ Y^T D Y = I \end{cases}$$

We only need to normalize either  $X$  or  $Y$ .Let us see what happens if we normalize  $X$ .

$$Z(X; \alpha) \triangleq \min \{ \phi(X; Y) \mid Y \text{ free} \}$$

minimum attained for  $Y = D^{-1} G^T X$  [or  $G^T X$ ]

$$Z(X; \alpha) = \text{tr } X^T [I - P_\alpha] X$$

now define

$$Z(\alpha; w) \triangleq \min \{ Z(X; \alpha) \mid X \text{ norm} \}$$

then

$$Z(\alpha; w) = P - \sum_{s=1}^P \lambda_s (P_\alpha)$$

$$\boxed{\lambda_s (P_\alpha) = 1}$$

Interpretation:  $X$  are eigenvectors of the average projector  $P_\alpha$ , and  $\lambda_j$  is the centroid of the corresponding  $X$ .

Now proceed dually

$Z(x, Y) \triangleq \min \{ G(x, Y) \mid x \text{ free} \}$   
 attained for  $x = \frac{1}{m} \sum G_j Y_j$

$L(x, Y) = \frac{1}{m} \sum Y_j' [D - \frac{1}{m} C] Y$

then [Theorem, not definition]

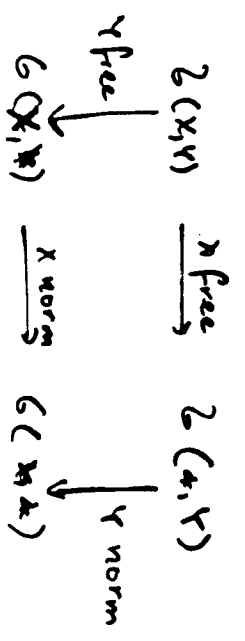
$Z(x, Y) = \min \{ G(x, Y) \mid Y \text{ norm} \} -$

$= p - \sum_{s=1}^p \lambda_s [ \frac{1}{m} D^{-1/2} C D^{-1/2} ]$

Proof. let  $H = m^{-1/2} D^{-1/2} G'$ . Then

$P_k = HH'$  and  $m^{-1} D^{-1/2} C D^{-1/2} = HH'$ .

The <sup>nonzero</sup> eigenvalues of  $HH'$  and  $HH$  are the same.  $\square$



- Thus nonnegativity constraints, computationally, is finding free p largest eigenvalues of either
- the Burt table C [corrected for marginals]
  - the average projector  $P_k$

Pictures of loss

[Example from Beyond]  
 [Also illustrating the MDS - algorithm]



	Price	Gas	Weight
Chevette	5.6	6.9	9.7
Dodge Colt	5.7	5.1	8.8
Plymouth Horizon	6.3	5.5	9.9
Fort Mustang	7.6	6.7	12.0
Pontiac Phoenix	8.6	6.9	12.1
Dodge Diplomat	9.4	10.2	13.5
Chevrolet Impala	10.1	7.5	16.9
Buick Regal	10.5	7.8	15.0
AMC Eagle	10.7	11.7	15.7
Oldsobile 98	13.3	8.7	18.3

table 1: Car data.

Chevette	1	1	1
Dodge Colt	1	1	1
Plymouth Horizon	1	1	1
Fort Mustang	2	1	2
Pontiac Phoenix	2	1	2
Dodge Diplomat	2	3	2
Chevrolet Impala	3	2	3
Buick Regal	3	2	2
AMC Eagle	3	3	2
Oldsobile 98	4	2	3

table 2: Car data, discrete.

Chevette	1	4	2
Dodge Colt	2	1	3
Plymouth Horizon	3	2	4
Fort Mustang	4	3	5
Pontiac Phoenix	5	3	7
Dodge Diplomat	6	9	9
Chevrolet Impala	7	6	6
Buick Regal	8	7	8
AMC Eagle	9	10	10
Oldsobile 98	10	8	10

table 3: Car data, ranked.

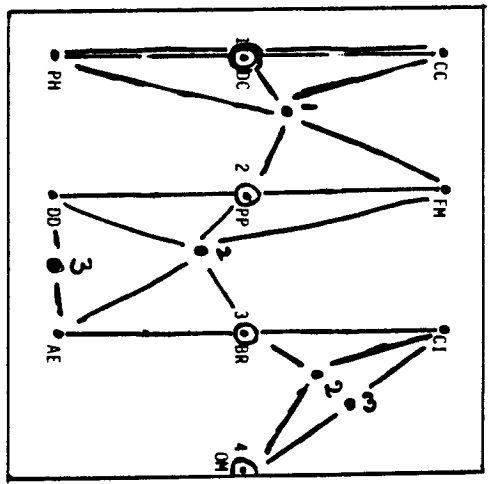


figure 1: loss variable 1, arbitrary solution

red  
variable 2  
variable 3

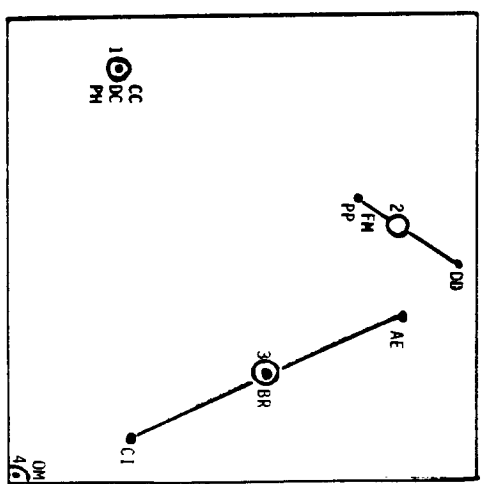
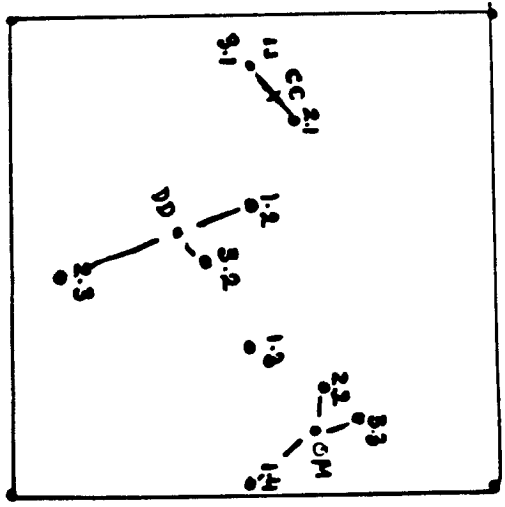


figure 2: loss variable 1, optimal solution



- it shrinks
- it is simple

Alg A  $X_0$  such that  $WX_0=0, X_0'X_0=I$

- Step 1:  $Y_j \leftarrow D_j^{-1} G_j' X$
- Step 2:  $X \leftarrow \frac{1}{n} \sum_j G_j Y_j$
- Step 3:  $X \leftarrow \text{GRAM}(X)$

Alg B:  $Y$  such that  $W'Y=0, Y_0'DY_0=I$ .

- Step 1:  $X \leftarrow \frac{1}{n} \sum_j G_j Y_j$
- Step 2:  $Y_j \leftarrow D_j^{-1} G_j' X$
- Step 3:  $Y \leftarrow \text{WGRAM}(Y)$

The algorithms both work by applying

- reciprocal averaging [Horn, 1935]
- le principe bodycentrique [Rehder]
- [the centroid principle]

Unfortunately one cannot have both  $Y$  as centroids of  $X$  and  $X$  of centroids of  $Y$ . One has to choose.

HOPMIS uses Alg A. Thus  $X$  is normalized, and category quantifications are centroids of scores of the objects in the categories.

Variables	Categories & codes
Thread	Yes - Y No - N. Flat - F Cup - U Cone - O Round - R Cylinder - Y.
Head	None - N Star - T Silt - L. Sharp - S Flat - F. (In half inches).
Head Indentation	Sharp - S Flat - F. (In half inches).
Bottom	Yes - Y No - N.
Length	
Brass	

Object	1	2	3	4	5	6
TACK	N	F	M	S	1	N
NAIL1	N	F	M	S	4	N
NAIL2	N	F	M	S	2	N
NAIL3	N	F	M	S	2	N
NAIL4	N	F	M	S	2	N
NAIL5	N	F	M	S	2	N
NAIL6	N	U	M	S	5	N
NAIL7	N	U	M	S	3	N
NAIL8	N	U	M	S	3	N
SCREW1	Y	O	T	S	5	N
SCREW2	Y	R	L	S	4	N
SCREW3	Y	R	L	S	4	N
SCREW4	Y	R	L	S	2	N
SCREW5	Y	R	L	S	2	N
BOLT1	Y	R	L	F	4	N
BOLT2	Y	O	L	F	1	M
BOLT3	Y	Y	L	F	1	M
BOLT4	Y	Y	L	F	1	M
BOLT5	Y	Y	L	F	1	M
BOLT6	Y	Y	L	F	1	M
TACK1	N	F	M	S	1	Y
TACK2	N	F	M	S	1	Y
NAILB	M	F	M	S	1	Y
SCREMB	Y	O	L	S	1	Y

Table 3.13 Hartigan's hardware

object	dim1	dim2
TACK	0.75	0.46
NAIL1	0.68	0.47
NAIL2	0.96	0.52
NAIL3	0.96	0.52
NAIL4	0.96	0.52
NAIL5	0.96	0.52
NAIL6	1.00	-1.69
NAIL7	1.25	-0.74
NAIL8	1.25	-0.74
SCREW1	-0.38	-3.96
SCREW2	-0.85	0.23
SCREW3	-0.91	0.26
SCREW4	-0.57	0.28
SCREW5	-0.63	0.31
BOLT1	-1.31	0.38
BOLT2	-1.10	-0.51
BOLT3	-1.30	0.40
BOLT4	-1.30	0.40
BOLT5	-1.30	0.40
BOLT6	-1.30	0.40
TACK1	0.93	0.67
TACK2	0.93	0.67
NAILB	0.93	0.67
SCREMB	-0.54	-0.44

Table 3.14A Hartigan's hardware object scores

category	dim1	dim2
1Y	-0.96	-0.15
1N	0.96	0.15
2U	0.90	0.52
2O	1.16	-1.06
2R	-0.70	-1.64
2Y	-0.91	0.30
3N	-1.12	0.36
3T	0.96	0.15
3L	-0.38	-3.96
4S	-1.02	0.19
4F	0.43	-0.08
5/1	-1.28	0.25
5/2	-0.34	0.31
5/3	0.44	0.44
5/4	1.25	-0.74
5/5	-0.60	0.24
6Y	0.31	-2.82
6T	0.36	0.39
6N	-0.11	-0.00

Table 3.14B Hartigan's hardware category quantifications

Variables	dim1	dim2
Thread	0.930	0.024
Head	0.951	0.635
Head Indentation	0.945	0.681
Bottom	0.546	0.020
Length	0.292	0.819
Brass	0.064	0.031
Eigenvalues	0.621	0.368

Table 3.14C Hartigan's hardware discrimination measures

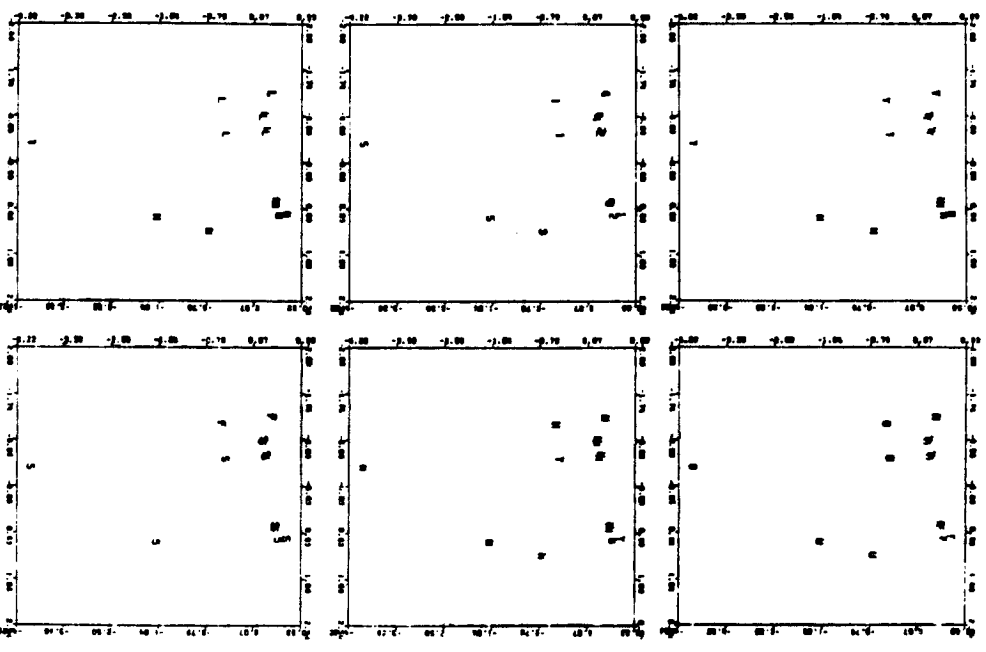


Figure 3.13. Hartigan's hardware object scores labeled by variables

dimensions. This gives the dis-  
 tances we get the total discrim-

(40)

best-of-fit of a HPMALS-repre-  
 sentation:

(41)

ation of the data. Classification of object points in figure 3.13 is an undifferentiated set of objects in 11 clusters and we need to identify clusters or regions in relation according to the criteria of high discrimination. By using an outside criterion, we can compare the old fit is used in our example which

ies, meaning the measurements description we refer to Vesica is a fit of 508, which is not better than the fit of 508. The map numbers given in the first two clusters: (5,6,7,8,9) (1,2) are nearer to the fish and 3-dimensional representation (6) however not far apart. It of 442 which is lower than the fit is the mean of the dimensions increases. The first dimension is similar to the two-dimensional representation of dolphins and sperm whales (8,4) from the HPMALS-representation is given in figure 3.13. The HPMALS-representation of the points on the map of figure 3.13, which is a balaen, grey and finback whales, the dolphins and porpoises (5,6) white whales (7). These relative solution. The sperm whale (4) HPMALS classification.

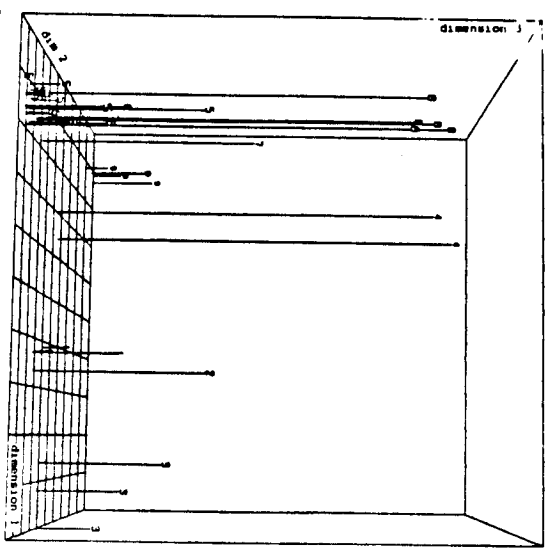


Figure 5  
 Three-dimensional HPMALS-representation. Objects are labeled by their group numbers.

Whales - example

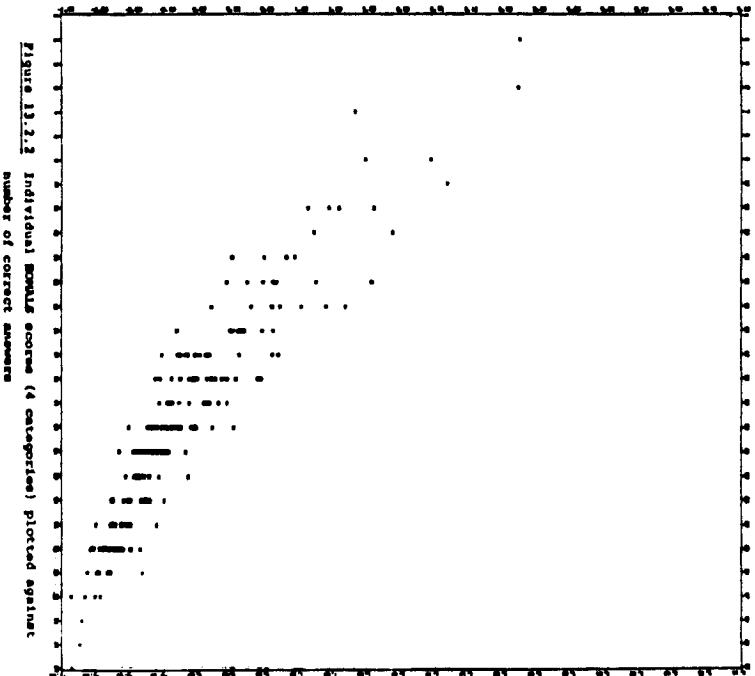
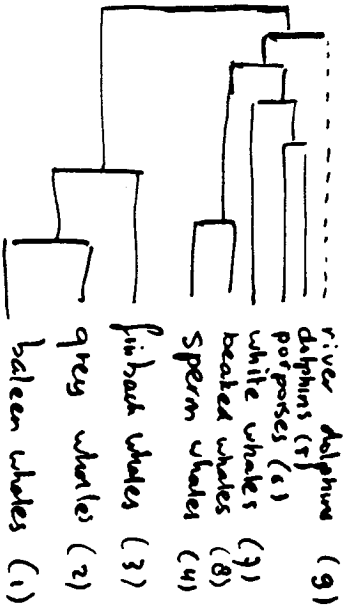
E. van der Burg

Hornik's classification of whales, porpoises, and dolphins

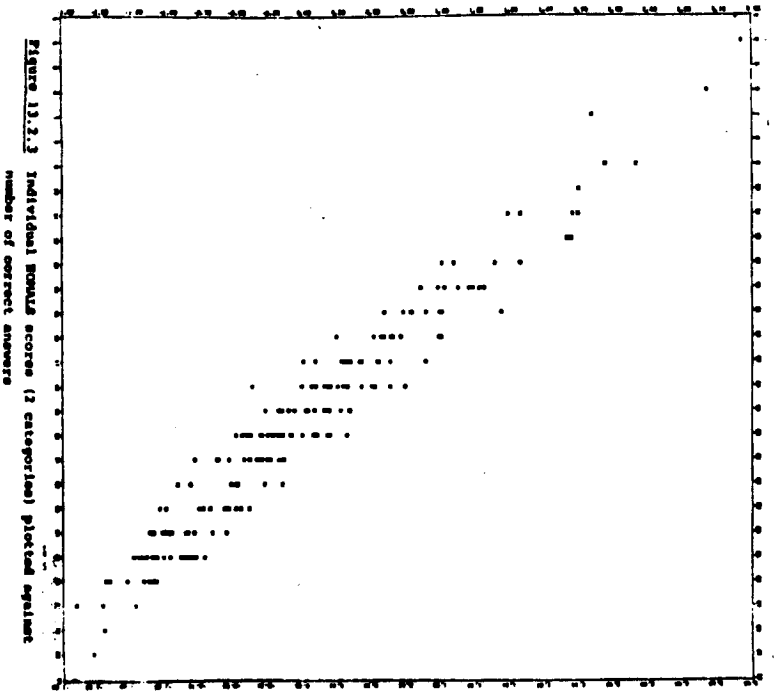
In J.F. Marcotorchino, J.H. Park, and J. Janssen (eds)

Data analysis in Real Life Environments

Elsevier Science Publishers, 1987



25 31



36

Our discussion of homogeneity analysis so far has been largely geometrical. What is the purpose of the technique? If  $p > 1$   $n$  is drawing a picture in which objects with similar characteristics pattern are close together, and categories containing the same objects are close together etc. If  $p=1$  the "picture" is really a "scale". How do you choose  $p$ ? That depends on the type of picture you want. Moreover the dimensions are "nested". Are we "estimating" something? Not necessarily. Neither are we testing anything.

But

There are connections with the chi-square test for independence.

variance of eigenvalues

$$\frac{1}{2} n \sum_i [m \lambda_i - 1]^2 = \sum_{j < e} \sum_{j \neq e} \chi^2_{je}$$

Thus, if all variables are independent, we have a random sample.

$$\frac{1}{2} n \sum_j [m \lambda_j - 1]^2 \rightsquigarrow \chi^2_{\sum_{j \neq e} (k_j - 1)(k_e - 1)}$$

Moreover, if variable  $j$  is independent of the others,

then

$$n \left( \frac{\sum_i y_i}{1 - \lambda_j} \right)^2 \xrightarrow{D} \chi^2_{k_j - 1}$$

discrimination measure

A less geometrical approach [  $p=1$ , already used in MC-example ].

We have seen that homogeneity analysis (if  $p=1$ ) can be formulated as

$$y^i C y \quad \max! \quad y^i D y = m \cdot \quad u^i D y = 0$$

Here  $y^i = (y_i^1, \dots, 1, y_i^m)$ . Now decompose

$$y_i = \alpha_j z_j \quad \text{with} \quad z_j^i D_j z_j = 1 \quad \text{and} \quad y_j^i D_j z_j = 0$$

Then

$$y^i D y = \sum_j y_j^i D_j y_j = \sum_j \alpha_j^2 z_j^i D_j z_j = \sum \alpha_j^2$$

and

$$y^i C y = \sum_j \sum_e \alpha_j \alpha_e z_j^i C_{je} z_e = \sum_j \sum_e \alpha_j \alpha_e z_j^i z_e^i$$

with  $R = \{ \alpha_j \}$  the correlation matrix induced by the scores  $y_j$  (or  $z_j$ ).



The homogeneity analysis problem can now be

written as

$$\max_{\alpha} \max_z \left\{ \alpha' R \alpha \mid \alpha' \alpha = 1 \quad z' D z = 1 \right\}$$

or

$$\max_z \left\{ \lambda_{\max}(R) \mid z' D z = 1 \right\}$$

In words: we want to find quantifications of the categories that maximize the largest eigenvalue of the induced correlation matrix.

Data for 15 [CRIME AND FEAR]

CORRELATION MATRIX BEFORE TRANSFORMATION

	1	2	3	4	5	6	7	8	9	10
1	1.000	0.101	0.227	0.430	0.047	0.316	-0.024	-0.029	-0.035	-0.019
2	0.101	1.000	0.408	0.164	0.165	0.170	0.042	0.006	-0.020	-0.128
3	0.227	0.408	1.000	0.265	0.274	0.217	0.091	-0.023	-0.000	-0.154
4	0.430	0.164	0.265	1.000	0.103	0.438	-0.007	-0.045	0.019	-0.142
5	0.047	0.165	0.274	0.103	1.000	0.061	0.115	-0.055	-0.051	-0.172
6	0.316	0.170	0.217	0.438	0.061	1.000	-0.021	-0.037	0.015	-0.004
7	-0.024	0.042	0.091	-0.007	0.115	-0.021	1.000	-0.072	-0.048	-0.004
8	-0.029	0.006	-0.023	-0.045	-0.055	-0.037	-0.072	1.000	0.016	-0.144
9	-0.035	-0.020	-0.000	0.019	-0.051	0.015	-0.048	-0.016	1.000	0.004
10	-0.019	-0.128	-0.154	-0.142	-0.172	-0.004	-0.004	-0.144	0.004	1.000

EIGENVALUES

1	2.274	1.545	1.159	1.003	0.974	0.853	0.747	0.645	0.520	0.444
---	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

COMPONENT LOADINGS

	1	2
1	0.545	-0.450
2	0.570	0.515
3	0.710	0.311
4	0.404	-0.591
5	0.394	0.449
6	0.412	-0.409
7	0.108	0.415
8	-0.055	0.174
9	-0.003	-0.253
10	-0.325	-0.400

DISCRIMINATION RESULTS

VARIABLES

1	2	3	4	5	6	7	8	9	10
0.310	0.348	0.508	0.491	0.161	0.407	0.021	0.028	0.038	0.039

CATEGORY QUANTIFICATIONS

CATEGORIES

VARIABLES

1	2	3	4	5	6	7	8	9	10	11	12	13	14
-1.505	-0.619	-0.828	0.288	0.942	0.942	0.288	0.942	0.942	0.288	0.942	0.942	0.288	0.942
-0.687	-0.001	0.043	0.766	1.408	1.408	0.766	1.408	1.408	0.766	1.408	1.408	0.766	1.408
-1.894	-0.916	-0.704	0.038	0.849	0.849	0.038	0.849	0.849	0.038	0.849	0.849	0.038	0.849
-1.997	-0.774	-0.575	0.572	0.822	0.822	0.572	0.822	0.822	0.572	0.822	0.822	0.572	0.822
-0.050	-0.301	0.128	0.171	0.019	0.019	0.171	0.019	0.019	0.171	0.019	0.019	0.171	0.019
0.181	-0.181	-0.094	-0.470	0.470	0.470	-0.181	0.470	0.470	-0.181	0.470	0.470	-0.181	0.470
0.088	0.170	-0.724	-0.041	0.181	-0.387	0.181	0.181	0.181	-0.387	0.181	0.181	-0.387	0.181
0.124	0.107	0.090	-0.078	-0.081	-0.100	-0.078	-0.081	-0.100	-0.078	-0.081	-0.100	-0.078	-0.081

CORRELATION MATRIX AFTER TRANSFORMATION

1	2	3	4	5	6	7	8	9	10
1.000	0.126	0.238	0.407	0.054	0.315	0.126	0.054	0.117	0.068
0.126	1.000	0.498	0.287	0.150	0.160	0.094	0.117	0.068	0.135
0.238	0.498	1.000	0.279	0.311	0.322	0.080	0.207	0.179	0.163
0.407	0.287	0.279	1.000	0.103	0.448	0.105	0.096	0.146	0.161
0.054	0.150	0.311	0.103	1.000	0.083	0.035	0.139	0.120	0.161
0.315	0.160	0.222	0.448	0.083	1.000	0.096	0.075	0.117	0.090
0.126	0.094	0.080	0.105	0.035	0.096	1.000	0.134	0.037	0.075
0.054	0.117	0.207	0.096	0.139	0.075	0.134	1.000	0.105	0.123
0.080	0.065	0.179	0.146	0.120	0.117	0.037	0.105	1.000	0.131
0.020	0.135	0.163	0.161	0.161	0.090	0.075	0.123	0.131	1.000

EIGENVALUES

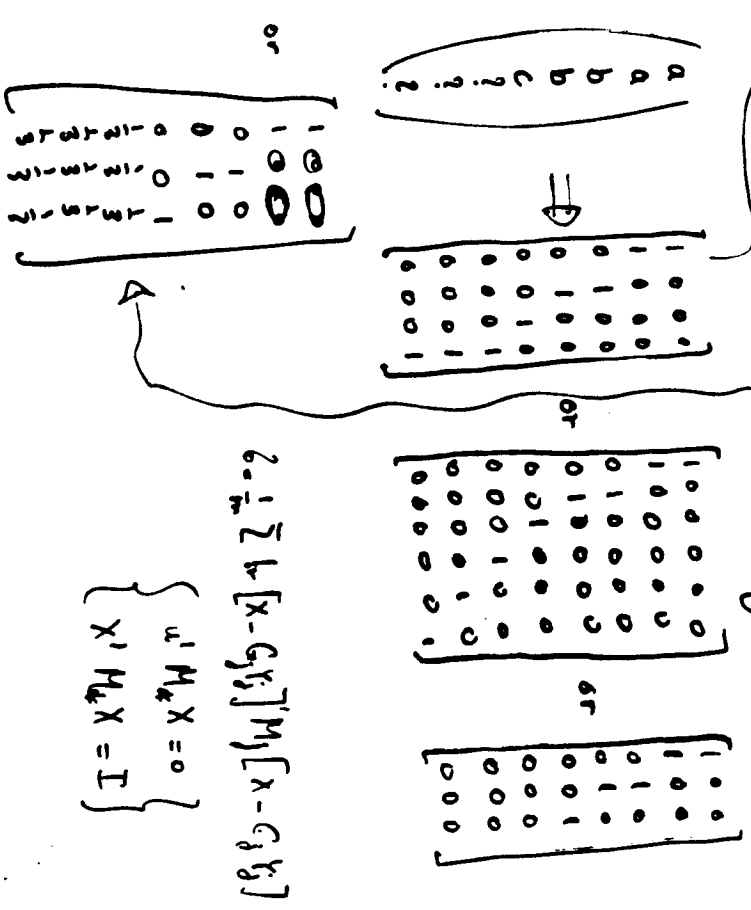
2.510	1.274	1.030	0.999	0.875	0.827	0.820	0.678	0.519	0.450
-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

COMPONENT LOADINGS

1	2
0.535	0.556
0.567	-0.292
0.707	-0.295
0.662	0.457
0.408	-0.471
0.592	0.443
0.285	0.040
0.362	-0.334
0.364	-0.138
0.368	-0.319

Missing data [refer to Meulman, in list].

- missing data single category
- missing data multiple categories
- missing data deleted
- [fuzzy coding]



It is clear that these choices can have important consequences.

- single category suggests more homogeneity than there usually is.
- multiple will lead to dimensions determined by missing values.
- deleted is standard [in HOWARD and PRUDENIS]

The options can also be used in a much more interesting and creative way.

- Example
- Munsinger - Rain
  - [Japanese]
  - [Skiing Resorts]

The example also illustrates the reordering [parallelgram analysis] used of homogeneity analysis. Also. multinomial gauge.

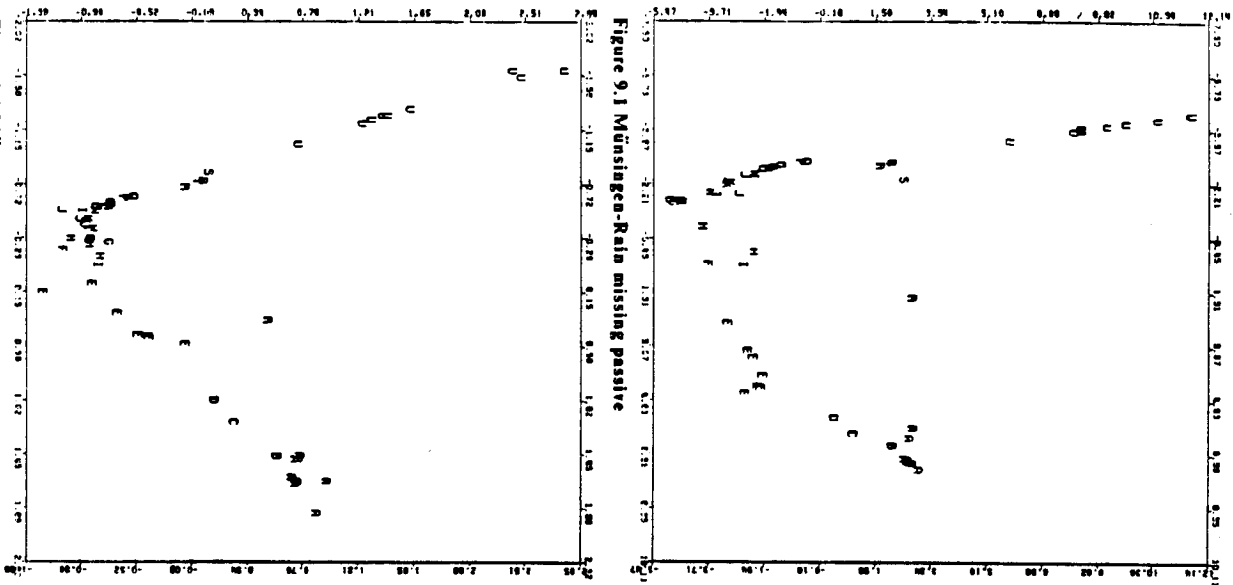


Figure 9.1 Münstingen-Rain missing passive

Figure 9.2 Münstingen-Rain missing multiple

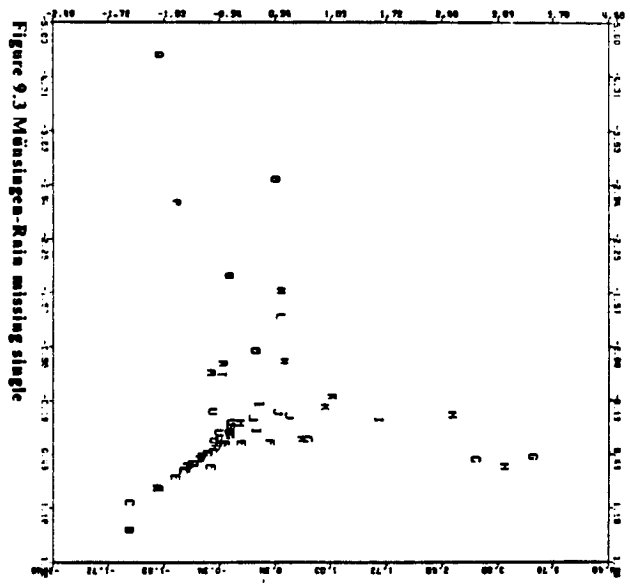


Figure 9.3 Münstingen-Rain missing single



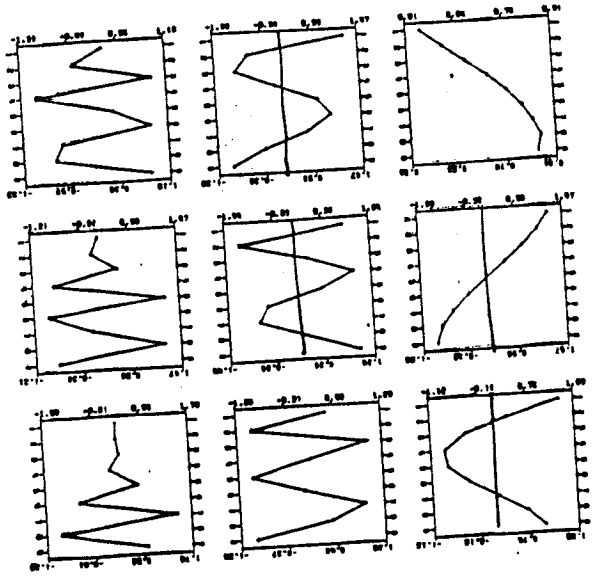








Gutman 1950



FIGUR 2. 2-eigenvektoren Gutman gespreid

-82-

At some later point we shall study applications of HORN'S to binary data & time series. But now ----- CORRESPONDENCE ANALYSIS.

[HISTORICAL REMARKS]

What is it? We can be brief. CA is MCA if  $m=2$ . That is a lot too brief, perhaps. CA is often [in France, for instance] presented in different ways, not on a special case of MCA. In fact .... the other way around.

(a) Bénédicti - distances [ $\chi^2$  distances].

Suppose  $F = \{f_{ij}\}$  is an  $m \times n$  table with <sup>relative</sup> frequencies. <sup>adding up to one.</sup> Let  $e_i = \sum_j f_{ij}$  and  $d_j = \sum_i f_{ij}$ .

Define the Berezin-distances between rows by

$$S_{ik}^2 = N \sum_{j=1}^n \frac{1}{d_j} \left[ \frac{f_{ij}}{c_i} - \frac{f_{kj}}{c_k} \right]^2$$

or

$$S_{ik}^2 = N (u_i - u_k)' E^{-1} F D^{-1} F' E^{-1} (u_i - u_k)$$

The idea is to approximate this by ordinary

Euclidean distances. we use the singular value

decomposition

$$N^{1/2} E^{-1/2} F D^{-1/2} = K \Psi L' \quad \left\{ \begin{array}{l} L'K = I \\ L'L = I \\ \Psi \text{ diag} \end{array} \right.$$

and define

$$X = E^{-1/2} K \Psi \quad (\text{thus } X'E X = \Psi K' K \Psi = \Psi)$$

Then

$$S_{ik}^2 = (u_i - u_k)' X X' (u_i - u_k) = d_{ij}^2(X)$$

If we truncate the SVD at p dimensions, then

$$d_1^2(X_1) \leq d_1^2(X_2) \leq \dots \leq d_1^2(X_r) = S_{11}^2$$

where  $r = \text{rank}(F)$ .

Example London (see 16)

Again, dually, we can define B-distances between columns and approximate them. This gives  $Y = D^{-1/2} L \Psi$ .

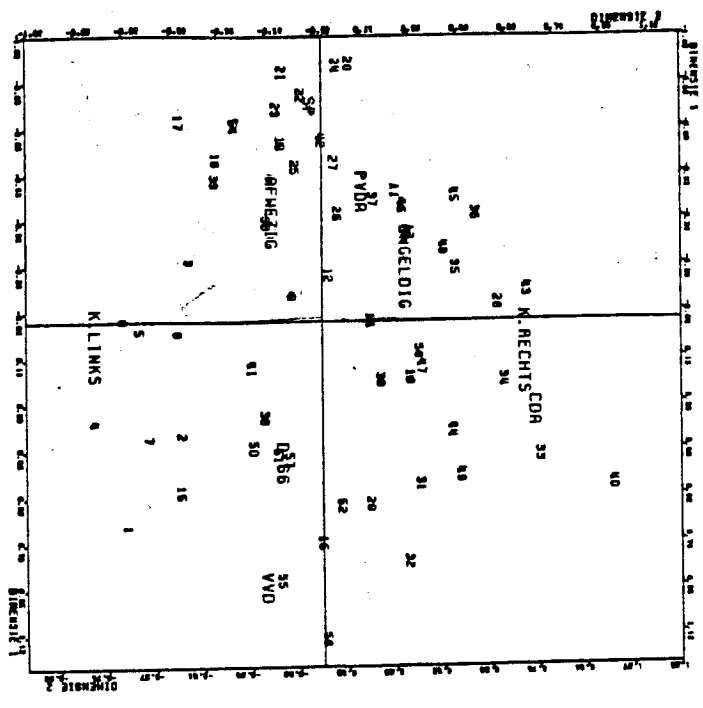
Transition equations [SVD in disguise]

$$X = E^{-1/2} F Y \Psi'$$

$$Y = D^{-1/2} F' X \Psi'$$

- $\alpha=0 \quad \beta=1 \rightarrow$  first centroid principle
- $\alpha=1 \quad \beta=0 \rightarrow$  Scree
- $\alpha=1/2 \quad \beta=1/2 \rightarrow$  biplot scaling [Gabriel]
- $\alpha=1 \quad \beta=1 \rightarrow$  French scaling.

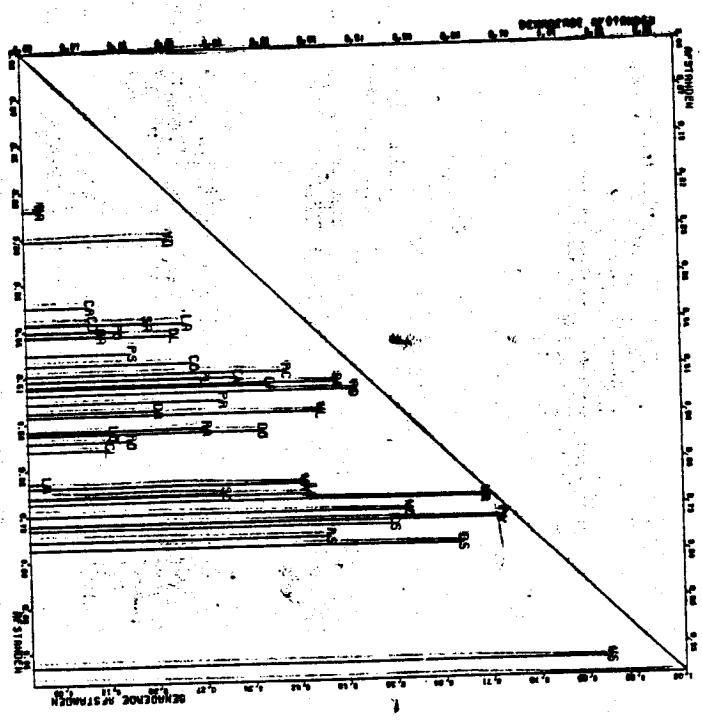
Figuur 2.2: Analyse van marginale matrix van stembedistricten x partijen  
 $\lambda_1 = .273$  ( $t_1 = .66$ ),  $\lambda_2 = .155$  ( $t_2 = .21$ )



relatief veel VVD en D'66 stemt. Het "relatief" bedoel ik hier "t.o.v. het gemiddelde profiel van Leiden". Het gemiddelde Leidse partijprofiel en stembedistrictprofiel liggen in de oorsprong. Ik roep hier in betrekking dat een partijpunt i dicht ligt bij een stembedistrict j als  $x_{ij} > e_{ij}$ , en dat een partijpunt i ver weg ligt van stembedistrict j als  $x_{ij} < e_{ij}$ . Door K en R te vergelijken, zou men eventueel kunnen zeggen of de oplossing wel juist getolereerd is. Het is immers mogelijk dat twee punten op de eerste twee dimensies dicht bij elkaar liggen, maar een grote afstand hebben op een hogere dimensie. De tweede dimensie maakt een onderscheid tussen stembedistricten uit de binnenstad, waar men veel klein links stemt (een verklaring hiervoor kan zijn dat hier veel studenten wonen) en stembedistricten waar men veel stemt op christelijke partijen. Men kan zien dat CA een helder

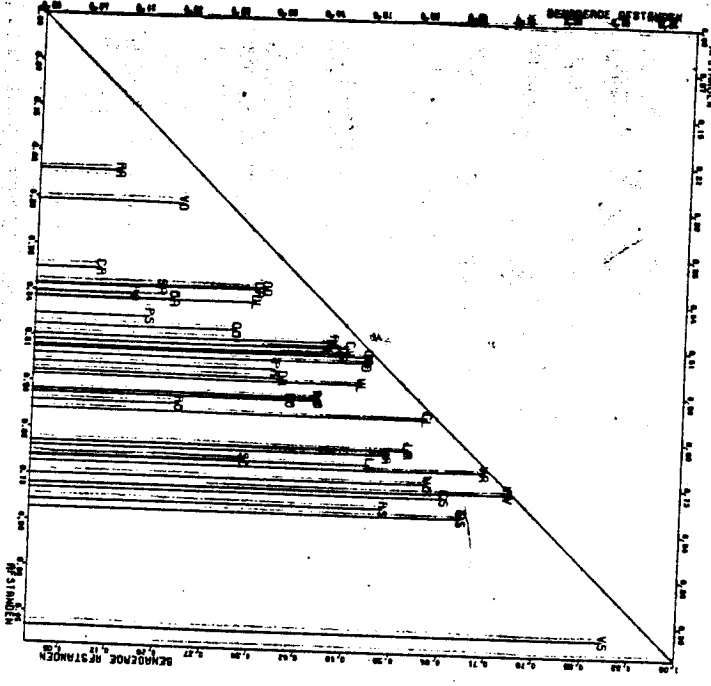
inzicht verschaft in de 9x9 matrix.  
 Het behulp van formule (2.14) worden de  $\chi^2$ -afstanden nader bestudeerd. Het behulp van formule (2.14) zijn de  $\chi^2$ -afstanden tussen de 9 partijen berekend, en deze zijn vergeleken met de in de eerste twee dimensies benaderde  $\chi^2$ -afstanden, berekend met het rechterlid van (2.14). In figuur 2.3 zijn voor de 36 paren rijpunten de ongebenederde  $\chi^2$ -afstanden (horizontaal) afgezet tegen de in de eerste dimensie benaderde  $\chi^2$ -afstanden (verticaal). Als labels zijn de eerste letters van de partijen genoemd (L voor klein links, en R voor klein rechts). De grootste afstanden zijn die tussen de BP- en resp. de VVD, D'66, klein rechts en het DA. Verder is de afstand groot tussen de VVD en resp. de PvdA, ONGERUDIG en AFWEZIG. In figuur 2.2 worden deze afstanden voor een groot deel op de eerste dimensie benaderd; de verticale lijnen

Figuur 2.3:  $\chi^2$ -afstanden (horizontaal) vs. benaderde afstanden voor een enkele dimensie (verticaal)



komen bijna tot de schuine lijn, waarvoor geldt dat de  $\chi^2$ -afstand gelijk is aan de benaderde  $\chi^2$ -afstand. Kleine afstanden zijn die tussen de paren PvdA en AVEZIG, VVD en D'66, CDA en Klein rechts. Deze paren halen relatief veel van hun stemmen uit dezelfde standplaatsen. Grote afstanden, die nauwelijks in de eerste dimensie worden benaderd, zijn bijvoorbeeld de afstanden tussen Klein links en resp. Klein rechts, CDA en ONGEIDIG. Deze benadering vindt plaats in de tweede dimensie. Figuur 2.4 laat zien in hoeverre de  $\chi^2$ -afstanden benaderd worden in twee dimensies. Alle verticale lijnen zijn hier langer dan in figuur 2.3, omdat een benadering in twee dimensies altijd beter is dan een benadering in één dimensie. Vooral de lijnen LR, CL, LO zijn een stuk langer geworden. Veel afstanden tot ONGEIDIG worden nog steeds slecht gerepresenteerd. Deze worden in hogere dimensies afgebeeld. Hogere

Figuur 2.4:  $\chi^2$ -afstanden (horizontaal) vs. benaderde afstanden voor eerste twee dimensies (verticaal)



	A	AB	AC	AD	AE	F	G
1	31	18	10	10	0	1164	111
2	0	15	54	18	27	113	14
3	0	40	82	10	142	160	0
4	19	0	35	19	28	163	52
5	14	7	35	22	17	0	56
6	69	72	38	55	184	48	0
7	74	0	86	14	0	184	48
8	78	0	80	5	17	0	0
9	74	19	33	12	26	0	0
10	80	68	67	15	29	0	0
11	108	48	108	4	10	108	0
12	109	13	5	17	39	0	46
13	16	35	69	24	0	26	41
14	26	86	60	6	48	48	28
15	290	10	6	0	0	0	0
16	184	48	82	42	134	0	0
17	29	0	0	0	41	211	32
18	0	19	56	0	39	75	0
19	0	22	45	42	60	230	69
20	178	0	90	28	68	0	0

Table 1.1: number of pages of RWA books devoted to several subjects.

1	ROLD	-1.1086	-0.6144	-0.3390
2	KEN1	0.0740	0.7625	-0.2526
3	KEN2	-0.2115	0.4605	-0.6923
4	ANDE	-0.7779	-0.1197	0.1556
5	COL1	0.0278	0.4065	1.0513
6	COL2	0.3578	0.6960	0.0928
7	MORI	-0.1641	-0.1572	0.4635
8	MOR2	-0.2502	-0.1863	0.3900
9	GEEL	0.7279	-0.1945	-0.0475
10	GEER	0.6840	0.2434	-0.1796
11	DEMP	0.0273	-0.3665	-0.4430
12	TATS	0.2680	-0.4475	0.2829
13	HARR	0.0219	0.4846	0.5172
14	DAGN	0.1205	0.5089	-0.1948
15	GREC	1.0831	-0.0321	0.0821
16	CAPA	0.6496	-0.0708	-0.1827
17	GIR1	-0.9635	-0.3927	-0.2502
18	GNAM	-0.4001	0.3292	-0.3889
19	KSMI	-0.7473	0.0810	-0.10051
20	THOR	0.5655	0.8145	-0.3526

Table 1.2: projections of books, professions of subjects and singular values from correspondence analysis on table 1.1.

LABDA	0.6038	-0.5103	0.3364
-------	--------	---------	--------



Carry over from MCA

$N \sum \psi_i^2 = X^2$  (decomposition)

of departure from independence, analysis of residuals. Van der Heyden.....

CA find scores for rows and columns that maximize the correlation. This makes it [perhaps] interesting as a scaling technique. Example Guilford, Bock, 1960; Guttman, 1946.

CA: linear regression; Example Rasen. of transition [irritable]. MCA.

Factor decomposition [Guttman] MCA

$F = \sum u_i u_i + X \psi_i \psi_i$

$F_i = d_j e_i \sum \psi_s X_{ij} \psi_s$

LL

cards in the same groups and equal category quantifications for all ten replications. The homogeneity approach has a natural interpretation in psychophysical contexts, in which we very commonly suppose that there 'is' a one-dimensional scale and that the different variables are merely replications.

Guilford's data matrix is reproduced in table 10.2.1. The first singular value of an ANCOR analysis on this table was .93, the remaining singular values conformed closely to  $\lambda_s = (.93)^s$ , which shows that we only have to pay attention to a single dimension. Table 10.2.2 contains optimal quantifications of spot patterns and of pilae (intervals). Both transformations are plotted in figures 10.11 and 10.2. It is clear from the transformation of the intervals that correspondence analysis does not follow the instruction to keep interpile distances equal, the intervals in the middle are larger, near the endpoints the distances are smaller. The transformation of the stimuli is fairly linear, deviations from linearity are in the direction of the stimuli. In this sense correspondence analysis, which produces the 'best' scale in a specific least squares sense, confirms Fechner's law, but this is obviously a very weak confirmation. It could very well be that a similar technique with a least absolute deviation loss function disconfirms Fechner's law.

S/R	1	2	3	4	5	6	7	8	9
15:	14	10	7	1					
16:	16	19	3	2					
17:	7	18	11	4					
19:	8	18	9	3	2				
20:	3	12	14	3	6	2			
22:	1	11	14	12	2				
24:			12	14	2				
26:			18	18	9	2			
28:			20	20	11	11			
30:			26	11	3				
32:		2	10	16	3	3			
35:			8	17	14	11			
37:			8	18	10	14			
40:			2	14	14	10			
43:				12	19	9			
46:				2	6	14			
48:					2	14			
53:						23			
58:						10			
56:						12			
60:						22			
64:						5			
69:						14			
74:						17			
						20			
						14			

Table 10.2.1 Guilford's Spot Pattern data

20163



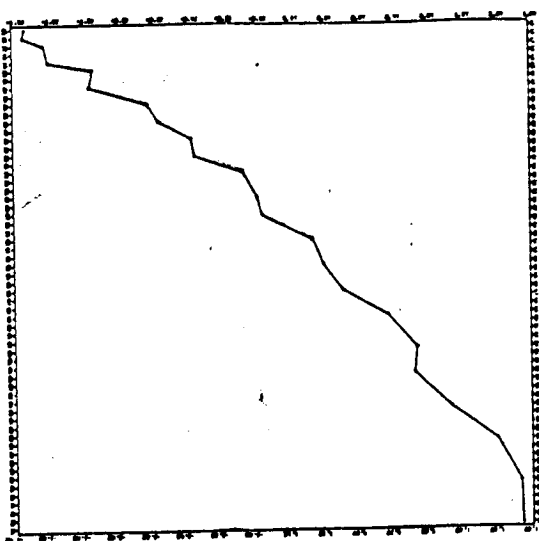


Figure 10.1 Gullford's spot patterns: transformations of the stimuli

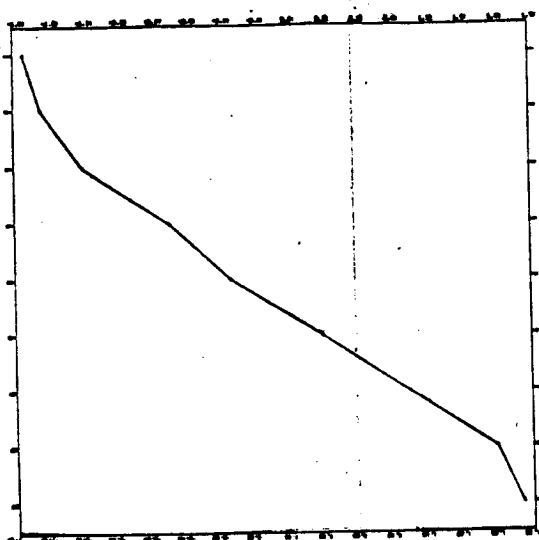
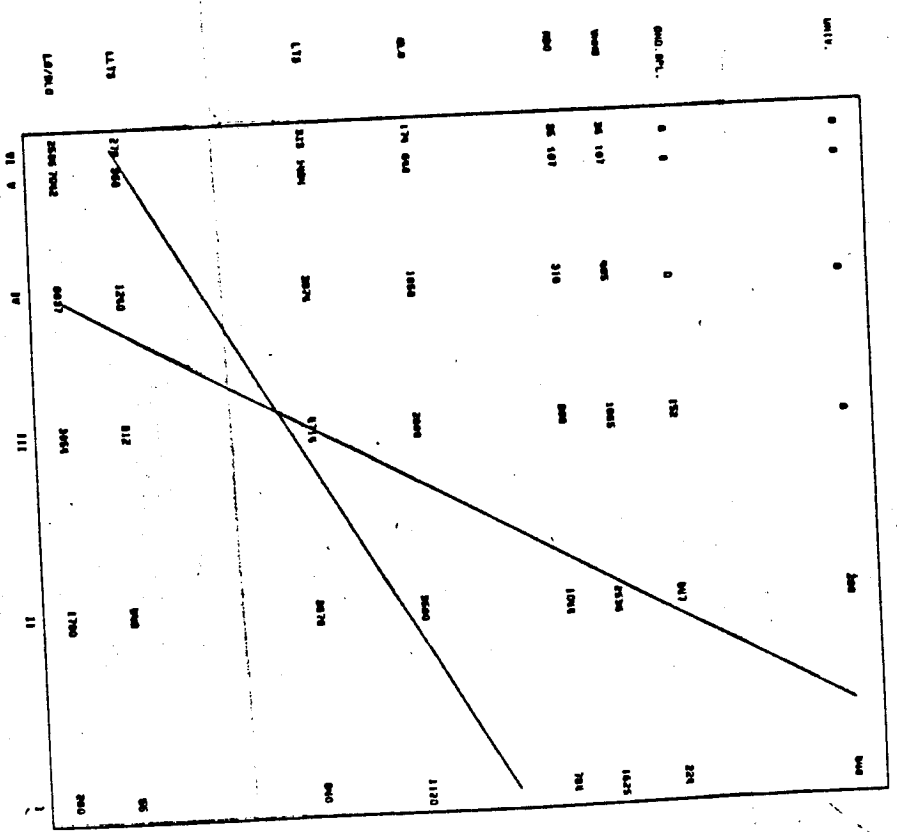


Figure 10.2 Gullford's spot patterns: transformations of the intervals

64

Figuur 3.4: Intelligentie (Raven-scores) en genoten schoolopleiding in 1952



Raven-scores langs de X-as in 6 klassen (1 zijn de hoogste scores)

65

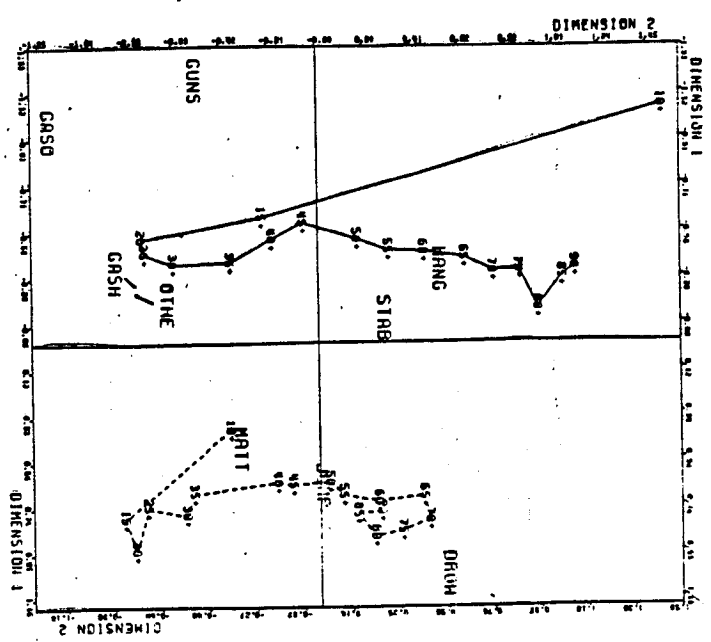


Table 1: Suicide behavior: age by sex by cause of death

- Labels for cause of death-categories:
1. Suicide by solid or liquid matter (MATT)
  2. Suicide by toxification of gas at home (GASH)
  3. Suicide by toxification with other gas (GASO)
  4. Suicide by hanging, strangling, suffocating (HANG)
  5. Suicide by drowning (DROW)
  6. Suicide with guns and explosives (GUNS)
  7. Suicide with knives etc. (STAB)
  8. Suicide by jumping (JUMP)
  9. Suicide with other methods (OTHE)

	MATT	GASH	GASO	HANG	DROW	GUNS	STAB	JUMP	OTHE	TOTAL
Men	4	0	0	247	1	17	1	6	9	285
10-15	348	7	47	578	22	179	11	74	175	1461
15-20	808	32	229	699	44	316	35	109	289	2561
20-25	808	32	229	699	44	316	35	109	289	2561
25-30	789	26	243	648	52	268	38	109	226	2399
30-35	916	17	257	825	74	291	52	123	281	2836
35-40	1118	27	313	1278	87	293	49	134	264	3567
40-45	926	13	250	1273	89	299	53	78	198	3179
45-50	855	9	203	1381	71	347	68	103	180	3227
50-55	684	14	136	1282	87	229	62	63	146	2703
55-60	502	6	77	972	49	151	46	66	77	1946
60-65	516	5	74	1249	83	162	52	92	122	2335
65-70	513	8	31	1360	75	164	58	113	95	2417
70-75	425	5	21	1268	90	121	44	119	82	2175
75-80	266	4	9	868	63	78	30	79	34	1429
80-85	159	2	2	479	39	18	9	46	19	782
85-90	70	1	0	259	16	10	9	18	10	393
90+	18	0	1	76	4	2	4	8	2	113
TOI	8917	176	1913	14740	946	2945	628	1340	2223	33828
Women	28	0	3	20	0	1	0	10	6	68
10-15	353	2	11	81	6	15	2	43	67	560
15-20	540	4	20	111	24	9	9	66	82	862
20-25	454	6	27	125	33	26	7	66	73	839
25-30	530	2	29	178	42	14	20	82	76	865
30-35	688	5	44	272	64	24	14	88	110	1319
35-40	568	4	24	343	78	18	22	103	86	1242
40-45	716	6	24	447	84	13	37	129	90	1504
45-50	942	7	26	691	184	21	21	131	92	1658
50-55	723	3	14	527	163	30	35	140	114	2083
55-60	820	8	8	702	245	14	38	156	90	2096
60-65	740	8	4	610	244	1	27	129	46	1691
65-70	624	6	4	420	161	1	29	129	35	1279
70-75	495	3	2	223	78	0	10	84	23	715
75-80	292	2	0	83	14	0	6	34	2	256
80-85	113	1	0	19	4	0	2	7	0	57
85-90	24	0	0	0	0	0	0	0	0	24
90+	24	1	0	19	4	0	2	7	0	57
TOT	8648	77	241	5637	1703	172	309	1505	1090	19382

Figure 1: Analysis of  $F_{jk}(S_{jk})$ , first two dimensions. The line for men is solid. Singular values with relative contributions:  $\lambda_1 = .312 (.519)$ ,  $\lambda_2 = .268 (.381)$ .



[M] [SA]

Author	Walter Heiden	DATE	10/20
Project	2	With	





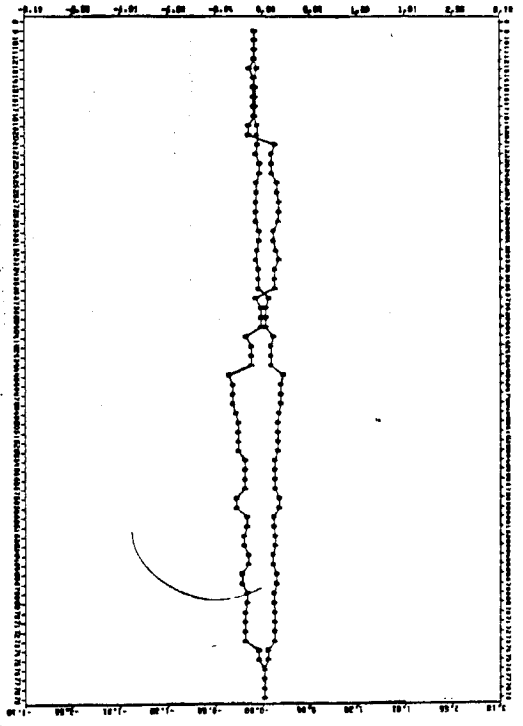


Figure 3: average induced curves for boys and girls

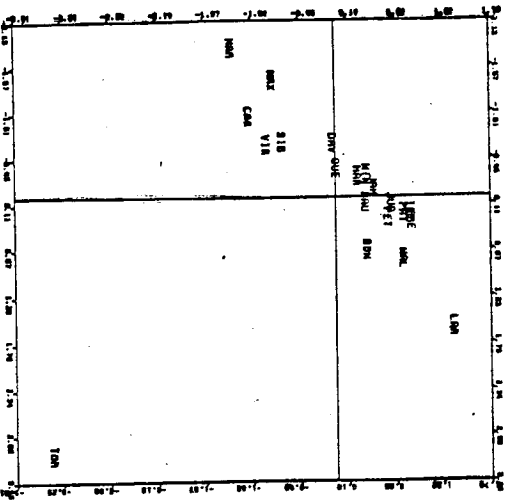


Figure 4: baby object scores in two dimensions

74

	1978	1979	1980	1981	1982
Ibo-	384	408	279	95	132
Ibo+	384	2	115	132	
havo 1	1695	725	64	2	
havo 2	1345	11	829	133	8
havo 3-			28	11	
havo 3+		11	377	262	276
havo 4-			589	460	
havo 4+	107				1
havo 1	1329	289	8		
havo 2			469	26	1
havo 3		1154	970	1011	472
havo 4				233	928
havo 5-				336	
havo 5+					
vwo 1	34				
vwo 2	1449	106			
vwo 3		1276			
vwo 4			175	5	22
vwo 5			1071	215	22
vwo 6			1	952	462
havo 6				1	792
havo 1			12	1	22
havo 2				426	608
havo 3				16	426
havo 4				5	18
havo 5					154
uitval	82	148	477	1142	1142
totaal	5464	5464	5464	5464	5464

Tabel 2  
schoolloopbaan categorieën, verdeling over de meetpunten.

75

	Leerjaar 2	Leerjaar 3	Leerjaar 4	Leerjaar 5	Leerjaar 6
lbo+	-1.598	-1.530	-1.440	-1.332	-1.319
mavo 1	-1.146	-1.085	-1.158	-1.051	-1.165
mavo 2	-0.932	-0.763	-0.868	-1.045	-0.972
mavo 3		-0.793	-1.021	-0.903	
mavo 3+			-0.691	-0.861	
mavo 4+			-0.698		
mavo 1	-0.384	-0.855	-0.257	-0.139	-0.495
mavo 2	+0.019	+0.124	-0.118	-0.145	-0.017
mavo 3			+0.346	+0.288	
mavo 4				+0.296	
mavo 5+					
mavo 1	+0.452	+0.718	+1.008	+1.271	+0.696
mavo 2	+1.396	+1.520	+1.633	+1.665	+1.167
mavo 3					+1.686
mavo 4					-1.080
mavo 5					-0.830
mavo 6					-0.697
mavo 1			-0.071	-0.863	-0.880
mavo 2				-0.156	-0.697
mavo 3				+0.035	+0.341
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					
mavo 5					
mavo 6					
mavo 1					
mavo 2					
mavo 3					
mavo 4					

(a) the 'elbow'-criterion indicates that the remaining singular values are approximately equal, (b) the remaining dimensions are not really 'common' factors but contrast either one time period or one activity with the rest, (c) three-dimensional plots are less attractive. Figure 1 shows the projections of the 940 individuals on the first two singular vectors. The very large cluster in the top right-hand section are the individuals who are mainly at home. The second large cluster, top left, are people who work full-time. Thus dimension one, the horizontal dimension, contrasts 'work' with 'being-at-home'. The second, vertical, dimension contrasts 'working' and 'being-at-home'. At the top, with other activities outside the house, at the bottom. These alternative activities are mainly in the category 'other', but also in 'shopping' and 'traveling'. This interpretation of the dimensions becomes beautifully clear if we plot the  $5 \times 22 = 110$  category quanti-

Figure 1: analysis 1, 940 object scores in two dimensions

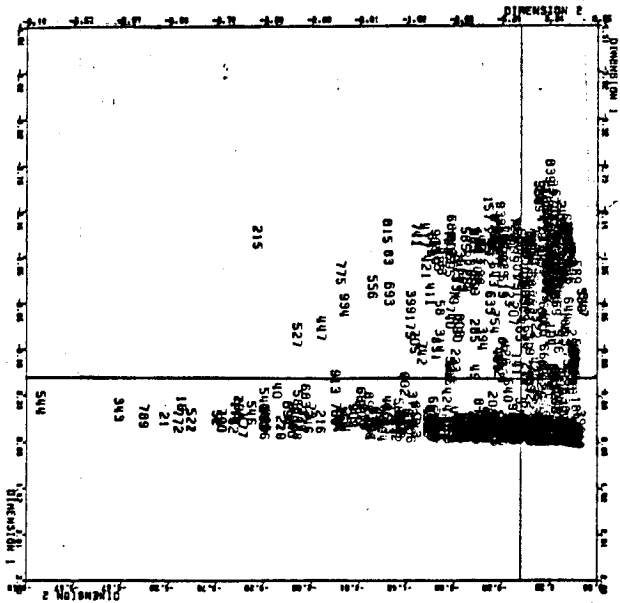
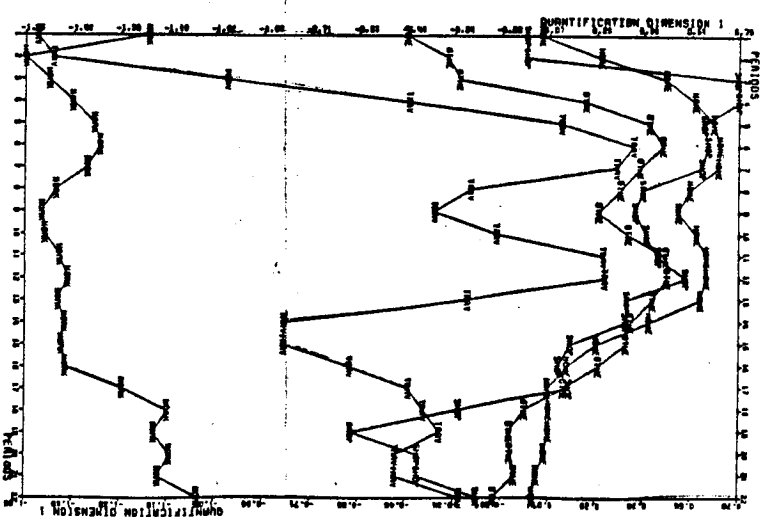


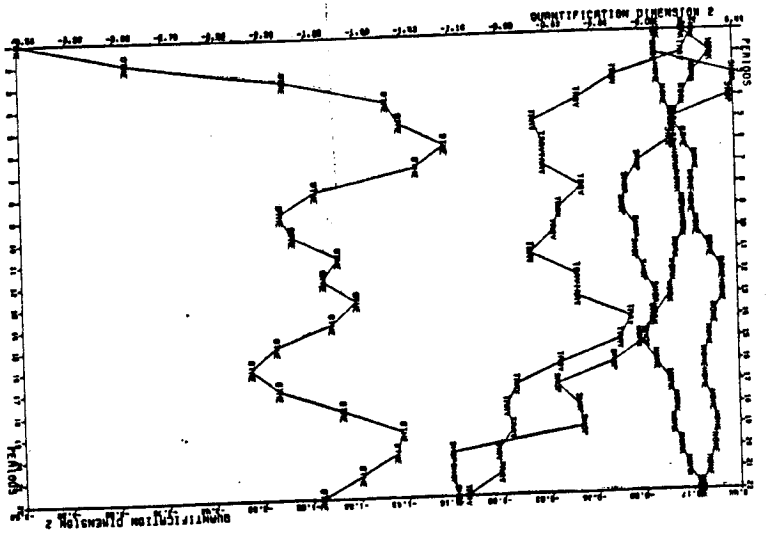
Figure 2a: analysis 1, category quantifications, dimension one against time



fications. In this analysis we have not plotted them in two dimensions, but we have plotted each dimension separately against time-period. Figure 2a is the plot for the first dimension. We have seen that these category quantifications make maximum discrimination of the individuals possible. They are related in a simple way to figure 1: the score for 'travel' in period 6 is the average score of all individuals who travel in period 6 (on dimension 1, and weighted with the number of minutes they travel). Individuals who work a lot are low on the dimension, individuals who are at home are high. During working hours the



Figure 2b: analysis 1, category quantifications, dimension two against time



average for travelling persons is close to the average for persons at home, during the lunch break. During the early morning and late evening hours it is much closer to the average of the persons who work. The same thing is true for shopping, although working people do not shop a great deal during lunch time, and do not even shop much during late afternoon. Shopping is done by those who stay at home. Figure 2b shows a similar plot for the second dimension. The social and cultural activities are concentrated in the morning, in the early afternoon, and in the early evening. The morning and early

afternoon shopping behaves like 'other', in the evening it is quite different. Travelling behaves in the opposite way. If you are going to visit somebody or something, then you have to travel before and after this visit. Thus travelling hours are just before and just after visiting hours. On the first dimension we can best discriminate people during working hours, on the second dimension we can best discriminate them outside working hours. More detailed interpretations of these plots are in Kreft and Mulder (1984). Another way to interpret the dimensions of the homogeneity analysis is by using passive variables (or supplementary variables). They play no role in the analysis, only in the interpretation. They are used afterwards to label the plots, and to compute centroids of groups of individuals. In our analysis we have used two passive variables: sex combined with work-situation of the family. Work situation has nine possible values: the head of the family can be employed full-time, part-time, or no-time, and the same thing is true for his/her partner. If we combine this with sex we have a new interactive passive variable

Figure 3: analysis 1, centroids for passive sex x work variable

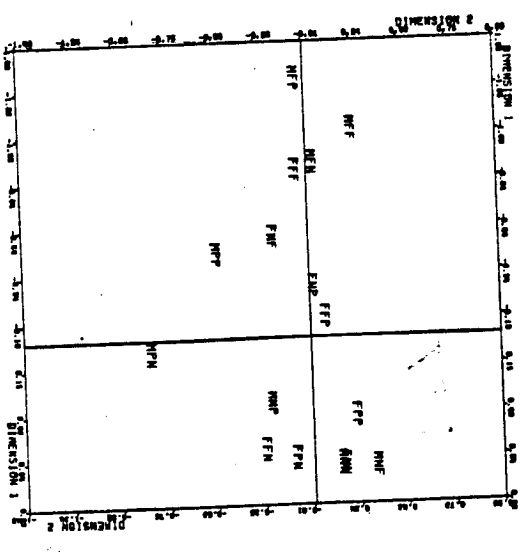


Figure 5a: analysis 2, 20680 x 5 table,  
average object scores for men and women,  
dimension one against time

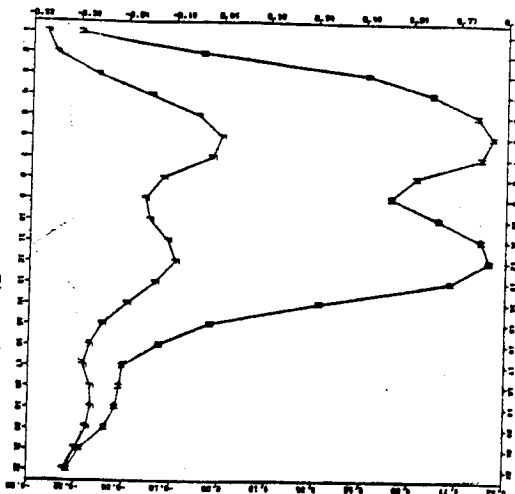
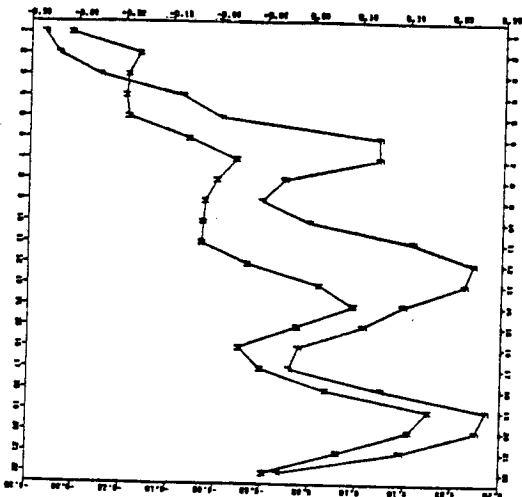


Figure 5b: analysis 2, 2068 x 5 table,  
average object scores for men and women,  
dimension two against time



Back to theory

Nonlinear PCA [program PRINONS]

NCA is related to PCA, but quite different.

If  $p=1$  then our analytical approach [near  $\lambda_{max}(R)$ ] shows that the relationship is quite close. But

if  $p>1$  then NCA finds a different transformation of the variables (and thus, implicitly, a new induced correlation matrix) for each dimension.

Since there are  $\sum_{j=1}^m (k_j-1)$  nontrivial PCA dimensions, this would mean as many correlation matrices [each of which could be used in a PCA!].

GIFI calls this DATA PRODUCTOR.

In order to create a link with RFA which is stronger, we shall now introduce various restrictions into the Giff systems. They are

- (a) rank restrictions
- (b) cone restrictions

For ease of reference we repeat our basic loss function

$$L(X; Y) = \frac{1}{n} \sum_{j=1}^m t [X - G_j Y]' [X - G_j Y]$$

and our normalization

$$u'Y = 0,$$

$$X'X = I.$$

### (a) Rank restrictions

$$Y_j = z_j a_j' \\ k_j x_p \quad k_j x_p$$

$$\left. \begin{array}{l} \text{IDENTIFICATION} \\ u_j' D_j z_j = 0 \\ z_j' D_j z_j = 1 \end{array} \right\}$$

### Geometrical [Figure CAR5]

$a_j$  defines a direction [through the origin] and  $z_j$  indicates the location on that direction. Thus we require

- (a) perfect homogeneity
- (b) all objects must be on parallel straight lines

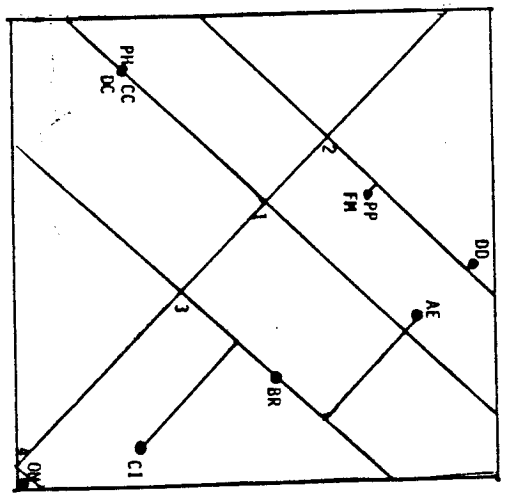


figure 3: single nominal loss, variable 1, arbitrary solution

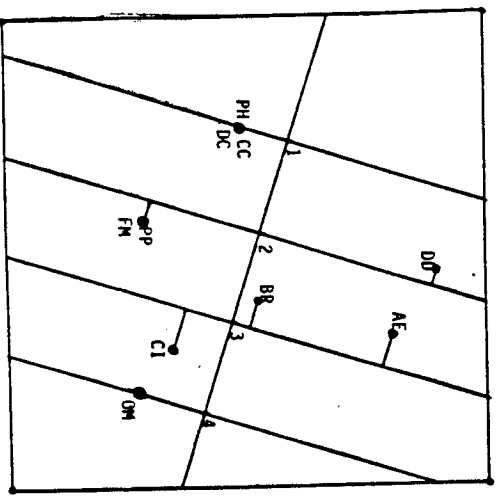


figure 4: single numerical loss, variable 1, optimal solution

ANALYTICAL

For the analytical interpretation we start from the dual formulation

max  $\text{tr } Y'CY$  with  $Y'DY = mI$ .

If all variables are single (i.e. not multiple, i.e. rank-restricted) then

$$\text{tr } Y'CY = \sum_j \text{tr } Y' C_j e_j e_j' Y - \sum_j \text{tr } a_j z_j' C_j z_j a_j'$$

$$= \sum_j \text{tr } a_j a_j' = \text{tr } R A A' = \text{tr } A' R A.$$

$$\text{tr } Y'DY = \sum_j \text{tr } Y' D_j Y = \sum_j \text{tr } a_j z_j' D_j z_j a_j' = \text{tr } A' A.$$

It follows (in the same way as before) that homogeneity analysis, when all variables single [nominal] is equivalent to

$$\max_{z_1, \dots, z_m} \lambda_1 [R] + \dots + \lambda_p [R]$$

where  $R$  depends on the  $z_j$  ( $r_{je} = z_j^i c_{je} z_e$ ).

It also follows that if we impose the further restriction that the  $z_j$  are known vectors, then the only unknown are in

$A_j$  and the technique becomes ordinary PCA for linear PCA.

Also: if  $p=1$  then rank-one is no restriction.

(b) Some restrictions

$$z_j \in K_j$$

- $K_j$  is  $\mathbb{R}^{k_j}$  [no restriction, single nominal]
- $K_j$  is a subspace of  $\mathbb{R}^{k_j}$  [for instance of polynomials]
- $K_j$  is a ray [linear vector, single numerical]
- $K_j$  is the case of monotone transformations [single ordinal]

In PRINCIPALS variables can be mixed

- multiple nominal [MCA]
- single nominal
- ordinal
- numerical [PCA]

TABLE 1. Table of social indicator statistics taken from statistical abstracts of the US (1977). U.S. Department of Commerce; Bureau of the census.

State	Popul	Income	Illit	Life	Homic	School	Freeze
Alabama	AB	3615	2.1	69.05	15.1	41.3	20
Alaska	AK	365	1.5	69.31	11.3	66.7	152
Arizona	AZ	2212	1.8	70.55	7.8	58.1	15
Arkansas	AR	2110	1.9	70.66	10.1	39.9	65
California	CA	21198	1.1	71.71	10.3	62.6	20
Colorado	CO	2541	0.7	72.06	6.8	63.9	166
Connecticut	CT	3100	1.1	72.68	3.1	54.6	103
Delaware	DM	579	0.9	70.06	6.2	52.6	11
Florida	FL	8277	1.3	70.66	10.7	52.6	11
Georgia	GA	4931	2.0	68.54	13.9	40.9	60
Hawaii	HA	868	1.9	73.68	6.2	59.5	0
Idaho	ID	813	0.6	71.87	5.3	59.5	126
Illinois	IL	11197	0.9	70.14	10.3	52.6	127
Indiana	IN	5313	0.7	70.88	7.1	59.0	140
Iowa	IA	2861	0.5	72.56	2.3	59.9	114
Kansas	KY	2280	0.6	72.58	4.5	59.9	114
Kentucky	KY	3387	1.6	70.10	10.6	38.5	95
Louisiana	LA	3806	2.8	68.76	13.2	42.2	12
Maine	ME	1058	0.7	70.39	2.7	54.7	161
Maryland	MD	4122	0.9	70.22	8.5	52.3	101
Massachusetts	MA	5814	1.1	71.83	3.3	58.5	103
Michigan	MI	9111	0.9	70.63	11.1	52.8	125
Minnesota	MN	3921	0.6	72.96	2.3	57.6	160
Mississippi	MS	2341	2.4	68.09	12.3	41.0	50
Missouri	MO	4767	0.8	70.69	9.3	48.8	108
Montana	MT	746	0.6	70.56	5.0	59.2	155
Nebraska	NE	1544	0.6	72.60	2.9	59.3	139
Nevada	NV	590	0.5	69.03	11.5	65.2	188
New Hampshire	NH	812	0.7	71.23	3.3	57.6	174
New Jersey	NJ	7333	1.1	70.93	5.2	52.5	115
New Mexico	NM	1144	2.2	70.32	9.7	55.2	120
New York	NY	18076	1.4	70.55	10.9	52.7	82
N. Carolina	NC	5461	1.8	69.21	11.1	38.5	80
N. Dakota	ND	637	0.8	72.78	1.4	50.3	186
Ohio	OH	10735	0.8	70.82	7.4	53.2	124
Oklahoma	OK	2715	1.1	71.42	6.4	51.6	82
Oregon	OR	2284	0.6	72.13	4.2	60.0	44
Pennsylvania	PA	11860	1.0	70.43	6.1	46.4	126
Rh. Island	RI	931	1.3	71.90	2.4	46.4	127
S. Carolina	SC	2816	2.3	67.96	11.6	37.8	65
S. Dakota	SD	681	0.5	72.08	1.7	53.3	172
Tennessee	TN	4173	1.7	70.11	11.0	41.8	70
Texas	TX	12237	2.2	70.92	12.2	47.4	35
Utah	UT	1203	0.6	71.64	4.3	67.3	137
Vermont	VT	472	0.6	71.72	5.5	57.1	168
Virginia	VA	4981	1.4	70.08	9.5	47.8	85
Washington	WA	3559	0.6	71.72	4.3	63.5	32
W. Virginia	WV	1799	1.4	69.48	6.7	41.6	100
Wisconsin	WI	4589	0.7	72.48	3.0	54.5	149
Wyoming	WY	376	0.6	70.29	6.9	62.9	173

School: Percent of the population over age 25 who are high school graduates  
 Freeze: Average numbers of days of the year in which temperature falls below freezing.

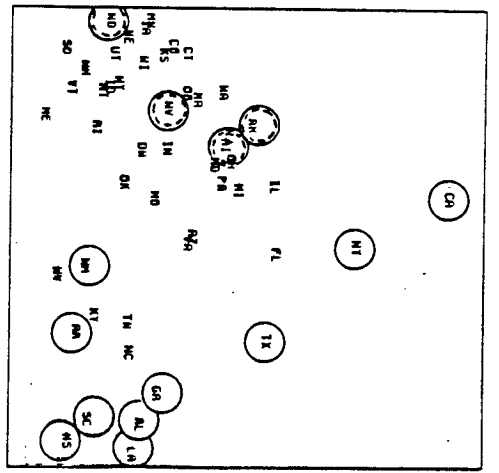


Figure 8. PCA solution for 50 states. Encircled points have dissimilarities larger than average. Dotted circles indicate more than average stress in addition.

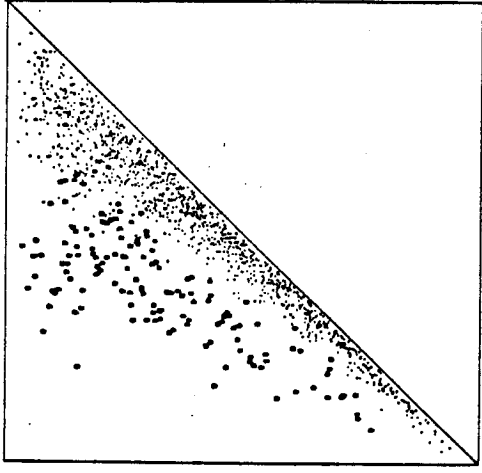


Figure 2. PCA solution for 50 states.  $d(Z)$  (horizontal axis) versus  $d(X)$  (vertical axis). Approximation from below. Ellipses refer to all pairs including AK, HI and NV.

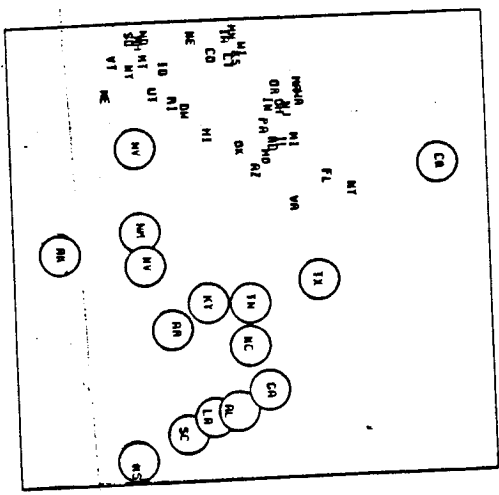


Figure 5. Nonlinear PCA solution for 50 states. Encircled points have dissimilarities larger than average.

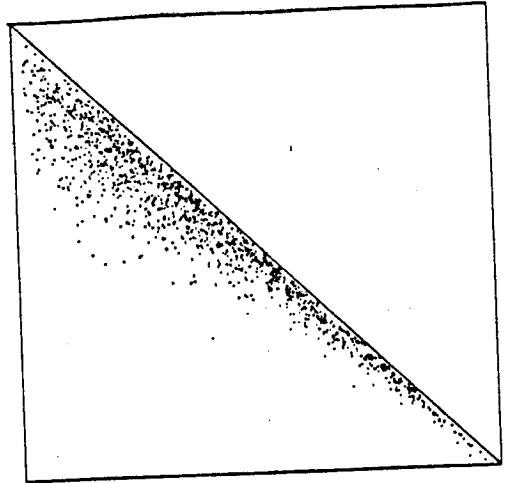


Figure 6. Nonlinear PCA solution for 50 states.  $d(Z)$  (horizontal axis) versus  $d(X)$  (vertical axis). Approximation from below.

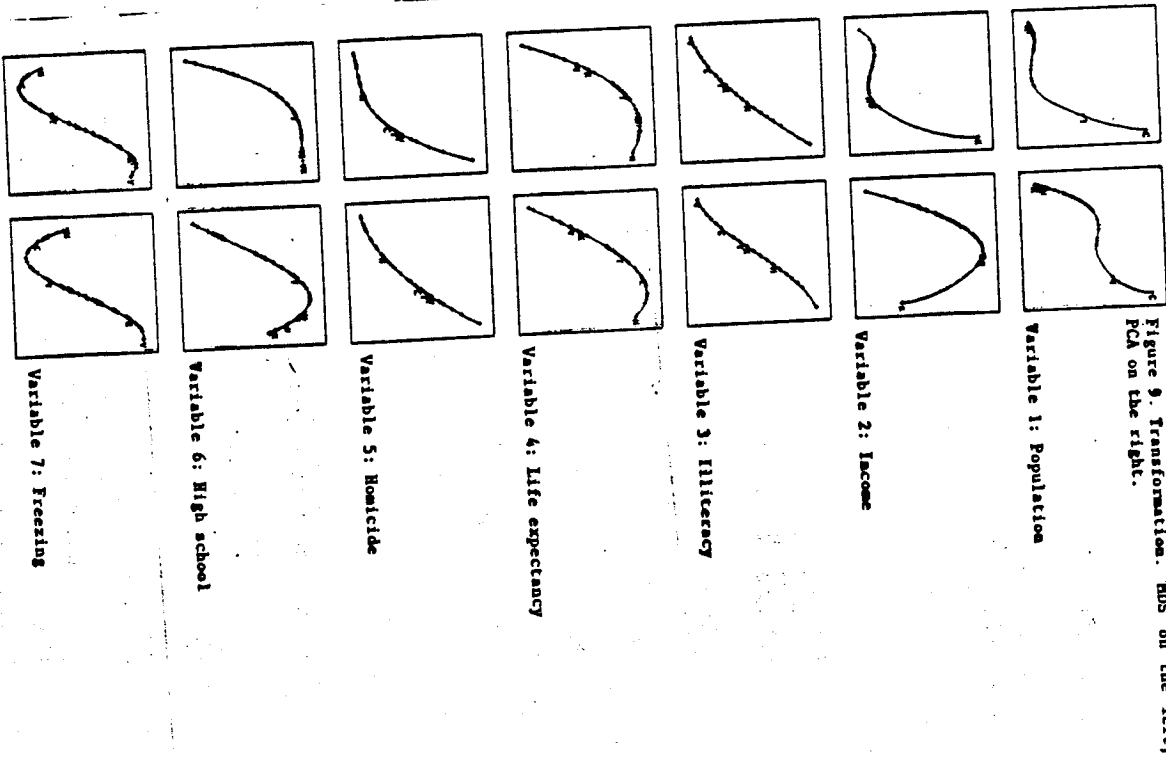


Figure 9. Transformation. MDS on the left, PCA on the right.

Variable 1: Population

Variable 2: Income

Variable 3: Illiteracy

Variable 4: Life expectancy

Variable 5: Homicide

Variable 6: High school

Variable 7: Freezing

1:	7	4	1	8	10	9	5	2	3	6	(S)
2:	7	6	2	9	3	8	10	1	4	5	(S)
3:	10	5	1	7	4	6	8	2	3	9	(S)
4:	6	5	3	7	4	8	9	2	1	10	(S)
5:	6	3	5	10	4	2	9	7	8	1	(D)
6:	8	7	4	9	2	5	10	6	3	1	(D)
7:	5	9	4	8	6	2	10	7	3	1	(D)
8:	6	7	4	9	5	3	10	8	2	1	(D)
9:	2	3	6	4	5	8	9	7	10	1	(D)
10:	5	8	2	9	1	7	10	6	4	3	(D)
11:	7	2	6	10	5	1	9	8	4	3	(D)
12:	8	7	2	9	1	6	10	5	3	4	(C)
13:	10	7	1	9	4	6	8	2	3	5	(C)
14:	5	2	3	4	1	8	7	9	6	10	(C)
15:	6	5	2	7	1	10	9	8	4	3	(C)
16:	4	7	5	2	8	9	1	6	3	10	(M)
17:	4	7	5	3	9	8	1	6	3	10	(M)
18:	5	4	7	3	9	8	1	10	2	6	(M)
19:	1	5	6	7	10	9	3	8	2	4	(E)
20:	1	5	8	7	9	3	6	10	2	4	(E)
21:	3	7	6	2	8	4	5	9	1	10	(E)
22:	1	3	8	6	9	7	4	10	2	5	(E)
23:	1	4	6	5	9	10	2	8	3	7	(E)
24:	1	7	5	4	10	9	3	8	2	6	(E)
25:	1	8	6	5	9	4	3	10	2	7	(E)
26:	1	2	5	6	10	4	7	9	3	8	(E)
27:	1	5	6	4	8	7	2	9	3	10	(E)
28:	4	6	5	1	7	10	3	8	2	9	(E)
29:	8	7	1	2	9	10	6	3	4	5	(R)
30:	7	4	1	2	9	10	6	3	4	5	(R)
31:	9	8	2	7	1	4	10	5	6	3	(R)
32:	7	1	5	8	2	6	3	9	4	10	(T)
33:	2	3	7	8	10	9	1	6	4	5	(T)
34:	10	4	2	9	3	5	6	8	1	7	(T)
35:	3	2	10	6	8	4	7	9	1	5	(T)
36:	6	1	3	9	4	7	10	2	5	8	(T)
37:	2	1	6	4	10	9	5	7	3	8	(T)
38:	2	3	6	5	7	8	4	9	1	10	(A)
39:	2	6	7	3	10	8	4	9	1	5	(A)

table 5.6: Roskam's journal data  
preference rank orders

### III.3. Roskam's journal preference data.

#### III.3.1.

Table III.3.1. gives preference rank orders of 39 psychologists for ten psychological journals (from Roskam, 1968). Columns of the table refer to the following journals:

- 1: JEXR: Journal of experimental psychology,
- 2: JAPP: Journal of applied psychology,
- 3: JPSF: Journal of personality and social psychology,
- 4: RUBR: Multivariate behavioral research,
- 5: JCLP: Journal of consulting psychology,
- 6: JEDP: Journal of educational psychology,
- 7: PWEK: Psychometrika,
- 8: HURE: Human relations,
- 9: BULL: Psychological bulletin,
- 10: Hude: Human development.

In addition, table III.3.1. contains a final column that identifies each psychologist with respect to the department he or she is affiliated. The codes are:

- S: social psychology (4),
  - D: educational and developmental psychology (7),
  - C: clinical psychology (4)
  - M: mathematical psychology and psychological statistics (3),
  - E: experimental psychology (10),
  - R: cultural psychology and psychology of religion (3),
  - T: industrial psychology (6),
  - A: physiological and animal psychology (2),
- Numbers between parentheses in the list above give the number of psychologists from the department.
- Table III.3.1. gives preferences in the usual way, from 1 (most preferred journal) to 10 (least preferred journal).

#### III.3.2.

A matrix of ranking such as given in table III.3.1. can be analyzed in two ways.

- (a) columns (journals) as variables, rows (psychologists) as objects,
- (b) rows (psychologists) as variables, columns (journals) as objects.





### The PRINCEPS algorithm

- Get  $\beta^1$
- (a) Start with  $X_0$ , where  $u^1 X_0 = 0$  and  $X_0^1 X_0 = I$ ,  
 and  $z_0$ , where  $u^1 z_0 = 0$  and  $z_0^1 z_0 = -1$ .
- step 1 :  $Y_j \leftarrow D_j G_j^1 X$
- step 2 : if variable  $j$  is multiple go to step 10
- step 3 :  $a_j \leftarrow Y_j^1 D_j z_j$
- step 4 :  $g$  variable  $j$  is single numerical  $g$  go to step 9
- step 5 :  $z_j \leftarrow Y_j a_j$
- step 6 : if variable  $j$  is single numerical go to step 8
- step 7 :  $z_j \leftarrow Proj(z_j)$
- step 8 :  $z_j \leftarrow z_j / (z_j^1 z_j)^{1/2}$
- step 9 :  $Y_j \leftarrow z_j a_j^1$
- step 10 : if  $j = h$  go to step 12
- step 11 :  $j \leftarrow j + 1$ , go to step 1
- step 12 :  $X = GRAM \left[ \sum_{j=1}^h G_j Y_j \right]$
- step 13 :  $\beta^1, g_0, h$  stop.

PRINCEPS is more general than computer programs in some respects, but less general in others. All transformations are discrete, ties remain fixed. But multiple and single can be mixed.

We do not have firsthand experience with the continuous normal / ordinal / numerical options. They are (a) slow and (b) unstable. In order to remedy this we use fuzzy coding again, this time in the form of splines.

Definition, example -  
 Cars  
 Public Spending  
 Cylinders

Program SPLINALS

= Yesterday we have discussed correspondence analysis as a special case of homogeneity analysis [with  $m=e$ ], and we have shown that Benzécri presents a geometrical derivation of the technique which makes multiple correspondence analysis a special case of ordinary correspondence analysis. We have also shown that CA maximizes correlation and linearizes regression.

= Homogeneity analysis gives multiple quantifications and thus multiple induced correlation matrices. This makes it somewhat dissimilar from PCA (except when pair). We make it more similar by introducing rank-one restrictions, in fact we make it identical by combining this with linear cone restrictions. The combination of homogeneity analysis, rank restrictions, and cone restrictions defines MPCA, program PRINDUALS.

Of course this is not the usual way in which PCA is introduced. Other derivations

- the inner product approximation [MDS]
- the linear approximation with unknown approximations [Real Rank Reg]  $[S \ll d]$
- the French derivation of typical variables (index numbers).

The one we have chosen has PCA to MCA, and thus is the centroid principle. We do not like inner products as a geometric "model". Some of the other approaches are perhaps covered by our "dual" approach via the Burt-table.

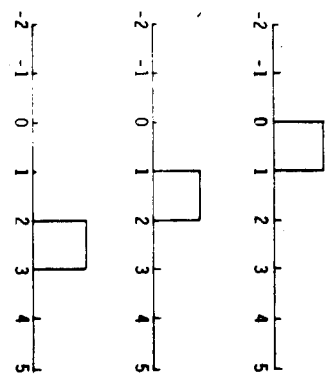


Figure 11.1.a  
zero degree B-splines

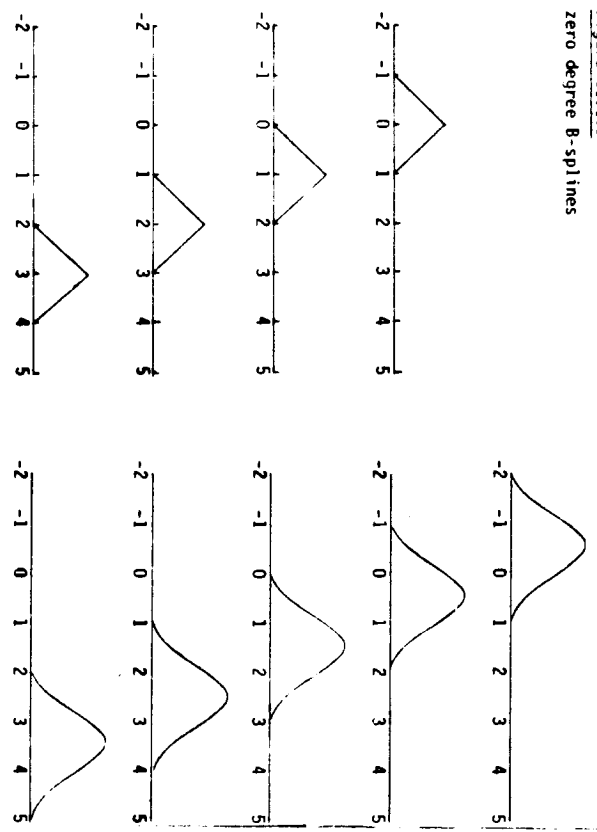


Figure 11.1.b  
first degree B-splines

Figure 11.1.c  
second degree B-splines

0.50	1 0 0	0.50 0.50 0.00 0.00	0.12500 0.75000 0.12500 0.00000 0.00000
0.75	1 0 0	0.25 0.75 0.00 0.00	0.03125 0.68750 0.28125 0.00000 0.00000
0.90	1 0 0	0.10 0.90 0.00 0.00	0.00500 0.59000 0.40500 0.00000 0.00000
2.30	0 0 1	0.00 0.00 0.70 0.30	0.00900 0.00000 0.24500 0.71000 0.04500

Table 11.1: B-splines of degree zero, one, two for four values of x.

var 1 = a <sub>1</sub>
var 2 = b <sub>1</sub>
var 3 = 2i(b <sub>1</sub> ) <sup>1/2</sup>
var 4 = 2a <sub>1</sub> (b <sub>1</sub> ) <sup>1/2</sup>
var 5 = a <sub>1</sub> b <sub>1</sub>
var 6 = (2a <sub>1</sub> ) <sup>-1</sup> a <sub>1</sub> b <sub>1</sub> <sup>2</sup>
var 7 = (2a <sub>1</sub> ) <sup>-1</sup> a <sub>1</sub> b <sub>1</sub> <sup>-1</sup>
var 8 = a <sub>1</sub> b <sub>1</sub> <sup>-1</sup>
var 9 = a <sub>1</sub> <sup>-1</sup> b <sub>1</sub>
var 10 = 2a <sub>1</sub> b <sub>1</sub> <sup>-2</sup>

Table 1. Thurstone's cylinder data

not accomplish this result because the second eigenvalue is not the second root of the correlation matrix R(θ). The fit in table 2 is defined as the sum of the first two eigenvalues of R(θ). The maximum fit is thus equal to 1.

Technique used	Type of transformations approximating φ	Total fit for p = 2	eigenvalues
matrix PCA	single linear	.87	.60 .27
ordinal PCA	single monotone	.94	.65 .28
splines-PCA k=1	single LS splines k=1	.89	.61 .28
splines-PCA k=2	single LS splines k=2	.98	.68 .30
splines-PCA k=3	single LS splines k=3	.99	.69 .29
Correspondence-analys+splines	multiple splines k=2	.77	.77 .66

Table 2. Fit, eigenvalues and transformation-types for several techniques.

The linear fit is surprisingly high compared with the ordinal fit, which on the other hand is clearly inferior to the single

TABLE 1  
Data on 1980 Automobiles

Make/Model	Price in \$	Engine Size	Miles/Gal.		Weight
			City	Highway	
Audi 5000	7100	98	19.3	35.8	2190
Chev Chevette	5100	98	18.6	35.9	2170
Datsun 210	4750	86	25.2	40.7	2000
Datsun 510	5950	119	21.9	42.7	2430
Dodge Colt	4800	98	24.2	41.2	1862
Ford Mustang	5800	141	15.6	29.8	2816
Ford Accord	4100	110	19.9	35.7	2250
Honda Civic EX	4100	91	22.7	39.7	1765
Kiaa GLC	4150	86	26.9	47.2	1960
Plymouth Horizon	5400	105	21.6	38.3	2156
Plymouth Sapparo	6500	159	13.8	28.8	2790
Pontiac Sunbird	4950	151	14.6	28.8	2700
Toyota Celica	6650	134	14.5	28.6	2410
Toyota Corolla	5700	134	16.6	34.9	2570
Toyota Tercel	4350	97	23.3	38.5	1970
VW Rabbit	6100	97	19.7	45.7	1870
VW Concord	8200	152	15.6	29.2	2910
ZK Eagle	7200	758	10.2	22.8	2660
Chev Citation	7200	173	14.4	31.8	3180
Chev Malibu	7200	232	13.9	22.8	3330
Dodge Aspen	6700	275	10.4	21.8	3530
Dodge Diplomat	7550	318	10.4	22.2	3530
Ford Fairmont	6550	140	13.1	29.8	2850
Mercury Monarch	7150	300	11.5	20.9	3510
Olds Cutlass	7550	331	11.9	25.5	3400
Pontiac LeMans	7400	231	11.9	25.5	2450
Pontiac Phoenix	7350	231	12.8	37.9	3480
Buick Regal	7400	231	9.8	22.4	3552
Chev Impala	8150	331	12.1	26.1	3360
Chev Monte Carlo	7900	305	10.3	22.2	3730
Dodge St. Regis	8100	318	10.3	21.5	3730
Mercury Marquis	8550	307	11.2	20.6	3720
Pontiac Catalina	8100	311	10.8	22.8	3610

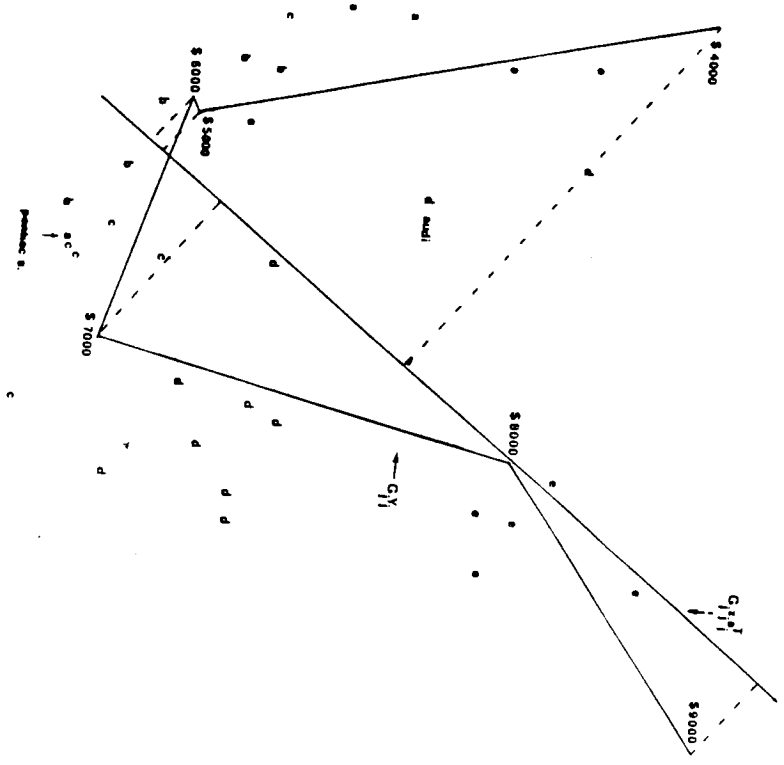


Figure 4 Multiple and single spline functions for price in the object scores plot, cars labeled for price level (a = between \$4000 and \$5000, etc.). Single hat analysis.

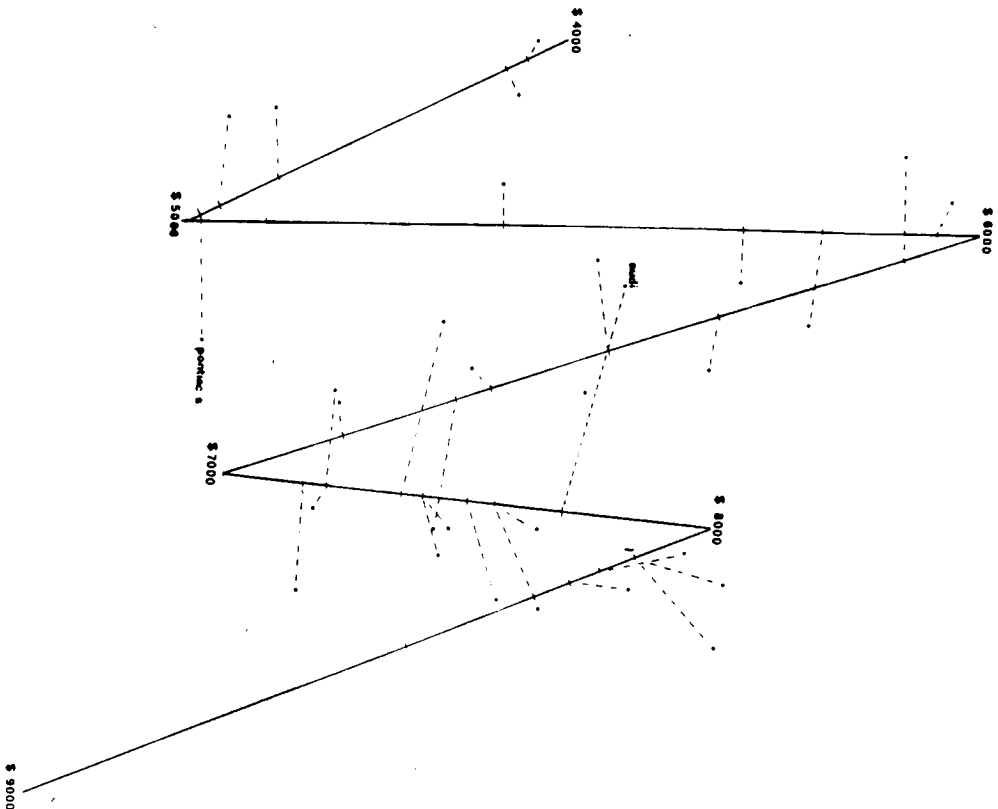
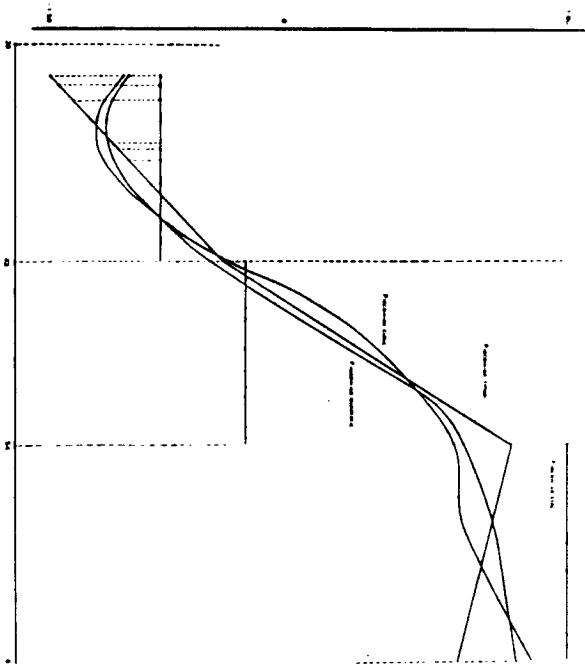


Figure 3 Multiple hat transformation of price with remaining variables. Linear (omitted) in the object scores plot. The object scores and the corresponding spline values are connected with dotted lines.

components accounted for 96% of the sum of all eigenvalues. From NCA with three components and with less restricted transformations we may expect the equal amount or more.

Figure 2: Transformations of public spending for varying degrees of smoothness.



**RESULTS: THE INTERPRETATION**

The results are reported concisely. An extensive, conceptually more detailed interpretation with further analysis is VAN KOOTEN & VAN RIJCKHOVEN (1985).

The eigenvalues are well separated: the first, .66, being three times as large as the second, .22. The third is half as important as the se-

Scaling criteria

We have seen that homogeneity analysis with  $p=1$  can be formulated as

$$\max_Z \lambda_1 \{R(z)\}$$

and NIPCA can be formulated as

$$\max_{Z \in \mathbb{R}^n} \sum_{s=1}^p \lambda_s \{R(z)\}$$

Are there other "PCA-like" criteria that could be optimized.

- criteria should be symmetric in the variables.

- concentrating variance on the largest eigenvalues is good. [shift from  $p \gg 2$  to nonlinear transformation]

Suppose

$$\mu(z) = \Phi \{R(z)\}$$

is such a criterion. We have a satisfying algorithm for the case in which

$\Phi$  is a convex function of  $R$ .

Special cases

- $\sum \tau_i$  HORST
- $\sum |\tau_i|^s$  KETTERUNG (s=2)
- $\det(R)$  NGV
- $\lambda_1 + \dots + \lambda_p$  PRINCIPALS etc

More satisfying  $\Phi$  in a Neoclassical sense) if

$\Phi$  is a norm on the eigenvalues of  $R$

- $\sum \tau_i^q$  not only for  $s=2$
- $\sum |\tau_i|^s$
- $\det(R)$
- $\lambda_1 + \dots + \lambda_p$
- $\sum |\lambda_i|^r$

In the meantime our activities have created two serious problems

- a proliferation of MLPCA techniques
- the fact that homogeneity analysis with P21 [let alone mixed PRINCIPALS] does not really fit in, and leads to Data Production.

We try to remedy both effects with a Research study, which (somewhat unfortunately) leads to two new forms of MPCA.

### COMPSTAT - paper.

We can test study these theoretical results if we remember how CA linearised regressions

Thus, for the cross product of the two variables,

$$C_{12} Y_2 = C_{12} D_1 Y_1 \Rightarrow D_1^{-1} C_{12} Y_2 = C_{12} Y_1$$

$$C_{21} Y_1 = C_{12} D_2 Y_2 \Rightarrow D_2^{-1} C_{21} Y_1 = C_{12} Y_2$$

If  $m > 2$  then homogeneity analysis (or any other technique) in general cannot linearise all  $2 \binom{m}{2}$  bivariate regressions

### Example SMOO

Now let us suppose that the data is such that all regressions can be linearised

[a model at last!]. Then here we can

$Z_{11} \dots Z_{1m}$  and  $C_{12}$  such that

$$C_{12} Y_2 = (C_{12} D_1) Y_1 \quad (A, B, C).$$



FIG. 7.5: Regressies PRE en KEUS Fig. 7.6: Regressies KEUS en BYA  
(.69) (.43)

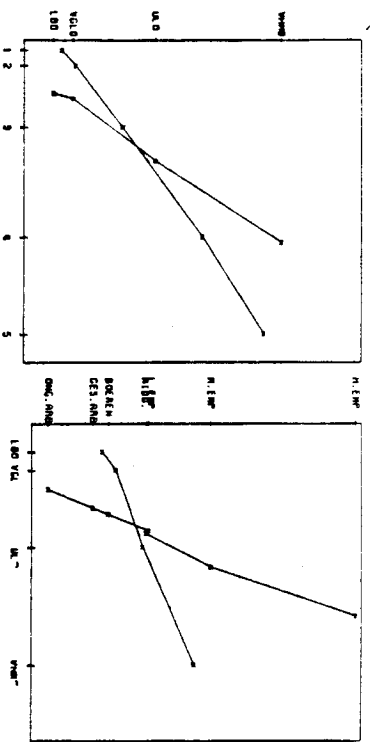
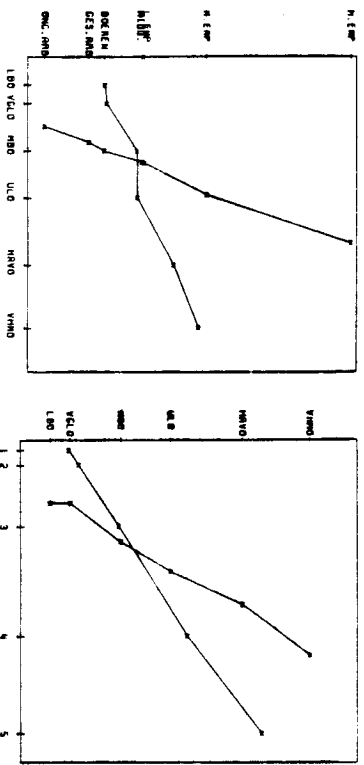


FIG. 7.7: Regressies EIN en BYA Fig. 7.8: Regressies PRE en EIN  
(.39) (.63)



Theorem : If such qualifications exist, homogeneity analysis will find them.

Proof :  $\sum_i g_i y_i \mathbf{e}_i = \sum_i x_i z_i \mathbf{D}_j \mathbf{z}_i = (\sum_i \alpha_i \rho_i \mathbf{e}_i) \mathbf{D}_j \mathbf{z}_i$

Now choose  $\alpha$  to be an eigenvector of  $R = f \rho_i \mathbf{e}_i$  with eigenvalue  $\lambda$ . Then

$$\sum_i \rho_i g_i (y_i \mathbf{e}_i) = \lambda \mathbf{D}_j (\alpha_i \mathbf{e}_i).$$

i.e.  $(\alpha_i \mathbf{e}_i), \alpha_2 \mathbf{e}_2, \dots, \alpha_m \mathbf{e}_m$  satisfies the equations of homogeneity analysis.  $\square$

Actually the construction above provides us with  $m$  solutions, because  $R$  has  $m$  eigenvectors, of the form  $\rho_i g_i \mathbf{z}_i$ . Collecting them in an  $k \times m$  matrix  $Y_i$  gives  $Y_i = \mathbf{z}_i \rho_i$ , i.e. there are single qualifications which all induce the same  $R$ .



Examples

- $m = 2$
- $k \hat{p} = 2$
- multinomial

Practical relevance

CMPSIAT paper [PREHOM]  
 [LINEARS]

Lineals  
 min  $\sum_i e [y_{ij}^i - r_{ij}^i]$

Prehom  
 take  $D^{-1/2} C D^{-1/2}$

minimize sum of "off-diagonal" elements by some a generalization of the Jacobi - method

Table 1 Separate classification by seven pathologists of most involved histological lesion of the uterine cervix

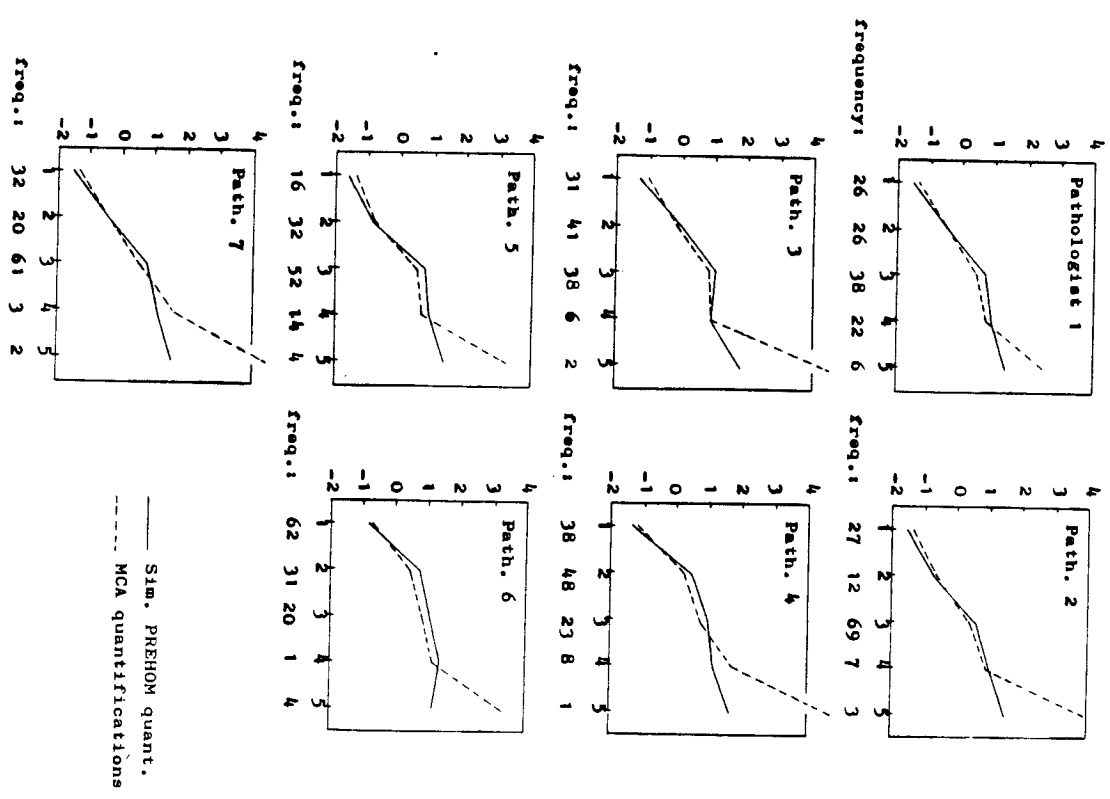
slide No.	pathologist							slide No.	pathologist						
	1	2	3	4	5	6	7		1	2	3	4	5	6	7
1	4	1	2	3	4	5	7	64	4	1	2	3	4	5	7
2	4	1	2	3	4	5	7	65	4	1	2	3	4	5	7
3	4	1	2	3	4	5	7	66	4	1	2	3	4	5	7
4	4	1	2	3	4	5	7	67	4	1	2	3	4	5	7
5	4	1	2	3	4	5	7	68	4	1	2	3	4	5	7
6	4	1	2	3	4	5	7	69	4	1	2	3	4	5	7
7	4	1	2	3	4	5	7	70	4	1	2	3	4	5	7
8	4	1	2	3	4	5	7	71	4	1	2	3	4	5	7
9	4	1	2	3	4	5	7	72	4	1	2	3	4	5	7
10	4	1	2	3	4	5	7	73	4	1	2	3	4	5	7
11	4	1	2	3	4	5	7	74	4	1	2	3	4	5	7
12	4	1	2	3	4	5	7	75	4	1	2	3	4	5	7
13	4	1	2	3	4	5	7	76	4	1	2	3	4	5	7
14	4	1	2	3	4	5	7	77	4	1	2	3	4	5	7
15	4	1	2	3	4	5	7	78	4	1	2	3	4	5	7
16	4	1	2	3	4	5	7	79	4	1	2	3	4	5	7
17	4	1	2	3	4	5	7	80	4	1	2	3	4	5	7
18	4	1	2	3	4	5	7	81	4	1	2	3	4	5	7
19	4	1	2	3	4	5	7	82	4	1	2	3	4	5	7
20	4	1	2	3	4	5	7	83	4	1	2	3	4	5	7
21	4	1	2	3	4	5	7	84	4	1	2	3	4	5	7
22	4	1	2	3	4	5	7	85	4	1	2	3	4	5	7
23	4	1	2	3	4	5	7	86	4	1	2	3	4	5	7
24	4	1	2	3	4	5	7	87	4	1	2	3	4	5	7
25	4	1	2	3	4	5	7	88	4	1	2	3	4	5	7
26	4	1	2	3	4	5	7	89	4	1	2	3	4	5	7
27	4	1	2	3	4	5	7	90	4	1	2	3	4	5	7
28	4	1	2	3	4	5	7	91	4	1	2	3	4	5	7
29	4	1	2	3	4	5	7	92	4	1	2	3	4	5	7
30	4	1	2	3	4	5	7	93	4	1	2	3	4	5	7
31	4	1	2	3	4	5	7	94	4	1	2	3	4	5	7
32	4	1	2	3	4	5	7	95	4	1	2	3	4	5	7
33	4	1	2	3	4	5	7	96	4	1	2	3	4	5	7
34	4	1	2	3	4	5	7	97	4	1	2	3	4	5	7
35	4	1	2	3	4	5	7	98	4	1	2	3	4	5	7
36	4	1	2	3	4	5	7	99	4	1	2	3	4	5	7
37	4	1	2	3	4	5	7	100	4	1	2	3	4	5	7
38	4	1	2	3	4	5	7	101	4	1	2	3	4	5	7
39	4	1	2	3	4	5	7	102	4	1	2	3	4	5	7
40	4	1	2	3	4	5	7	103	4	1	2	3	4	5	7
41	4	1	2	3	4	5	7	104	4	1	2	3	4	5	7
42	4	1	2	3	4	5	7	105	4	1	2	3	4	5	7
43	4	1	2	3	4	5	7	106	4	1	2	3	4	5	7
44	4	1	2	3	4	5	7	107	4	1	2	3	4	5	7
45	4	1	2	3	4	5	7	108	4	1	2	3	4	5	7
46	4	1	2	3	4	5	7	109	4	1	2	3	4	5	7
47	4	1	2	3	4	5	7	110	4	1	2	3	4	5	7
48	4	1	2	3	4	5	7	111	4	1	2	3	4	5	7
49	4	1	2	3	4	5	7	112	4	1	2	3	4	5	7
50	4	1	2	3	4	5	7	113	4	1	2	3	4	5	7
51	4	1	2	3	4	5	7	114	4	1	2	3	4	5	7
52	4	1	2	3	4	5	7	115	4	1	2	3	4	5	7
53	4	1	2	3	4	5	7	116	4	1	2	3	4	5	7
54	4	1	2	3	4	5	7	117	4	1	2	3	4	5	7
55	4	1	2	3	4	5	7	118	4	1	2	3	4	5	7
56	4	1	2	3	4	5	7	119	4	1	2	3	4	5	7
57	4	1	2	3	4	5	7	120	4	1	2	3	4	5	7
58	4	1	2	3	4	5	7	121	4	1	2	3	4	5	7
59	4	1	2	3	4	5	7	122	4	1	2	3	4	5	7
60	4	1	2	3	4	5	7	123	4	1	2	3	4	5	7
61	4	1	2	3	4	5	7	124	4	1	2	3	4	5	7
62	4	1	2	3	4	5	7	125	4	1	2	3	4	5	7
63	4	1	2	3	4	5	7	126	4	1	2	3	4	5	7

Statistica Neerlandica 36 (1982), nr. 2

PREHOM  
HOM.

R	Comp			
1	1	5.50	1	3.56
2	1	5.24	2	5.14
3	1	2.50	3	2.75
4	1	2.21	4	2.10
3	2	1.53	5	1.68
4	2	1.15	6	1.42
3	3	1.06	7	1.17
3	3	0.86	8	0.89
5	4	0.81	9	0.81
4	4	0.72	10	0.71
2	5	0.63	11	0.62
3	5	0.63	12	0.62
3	6	0.62	13	0.61
4	6	0.61	14	0.60
4	6	0.61	15	0.60
2	7	0.59	16	0.58
3	7	0.57	17	0.56
3	7	0.57	18	0.56
2	8	0.52	19	0.51
2	8	0.52	20	0.51
1	9	0.46	21	0.45
1	9	0.46	22	0.45
1	9	0.46	23	0.45
1	9	0.46	24	0.45
1	9	0.46	25	0.45
1	9	0.46	26	0.45
1	9	0.46	27	0.45
1	9	0.46	28	0.45
1	9	0.46	29	0.45
1	9	0.46	30	0.45
1	9	0.46	31	0.45
1	9	0.46	32	0.45
1	9	0.46	33	0.45
1	9	0.46	34	0.45
1	9	0.46	35	0.45
1	9	0.46	36	0.45
1	9	0.46	37	0.45
1	9	0.46	38	0.45
1	9	0.46	39	0.45
1	9	0.46	40	0.45
1	9	0.46	41	0.45
1	9	0.46	42	0.45
1	9	0.46	43	0.45
1	9	0.46	44	0.45
1	9	0.46	45	0.45
1	9	0.46	46	0.45
1	9	0.46	47	0.45
1	9	0.46	48	0.45
1	9	0.46	49	0.45
1	9	0.46	50	0.45

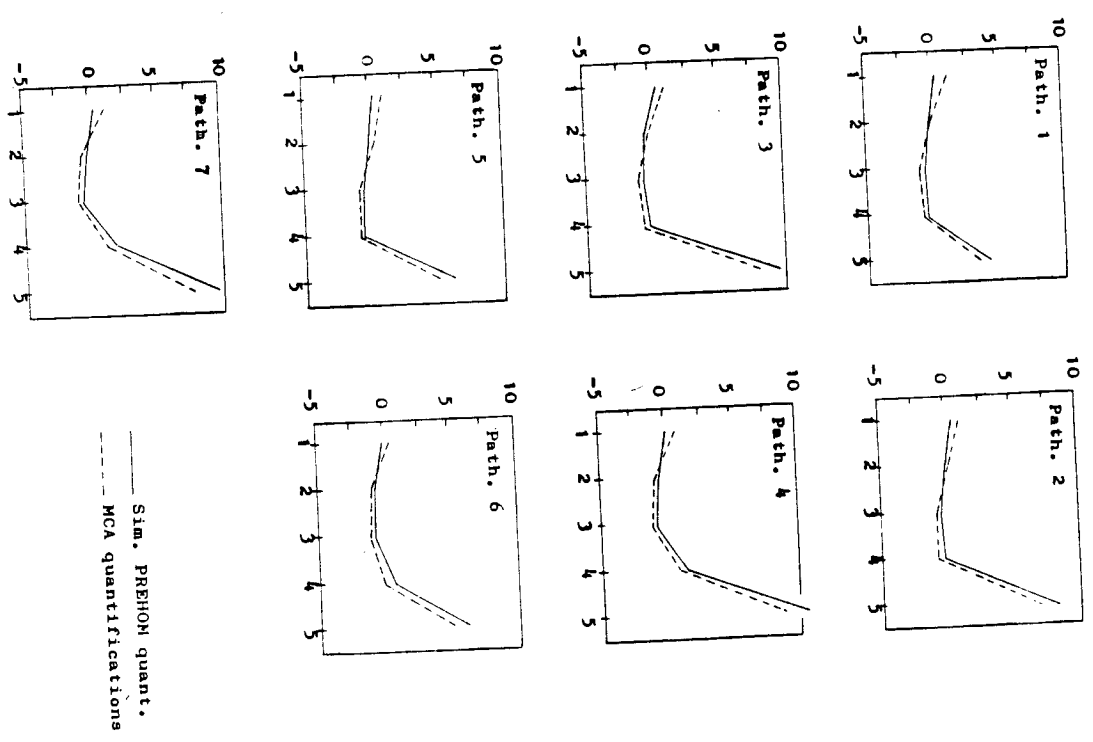
Figure 5.1: Category quantifications of the first dimension



98 17

Figure 5.2: Category quantifications of the second dimension

118



Example last  
 no 9

112

SOC	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
LBO	0.59	1.00	0.10	0.24	0.54	0.29	0.15	0.65	1.00
WRK	0.26	0.10	1.00	0.22	0.26	0.20	0.15	0.65	0.18
BIB	-0.15	-0.15	0.22	1.00	0.26	0.15	0.65	1.00	1.00
TYP	-0.28	-0.15	0.22	0.26	1.00	0.15	0.65	0.18	1.00
ONG	0.25	0.15	0.82	0.20	0.20	1.00	0.14	0.18	1.00
JGD	0.30	0.17	0.82	0.20	0.15	0.65	1.00	1.00	1.00
ABO	0.47	0.29	0.19	-0.10	-0.17	0.14	0.18	1.00	1.00

TABEL 3: Korrelaties ruwe data gediskretiseerd

SOC	0.50	-0.69
LBO	0.35	-0.62
WRK	0.90	0.17
BIB	0.30	0.64
TYP	0.25	0.72
ONG	0.86	0.20
JGD	0.86	0.08
ABO	0.37	-0.54

TABEL 4: Komponent-ladingen uit tabel 3

SOC	1.00	0.62	0.34	0.14	0.15	0.32	0.33	0.37	0.51
LBO	0.60	1.00	0.14	0.16	0.17	0.19	0.21	0.33	0.33
WRK	0.29	0.16	1.00	0.24	0.23	0.83	0.82	0.24	0.24
BIB	0.16	0.20	0.28	1.00	0.55	0.28	0.25	0.14	0.21
TYP	0.32	0.22	0.25	0.54	1.00	0.34	0.18	0.21	0.22
ONG	0.30	0.21	0.83	0.26	0.29	1.00	0.66	0.22	0.22
JGD	0.32	0.21	0.82	0.20	0.16	0.67	1.00	0.30	0.30
ABO	0.48	0.30	0.24	0.10	0.18	0.19	0.22	1.00	1.00

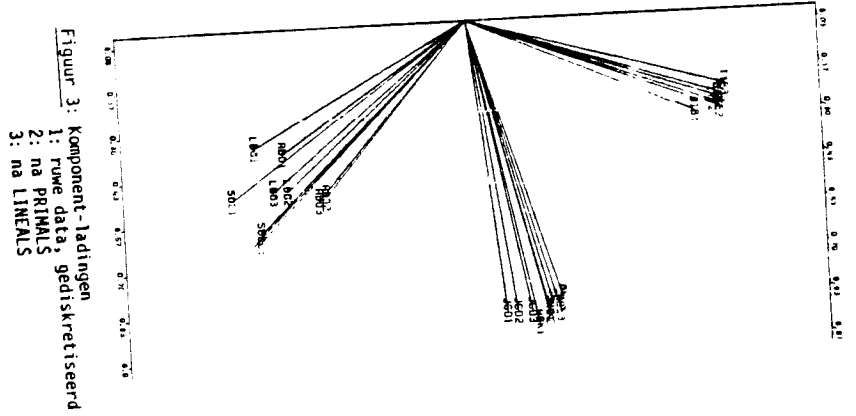
TABEL 5: Korrelatie-ratios ruwe data gediskretiseerd

119

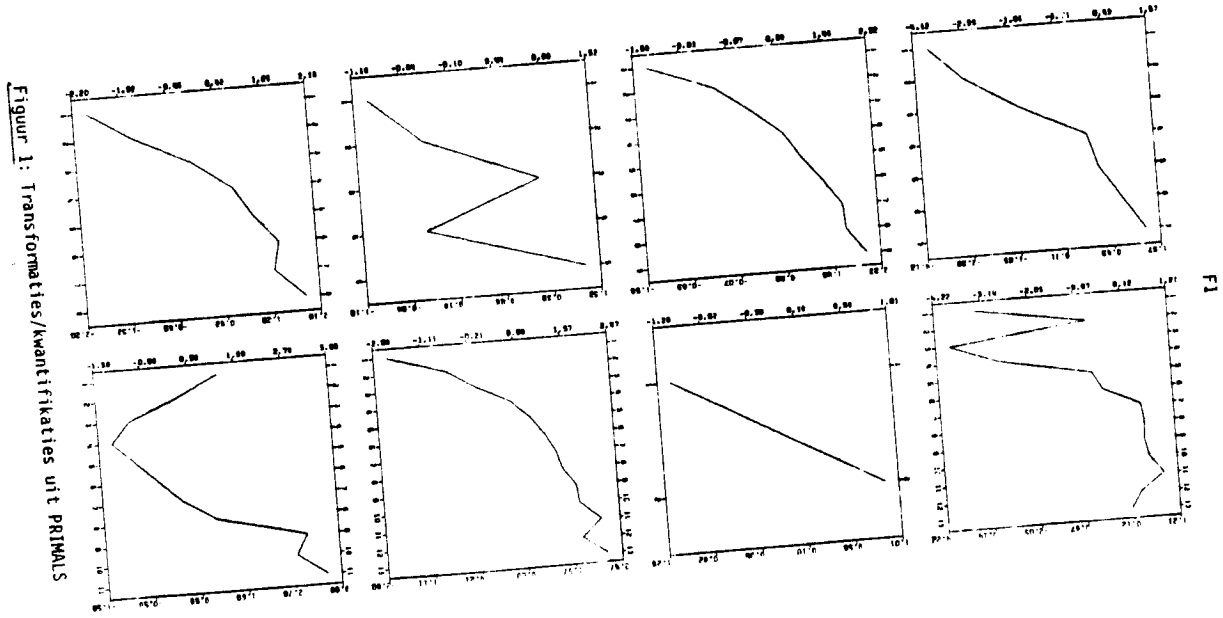


F3

122



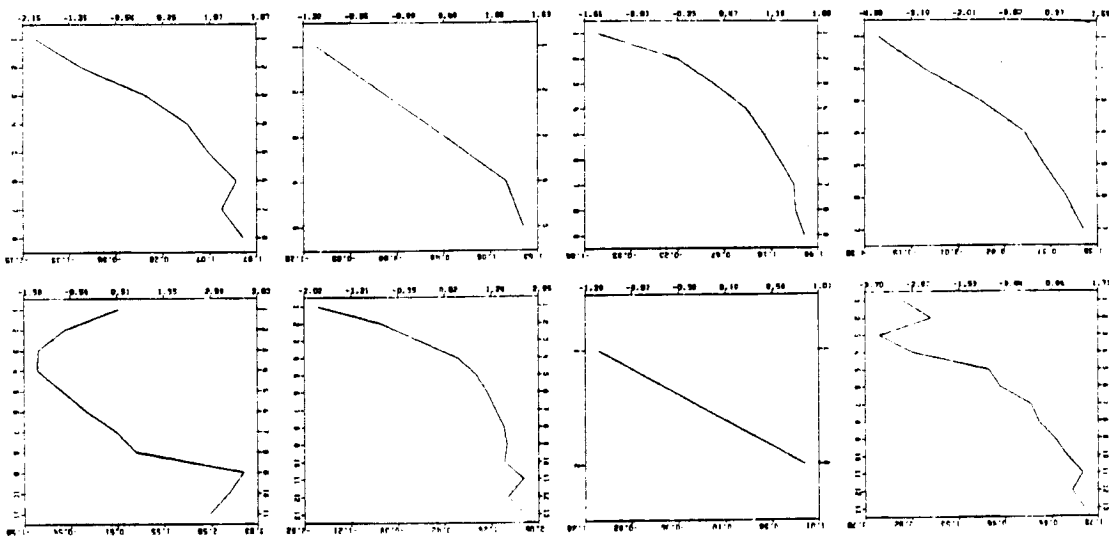
Figuur 3: Komponent-ladingen  
 1: ruwe data, gediskretiseerd  
 2: na PRIMALS  
 3: na LINEALS



Figuur 1: Transformaties/kwantifikasies uit PRIMALS

123

Figur 2: Transformaties/kwantificaties uit LINEALS

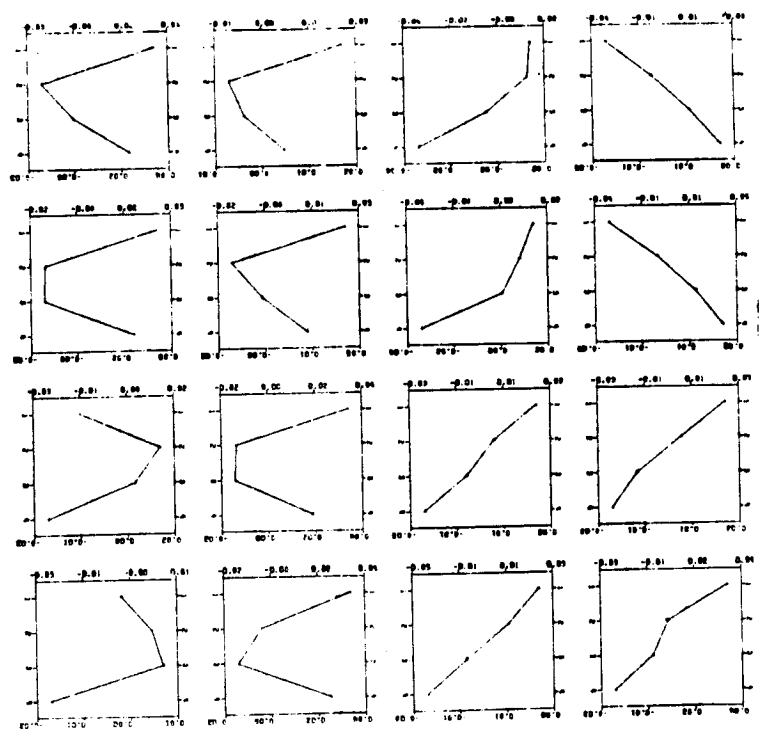


F2

124

MULTINOMIAL GAUSS

Figur 3: Four comparisons (rows), four variables (columns).



125





Of course adding restrictions can be combined with some restrictions and rank restrictions. This combination defines the SURFATS algorithm.

Start with

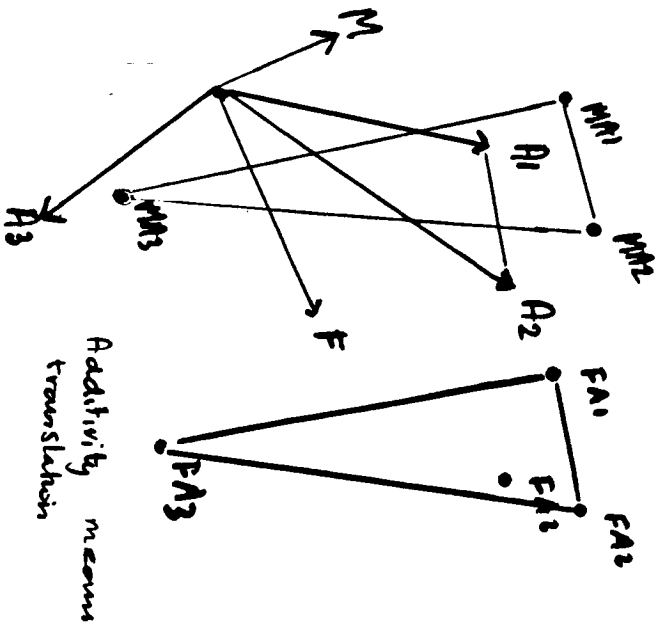
$$b(X;Y) = \frac{1}{n} \sum_{j=1}^n (X - G_j Y)' [X - G_j Y]$$

Now suppose there are additivity restrictions on all  $Y_j$ . We can write

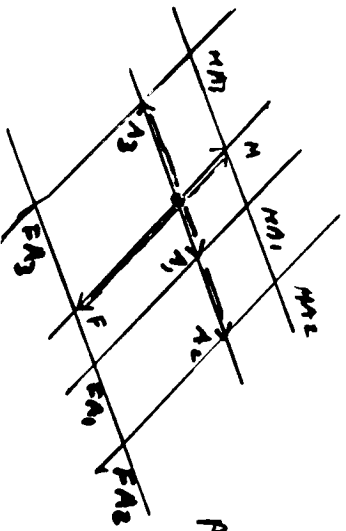
$$b(X;Y) = \frac{1}{n} \sum_{j=1}^n (X - \sum_s G_{js} Y_s)' [X - \sum_s G_{js} Y_s]$$

and then if course we can have  $Y_{js} = \sum_i a_{ij} a_i'$  for all or some  $(j,s)$ .

Geometry



Additivity means translation



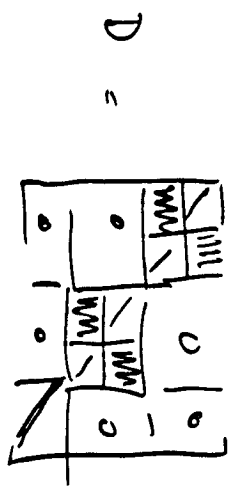
Additivity + rank means gnd.

Analytically

Using the same methods as before

$$t \quad Y_1 C Y \quad \max \quad Y_1 D Y = m I$$

where now



In the case of single quantifications we can again try to formulate the problem in terms of correlations, and a property of them that is optimized. But what if some variables are multiple?

To integrate multiple variables we need the notion of copies. Suppose  $G_j$  is the indicator of a ~~single~~ variable, and we use this same  $G_j$  more than once in the set, in all cases treating it as single. It will occur in the loss function in the form

$$G_j z_j a_j + \dots + G_j z_j m a_j + \dots = G_j Y_j \quad \text{with} \quad Y_j = z_j A_j$$

But if  $\mu \geq \min [p, k_j - 1]$  then any  $Y_j$  can be represented in the form  $Y_j = z_j A_j$ . Thus multiple nominal  $\equiv$  p copies of a single variable in the set.

Actually the notion of copies is more fruitful still.

If we have  $q < p$  copies, we have rank- $q$  restrictions (between single and multiple). If we impose orthonormality restrictions on the  $Z_j$  of the copies we have multiple ordinal. We can also make the first copy linear and the second quadratic, etc.

It also means that if we have found a multiple  $Y_j$  [from some homogeneity analysis] we could look at decompositions  $Y_j = Z_j A_j$  which translate it to a set of single variables.

- use orth. polys. in  $Z_j$
- use the SVD  $Z_j' D_j Z_j = I$
- use  $Z_j = D_j^{-1/2}$

We can now make a familiar list

	number of sets	number of vars in set
PCA	1	1
HR	2	$m_1 > 1, m_2 = 1$
DISCR	2	$m_1 > 1, m_2 = 1$ [null]
CANCOR	2	$m_1 > 1, m_2 > 1$
MANOVA	2	$m_1 > 1, m_2 > 1$ [Def]
ANCOV	2	$m_1 > 1, m_2 = 1$
MANCOVA	3	[Def] $m_1 > 1, m_2 = 1$

Clearly we can add many "new" species across to this list.

All these techniques (as techniques) as special cases of homogeneity analysis. Of course here statistical properties depend on the additional assumptions we want to make.

Selected Illustrations

It will become obvious now, why there was so much emphasis on MCA/CA/PCA in this course.

MORALS (example Rut n<sub>2</sub> 29) } M<sub>1</sub> > 1  
 Wilson, Science, 64, 1920, 47-57. } M<sub>2</sub> = 1

~~MORALS~~ (example Rut n<sub>2</sub> 29)

CANALS (example Rut n<sub>2</sub> 30, in Giff.)

(example in paper —) (n<sub>2</sub> 29)

(physical example) (n<sub>2</sub> 31)

REDUCTION ANALYSES

ESRHEIS, PM 1984, 331-344

	r <sub>xy</sub>	r <sub>xz</sub>	r <sub>yz</sub>	b <sub>x</sub>	b <sub>y</sub>	R <sup>2</sup>
linear	-0.3804	0.8156	0.1534	1.0219	0.5421	0.3166
rational	-0.4008	0.8286	0.1663	1.0666	0.5938	0.9826
step 1	-0.2302	0.7248	0.2189	0.8186	0.4073	0.6825
step 2	-0.3847	0.8475	0.1406	1.0582	0.5477	0.9739
spline 1	-0.3891	0.8394	0.1663	1.0654	0.5809	0.9910
spline 2	-0.3884	0.8413	0.1655	1.0664	0.5797	0.9931

Table 11.2 Regression statistics for ten different analyses of Gibbs data.

a	temperature	16	39	10																
	pressure	34	17	13																
	density	28	21	16																
b	temperature	6	10	7	8	9	8	7	8	7	2	5	4	4	5					
	pressure	9	12	13	8	7	6	6	6	6	6	6	6	6	6					
	density	9	12	7	9	6	6	6	6	6	6	6	6	6	6					
c	temperature	0.2	16.2	35.0	12.8	0.8														
	pressure	2.9	29.4	19.2	11.0	2.6														
	density	3.4	24.7	21.1	14.0	1.8														
d	temperature	6.0	10.0	7.0	8.0	9.0	8.0	9.0	8.0	7.0	2.0	2.0	8.0	8.0	8.0					
	pressure	0.2	8.3	12.2	12.6	9.3	5.6	2.5	4.6	4.6	3.1	1.1	0.8							
	density	1.5	7.2	11.0	8.4	8.4	6.3	6.3	6.6	4.6	4.0	0.7								

Table 11.3 Marginals for:  
 a: step 1: step-functions crude.  
 b: step 2: step-functions fine.  
 c: spline 1: splines crude.  
 d: spline 2: splines fine.

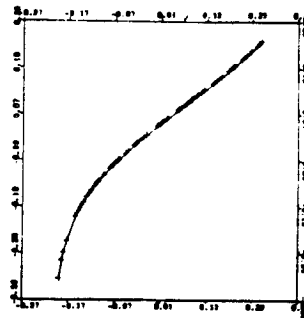
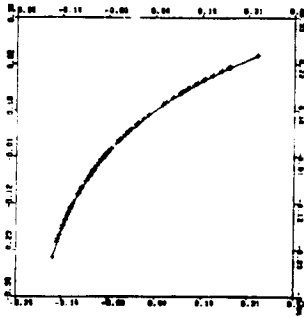
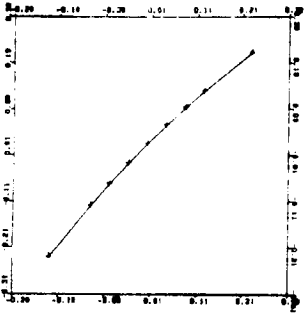


Figure 11:  $\lambda$  Cyclic Wilson example  
rational form functions.

136

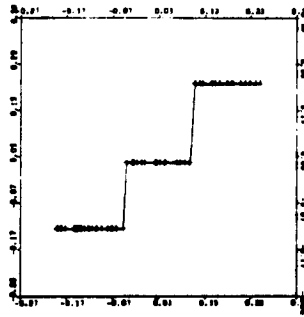
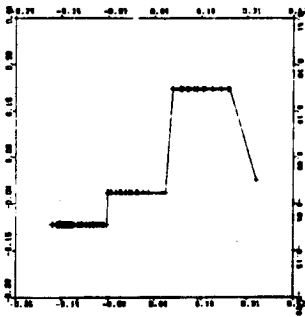
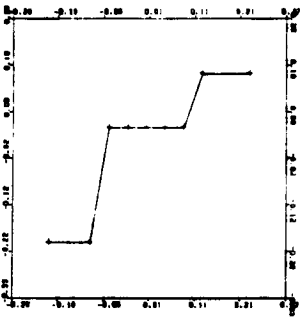


Figure 12:  $\lambda$  Cyclic Wilson example  
rational form functions.

298

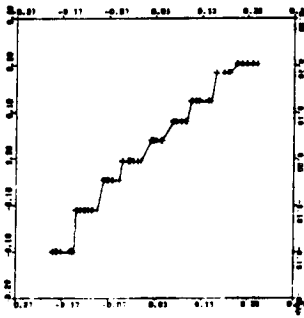
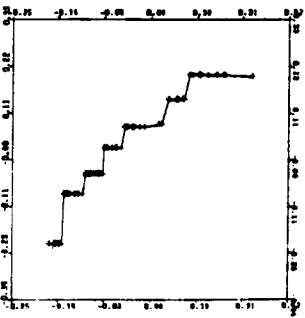
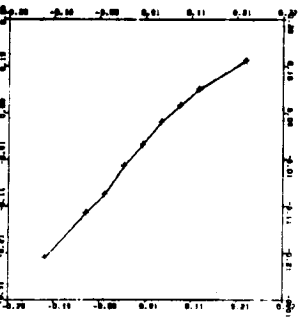


Figure 13:  $\lambda$  Cyclic Wilson example  
rational form functions.

137

254  
138

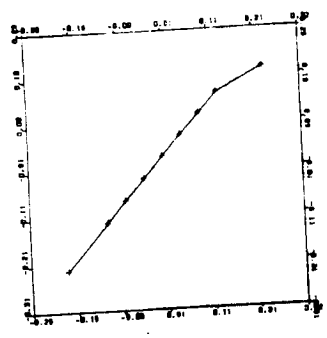
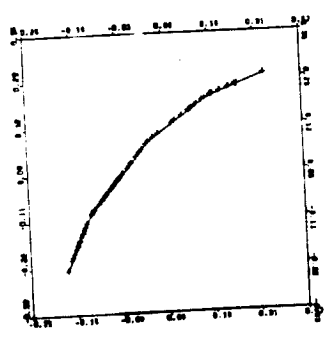
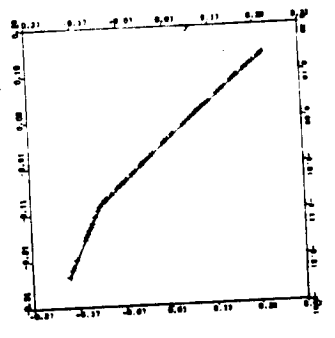


Figure 11.5 (auto-calculation example)  
from solution

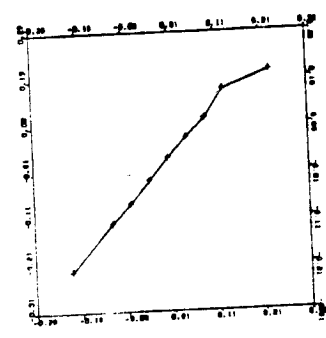
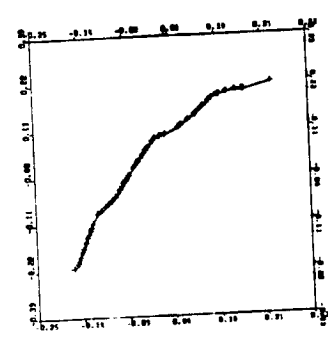
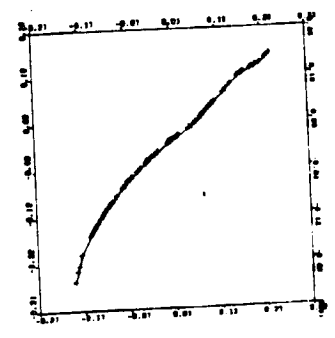
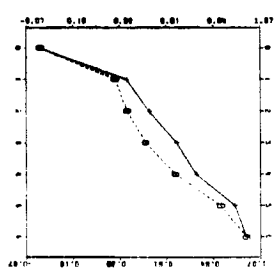
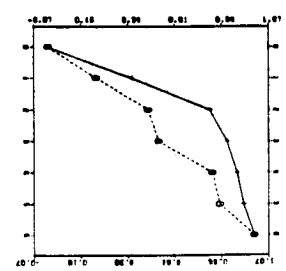
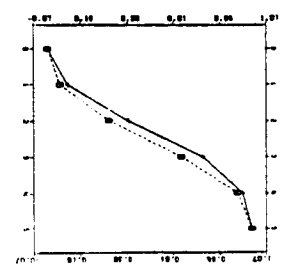
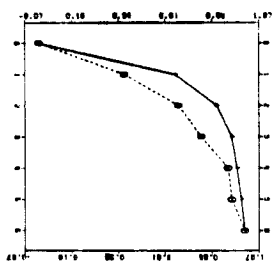
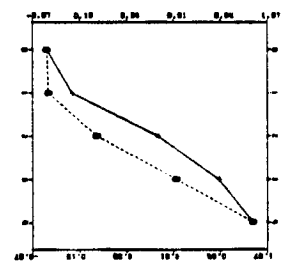
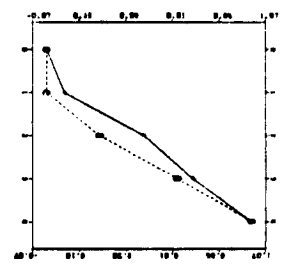
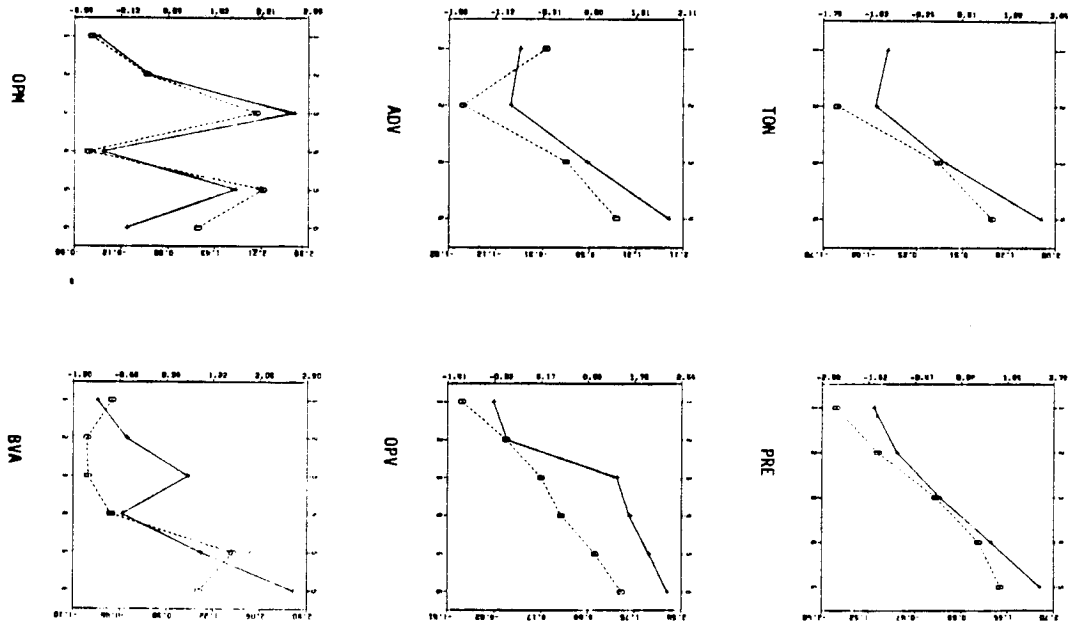


Figure 11.3 (auto-calculation example)  
from solution



Figur 1: kumulatieve verdelingen VJJ(+) en SWV(0).



figuur 2: CANALS-transformaties VUTJ(+) en SMVO(0).

tabel 3a: korrelaties VUTJ uit De Jong e.a.

	TON	PRE	ADV	OPV	OPM	BVA
TON	1.00					
PRE	.63	1.00				
ADV	.74	.68	1.00			
OPV	.37	.30	.34	1.00		
OPM	.24	.19	.24	.45	1.00	
BVA	.31	.28	.29	.60	.35	1.00

tabel 3b: korrelaties SMVO uit De Jong e.a.

	TON	PRE	ADV	OPV	OPM	BVA
TON	1.00					
PRE	.63	1.00				
ADV	.80	.67	1.00			
OPV	.37	.27	.36	1.00		
OPM	.36	.21	.31	.48	1.00	
BVA	.33	.26	.33	.59	.40	1.00

tabel 3b: korrelaties SMVO uit De Jong e.a.

tabel 4: beta-gewichten en multiple korrelaties uit tabel 3.

	VUTJ	SMVO
PRE	.22	.17
ADV	.55	.64
OPV	.09	.04
OPM	.01	.10
BVA	.03	.01
R <sup>2</sup>	.59	.67

tabel 4: beta-gewichten en multiple korrelaties uit tabel 3.



TON	PRE	ADV	OPV	OPM	BVA
1.00					
PRE	1.00				
ADV	.68	1.00			
OPV	.80	.69	1.00		
OPM	.44	.30	.36	1.00	
BVA	.27	.18	.22	.34	1.00
	.42	.31	.32	.56	.27
					1.00

tabel 5a: korrelaties VJTJ uit CANALS.

TON	PRE	ADV	OPV	OPM	BVA
1.00					
PRE	1.00				
ADV	.61	1.00			
OPV	.38	.27	.35	1.00	
OPM	.23	.14	.22	.26	1.00
BVA	.28	.20	.25	.47	.21
					1.00

tabel 5b: korrelaties SMVO uit CANALS.

VJTJ	SMVO
PRE	.21
ADV	.57
OPV	.10
OPM	.05
BVA	.10
R <sup>2</sup>	.74

tabel 6: beta-gewichten en multipele korrelaties uit CANALS.

TON	VJTJ	BS1	BS2	BS3	BS4	BS5
VGL0:	-0.74	-0.79	-0.76	-0.75	-0.67	-0.88
LBO:	-0.93	-0.89	-0.89	-0.83	-0.88	-0.98
ULO:	0.24	0.20	0.23	0.13	0.12	0.32
VHMO:	1.79	1.88	1.84	1.92	1.83	1.64
PRE: 1	-1.39	-1.32	-1.49	-1.32	-1.40	-1.40
2	-0.89	-0.86	-0.89	-0.82	-0.90	-0.86
3	0.09	-0.01	0.21	0.06	0.13	0.04
4	1.27	1.42	1.31	1.30	1.14	1.24
5	2.34	2.26	2.02	2.43	2.58	2.34
ADV: VGL0:	-0.69	-0.72	-0.64	-0.70	-0.62	-0.95
LBO:	-0.86	-0.80	-0.83	-0.76	-0.79	-0.89
ULO:	0.45	0.39	0.41	0.25	0.25	0.60
VHMO:	1.85	1.99	1.94	2.04	1.93	1.60
OPV: LO	-0.64	-0.83	-0.49	-0.49	-0.58	-0.71
LBO:	-0.38	-0.09	-0.33	-0.55	-0.40	-0.19
ULO:	1.23	1.51	1.12	1.61	1.56	0.23
MBO:	1.72	0.31	0.65	1.15	2.20	2.60
VHMO:	2.20	2.20	0.25	2.61	3.02	2.27
HBO:	2.42	2.68	3.93	2.22	1.65	1.98
OPM: LO	-0.53	-0.41	-0.37	-0.52	-0.37	-0.47
LBO:	0.36	0.01	-0.28	0.03	0.61	0.08
ULO:	2.72	0.53	2.97	2.36	2.93	2.39
MBO:	0.32	4.11	2.11	2.32	-2.35	-0.21
VHMO:	1.83	3.34	2.11	2.85	-0.27	3.43
HBO:	0.09	2.34	-0.09	1.95	-1.21	1.23
BVA: 1	-0.90	-0.88	-0.83	-0.88	-0.76	-0.93
2	-0.35	-0.03	0.01	-0.37	-0.48	-0.27
3	0.96	1.08	0.63	1.73	1.20	0.71
4	-0.40	-0.60	-1.04	-0.12	-0.17	-0.40
5	1.00	1.30	1.54	0.83	0.16	1.41
6	2.43	1.54	1.48	1.34	2.94	1.79
B	PRE	.24	.26	.21	.22	.20
ADV	.69	.68	.70	.72	.74	.68
OPV	.13	.13	.14	.12	.12	.10
OPM	.06	.06	.08	.07	.07	.07
BVA	.11	.11	.13	.07	.12	.14
Rmult	.84	.84	.85	.86	.86	.85

tabel 7: CANALS analyse op VJTJ en op vijf bootstrap-samples uit VJTJ.

	GINI	FARM	RENT	CMPR	LA	IMST	ECK	DEMT	DEB
ARGENTINA	86.3	98.2	22.9	374	23	13.6	57	217	2
AUSTRALIA	82.9	99.2	22.2	1212	16	11.3	0	0	1
BELGIUM	58.7	85.8	62.3	1012	16	11.3	0	0	1
BOLIVIA	93.4	97.7	20.0	68	72	12.5	8	463	1
BRASIL	81.7	98.5	9.1	262	61	15.5	49	0	1
CANADA	49.7	82.9	7.2	1667	12	11.3	22	0	1
CHILE	93.8	99.7	13.4	180	30	14.2	21	2	2
COLOMBIA	88.9	98.1	12.1	330	55	14.6	47	316	2
CUBA	88.1	99.1	5.4	307	55	14.6	19	24	2
CYPRUS	88.1	99.1	5.4	307	55	14.6	19	24	2
DEMOCRACY REF.	78.5	98.5	20.5	202	62	14.6	100	2900	3
EGYPT	86.4	99.3	14.4	204	53	15.3	4	18	3
FINLAND	74.0	98.1	11.6	133	64	15.4	41	0	1
FRANCE	82.8	98.8	15.1	244	63	15.1	9	2	3
GERMANY	56.4	88.2	37.3	201	59	14.0	15	292	2
HONGKONG	59.9	86.3	2.4	941	46	15.6	4	0	2
INDIA	80.3	98.0	23.8	442	29	15.5	51	1	2
INDONESIA	43.7	79.8	0.0	297	67	0.0	9	0	3
ITALY	82.8	87.7	18.8	1184	23	12.8	8	0	1
JAPAN	60.8	82.0	8.5	788	75	14.8	2	0	3
KOREA	40.8	82.0	8.5	788	75	14.8	2	0	3
NETHERLANDS	75.7	96.2	51.3	208	11	13.6	16	163	1
NICARAGUA	77.3	95.5	22.3	1258	16	12.8	1	0	1
NORWAY	66.9	87.5	7.5	969	26	12.8	1	0	1
PAKISTAN	73.7	95.0	12.3	350	54	15.6	29	35	3
PERU	87.5	96.9	0.0	140	60	14.6	23	26	3
POLAND	45.0	77.7	0.0	468	57	8.5	19	5	3
ROMANIA	81.0	93.5	43.7	254	50	0.0	22	1	3
RUSSIA	81.7	96.1	59.0	102	50	0.0	3	0	3
SPAIN	81.7	96.1	59.0	102	50	0.0	3	0	3
UNITED STATES	90.9	99.3	20.6	782	27	14.6	16	11	1
WEST GERMANY	67.4	93.0	5.7	762	14	3.0	4	0	2
SOUTH VIETNAM	57.7	87.2	18.9	1165	13	0.5	0	1000	3
SWEDEN	49.8	81.5	18.9	1229	10	0.5	0	0	1
SWITZERLAND	49.8	81.5	18.9	1229	10	0.5	0	0	1

Table 7.2.a Original data

	ECON. VAR	POL. VAR
GINI	5	5
FARM	4	4
RENT	4	4
CMPR	7	1
LA	3	6
IMST	5	6
ECK	1	2
DEMT	2	3
DEB	3	3
ARGENTINA	1	1
AUSTRALIA	1	1
BELGIUM	1	1
BOLIVIA	1	1
BRASIL	1	1
CANADA	1	1
CHILE	1	1
COLOMBIA	1	1
CUBA	1	1
CYPRUS	1	1
DEMOCRACY REF.	1	1
EGYPT	1	1
FINLAND	1	1
FRANCE	1	1
GERMANY	1	1
HONGKONG	1	1
INDIA	1	1
INDONESIA	1	1
ITALY	1	1
JAPAN	1	1
KOREA	1	1
NETHERLANDS	1	1
NICARAGUA	1	1
NORWAY	1	1
PAKISTAN	1	1
PERU	1	1
POLAND	1	1
ROMANIA	1	1
RUSSIA	1	1
SPAIN	1	1
UNITED STATES	1	1
WEST GERMANY	1	1
SOUTH VIETNAM	1	1
SWEDEN	1	1
SWITZERLAND	1	1

Table 7.2.b Discretized data

144

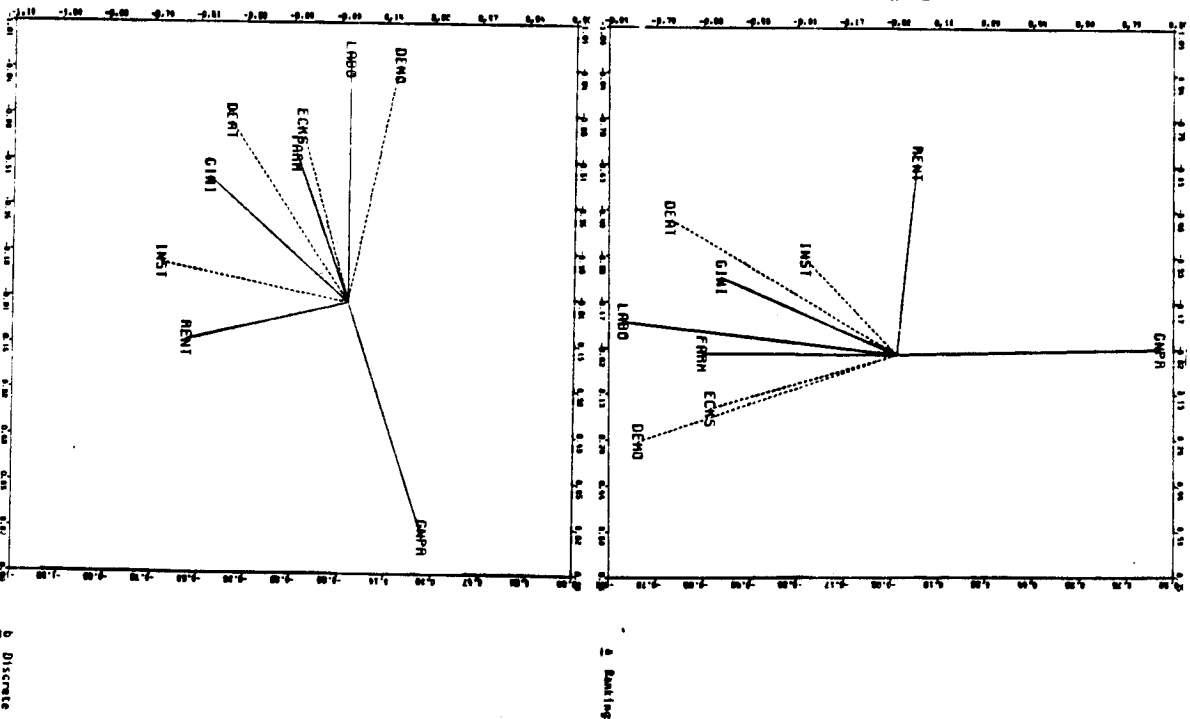
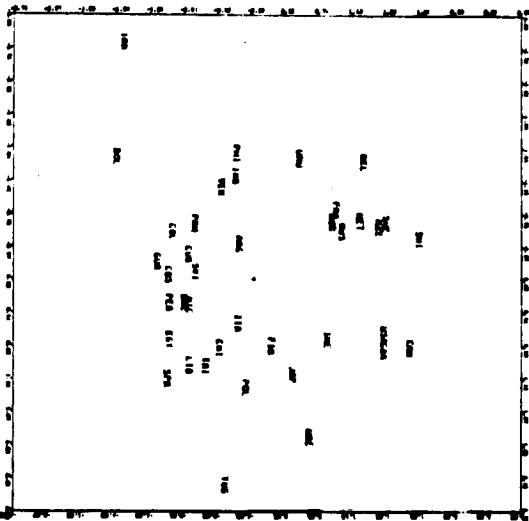
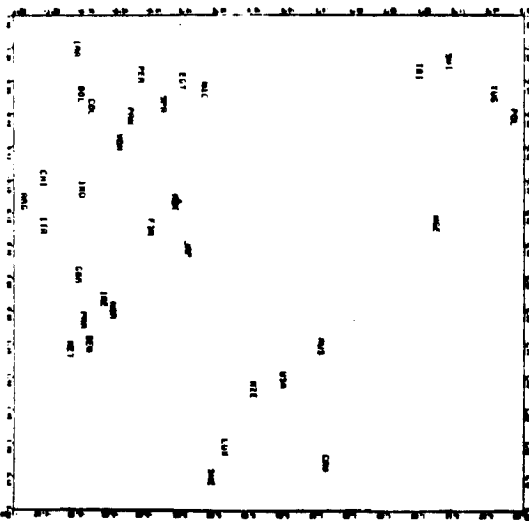


Figure 7.2 Correlations in the canonical space of economic variables

145

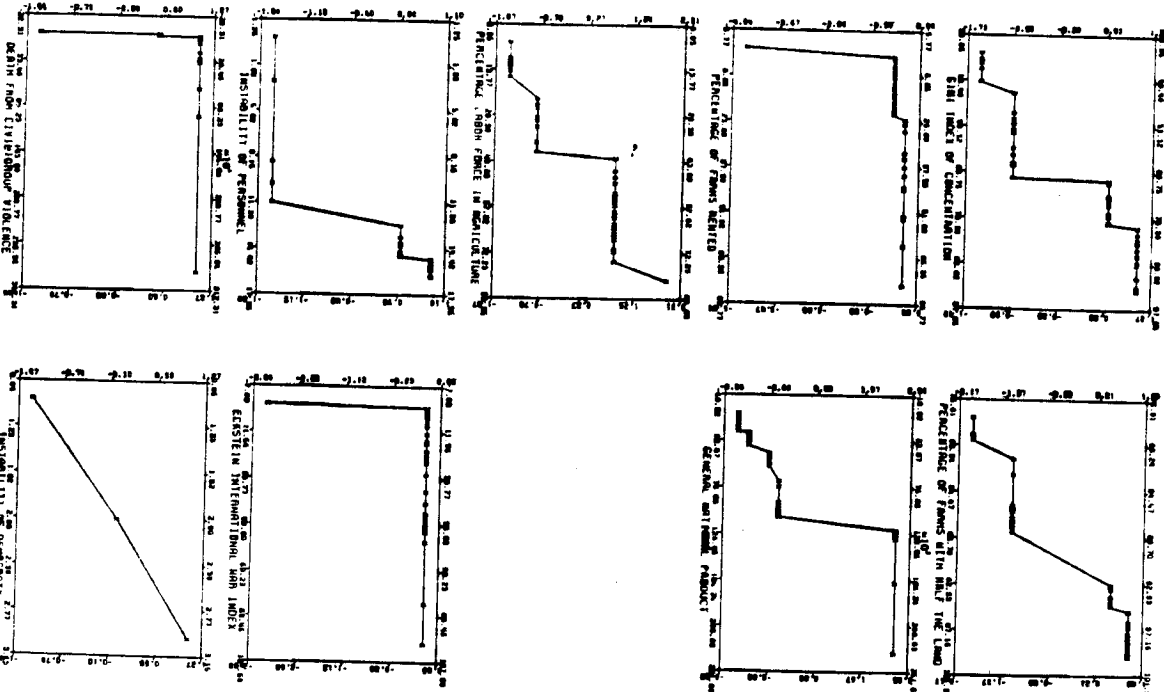


2. Montenegro



2. Discrete

Figure 2.3 Canonical scores in the space of economic variables



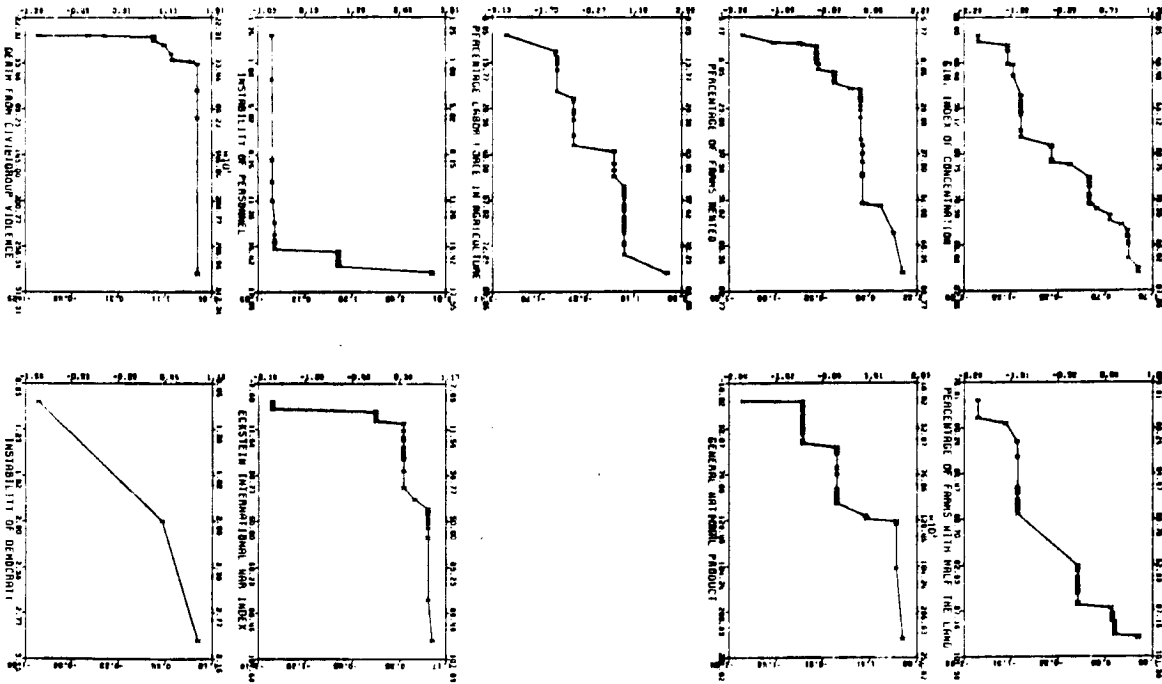


Figure 7.4.8 Transformations of the rankings

ORDERING OF MEFV-CURVES RELATIVE TO RESPIRATORY SYMPTOMS

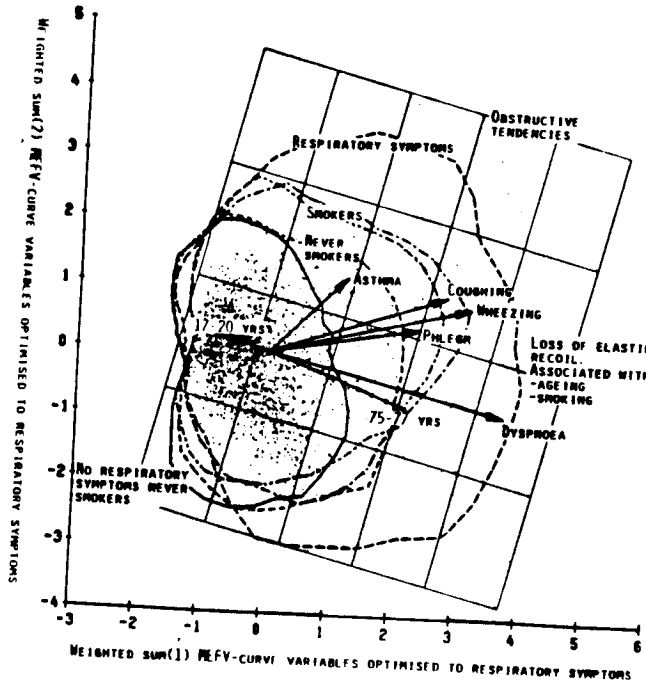


Figure 1. SCATTERGRAM

- Most parametric ordering of MEFV-curves with respect to chronic respiratory symptoms. The ordering is expressed in form of a specific weighted sum of the MEFV-curve variables (abscissa and ordinate). In this way each MEFV-curve gets its position in the scattergram.
- ARRIMS
  - The position of a curve in the scattergram is related to the probability on chronic respiratory symptoms. These increase in the direction of the full arrows.
  - Curves having a high probability for ASTHMA are more upwards and somewhat to the right; indications for DYSPNOEA mainly to the right; lungfunctions are found at the bottom left.
- CONTOURS
  - Subgroups (having identical age distributions) indicated by way of percentile contours contain 95% of their values.
  - SMOKERS (---) occupy an intermediate position between the subgroup NO RESPIRATORY SYMPTOMS (---) and RESPIRATORY SYMPTOMS (---). They contain more curves with a larger probability for one or more respiratory symptoms than NO SMOKERS (---).
  - KNOTTED LINE (---).
  - The knotted line for age shows that with increasing age, curves are more probably associated with respiratory symptoms, especially DYSPNOEA.
- INTERPRETATION
  - The association with smoking and especially aging in direction of DYSPNOEA, suggests that loss of elastic recoil may be involved in the change in configuration towards the right.
  - Relatively independent of the former trend, direction of the vector ASTHMA suggests that in more vertical direction obstructive tendency increases from other causes than loss of elastic recoil.

149

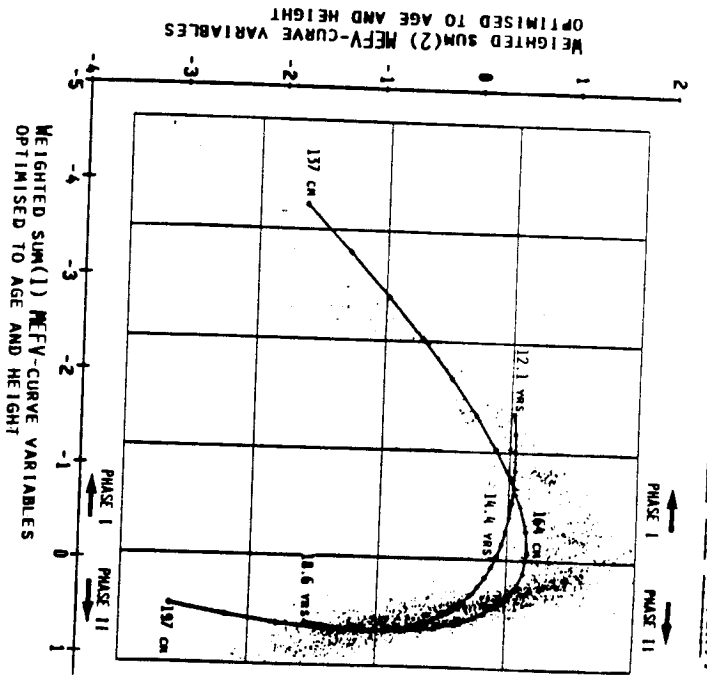


Figure 1  
SCATTERGRAM  
Most parsimonious ordering of MEFV-curves with respect to age and height.  
The ordering is expressed in terms of 2 specific weighted sums of the MEFV-curve variables. In this way each MEFV-curve gets its position in the scattergram.  
PHASES  
The scattergram has a boomerang-like shape suggesting two phases.  
Phase 1 covers mostly the age range 12-14.5 yrs. Curve variability is mainly related to differences in height, mostly 137-164 cm.  
Phase 2 concerns a larger age range mostly 14.5-18 yrs. Curve variability is also mainly related to differences in height, mostly 165-197cm.

of nomination and renomination on internal party processes could give central party organs a strong opinion with which to disagree. The electorate gives a mandate to a party and to make what members of the parliament do but members cannot bring their own to individual constituents but must satisfy themselves with such bearing as they can get within the party or leave.

5.1. Description of the data

In 1972 the members of the Dutch Parliament (MPs) were interviewed. Among other things, the MPs gave their opinions on a number of issues and their preference votes for the political parties. The issues concerned development and abortion, law and order, income differences, worker participation, taxation and defence. The opinions were measured on a nine-point scale of which the lowest and the highest category were described (Table 1). The party preferences were recorded in a table of rank

Table 1. The issues and the meanings of the lowest and the highest category

1 DEVELOPMENT AID	(1) the government should spend less money on aid to developing countries	(9) the government should spend less money on aid to developing countries
2 ABORTION	(1) the government should prohibit abortion completely	(9) a woman has the right to decide for herself about abortion
3 LAW AND ORDER	(1) the government takes too strong action against public disturbances	(9) the government should take stronger action against public disturbances
4 INCOME DIFFERENCES	(1) income differences should remain as they are	(9) income differences should become much less
5 PARTICIPATION	(1) only management should decide important matters in industry	(9) workers must also have participation in decisions important for industry
6 TAXATION	(1) taxes should be increased for general welfare	(9) taxes should be decreased so that people can decide for themselves how to spend their money
7 DEFENCE	(1) the government should insist on shrinking the Western armaments	(9) the government should insist on maintaining strong Western armaments

orders. The scores in this table tell us the rank order each member of the parliament gave to the different parties (2 = highest preference, 15 = lowest preference). The lowest score (2) was always used for the MP's own party. For our illustration we only consider the preferences for the four largest parties, which are:

- PvdA - labour party (socialists) (39)
- ARP - Anti Revolutionary Party (christian democrats) (13)
- KVP - catholic party (christian democrats) (35)
- VVD - liberal party, referred to as conservatives (16)

The figure in parentheses is the number of MPs. The other parties in 1972 were:

- CHU - Christian Historical Union (christian democrats) (10)
- D'66 - democrats (66) (liberals) (11)

	GINI	FARM	RENT	GNPR	LA	INST	ECR	DEBT	DEPR	ECON	VAR	POL	VAR
ARGENTINA	86.3	98.2	32.9	374	25	13.6	37	217	2	5	4	4	2
AUSTRALIA	92.9	99.6	---	1215	14	11.3	0	0	1	6	5	9	7
BELGIUM	58.7	85.0	62.3	1015	12	15.5	8	1	1	2	3	5	6
BOLIVIA	93.8	97.7	20.0	68	72	11.3	53	463	3	6	2	1	4
BRASIL	83.7	98.5	7.3	287	61	13.3	29	1	1	7	2	3	4
BURUNDI	91.8	94.7	13.4	180	30	14.2	21	2	2	6	5	3	2
CHINA	84.9	98.1	12.1	330	55	14.6	47	316	2	5	2	4	3
COLUMBIA	80.1	99.1	5.4	307	55	14.6	19	242	2	5	2	4	3
COSTA RICA	79.2	97.0	53.8	361	42	13.6	100	3900	3	4	5	6	2
CUBA	45.8	79.3	3.5	913	23	14.6	0	0	1	1	2	3	4
DENMARK	78.5	98.5	20.8	205	53	11.3	6	31	3	5	3	3	3
DOMINICAN REP.	86.4	99.3	14.6	204	53	11.3	4	18	3	5	3	3	3
EQUADOR	84.0	98.1	11.6	133	63	11.6	45	2	3	4	3	2	3
EGYPT	82.8	98.8	13.1	240	53	11.6	12	292	3	4	3	2	3
EL SALVADOR	59.9	86.3	2.4	841	46	15.6	4	0	2	2	3	2	3
FINLAND	58.3	86.1	26.0	1046	26	16.3	46	1	2	2	3	2	3
FRANCE	86.0	99.7	17.0	174	68	14.9	45	57	3	5	3	3	3
GUATEMALA	74.7	99.4	17.7	239	48	15.8	9	2	2	4	5	3	3
GREECE	71.0	93.4	44.5	988	5	13.6	12	111	3	4	5	2	1
HONGKONG	75.7	97.4	16.7	137	66	13.6	45	1	1	4	5	2	1
HONDURAS	59.8	85.9	2.5	509	40	14.2	9	0	1	2	3	2	3
IRELAND	52.2	86.9	53.0	172	71	2.0	93	348	2	2	3	1	4
INDIA	80.1	99.3	23.0	143	21	15.2	54	1	1	4	5	2	1
INDONESIA	47.0	81.5	2.9	240	40	15.7	22	1	2	1	2	3	2
ITALY	43.7	79.8	0.0	287	67	0.0	9	0	3	1	1	4	1
YUGOSLAVIA	63.8	87.7	18.8	1194	23	12.8	0	0	1	3	3	3	2
LUXEMBOURG	70.0	93.0	8.5	708	11	13.6	2	0	1	3	4	2	1
LIBYA	60.5	86.2	53.3	708	11	13.6	2	0	1	3	3	3	2
NETHERLANDS	75.7	96.4	---	254	68	12.8	16	16	3	4	5	4	1
MICROGASIA	77.3	95.5	22.3	1259	16	12.8	0	0	1	4	5	4	1
NEW ZEALAND	64.9	87.5	7.5	965	76	12.8	1	0	1	4	5	2	1
NORWAY	74.6	87.4	10.7	350	52	15.8	29	25	3	4	4	1	4
AUSTRIA	43.5	95.9	---	140	60	14.6	23	28	3	5	5	9	2
PANAMA	87.5	95.9	---	468	57	8.5	19	5	3	1	1	5	3
POLAND	45.0	77.7	0.0	468	57	8.5	19	5	3	1	1	5	3
SPAIN	78.0	99.5	43.7	254	50	0.0	22	3	3	4	5	4	3
TAIWAN	65.2	94.1	40.0	102	50	0.0	3	0	3	3	4	4	3
URUGUAY	81.7	96.6	34.7	569	37	14.6	1	1	1	5	5	3	2
VENEZUELA	90.9	99.3	20.6	762	42	14.9	36	118	3	6	5	3	6
UNITED STATES	70.5	95.4	20.4	243	10	12.8	22	0	2	3	3	2	3
WEST GERMANY	67.4	93.0	20.7	743	44	10.0	50	1000	3	4	4	3	3
SOUTH VIETNAM	87.1	87.5	18.9	1165	13	8.5	0	0	1	2	3	3	3
SWEDEN	49.6	81.5	18.9	1229	10	8.5	0	0	1	1	2	3	3
SWITZERLAND	49.6	81.5	18.9	1229	10	8.5	0	0	1	1	2	3	3

Table 7.2.a Original data

Table 7.2.b Discretized data

144

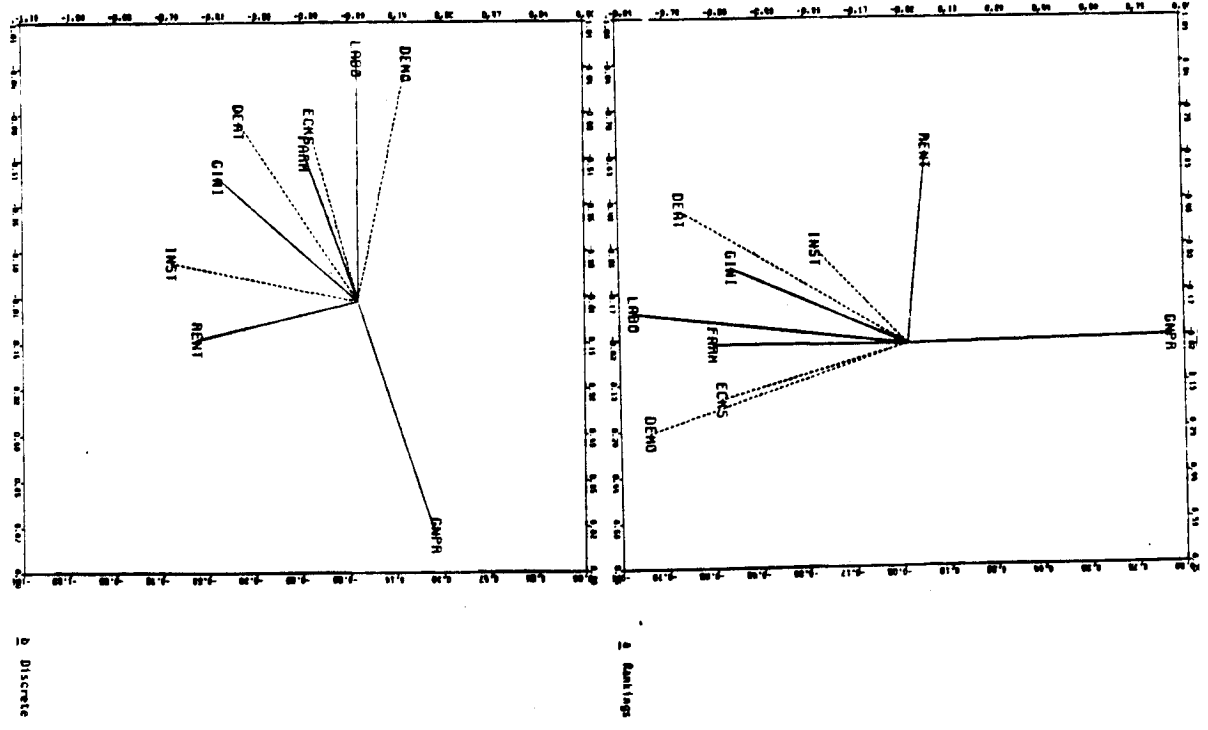


Figure 7.2 Correlations in the canonical space of economic variables

145

2 Discrete

28/146

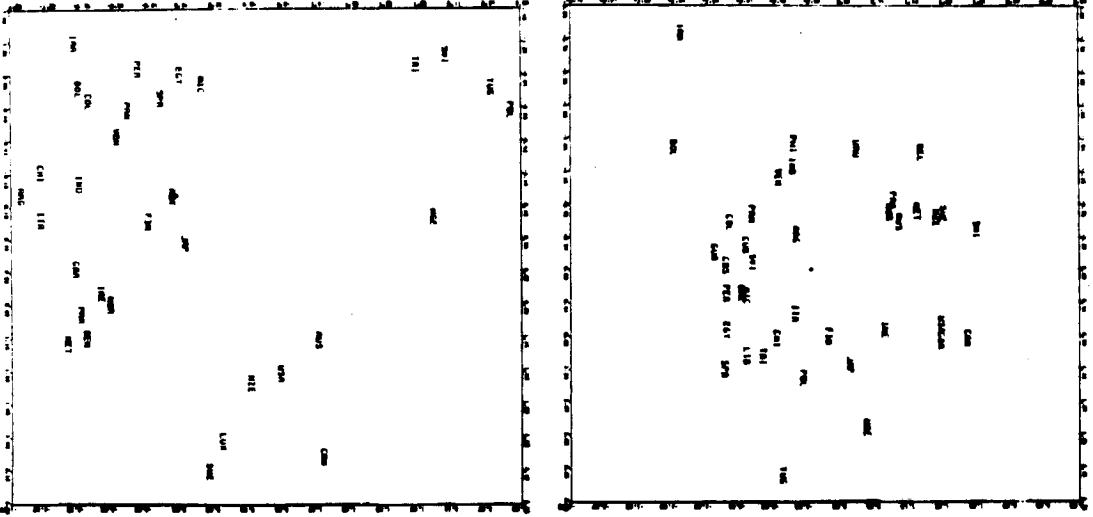
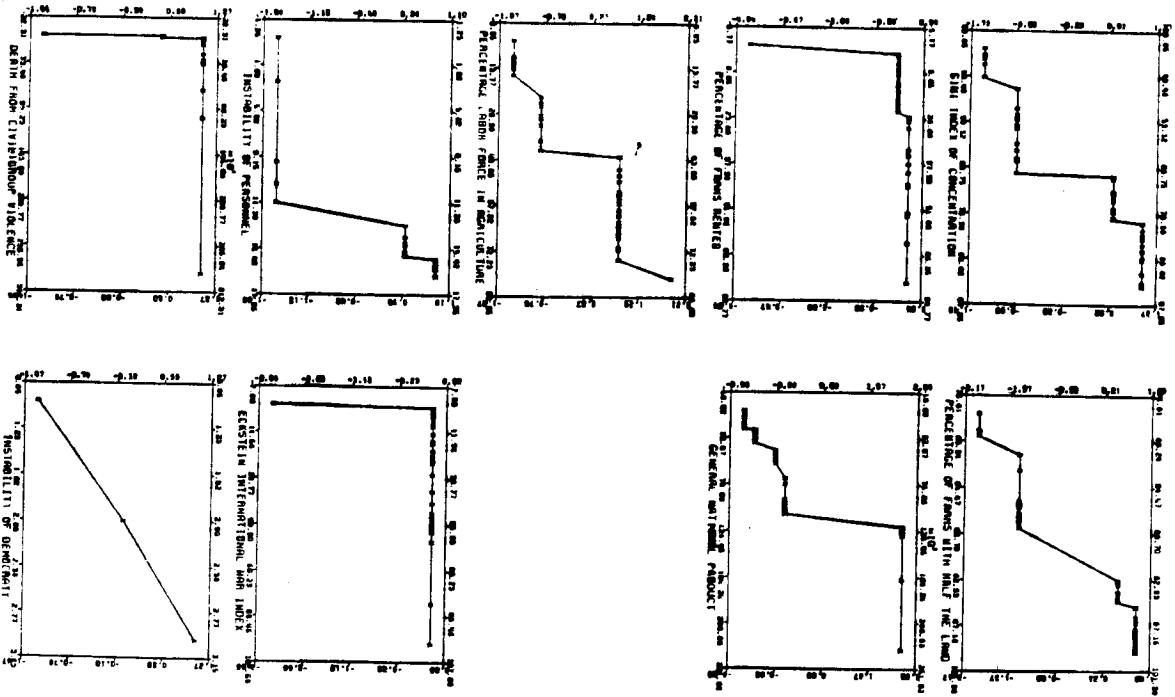


Figure 7.3 Canonical scores in the space of economic variables



147

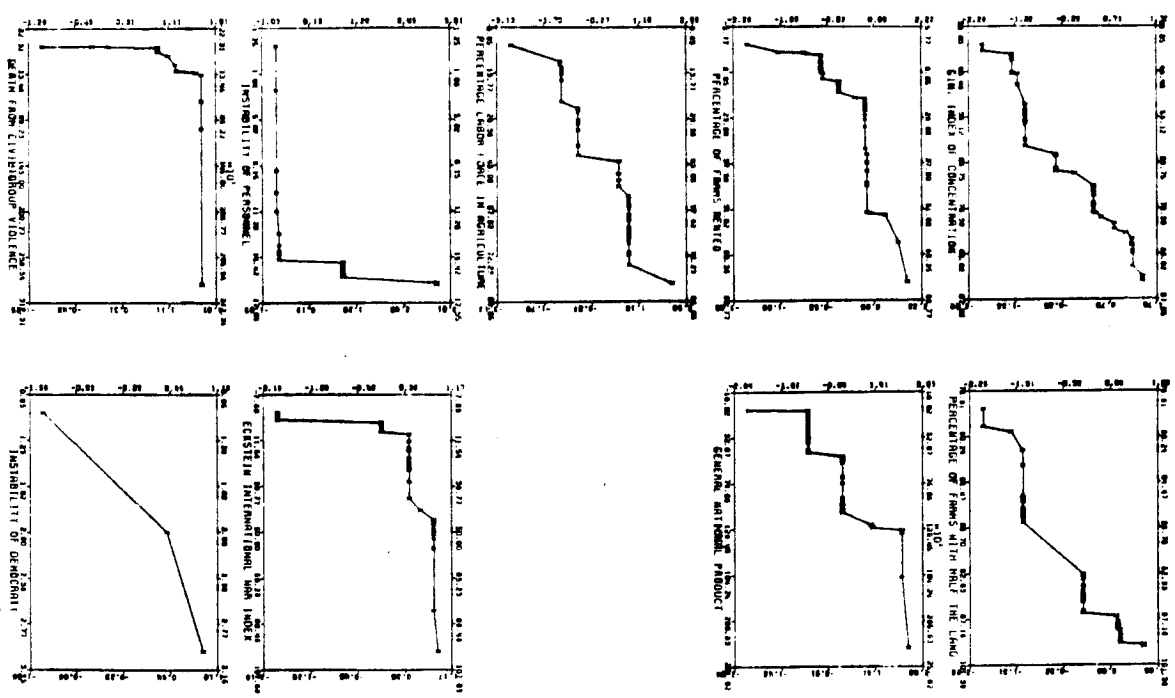


Figure 7.4.2 Transformations of the rankings

ORDERING OF MEFV-CURVES RELATIVE TO RESPIRATORY SYMPTOMS

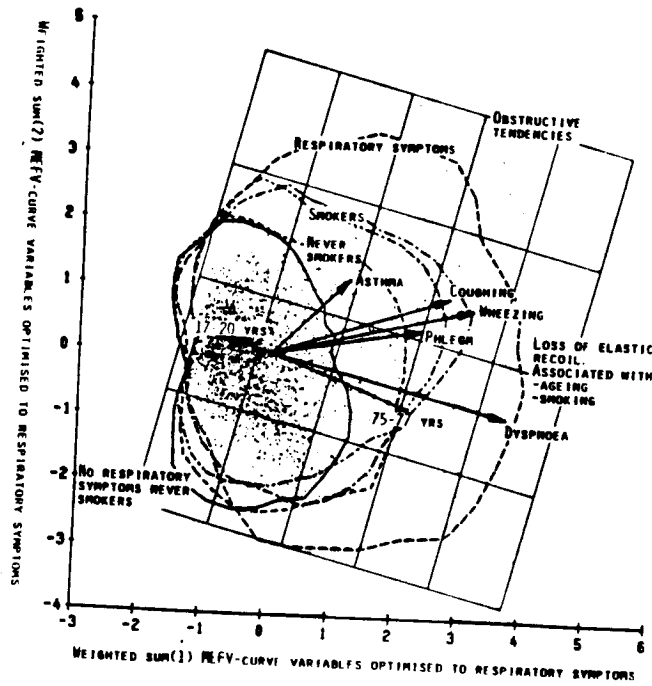


Figure 1. SCATTERGRAM

- Most judicious ordering of MEFV-curves with respect to chronic respiratory symptoms. The ordering is expressed in terms of a specific weighted sum of the MEFV-curve variables (abscissa and ordinate). In this way each MEFV-curve gets its position in the scattergram.

ARROWS

- The position of a curve in the scattergram is related to the probability of chronic respiratory symptoms. These increase in the direction of the two arrows.
- Curves having a high probability for ASTHMA are more upwards and somewhat to the right. Those indicative for DYSPNOEA mainly to the right. Lung functions are found at the bottom left.

CONTOURS

- Subgroups (having identical age distribution) indicated by way of percentile contours contain 95% of their values.
- SMOKERS (---) occupy an intermediate position between the subgroups NO RESPIRATORY SYMPTOMS (---) and RESPIRATORY SYMPTOMS (---). They contain more curves with a larger probability for one or more respiratory symptoms than NO SMOKERS (---).

KNOTTED LINE (---)

- The knotted line for age shows that with increasing age, curves are more probably associated with respiratory symptoms, especially DYSPNOEA. Independent of age, height is hardly related to curve positions.

INTERPRETATION

- The association with smoking and especially ageing in the direction of DYSPNOEA, suggests that loss of elastic recoil may be involved in the change in the configuration towards the right.
- Relatively independent of the former trend, direction of the vector to ASTHMA suggests that in more vertical direction obstructive tendency increases from other causes than loss of elastic recoil.



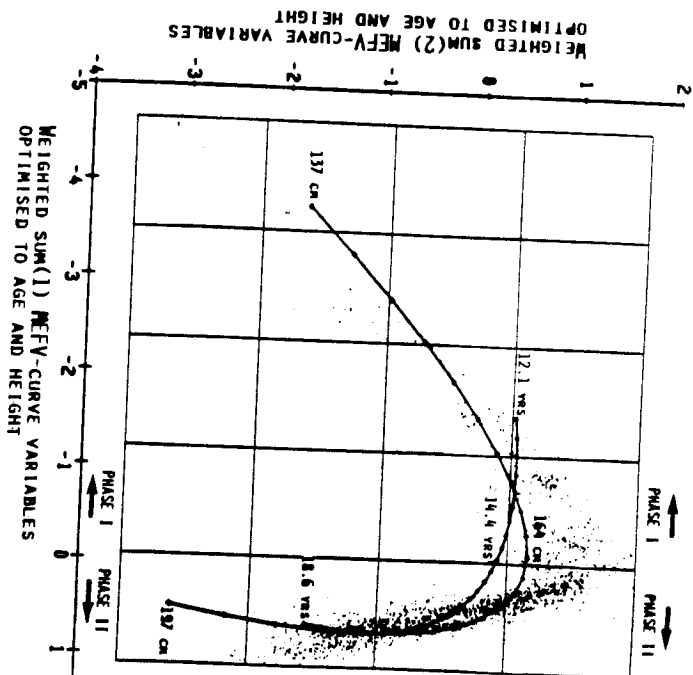


Figure 1  
SCATTERGRAM  
- Most parsimonious ordering of MEFV-curves with respect to age and height.  
The ordering is expressed in terms of 2 specific weighted sums of the MEFV-curve variables. In this way each MEFV-curve gets its position in the scattergram.  
- PHASES  
- The scattergram has a boomerang-like shape suggesting two phases.  
- Phase 1 covers mostly the age range 12-14.5 yrs. Curve variability is mainly related to differences in height, mostly 137-164 cm.  
- Phase 2 concerns a larger age range mostly 14.5-18 yrs. Curve variability is also mainly related to differences in height, mostly 165-197cm.

of nomination and re-nomination on internal party processes could give central party organs a strong weapon with which to discipline deviant behaviour. The electorate gives a mandate to a party, not to individual members of the parliament. Individual members cannot bring their own individual contribution but must submit themselves with such hearing as they can get within the party or house.

5.1 Description of the data

In 1972 the members of the Dutch Parliament (MPs) were interviewed. Among other things the MPs gave their opinions on a number of issues and their preference votes for the political parties. The issues concerned development and abortion law and order, income differences, worker participation, taxation and defence. The opinions were measured on a nine-point scale of which the lowest and the highest category were described (Table 1). The party preferences were recorded in a table of rank

Table 1. The issues and the meanings of the lowest and the highest category

1 DEVELOPMENT AID	the government should spend more money on aid to developing countries	(1) ..... (9) the government should spend less money on aid to developing countries
2 ABORTION	the government should prohibit abortion completely	(1) ..... (9) a woman has the right to decide for herself about abortion
3 LAW AND ORDER	the government takes too strong action against public disturbances	(1) ..... (9) the government should take stronger action against public disturbances
4 INCOME DIFFERENCES	income differences should remain as they are	(1) ..... (9) income differences should become much less
5 PARTICIPATION	only management should decide important matters in industry	(1) ..... (9) workers must also have participation in decisions important for industry
6 TAXATION	taxes should be increased for general welfare	(1) ..... (9) taxes should be decreased so that people can decide for themselves how to spend their money
7 DEFENCE	the government should insist on striking the Western enemy	(1) ..... (9) the government should insist on maintaining strong Western virtues

order. The scores in this table tell us the rank order each member of the parliament gave to the different parties (2 = highest preference, 15 = lowest preference). The lowest score (2) was always used for the MP's own party. For our illustration we only consider the preferences for the four largest parties, which are:

- PvdA - labour party (socialists) (39)
- ARP - Anti Revolutionary Party (christian democrats) (13)
- KVP - catholic party (christian democrats) (35)
- VVD - liberal party, referred to as conservatives (10)

The figure in parentheses is the number of MPs. The other parties in 1972 were:

- CHU - Christian Historical Union (christian democrats) (10)
- DV66 - democrats 66 (liberals) (11)



154

5.4. The category quantifications

To get an impression of the transformations of the variables, we give a plot of the original category numbers against the category quantifications (Fig. 7). The points are connected to allow the monotone transformation of each variable (see Section 4.2). For the issue 'development aid' we see for instance that the MPs who agree with

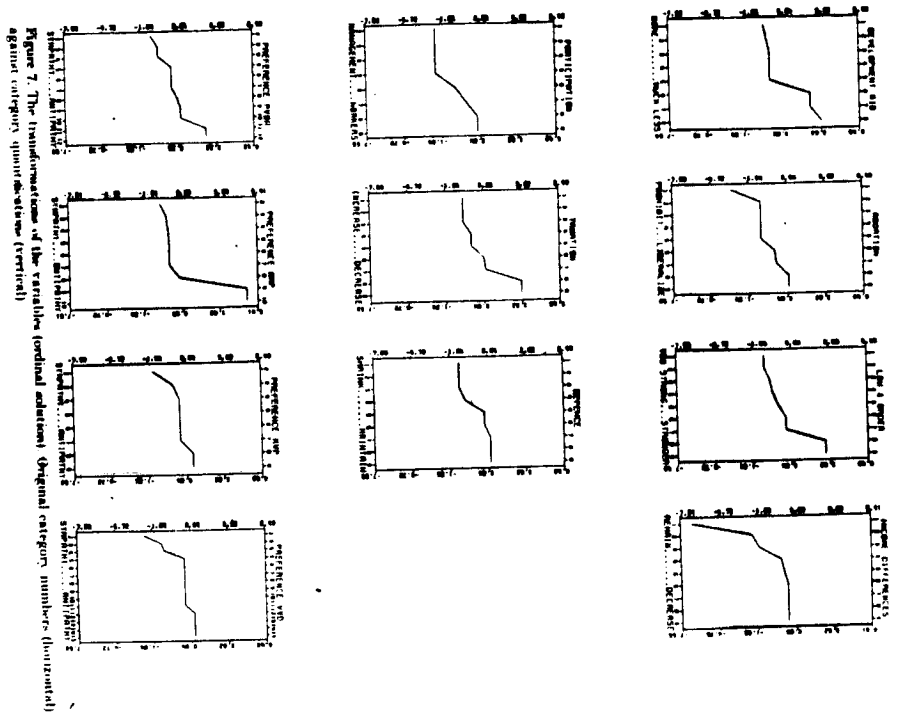


Figure 7. The transformations of the variables (ordinal relations). Original category numbers (horizontal) against category quantifications (vertical)

Summary of Thursday

Yesterday we discussed the use of B-splines as a

form of fuzzy coding, in connection with MPCA.

We also discussed various alternative criteria of  $R^2$

form  $\mu = \Phi(R)$  that can be used to

define other types of MPCA. Theory was developed

which shows how the various forms were related.

Finally we discussed the use of additivity restrictions on

the quantifications to introduce sets of variables, and

thus other forms of MVA such as regression, (M)ANOVA,

canonical correlation analysis, and so on.

We ended the talk yesterday by giving various examples

of the use of additivity restrictions. We shall continue

by giving some more.

Example 28 [ISRAELI, PM, 49, 1984, 311-342]

15c

and Activity. Also the explanatory variables are considered to be qualitative variables, corresponding with the way these variables are measured in the survey. In the first column of Table 2 the meaning of the categories is given, and in the last column the number of people falling in each category. The total sample consisted of 4108 people, all being 15 years or older.

In the second column the scale values  $c^i$  for Satisfaction are given as well as the regression coefficients for the dummy-variables, i.e. the vector  $d$ . The scale values for Satisfaction are determined in such a way that Satisfaction becomes a quantitative variable with mean zero and variance one. For the explanatory variables neither holds. By

TABLE 2

Qualitative Regression Analysis of Satisfaction on Marital status, Schooling, Income and Activity

Variable	Coefficients (c and d)	Coefficients after centeration	Standardized scale values	Number of people
<b>Satisfaction</b>				
Not too satisfied	-3.43	idem	idem	264
Rather satisfied	-0.81			539
Satisfied	0.25			1 857
Very satisfied	0.63			1 088
Extremely satisfied	0.53			360
<b>Marital status</b>				
Married	-0.21	0.04	0.26	2 940
Widowed	-0.56	-0.32	-2.34	268
Divorced	-1.17	-0.92	-6.72	53
Single	-0.21	0.03	0.24	847
<b>Schooling</b>				
Low	0.13	0.03	0.68	2 492
Medium	-0.09	-0.02	-0.42	1 027
High	0.03	-0.07	-1.88	401
Unknown	—	-0.10	-2.69	188
<b>Income (= Difl. 1,000)</b>				
< 21	-0.11	-0.12	-1.24	1 356
21 - < 40	0.07	0.06	0.55	1 689
> 40	0.19	0.18	1.81	454
Unknown	—	-0.01	-0.11	609
<b>Activity</b>				
Employed	0.17	0.04	0.17	1 987
Unemployed	-1.00	-1.13	-4.96	65
Not able to work	-0.86	-0.99	-4.35	118
Retired	0.26	0.13	0.59	314
Student	0.17	0.04	0.19	336
Housewife	0.18	0.06	0.25	1 203
Unknown	—	-0.13	-0.56	85

N = 4108

$R^2 = .082$  [int. .0624, ridit. .0538]

17

TABLE 3  
Qualitative Regression Analysis of Happiness on Marital status, Schooling, Income and Activity

Variable	Coefficients (c and d)	Coefficients after centeration	Standardized scale values	Number of people
<b>Happiness</b>				
Unhappy	-5.41	-4.17	-4.17	11
Not too happy	-1.66	-1.66	-1.66	92
Happy, not unhappy	0.28	0.28	0.28	561
Happy	0.73	0.73	0.73	2 555
Very happy	—	—	—	889
<b>Marital status</b>				
Married	-0.14	0.09	0.46	2 940
Widowed	-0.84	-0.61	-2.94	268
Divorced	-1.16	-0.93	-4.49	53
Single	-0.31	-0.08	-0.37	847
<b>Schooling</b>				
Low	0.03	0.00	0.28	2 492
Medium	0.04	0.01	0.71	1 027
High	-0.00	-0.03	-2.48	401
Unknown	—	-0.03	-2.24	188
<b>Income (= Difl. 1,000)</b>				
< 21	-0.08	-0.12	-1.22	1 356
21 - < 40	0.10	0.07	0.68	1 689
> 40	0.20	0.16	1.63	454
Unknown	—	-0.04	-0.39	609
<b>Activity</b>				
Employed	0.15	-0.01	-0.08	1 987
Unemployed	-0.12	-0.28	-2.17	65
Not able to work	-0.50	-0.66	-5.10	118
Retired	0.24	0.07	0.57	314
Student	0.22	0.10	0.80	336
Housewife	0.23	0.06	0.47	1 203
Unknown	—	-0.16	-1.27	85

$c^i, c^j = 1$

Because  $E_{ij}$  is a diagonal matrix we can construct a  $(2 \times 2) \times (2 \times 2)$  diagonal matrix  $A_0$  with block-diagonal submatrices  $E_{ij}$ . The diagonal of  $A_0$  contains the probabilities of falling in the various categories. Defining  $c = (c^1, \dots, c^k)$ ,  $(k=2)$  terms into

(4.3).

For the time being, we do not put separate normalizations on each of the  $c^i$ -vectors. This kind of normalizations can be done afterwards. Notice that the parameters  $c^i$  can be considered as a product of a scaling parameter  $c^{i'}$ , with norm 1 for each variable, and a

$R^2 = .0754$

TABLE 4

Qualitative Redundancy Analysis of Satisfaction and Happiness on Marital Status, Schooling, Income and Activity

Variable	Coefficients (c and d)	Coefficients after centeration	Standardized scale values
<u>Satisfaction</u>			
Not too satisfied	-2.46		-3.34
Rather satisfied	-0.66		-0.90
Satisfied	0.16	idem	0.22
Very satisfied	0.51		0.69
Extremely satisfied	0.41		0.55
<u>Happiness</u>			
Unhappy	-3.71		-5.49
Not too happy	-2.79	idem	-4.14
Not happy, not unhappy	-1.11		-1.65
Happy	0.19		0.28
Very happy	0.49		0.73
<u>Marital status</u>			
Married	-0.73	0.09	0.39
Widowed	-0.96	-0.65	-2.76
Divorced	-1.63	-1.30	-5.52
Single	-0.36	-0.03	-0.12
<u>Schooling</u>			
Low	0.12	0.02	0.61
Medium	0.10	-0.01	-0.16
High	0.02	-0.08	-2.10
Unknown	—	-0.10	-2.71
<u>Income (* Df1. 1,000)</u>			
< 21	-0.14	-0.18	-1.24
21 - < 40	0.12	0.09	0.61
> 40	0.28	0.24	1.73
Unknown	—	-0.03	-0.23
<u>Activity</u>			
Employed	0.21	0.02	0.08
Unemployed	-0.81	-1.00	-4.07
Not able to work	-0.98	-1.17	-4.75
Retired	0.35	0.15	0.62
Student	0.30	0.10	0.42
Housewife	0.27	0.08	0.33
Unknown	—	-0.19	-0.78

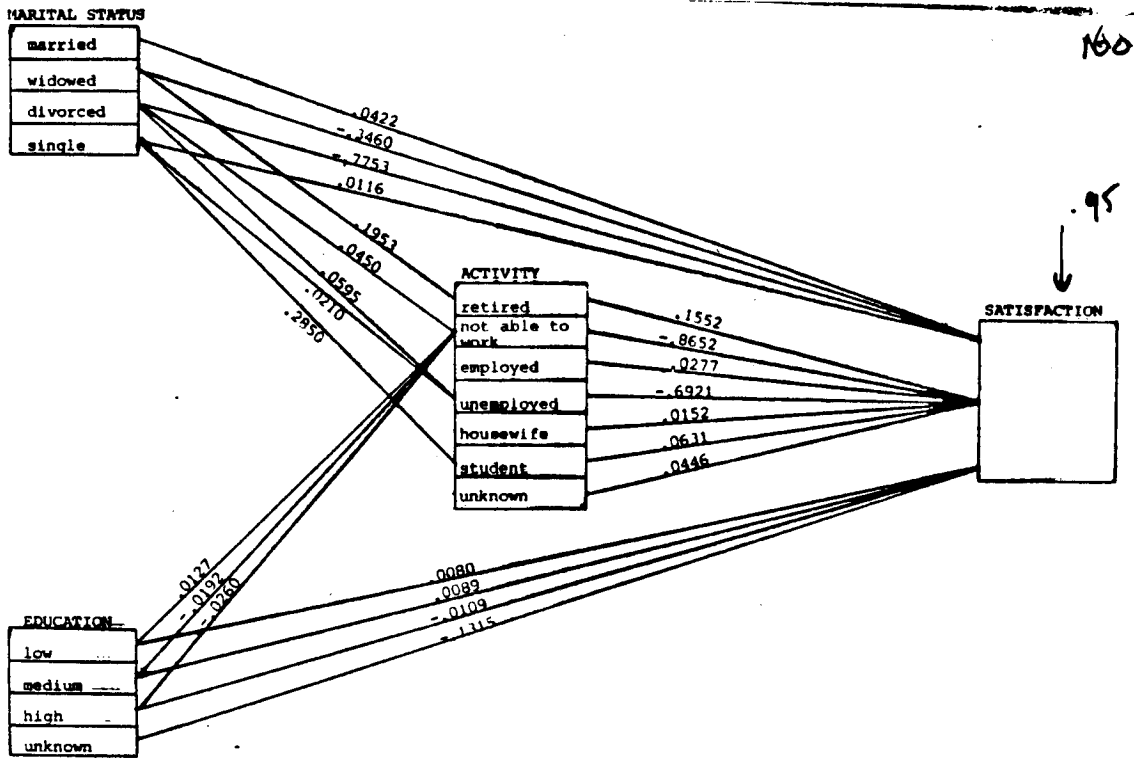
Sample

$$\text{Satisfaction} \approx .137 \text{ Marital Status} + .039 \text{ Schooling} + .100 \text{ Income} + .227 \text{ Activity}$$

$$\text{Happiness} \approx .261 \text{ Marital Status} + .014 \text{ Schooling} + .098 \text{ Income} + .130 \text{ Activity}$$

$$\begin{aligned} &.737 \text{ Satisfaction} + .676 \text{ Happiness} \approx \\ &.236 \text{ Marital Status} + .037 \text{ Schooling} + \\ &+.141 \text{ Income} + .247 \text{ Activity} \end{aligned}$$

Figure 2. Path diagram of Satisfaction, Activity, Marital status and Education



*variables 7.*  
*dependent*

TABLE 1

Variables from the study of chronic lung disease

set 1	VLA: Residence, (1) Vlagtwedde, (2) Vlaardingen.
set 2	SNO: Smoking, (1) never smoker, (2) ex-smoker, (3) current smoker. RATE: Rate of smoking (amount of tobacco), (1) never smoker, (2) low rate, ..... (9) high rate. PER: Time period smoked, (1) never smoker, (2) short period, ..... (13) long period. TIME: Time since last cigarette, (1) never smoker, (2) long ago, ..... (5) recently, (6) current smoker.
set 3	AGE: Age discretized into periods of 3.5 years, (1) age 19-22.5, ..... (10) age 52.5-56. SEX: Sex, (1) male, (2) female.
set 4	COU: Coughing, (1) no, (2) persistent. PHE: Phlegm, (1) no, (2) persistent. DIS: Dyspnoea or shortage of breath, (1) no, (2) slight/moderate, (3) severe. HHE: Wheezing, (1) never, (2) ever, (3) severe. AST: Asthma, (1) ever, (2) never.

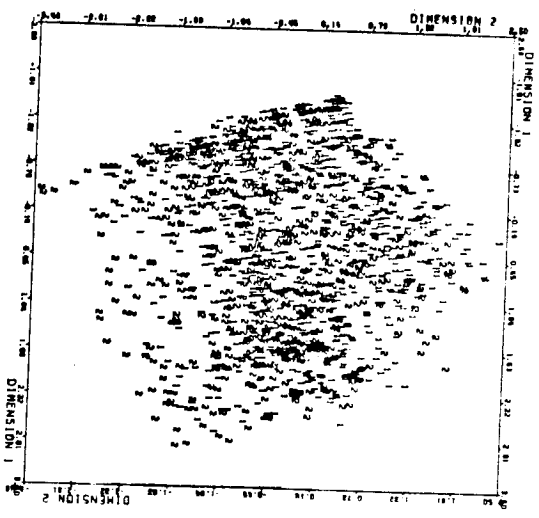


Figure 2  
 Object scores 1 Wageningen  
 2 Vlaardingen

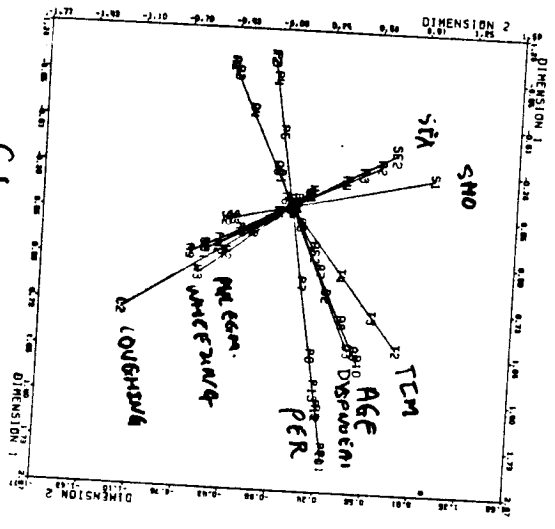


Figure 3  
 Category quantifications

The OVERALS algorithm (by implication this algorithm can also be used for the various special cases).

- ① Start with  $X_0, Z_0$
- ② Loop over variables

$$\min_k \left\{ [X - \sum_{j \in J} G_j Y_j]' [X - \sum_{j \in J} G_j Y_j] \right\}$$

$$\min_k \left\{ [X - \sum_{j \in J} G_j Y_j] - G_k Y_k \right\}' \left\{ \dots \right\}$$

- ③ first multiple quantification

$$Y_k \leftarrow D_k^{-1} \left[ X - \sum_{j \in J} G_j Y_j \right]$$

Observe that the centroid principle is not satisfied in the same way as before!

① Then  
adjust for cone & want one restriction of  
necessary

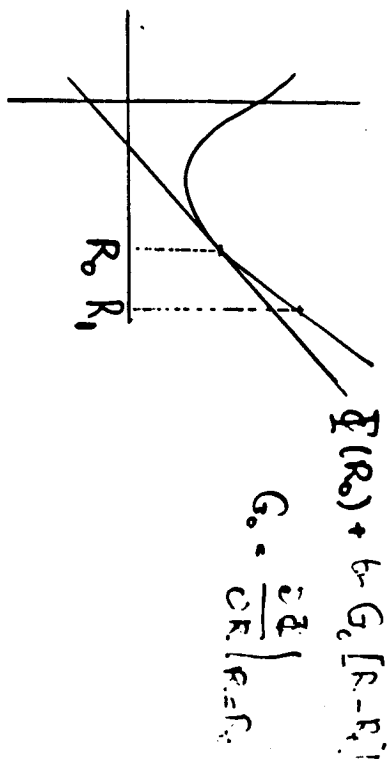
② Find new X at end of variable loop.

This is, of course, very similar to the PRDUALS  
algorithm, with just one additional complication.

We shall also briefly discuss the <sup>future</sup> alternative algorithm  
based on the Burt table. We want to maximize

$$\mu = \Phi[R(y)]$$

over quantifications. All variables are assumed to  
be single [table-one], multiple variables and  
ads must be introduced through copies and  
additivity restrictions.



$$\text{Thus } \Phi(R) \geq \Phi(R_0) + t G_0 (R - R_0)$$

Now suppose we maximize to  $G_0 R$  over  $R$   
(i.e. over  $y_1, \dots, y_m$ ), and this gives  $R_1$ . Then

$$\begin{aligned} \Phi(R_1) &\geq \Phi(R_0) + t G_0 (R_1 - R_0) \\ &\geq \Phi(R_0) \end{aligned}$$

Thus by maximizing to  $G_0 R$  we increase  $\Phi$ .  
This gives a convergent algorithm.



How does  $\text{tr } G_0 R$  look. We can write  
 $\text{tr } G_0 R = y_i' G_{ij} y_j$ ,  $y_i' D_j y_i = 1$ . Thus  
 $\text{tr } G_0 R = \sum_i \sum_j y_i' G_{ij} y_j$ . Thus if  
 we define  $\tilde{C}_0$  to be the Burt table with  
 each subtable multiplied by  $g_{ij}$ , then  
 $\text{tr } G_0 R = y_i' \tilde{C}_0 y$ , where must be  
 maximized over all  $y$  satisfying  $y_i' D_j y_i = 1$   
 (and cone + additivity restrictions, i.e. cone  
 restrictions).

If  $G_0$  is positive - semi-definite additional  
 simplification is possible. Then  $\tilde{C}_0$  is PSD,  
 and  $y_i' \tilde{C}_0 y$  is convex. We now apply  
 the same trick (majorization)

$$y_i' \tilde{C}_0 y \geq y_i' \tilde{C}_0 y_0 + 2 y_i' \tilde{C}_0 [y - y_0]$$

As a consequence we need to maximize the  
 linear function  $y_i' \tilde{C}_0 y_0$  over  $y$  satisfying  
 the restrictions. This can be done for each  
 $i$  separately by minimizing, setting  $\tilde{y}_i = \tilde{C}_0 y_0$ ,

$$[y_i - \tilde{y}_i] D_j [y_i - \tilde{y}_i]$$

over  $y_i \in K_i$  (the cone), and then  
 setting the new  $y_i$  equal to the normalized  
 projection.

Some examples of  $G_0$ .

$$\text{if } \Phi(R) = \lambda_1 + \dots + \lambda_p \quad \text{then } g_{ij}' = \sum_{s=1}^p \alpha_{ij}^s r_{ij}^s$$

$$\text{if } \Phi(R) = \sum_i \sum_j r_{ij} \quad \text{then } g_{ij}' \equiv 1$$

$$\text{if } \Phi(R) = \sum_i \sum_j r_{ij} \quad \text{then } g_{ij}' = r_{ij}$$

$$\text{if } \Phi(R) = -|R| \quad \text{then } g_{ij}' = r_{ij}$$

### Analysis of statistical stability

There are quite a few results on analytic and algebraic stability in Giff:

Techniques used - delta method

- Bootstrap

- Jostensnik.

We shall first explain these techniques briefly.

They are essentially nonparametric, and apply to all functions of counts (frequencies).

Suppose we have  $n$  repeated independent replications of a discrete random variable, taking  $m$  possible values. This defines a vector

$P_n$  of relative frequencies (proportions). We know that

$$MLN (a) \quad P_n \xrightarrow{P} \pi$$

$$CLT (b) \quad n^{1/2} [P_n - \pi] \xrightarrow{L} N(0, \pi - \pi\pi')$$

Now suppose  $\psi$  is a function of  $P$  [in the unit simplex of  $\mathbb{R}^m$ ] with values in  $\mathbb{R}^t$ .

If  $\psi$  is continuous (at  $\pi$ ), then

$$\psi(P_n) \xrightarrow{P} \psi(\pi).$$

If  $\psi$  is differentiable (at  $\pi$ ), with derivative  $G_\pi$ , then

$$n^{1/2} [\psi(P_n) - \psi(\pi)] \xrightarrow{L} N(0, G_\pi V G_\pi')$$

Another result which has been known at least since Laplace.

Thus: any diff function of proportions is asymptotically normal

If  $\psi$  is continuously differentiable in a neighborhood of  $\pi$ , then  $\hat{\sigma}_{\psi}^{VGR}$  can be estimated.

$$\hat{\sigma}_{\psi}^{VGR} = \frac{1}{n} \sum_{i=1}^n \left[ \frac{\partial \psi}{\partial \pi} \right] G_i' \frac{1}{n} \sum_{i=1}^n G_i \left[ \frac{\partial \psi}{\partial \pi} \right]' G_i'$$

This defines the Delta Method.

Everything computed in Gifi is a function of the profile frequencies [even if splines are used] and consequently the delta method can be used. We merely need

- expressions for the derivatives
- programs for computing them.

This means that there is a routine way in which standard errors and confidence ellipsoids can be computed for our terminology [at least if you are willing to assume that the data are some random sample]. If the data are nonrandom the confidence regions can still be interpreted as measures of stability.

Most of the things computed in Gifi are functions of the Part table C. We can write

$$C = \sum_{j=1}^m p_{ji} g_{ji} g_i'$$

with  $g_{ji}$  the profile-indicators. We can write  $\psi(p)$  in the form  $\psi(C(p))$  and

obtain (by the chain rule)

$$\frac{\partial \psi}{\partial p_{jk}} = \sum_{je} \frac{\partial \psi}{\partial c_{je}} g_{je} g_{jk}$$

This often simplifies matters. Many other things are functions of the reduced correlations  $r_{je} = y_j' C_{je} y_e$ .

Then

$$\frac{\partial \bar{Y}}{\partial P_n} = \sum_{je} \frac{\partial \bar{Y}}{\partial r_{je}} \frac{\partial r_{je}}{\partial P_n} = \sum_{je} \frac{\partial \bar{Y}}{\partial r_{je}} \left[ y_j' C_{je} \frac{\partial y_e}{\partial P_n} + y_j' \frac{\partial C_{je}}{\partial P_n} y_e + \frac{\partial y_j'}{\partial P_n} C_{je} y_e \right]$$

If (model!) we have (in the population) linear

bivariate regression then  $C_{je} y_e = r_{je} D_j y_j$  can

$C_{ij} y_j = r_{je} D_e y_e$ . Thus

$$\frac{\partial \bar{Y}}{\partial P_n} = \sum_{je} \sum_{ie} \frac{\partial \bar{Y}}{\partial r_{je}} \left[ r_{je} y_j' D_j \frac{\partial y_i}{\partial P_n} + r_{je} y_e' D_e \frac{\partial y_e}{\partial P_n} + (y_j' D_{je}) (y_e' D_{ie}) \right]$$

But we use the restrictions  $y_j' D_j y_j = y_e' D_e y_e = 1$ .

Thus

$$y_j' D_j \frac{\partial y_j}{\partial P_n} + \left( \frac{\partial y_j'}{\partial P_n} \right)' D_j y_j + y_j' \frac{\partial D_j}{\partial P_n} y_j = 0$$

or

$$y_j' D_j \frac{\partial y_j}{\partial P_n} = -\frac{1}{2} y_j' \frac{\partial D_j}{\partial P_n} y_j$$

It follows that

$$\frac{\partial \bar{Y}}{\partial P_n} = \sum_{je} \frac{\partial \bar{Y}}{\partial r_{je}} \left[ y_j' \frac{\partial C_{je}}{\partial P_n} y_e - \frac{1}{2} r_{je}' \frac{\partial D_j}{\partial P_n} y_j - \frac{1}{2} r_{je} y_e' \frac{\partial D_e}{\partial P_n} y_e \right]$$

But this is the same derivative as the one we find if the  $y_j$  are known scores (not computed).

Implication for - PRIMALS [first step]

or any other first-step method.

Implications in combination with LISREL.  
(for instance).

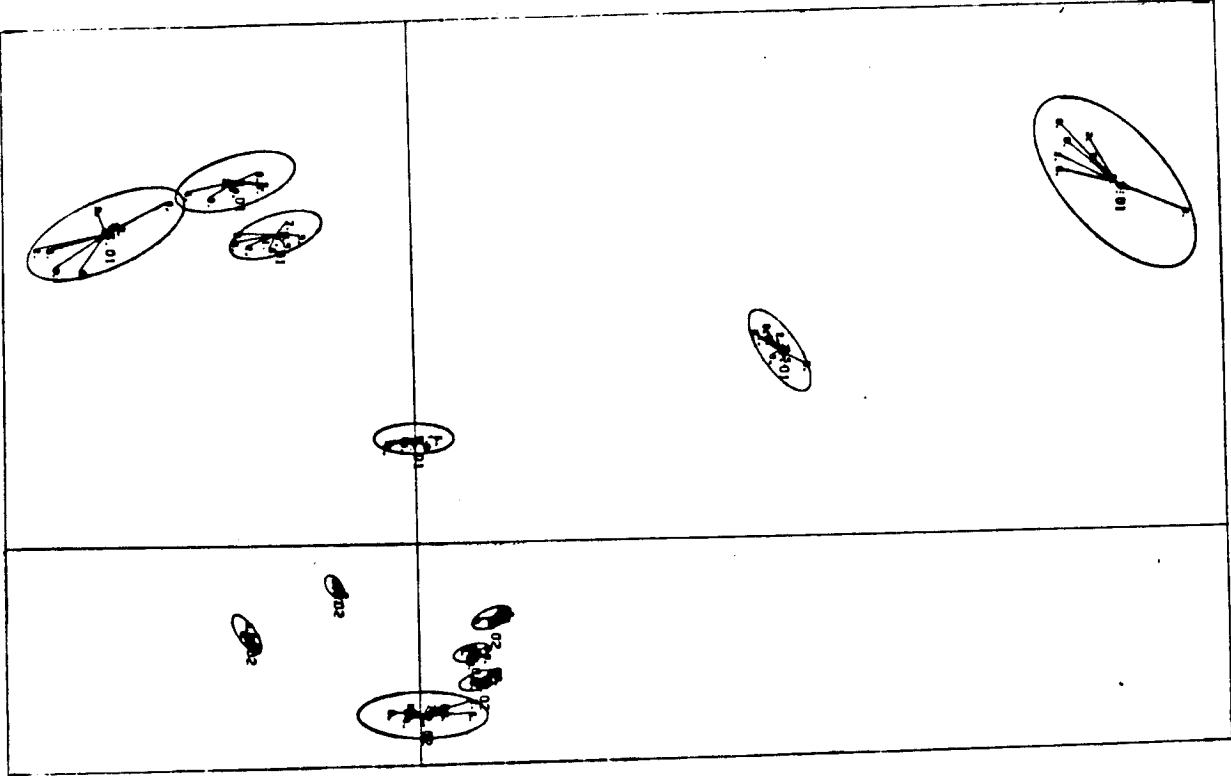


Figure 12.4 Stability of the ANNEX solution of the Sigman det. 95% ellipses and 10 bootstraps

Future = Derive the necessary derivatives for an application of the delta method to all techniques

[done for HODRHS  
 RWALOR → option  
 PRISMS  
 ANN PROF ← option]

= apply two-step model to LISREL etc  
 [and improve two-step estimates in order to construct efficient ones]. Construct test for linearity of all regressions.

Next Bootstrap. This is very easily explained. We know that

$$P_n | \pi \xrightarrow{L} \mathcal{D}(0, V(\pi))$$

Now suppose we draw a random sample, with replacement, from the multivariate with

parameters  $(p_n, n)$ . This gives  $\tilde{p}_n$  Now

$$\tilde{p}_n | p_n \xrightarrow{L} \mathcal{D}(0, V(p_n))$$

and because  $V(p_n) \xrightarrow{P} V(\pi)$  we can use

$\tilde{p}_n | p_n$  to approximate  $p_n | \pi$ . This amounts

to drawing random samples from your own data and making "inferences" to the "population" on the basis of this. In many recent publications (Efron, Beran, Freedman, Diaconis) it has been shown

that Bootstrap methods are robust and efficient in many circumstances.

Examples

- Back to Japan
- Smoothing
- Bias correction  $2\hat{F}(e_n) - \hat{F}(\hat{e}_n)$ .

Future research (in progress)

- how many samples are needed.
- how best implemented.

Jackknife We can be brief.

$$\hat{\Phi}_\mu(\hat{e}_n) \stackrel{\Delta}{=} \Phi \left[ \frac{n p_n - e_\mu}{n-1} \right] = \Phi \left[ p_n + \frac{1}{n-1} [p_n - e_\mu] \right]$$

$$\sim \Phi(p_n) + \frac{1}{n-1} \frac{\partial \Phi}{\partial p_n}(p_n - e_\mu)$$

$$n \cdot \Phi(p_n) - (n-1) \Phi_\mu(p_n) \approx \Phi(p_n) - \frac{\partial \Phi}{\partial p_n}(p_n - e_\mu)$$

Pseudo - values  $\tilde{\Phi}_\mu(p_n)$

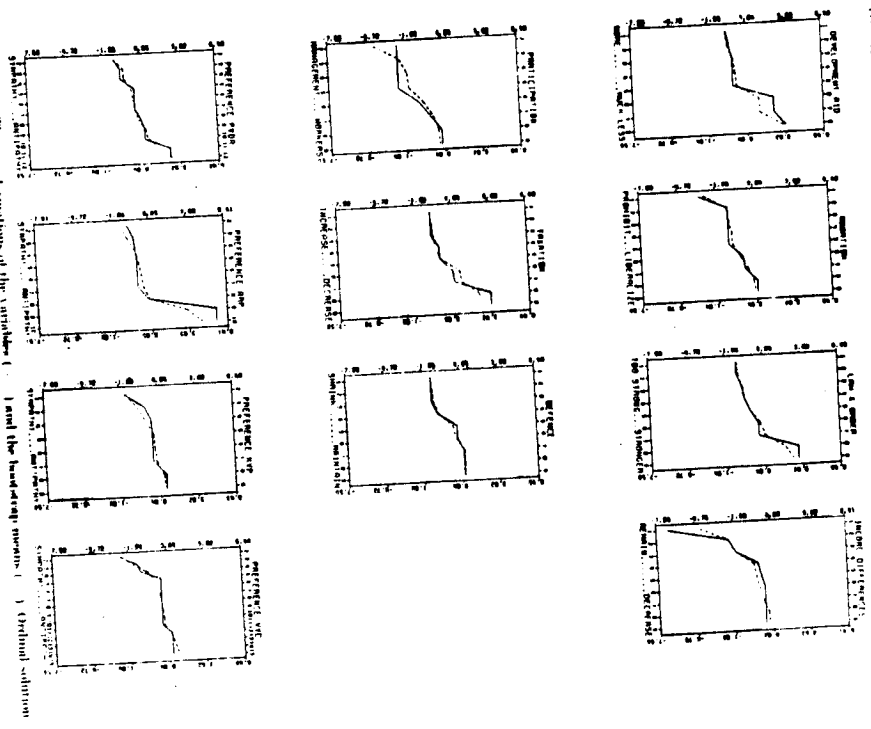


Figure 8. The transformation of the variables (1) and the best-order means (1) (orthogonal solution). Advantage of this technique is the fact that we are free to choose the measurement level of each variable separately, so that we do not have to impose strong restrictions on the analysis unnecessarily. The CANALS technique is a real extension and improvement of the technique proposed by Young *et al.* (1976). The alternative least squares algorithm of the CANALS program is a very nice way of avoiding expensive computation and computing results rather quickly and precisely. Another advantage

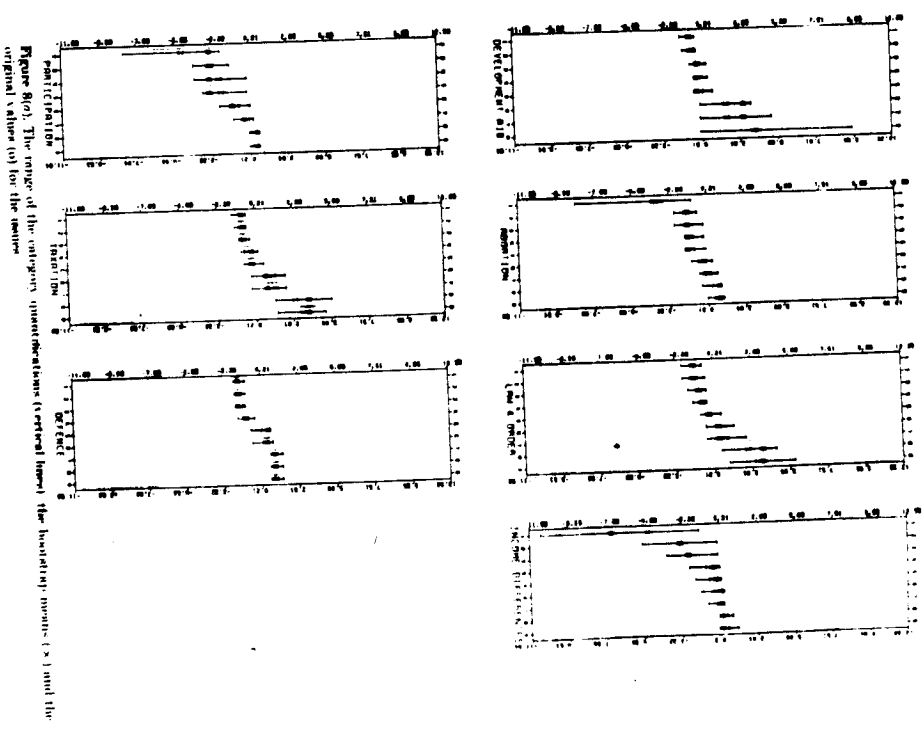


Figure 9(a). The range of the category observations (vertical lines) the best-order means (s) and the original values (o) for the issues.

$$\sum_p p_r \hat{f}_r(x_n) \approx \hat{F}(x_n)$$

$$\sum_p p_r [\hat{F}_r(x_n) - \hat{F}(x_n)]^2 \approx$$

$$g'(x_n) V(x_n) g(x_n)$$

which is the delta method variance estimate.

Thus, aside from profile, compute the pseudo-values and their variance.

Future: Comparison Bootstrap & Jackknife.

BS & JK are both useful as data analysis techniques, even if there is no random sampling.

### Significance testing of generalised correlation

De Leeuw & Van Buuren, Data analysis & hypothesis testing

We use permutation methods. If variables are partitioned into sets we leave set 1 intact.

Then permute all observations in set 2 randomly, same in set 3, etc. This creates the null-distribution. There are various ways to assess it

- enumeration (Fisher exact test)
- use asymptotic normality of the Pearson table (as in Lebart, O'Neill for CA)
- Monte Carlo
  - Via C
  - directly



TABLE 3: Generalized canonical correlations: empirical values (ev) and generated values at 5, 25, 50, 75 and 95 percent. Random permutation method. From Year to Year data, multiple nominal, single nominal and numerical.

	mult nominal		single nom		numerical	
	dim1	dim2	dim1	dim2	dim1	dim2
ev	.744	.570	.735	.530	.714	.496
5	.549	.532	.533	.511	.501	.482
25	.554	.542	.543	.521	.511	.494
50	.559	.545	.549	.527	.518	.499
75	.564	.550	.554	.532	.525	.504
95	.570	.557	.565	.541	.532	.514

TABLE 4: Generalized canonical correlations: empirical values (ev) and generated values at 5, 25, 50, 75 and 95 percent. C-matrix method. From Year to Year data, multiple nominal and numerical.

	multiple nom		numerical	
	dim1	dim2	dim1	dim2
ev	.744	.570	.714	.496
5	.553	.545	.509	.500
25	.560	.550	.514	.503
50	.565	.554	.521	.506
75	.569	.557	.525	.504
95	.574	.562	.535	.514

canonical correlation problems with  $p=2$  were computed separately. The eigenvalues are in the first row of table 5, while the rest of the table shows percentiles of the BD (all estimated by the random permutation method). In the S & M analysis probability plots [given in figure 2] show a large deviation from normality for the first eigenvalue, all eigenvalues are very clearly significant, however. All estimated PD's have small variance, we see that  $\lambda(95) - \lambda(5) < .03$ , except for dimension 1 of S & M (which has a very light tail on the right,  $\lambda(95) - \lambda(75) = .035$ ).

TABLE 6: Generalized canonical correlations: empirical values (ev) and generated values at 5, 25, 50, 75 and 95 percent. Random permutation method and C-matrix method. Russett data, multiple nominal and numerical

	Random permutation method		C-matrix method	
	mult nom	numerical	mult nom	numerical
	dim1	dim2	dim1	dim2
ev	.815	.737	.687	.462
5	.724	.695	.462	.397
25	.754	.714	.499	.426
50	.773	.731	.525	.445
75	.794	.752	.545	.467
95	.830	.778	.576	.488

CONCLUSIONS

It appears from our examples that both the random permutation method and the C-matrix method give a fairly good approximation to significance probabilities and permutation distributions. In order to be sure, we shall have to apply enumeration somewhat more extensively. We also need to generalize the exact computation of  $G_s$  to the case of more than two sets. But from our experience so far, we can say that the Monte Carlo methods give the correct indication of the order of significance, and they show that with multiple options and not too many individuals the eigenvalues have to be very high to be significant. It is not surprising, for instance, that Gifi (1981) had great difficulty in interpreting the second dimension in the solution for the Russett data. Our approximations to the permutation distribution guard us against chance capitalization, and against trying to interpret effects which are not really there.

184

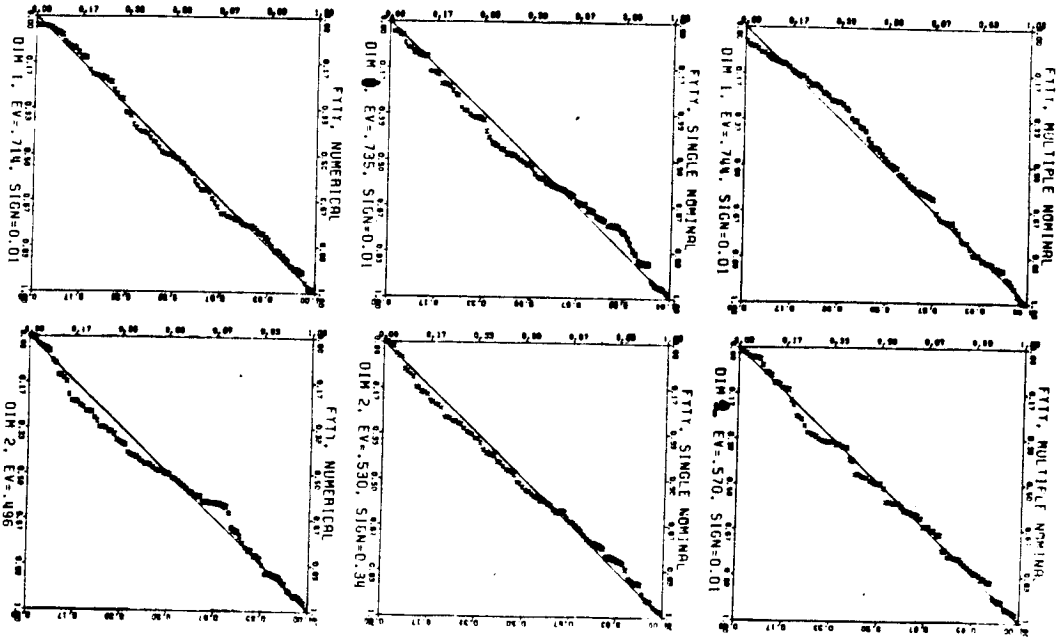


FIGURE 1: From Year to Year random permutations. Probability plot of eigenvalues.

185

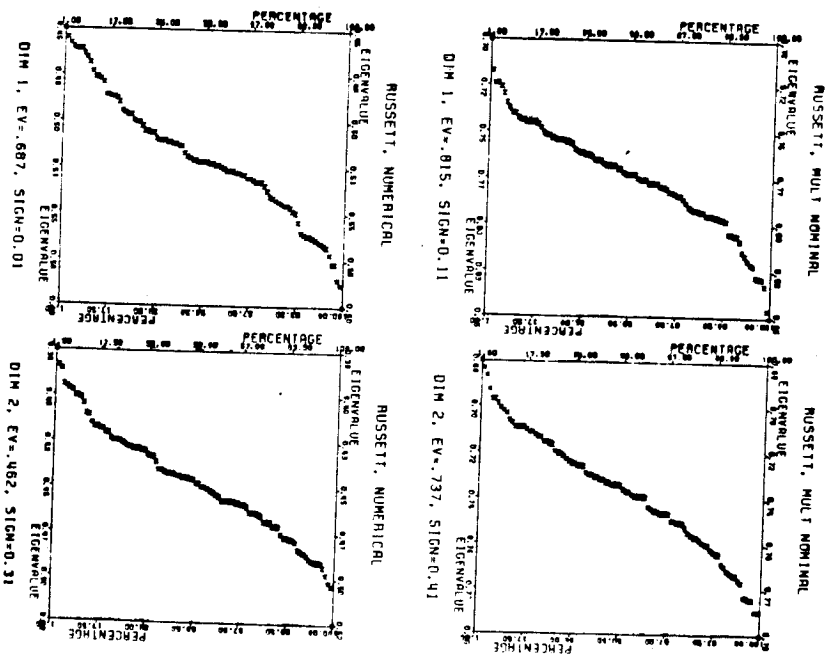


FIGURE 3: Russett random permutations. Eigenvalues (ordered) against percentages.

Possibilities for partial canonical correlation analysis

It can be incorporated in OVERALS by using copies of a variable in different sets [in all sets]. Using p copies again means multiple elimination, using one copy means single elimination. We must take care, however, that the variable has the same quantification in all sets in which it occurs.

Possibilities for "causal" analysis

We depart from the Gif-system, with its heavy emphasis on component analysis. The more general loss function we now study is

$$Z(Y) = \sum_{j=1}^m \lambda_j [G_j Y - \sum_{i \in I_j} G_i Y] [I \dots I]$$

Here  $j=1, \dots, m$  are endogenous variables and  $I_j$  is the set of "causes" of variable  $j$ . For normalization purposes it is convenient to require  $Y' G_j Y = I$  for all  $j$ , and to write

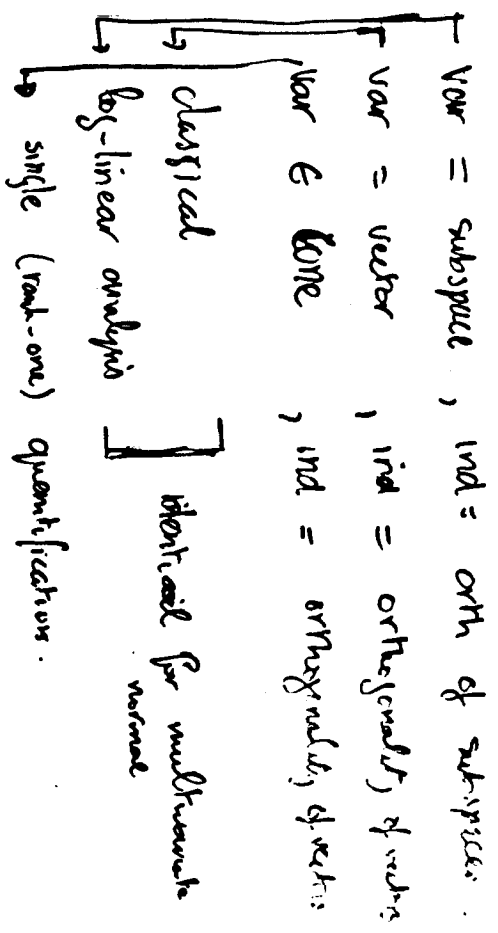
$$G(Y, B) = \sum_{j=1}^m \lambda_j [G_j Y - \sum_{i \in I_j} G_i Y e_i B_{ij}] [I \dots I]$$

Let us explain the idea behind this loss function.

- a) A causal model is a digraph.
- b) Exogenous variables are sources, or transmitters.



This is still not entirely quantitative (we have to specify how we represent "variables" and how we mean "independence".



The production-oriented approach checks (by using, in our case, a least squares loss function) if, and in how far, a variable is in the space spanned by its direct causes.

Thus we want residuals to be zero, not orthogonal.  
Key result: In recursive path models the predictive and model-fitting approaches for numerical variables are equivalent [Udd, 1966].  
If there is missing information this is no longer true.  
Now missing information can be missing for various reasons

- (a) proper
- (b) non-numerical (ordinal, nominal)
- (c) latent

The basic idea behind our approach is latent variables is that they are variables with a very primitive measurement level (they are certainly single nominal! No restriction at all).

This means now we can simply incorporate "latent" info. our hierarchy of measurement levels.

Future: - Compare here various ways of making qualitative path models quantitative.

- Contrast PATHALS algorithm, compare with LISREL and ALS. \*

\* Already been compared in detail, but not very well. The main difference is how the missing information is estimated. In PLS / GFI by LS and DS  $\Rightarrow$  inconsistency

IN LISREL / ML by  
constraint equations (EM - algorithm)  $\Rightarrow$   
consistency

New version of GFI [published by DSW-press, estimated early 1997]  
in three volumes.

Nonlinear Multivariate Data Analysis  
Annotated References

- Breiman, L., & Friedman, J.H. Estimating Optimal Transformations for Multiple Regression and Correlation. J. Amer. Statistical Association, 80, 1985, 580-619.
- Advanced mathematical treatment of regression with optimal (smoothed) scaling by alternating least squares (called alternating conditional estimation).
- Gifi, A. Non-linear Multivariate Analysis. Leiden, University of The Netherlands: Dept. of Data Theory, 1981.
- Preliminary version of a book covering most nonlinear multivariate techniques. Uses ALS throughout.
- Gittins, R. Canonical Analysis. Berlin: Springer-Verlag, 1985.
- Reviews statistical and algorithmic theory for Canonical Analysis. Contains many ecological examples worked out in detail.
- Greenacre, M.J. Theory and Applications of Correspondence Analysis. New York: Academic Press, 1984.
- Discusses Correspondence Analysis from the French point of view. Comprehensive, including algorithmic and statistical aspects, with examples.
- Heiser, W. Unfolding Analysis of Proximity Data. Leiden, The Netherlands: Dept. of Data Theory, University of Leiden, 1981.
- Discusses Correspondence Analysis as a way of performing Multidimensional Unfolding (i.e., from the Multidimensional Scaling viewpoint).
- Lebart, L., Morineau, A., & Warwick, K.M. Multivariate Descriptive Statistical Analysis: Correspondence, and Related Techniques for Large Matrices. New York: Wiley & Sons, 1984.
- Translation of a French book on Correspondence Analysis (and cluster analysis). Elementary.
- de Leeuw, J. Canonical Analysis of Categorical Data. Leiden, The Netherlands: DSWO-Press, 1985 (reprint of 1973 dissertation)
- The progenitor of the current work on Nonlinear Multivariate Data Analysis. Algebraic.
- de Leeuw, J. Nonlinear Principal Components Analysis. In: Causinus, H. (ed.), COMPSTAT 1982. Vienna, Physica-Verlag, 1982.
- Concise review of various approaches to nonlinear principal components analysis, including correspondence analysis.
- de Leeuw, J. The Gifi system for nonlinear multivariate analysis. In: Diday, E. (ed.) Data Analysis and Informatics, Vol 3. Amsterdam: North Holland Publishing Company, 1983.
- Very concise review paper covering most of the material in this course. Discursive overview. Difficult but not mathematical.
- Mardia, K.V., Kent, J.T., & Bibby, J.M. Multivariate Analysis. New York: Academic Press, 1979.
- Covers nearly all the models discussed in this course, but without nonlinear data transformations. Advanced. Mathematical.
- Marrlott, F.H.C. The Interpretation of Multiple Observations. New York: Academic Press, 1974.
- Beautifully written concise (117pp) overview of linear multivariate analysis. Requires elementary matrix algebra only.
- Meulman, J. Homogeneity analysis of incomplete data. Leiden, The Netherlands: DSWO-press, 1982.
- Clear introduction to homogeneity analysis (multiple correspondence analysis). Discusses missing data.
- Perrault, M.D., & Young, F.W. Alternating Least Squares Optimal Scaling: Analysis of Nonmetric Data in Marketing Research. J. Marketing Research, 17, 1980, 1-13.
- Overview of ALS methods as they apply to marketing research.

Stone, C.J. Additive Regression and other nonparametric models.  
Annals of Statistics, 13, 1985, 689-705.

Discusses nonlinear regression and conjoint analysis methods in a nonparametric statistical framework. Difficult and advanced.

Tenenhaus, M., & Young, F.W. An Analysis and Synthesis of Multiple Correspondence Analysis, Optimal Scaling, Scaling, Homogeneity Analysis, and other methods for Quantifying Categorical Multivariate Data. Psychometrika, 50, 1985, 91-120.

Compares the many different proposals for analysis methods which are all equivalent to Correspondence Analysis.

Young, F.W., & Sarle, W.S. Exploratory Multivariate Data Analysis. Cary, NC: SAS Institute, Inc., 1983.

Notes for a course presented to those who wish to apply linear exploratory multivariate data analysis. Detailed examples. Little mathematics. Some nonlinear treatment.

Young, F.W. Quantitative Analysis of Qualitative Data. Psychometrika, 46, 1981, 357-388.

Overview of the first five years of developments in nonlinear multivariate data analysis. Algorithmic. Intermediate mathematics.

APPENDIX A:  
 ILLUSTRATIVE DATA SETS  
 USED IN NLWVA COURSE



(1) 1980 Car data, taken from consumer reports, analyzed by Winsberg and Ramsay, our analysis by Van Rijckevorsel.

33 objects (cars)  
5 variables: price in \$

- engine size
- gas consumption city
- gas consumption highway
- weight

(2) baby data, taken from Shirley, The first two years, 1931. Analyzed by De Leeuw.

20 objects (babies)  
71 variables (weeks) each with the categories stepping - standing - walking with help - walking alone - not yet stepping.

(3) skiing resort data, taken from tourist information, analyzed by Heiser.

115 objects (skiing resorts)  
6 binary variables: suitable for beginners

- good for walking
- lots of amusement
- for advanced skiers
- round trips on skis
- cross country skiing

(4) banks and industries, analyzed earlier by Levine, our analysis by Heiser.  
cross table of 70 industries and 14 banks, counting the number of board members they have in common.

(5) japanese religion data, taken from Sugiyama. Our analysis by Giffi, Heiser.

4243 objects  
6 variables (religious practices): attending services

- visiting graves
- read religious books
- visit shrines and temples
- keep talisman
- draw fortune

(6) American states data, from social indicator statistics, previously analyzed by Mainer, our analysis by De Leeuw and Neuman.

50 objects (states)  
7 variables: population

- income
- illiteracy rate
- life expectancy
- homicide rate
- percent high school graduates
- number of days below freezing point

(7) Vlaardingen-Vlagtwedde data. Analysis by Van der Burg, Van Pelt.

4000 objects  
variables: set 1: Vlaardingen - Vlagtwedde

- set 2: smoking habits
  - ever smoked
  - rate of smoking
  - period of smoking
  - how long ago
- set 3: background
  - sex
  - age

- set 4: symptoms
  - coughing
  - phlegm
  - dyspnoea
  - wheezing
  - asthma
- set 5: maximum expiratory flow volume curve measurements
  - PEF
  - FVC
  - etc

(8) uterine cervix data

126 objects (slides)  
7 variables (pathologists) with categories

- negative
- atypical squamous hyperplasia
- carcinoma in situ
- squamous carcinoma with early

stromal invasion  
- invasive carcinoma

300

- (9) Cities, villages, etc. Analyzed by De Leeuw.  
792 objects (towns, cities, villages in Netherlands)  
8 variables: social backwardness,

percentage lower technical education  
percentage unemployed  
public library present  
number of types of schools  
percentage of unemployed without proper schooling  
percentage of young unemployed  
number of journal/magazine subscriptions per family

- (10) Event history data. Analyzed by De Leeuw, Van der Heijden, Kreft.  
940 objects (individuals)  
1440 variables (minutes) with categories - work, including school

- being at home  
- shopping  
- travelling  
- other, including visits, sports,  
culture.

- (11) Whale data. Analyzed by Van der Burg.  
30 objects (whale, porpoise, dolphin species)  
15 variables (neck, snout, head, beak, dorsal fin, flippers, teeth,  
feeding, spout/hole, colour, cervical vertebrae, lachrymal & jugal  
bones, habitat, throat, head bones).

- (12) Suicide data. Taken from German social statistics, analyzed by  
Van der Heijden and De Leeuw.  
about 50000 objects (successful suicides)  
three variables (sex, age, cause of death)

- (13) Hartigian's Hardware. From the book by Hartigian, analyzed by Giffi.  
24 objects (nails, tacks, screws, bolts).  
6 variables (thread, head, head indentation, bottom, length, brass).

- (14) Multiple choice. Data from Introductory Psychology Course, analysis by Giffi.  
190 objects (students)  
30 variables (items)

301

- (15) Crime and fear. Data from Cozijn and Van Dijk, analyzed by Giffi, and by  
Van de Geer and Meulman.  
1216 objects (respondents)  
10 variables, first six from a questionnaire on crime prevention methods,  
last four background

1. Re-education of offenders
2. Locking-up offenders
3. More severe punishments
4. Social work and rehabilitation
5. Labor camps
6. Better employment for potential offenders

First six variables five-point rating scaling from 5 = 'very effective'  
to 1 = 'very ineffective'.

7. Religion (5 cats)
8. Voting (10 cats)
9. Occupational status (high to low, 7 cats)
10. Age (low to high, 6 cats)

- (16) Votes in Leiden. Analyzed by Van der Heijden.  
Contingency table of 58 districts x 9 political parties, number of votes.

- (17) Books. Constructed and analyzed by Giffi.  
Contingency table of 20 books x 7 topics, number of pages.

- (18) Intelligence vs schooling. Dutch Army Data, analyzed by Meester and  
De Leeuw.  
Contingency table of 8 levels of schooling x 6 levels of intelligence,  
cells number of conscripts.

- (19) Shoplifting data. Dutch Central Bureau of Statistics Data, analyzed  
by Israels and Sikkal.  
Compound contingency table of (2 x 9) sex-age combinations by 13  
types of goods stolen. Cells number of suspected persons.

(20) Spot patterns. Psychophysical Data from Guilford. Analyzed by Giffi. Contingency table with 23 cards with varying number of spots by nine 'equal-appearing' response categories.

(21) Munsingen-Rain data. Archeological data, due to Hodson. Analyzed by Meulman. Presence-absence data. 70 varieties of objects are/are not in each of 59 graves.

(22) School-career data. From the SMW0-survey of the Central Bureau of Statistics.

5464 objects (pupils in secondary education)

5 variables describing school career, all with 25 categories

6 passive background variables

- province
- sex
- advice teacher
- occupational status father
- score arithmetic test
- score language test

(23) Journal preference data. Collected by Roskam. Analyzed by Giffi. Ranking data. 39 psychologists rank 10 psychology journals.

(24) Cylinder data. Using an idea of Thurstone. Analyzed by Giffi, and by Van Rijckevorsel.

20 objects (cylinders)

10 variables, all functions of height and radius

(25) macroeconomical variables. Analyzed by Van Kooten and Van Rijckevorsel.

31 objects (years),

6 variables (indices) - factor shares,

- unemployment
- profits
- institutional constraints
- prices
- public sector

(26) GL0 64-65 data. Collected by Dutch Central Bureau of Statistics. Analyzed by Meester and De Leeuw.

10455 objects (pupils)

6 variables - test score 6th grade

- teachers advice 6th grade
- choice of secondary education
- final form of secondary education
- occupational status father
- sex

(27) Gibbs-Wilson data. From the work of Willard Gibbs, analyzed by Giffi, analyzed earlier by E.B. Wilson.

65 objects (experiments)

3 variables - temperature

- pressure
- density

(28) Well-being research. Data collected by Dutch Central Bureau of Statistics. Analysis by Israels.

?? Objects

6 variables - satisfaction (5 cats)

- happiness (5 cats)
- marital status (4 cats)
- schooling (4 cats)
- income (4 cats)
- activity (7 cats)

(29) Comparative school careers. Two subsamples from the GL0 64-65 and the SMW0 1977 cohorts, data by Netherlands Central Bureau of Statistics. Analysis by De Leeuw, Van der Burg, and Bettonvil.

1788 + 1519 objects (pupils)

6 variables: choice of secondary education (5 cats)

- test score 6th grade (6 cats)
- advice teacher 6th grade (5 cats)
- school level father (7 cats)
- school level mother (7 cats)
- fathers profession (7 cats)

(30) Russett data. Analyzed by Giffi.

47 objects (countries)

9 variables - Gini index of income inequality

- Percentage of farmers owning half of the land, starting from below
- Percentage of farmers that rent all their land
- gross national product per capita
- percentage of labour force in agriculture
- average number of years in office of first executive
- total number of violent internal war accidents in 1946-1961.
- total number of people killed due to internal conflict
- stable democracy, unstable democracy, dictatorship

(31) POLP data, interview of members of Dutch parliament. Analyzed by Van der Burg and De Leeuw.

138 objects (MP's)

7 + 4 variables.

First seven nine-point rating scales, measuring attitude towards development aid, abortion, law and order, income differences, participation of workers in industrial decision making, tax increases, money for defence. Next four: rank number of PvdA (socialists), ARP (christian democrats, protestant), KVP (christian democrats, catholic), VVD (liberals, conservative) in a rank order of all 15 parties.

(32) Population and emancipation. Data of the Netherlands Social Cultural Planning Bureau, analyzed by Kreft en De Leeuw.

4693 objects (random sample)

4 + 2 + 2 variables.

First four attitude questions concerning birth control, and the role of government in giving information about birth control. Then the background variables sex x political choice (2 x 6 categories) and age (7 categories), and two attitude questions concerning emancipation.

