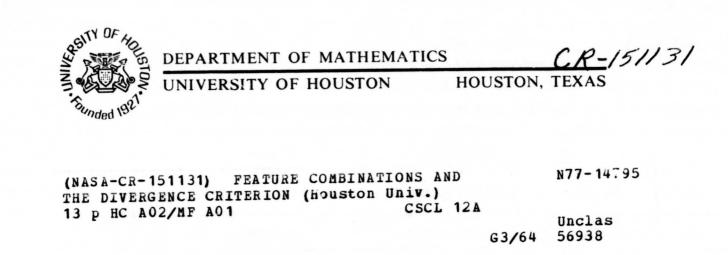
General Disclaimer

One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.
- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.
- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.
- This document is paginated as submitted by the original source.
- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.

Produced by the NASA Center for Aerospace Information (CASI)



FEATURE COMBINATIONS AND THE DIVERGENCE CRITERION BY H.P. DECELL AND S.M. MAYEKAR REPORT #56 JUNE 1976



1

PREPARED FOR EARTH OBSERVATION DIVISION, JSC UNDER CONTRACT NAS-9-15000

HOUSTON, TEXAS 77004

FEATURE COMBINATIONS AND THE DIVERGENCE CRITERION

ł

by

Henry P. Decell, Jr. and Shailesh M. Mayekar Department of Mathematics University of Houston Houston, Texas

> Report #56 NASA Contract NAS-9-15000

FEATURE COMBINATIONS AND THE DIVERGENCE CRITERION

Henry P. Decell, Jr. and Shailesh M. Mayekar

Department of Mathematics University of Houston

ABSTRACT

Classifying large quantities of multidimensional data (e.g., remotely sensed agricultural data)(Remote, 1968) requires efficient and effective classification techniques and the construction of certain transformations of a dimension-reducing, informationpreserving nature. This paper will deal with the construction of transformations that minimally degrade information (i.e., class separability). We will only consider the construction of linear dimension-reducing transformations for multivariate normal populations and information content will be measured by divergence (Kullback, 1968).

1. INTRODUCTION

For n-dimensional normal classes $N(m_i, V_i)$ i = 1,...,m, the divergence between class i and j (Kullback, 1968) is given by

$$D_{ij} = \frac{1}{2} \operatorname{tr}[(V_i - V_j)(V_j^{-1} - V_i^{-1})] + \frac{1}{2} \operatorname{tr}[(V_i^{-1} + V_j^{-1})(m_i - m_j)(m_i - m_j)^T]$$
Let $\delta_{ij} = m_i - m_j$. Then
$$D_{ij} = \frac{1}{2} \operatorname{tr}[(V_i - V_j)(V_j^{-1} - V_i^{-1})] + \frac{1}{2} \operatorname{tr}[(V_i^{-1} + V_j^{-1})(\delta_{ij})(\delta_{ij})^T]$$

$$= \frac{1}{2} \operatorname{tr}[V_i^{-1}(V_j + \delta_{ij}\delta_{ij}^T)] + \frac{1}{2} \operatorname{tr}[V_j^{-1}(V_i + \delta_{ij}\delta_{ij}^T)] - n.$$

The <u>interclass</u> <u>divergence</u> (Decell and Quirein, Oct. 1973) for m populations is given by

$$D = \sum_{i=1}^{m-1} \sum_{j=1}^{m} D_{ij}$$
$$i \neq j$$

and it follows that

:

$$D = \frac{1}{2} \operatorname{tr} \left[\sum_{i=1}^{m} V_i^{-1} \left(\sum_{\substack{j=1 \ i \neq j}}^{m} (V_j + \delta_{ij} \delta_{ij}^T) \right) \right] - \frac{m(m-1)}{2} n$$
$$= \frac{1}{2} \operatorname{tr} \left[\sum_{i=1}^{m} V_i^{-1} S_i \right] - \frac{m(m-1)}{2} n,$$

where

$$\mathbf{s_i} = \sum_{\substack{\mathbf{j=1}\\\mathbf{i\neq j}}}^{m} (v_{\mathbf{j}} + \delta_{\mathbf{ij}} \delta_{\mathbf{ij}}^{T}).$$

If B is a $k \times n$ rank k matrix, the <u>B-interclass diver</u><u>gence</u> (Decell and Quirein, Oct. 1973) is given by

•.

$$\mathbf{D}_{\mathbf{B}} = \sum_{\mathbf{i}=1}^{m-1} \sum_{\substack{\mathbf{j}=1\\\mathbf{i}\neq\mathbf{j}}}^{m} \mathbf{D}_{\mathbf{B}}(\mathbf{i},\mathbf{j})$$

$$D_{B} = \frac{1}{2} tr \left[\sum_{i=1}^{m} (BV_{i}B^{T})^{-1} (BS_{i}B^{T})\right] - \frac{m(m-1)}{2} k.$$

As in the case of average interclass divergence, the B-interclass divergence is a measure of the "separation" in the classes $N(Bm_i, BV_iB^T)$ i = 1,...,m, and is a useful tool for constructing rank k linear transformations that preserve "class separability". It has been shown (Decell and Quirein, Oct. 1973) that whenever $D = D_B$, the probability of misclassification (Anderson, 1958) for the classes $N(Bm_i, BV_iB^T)$, i = 1,...,m is the same as the probability of misclassification for the classes $N(m_i, V_i)$, i = 1,...,m.

2. THEORETICAL PRELIMINARIES

We will assume that k is an integer (k < n) and develop a procedure for selecting a k x n rank k matrix B such that D_B is maximum. The procedure will be based upon the following theorem (Decell and Smiley, to appear). We will let $C = \{u \in R^n : ||u||=1\}$ and $T(H) = \{H = I-2uu^T : u \in C\}$ denote the set of Householder transformations defined on R^n (Householder, 1968).

Theorem. For each positive integer i let $H_i \in T(H)$ be inductively chosen such that

$${}^{\mathrm{D}}(\mathbf{I}_{k}|\mathbf{Z})\mathbf{H}_{i}\mathbf{H}_{i-1}\cdots\mathbf{H}_{1} = \frac{1.u.b.[D}{\mathrm{HeT}(\mathbf{H})}(\mathbf{I}_{k}|\mathbf{Z})\mathbf{H}\mathbf{H}_{i-1}\cdots\mathbf{H}_{1}]$$

where

. . .

$$D(I_k|Z)H_1 = 1.u.b. D(I_k|Z)H$$

I,

The following hold:

- (1) $D(\mathbf{I}_{\mathbf{k}}|\mathbf{Z})\mathbf{H}_{\mathbf{i}}\mathbf{H}_{\mathbf{i}-1}\cdots\mathbf{H}_{\mathbf{1}} \stackrel{\leq}{=} D(\mathbf{I}_{\mathbf{k}}|\mathbf{Z})\mathbf{H}_{\mathbf{i}+1}\mathbf{H}_{\mathbf{i}}\cdots\mathbf{H}_{\mathbf{1}}$
- (2) $D(\mathbf{I}_{k}|Z)H_{i}H_{i-1}\cdots H_{1}H \stackrel{\leq}{=} D(\mathbf{I}_{k}|Z)H_{i+1}H_{i}\cdots H_{1}$, for every $H \in T(H)$.

(3)
$${}^{D}(I_{k}|Z)HH_{i}H_{i-1}\cdots H_{1} \stackrel{\leq}{=} {}^{D}(I_{k}|Z)H_{i+1}H_{i}\cdots H_{1}$$
, for every $H \in T(H)$.
(4) ${}^{D}(I_{k}|Z)H_{i}H_{i-1}\cdots H_{i-(p-1)}HH_{i-(p+1)}\cdots H_{1} \stackrel{\leq}{=} {}^{D}(I_{k}|Z)H_{i+1}\cdots H_{1}$,
for every $H \in T(H)$, $\mu = 0, 1, \dots, i-2$.

(5) The monotone sequence

$$\{D_{B_{i}}\}_{i=1}^{\infty} \equiv \{D_{(I_{k}|Z)H_{i}\cdots H_{i}}\}_{i=1}^{\infty}$$
 is bounded above,

and hence

$$\lim_{i\to\infty} D(\mathbf{I}_{k}|\mathbf{Z})\mathbf{H}_{i}\cdots\mathbf{H}_{1} = 1.u.b. \{D(\mathbf{I}_{k}|\mathbf{Z})\mathbf{H}_{i}\cdots\mathbf{H}_{1}\}.$$

We would, of course, be pleased if it were the case that 1.u.b. $\{D_{(I_k \mid Z)H_1} \cdots H_1\} = D$. This, unfortunately, is not always i the case for some choice of k < n and is not possible, in general, for any k < n. We do know that there is some $k \times n$ rank kmatrix B for which D_B is maximum and, in general, that $D_B \leq D$ (Decell and Quirein, Oct. 1973). It follows, moreover, that since the matrices of the form $(I_k \mid Z)H_1 \cdots H_1$ have rank k, ...

 $D(I_k|Z)H_1 \cdots H_1 \stackrel{\leq}{=} D_B \stackrel{\leq}{=} D$ for every integer i.

We will call the sequence ${D(I_k|Z)H_1\cdots H_1}^{\infty}$ suboptimal whenever

1.u.b.
$$\{D_{[\mathbf{I}_k|Z]H_1} \cdots H_1\} < D_B$$

(and optimal in the case of equality).

There are several open theoretical questions that deal with the conjecture that the sequence is, in general, optimal and cofinally constant beyond the index $i = \min\{k, n-k\}$ (Decell and Smiley, to appear). In what follows we will develop a procedure for constructing the subject sequence and demonstrate its application to agricultural data.

3. THE GRADIENT OF D

It has been shown (Quirein, Nov. 1972) that the differential dD_B of D_B (regarded as a function of the k × n matrix 3) can be expressed in the form $dD_B = F + G$, where, when the indicated inverses exist,

4

1

$$F = \frac{1}{2} tr[\sum_{i=1}^{m} (BV_{i}B^{T})^{-1} (dB S_{i}B^{T} + BS_{i}dB^{T})]$$

$$= \frac{1}{2} tr[\sum_{i=1}^{m} (dB S_{i}B^{T}) (BV_{i}B^{T})^{-1}]$$

$$+ \frac{1}{2} tr[\sum_{i=1}^{m} (BS_{i} dB^{T}) (BV_{i}B^{T})^{-1}]$$

$$= tr[\sum_{i=1}^{m} (dB S_{i}B^{T}) (BV_{i}B^{T})^{-1}]$$

and

$$G = -\frac{1}{2} tr \left[\sum_{i=1}^{m} (BV_{i}B^{T})^{-1} (dB \ V_{i}B^{T} + BV_{i}dB^{T}) (BV_{i}B^{T})^{-1} (BS_{i}B^{T})\right]$$

$$= -\frac{1}{2} tr \left[\sum_{i=1}^{m} (dB \ V_{i}B^{T}) (BV_{i}B^{T})^{-1} (BS_{i}B^{T}) (BV_{i}B^{T})^{-1}\right]$$

$$-\frac{1}{2} tr \left[\sum_{i=1}^{m} (BV_{i}B^{T})^{-1} (BS_{i}B^{T}) (BV_{i}B^{T})^{-1} (BV_{i}dB^{T})\right]$$

$$= - tr \left[\sum_{i=1}^{m} (dB \ V_{i}B^{T}) (BV_{i}B^{T})^{-1} (BS_{i}B^{T}) (BV_{i}B^{T})^{-1}\right].$$

Thus,

$$dD_{B} = tr[\sum_{i=1}^{m} dB\{S_{i}B^{T} - V_{i}B^{T}(BV_{i}B^{T})^{-1}(BS_{i}B^{T})\}(BV_{i}B^{T})^{-1}]$$

= tr $\sum_{i=1}^{m} dB Q_{i}$

where

$$Q_{i} = [\{S_{i}B^{T} - V_{i}B^{T}(BV_{i}B^{T})^{-1}(BS_{i}B^{T})\}(BV_{i}B^{T})^{-1}].$$

We are, of course, interested in extremizing D_B over the particular subclass of $k \times n$ rank k matrices of the form $(I_k|Z)H$ where $H \in T(H)$ (e.g., for i = 1 we find H_1 that maximizes $D_{(I_k|Z)H}$). Actually, one need only consider what is required to compute H_1 . The computation of H_2 is accomplished by the same procedure as that for H_1 . It is simply a matter of, after selecting H_1 , redefining the m classes to be $N(H_1m_1, H_1V_1H_1)$, $i = 1, \ldots, m$ and proceeding as in the selection of H_1 .

With these facts in mind we will simply calculate the gradient of D_B where B is restricted to having the form $B = (I_k | Z)H$, H \in T(H). The restrictions H \in T(H) can be accomplished by considering those $k \times n$ rank k matrices of the form

$$B = (I_k^{\dagger}Z)(I - 2\frac{ww^{\dagger}}{w^{T}w}), \quad w \in \mathbb{R}^n (w \neq \theta)$$

It follows that

-

$$dB = d[(I_k|Z) (I - 2 \frac{ww^T}{w^Tw})] = -2(I_k|Z) d(ww^T/w^Tw)$$
$$= -2(I_k|Z) [\frac{w^Twd(ww^T) - ww^Td(w^Tw)}{(w^Tw)^2}]$$

$$= -\frac{2(I_{k}|Z)}{(w^{T}w)^{2}}[w^{T}w(dw w^{T} + wdw^{T}) - ww^{T}(w^{T}dw + dw^{T} w)]$$

$$= -\frac{2(I_{k}|Z)}{(w^{T}w)^{2}}[dw w^{T}w w^{T} + w w^{T}w dw^{T} - w w^{T}dw w^{T} - w dw^{T}w w^{T}]$$

$$= -\frac{2(I_{k}|Z)}{(w^{T}w)^{2}}[(dw w^{T} - wdw^{T})ww^{T} - ww^{T}(dw w^{T} - wdw^{T})]$$

Substituting the latter in the expression for $dD_B^{}$,

$$dD_{B} = tr \sum_{i=1}^{m} \left\{ -\frac{2(I_{k}|Z)}{(w^{T}w)^{2}} \left\{ (dw w^{T} - wdw^{T})ww^{T} - ww^{T}(dw w^{T} - wdw^{T}) \right\} Q_{i} \right\}$$

$$= tr \sum_{i=1}^{m} \left[-\frac{2Q_{i}(I_{k}|Z)}{(w^{T}w)^{2}} \left\{ (dw w^{T} - wdw^{T})ww^{T} - ww^{T}(dw w^{T} - wdw^{T}) \right\} \right]$$

$$= tr \sum_{i=1}^{m} \frac{-2}{(w^{T}w)^{2}} [ww^{T} Q_{i} (I_{k}|Z) (dw w^{T} - wdw^{T})$$

$$- Q_{i}(I_{k}|Z)ww^{T}(dw w^{T} - wdw^{T})]$$

$$= \frac{-2}{(w^{T}w)^{2}} tr \sum_{i=1}^{m} [M_{i}dw w^{T} - M_{i}wdw^{T} - N_{i}dw w^{T} + N_{i}wdw^{T}]$$

$$Where M_{i} = ww^{T}Q_{i}(I_{k}|Z) and N_{i} = Q_{i}(I_{k}|Z)ww^{T}.$$

$$\vdots$$

$$dD_{B} = \frac{-2}{(w^{T}w)^{2}} tr [\sum_{i=1}^{m} \{w^{T} M_{i} dw - w^{T} N_{i} dw + N_{i} w dw^{T} - M_{i} w dw^{T}\}]$$

$$= \frac{-2}{(w^{T}w)^{2}} tr [\sum_{i=1}^{m} \{dw^{T} M_{i}^{T} w - dw^{T} N_{i}^{T} w + N_{i} w dw^{T} - M_{i} w dw^{T}\}]$$

$$dD_{B} = \frac{-2}{(w^{T}w)^{2}} \operatorname{tr} \left[\sum_{i=1}^{m} \{ M_{i}^{T} w dw^{T} - N_{i} w dw^{T} + N_{i} w dw^{T} - M_{i} w dw^{T} \} \right]$$
$$= \frac{-2}{(w^{T}w)^{2}} \operatorname{tr} \left[\sum_{i=1}^{m} \{ (M_{i} - N_{i})^{T} - (M_{i} - N_{i}) \} w dw^{T} \} \right].$$

The necessary condition that w be extremal is then,

$$G(w) = \frac{-2}{(w^{T}w)^{2}} \sum_{i=1}^{m} \{ (M_{i} - N_{i})^{T} - (M_{i} - N_{i}) \} w = \theta \text{ (the zero vector).}$$

We note that G(w) is the gradient of $D(I_k|Z)(I - 2\frac{ww^T}{w^T})$ and

use a steepest descent procedure for finding the extremal w. The process is repeated for each sequential index until corresponding values of divergence "ctabilize." Test results are presented in the following tables. The C-1 flight line data is twelve channel data for nine agricultural classes: soybeans, corn, oats, redclover, alfalfa, rye, bare soil, and two types of wheat. The Hill County data is sixteen-channel data for five agricultural classes: winter wheat, fallow crop, barley, grass, and stubble.

The starting value w_0 for the steepest descent procedure for selecting each successive Householder transformation H_1, H_2, H_3, \dots was arbitrarily chosen to be $w_0 = (\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}})^T$. Choosing starting values in this arbitrary fashion is certainly not the most clever thing to do in the presence of the monotone behavior of the sequence $D_{(I_k | Z)H_1 \cdots H_1}$. One would expect, for example, that the starting values for the selection of H_{i+1} should depend upon the unit vectors previously selected as generators of H_i, H_{i-1}, \dots, H_1 in such a way as to guarantee that the starting value w_0 , for the descent procedure for selecting H_{i+1} , satisfies

$$\mathbf{D}(\mathbf{I}_{\mathbf{k}}|\mathbf{Z})\mathbf{H}_{\mathbf{i}}\cdots\mathbf{H}_{\mathbf{1}} \stackrel{\leq \mathbf{D}}{=} (\mathbf{I}_{\mathbf{k}}|\mathbf{Z})(\mathbf{I} - 2\frac{\mathbf{w}_{\mathbf{0}}\mathbf{w}_{\mathbf{0}}^{\mathsf{T}}}{\frac{\mathbf{T}_{\mathbf{0}}\mathbf{w}_{\mathbf{0}}}{\mathbf{T}_{\mathbf{0}}})\mathbf{H}_{\mathbf{i}}\cdots\mathbf{H}_{\mathbf{1}}.$$

This rather arbitrary selection of the starting vector does, as the examples demonstrate, violate the latter inequality. The question about how to choose starting vectors, according to the latter inequality, is still an open one and its answer would certainly decrease computation time.

C-1 Flight Line Date n=12, k=6, m=9, D=10,660

Iteration for H

| No * | Divergence D _B |
|------|---------------------------|
| 1 . | 1982 |
| 2 | 3536 |
| 3 | 4533 |
| 4 | 5781 |
| 5 | 6910 |
| 6 | 7522 |
| 7 | 7710 |
| 8 | 7790 |
| 9 | 7838 |
| 10 | 7865 |
| 11 | 7881 |
| 12 | 7892 |

Hill County Data n=16, k=8, m=5, <u>D=636</u>

Iteration for H1

| No * | Divergence | DB |
|------|------------|-----|
| 1 | 114.58 | |
| 2 | 136.66 | |
| 3 | 152.27 | |
| 4 | 179.69 | |
| 5 | 223.81 | • • |
| 6 | 247.42 | |
| 7 | 252.78 | |
| 8 | 257.12 | |
| 9 | 260.74 | |
| 10 | 263.95 | |

.

*Iteration counter

C-1 Flight Line Data (cont.)

1

Iteration for H₂

| No* | Divergence D _B | |
|-----|---------------------------|--|
| 1 | 7815 | |
| 2 | 8797 | |
| 3 | 9542 | |
| 4 | 9785 | |
| 5 | 9901 | |
| 6 | 9966 | |
| 7 | 10,005 | |
| 8 | 10,031 | |
| 9 | 10,048 | |

Hill County Data (cont.)

1

1

Iteration for H2

| No [*] | Divergence D _B |
|-----------------|---------------------------|
| 1 | 269.00 |
| 2 | 280.48 |
| 3 | 293.32 |
| 4 | 300.68 |
| 5 | 304.07 |
| 6 | 306.19 |
| 7 | 307.74 |
| 8 | 308.95 |
| 9 | 309.93 |

Iteration for H₃

Iteration for H_3

| No* | Divergence | DB |
|-----|------------|---------------|
| 1 | 7582 | |
| 2 | 8705 | |
| 3 | 9809 | |
| 4 | 9947 | |
| 5 | 9995 | Sector Sector |
| 6 | 10,020 | |
| 7 | 10,037 | |
| 8 | 10,049 | |
| 9 | 10,058 | |

| No* | Divergence | DB |
|-----|------------|----|
| 1 | 312.18 | •• |
| 2 | 344.52 | |
| 3 | 380.83 | |
| 4 | 387.20 | |
| 5 | 391.70 | |
| 6 | 392.96 | |
| 7 | 394.58 | |
| 8 | 399.47 | |

Iteration for H_4

| No* | Divergence D _B | |
|-----|---------------------------|------|
| 1 | 371.12 | |
| 2 | 394.75 | |
| 3 | 398.62 | |
| 4 | 400.69 | |
| 5 | 402.03 | |
| 6 | 402.98 | |
| 7 | 403.74 | 19-1 |

:

*Iteration counter

BIBLIOGRAPHY

Anderson, T.W. (1958). An Introduction to Multivariate Statistical Analysis. New York: John Wiley & Sons, Inc.

Decell, H.P. and Quirein, J.A. (October 1973). "An Iterative Approach to the Feature Selection Problem". IEEE Cat. #CH0834-2, pp. 3B-1-3B-12.

- Decell, H.P. and Smiley, W. (to appear). "Householder Transformations and Optimal Linear Combinations." Comm. in Stat.
- Householder, A.S. (1968). "Unitary Triangularization of a Non-Symmetric Matrix." J. Assoc. Comput. Mach., pp. 339-342.
- Kullback, S. (1968). Information Theory and Statistics. New York: Dover Publications.
- Quirein, J.A. (November 1972). "Some Necessary Conditions for an Extremum." Report #12 NAS-9-12777. Dept. of Mathematics, Univ. of Houston, Texas
- "Remote Multispectral Sensing in Agriculture." (1968). <u>Report of</u> <u>the Laboratory for Agricultural Remote Sensing</u> Vol. 3, <u>Research Bulletin #844.</u> Purdue Univ., Lafayette, Indiana.