

PUC

Series: Monographs in Computer Science
and Computer Applications

Nº 9/69

IDENTIFICATION OF MINOR SUBSETS PERTAINING TO A
MAJOR SET THROUGH SIMULTANEOUS ANALYSIS OF
RANDOM SAMPLES

by

Fernando Olavo Franciss

Abelardo L. Puccini

Computer Science Department - Rio Datacenter

CENTRO TÉCNICO CIENTÍFICO
Pontifícia Universidade Católica do Rio de Janeiro
Rua Marquês de São Vicente, 209 — ZC-20
Rio de Janeiro — Brasil

IDENTIFICATION OF MINOR SUBSETS PERTAINING TO A
MAJOR SET THROUGH SIMULTANEOUS ANALYSIS OF
RANDOM SAMPLES

CONTENTS

1. Introduction
2. Recommended Solution
3. Conclusion
4. FORTRAN IV program

Fernando Olavo Franciss
Associate Professor
Civil Engineering Dept.

Abelardo L. Puccini
Associate Professor
Computer Science Dept.

1. Introduction

Knowing m random samples X_i of a set R , Being each sample referred to a non-Euclidian space by a function of n parameters X_{ij} , corresponding to its various characteristics, the recognition of subsets $S_I, S_{II}, S_{III}, \dots$ of the set R , is a difficult problem not so far solved statistical analysis (1), because the complexity grows as the number of parameters x_{ij} increases.

2. Recommended Solution

Assuming that the parameters x_{ij} have properties of equivalence, order, interval and value relationships, and are normally distributed, the recommended solution has the following steps:

Step I: Normalization of the Parameters x_{ij}

The normalization of the parameters x_{ij} is necessary to cancel the effects of the parameters being measured by different system of units, and/or when they have different amplitudes of variation.

It is done as following:

a) Define the matrix of $m \times n$ elements x_{ij}

$$\begin{matrix}
 x_{11} & x_{12} & \dots & x_{1n} \\
 x_{21} & x_{22} & \dots & x_{2n} \\
 x_{31} & x_{32} & \dots & x_{3n} \\
 \dots & \dots & \dots & \dots \\
 x_{m1} & x_{m2} & \dots & x_{mn}
 \end{matrix}$$

Which corresponds to the column of m random samples x_i :

$$x_1 \ x_2 \ x_3 \ \dots \ x_m$$

Now determine for each column j , the values \bar{x}_j and s_j , using the formulas:

$$\bar{x}_j = \frac{1}{m} \sum_{i=1}^m x_{ij}$$

$$s_j = \left\{ m \sum_{i=1}^m x_{ij}^2 - \left(\sum_{i=1}^m x_{ij} \right)^2 / m(m-1) \right\}^{1/2}$$

- b) Replace the former values of the elements x_{ij} , by the normalized elements x'_{ij} , as defined below:

$$x'_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j}$$

Step II: Equalization of the Normalized Parameters x'_{ij}

The equalization of the normalized parameters x'_{ij} is necessary to cancel the effects of the parameters corresponding usually to characteristics not equally relevant.

It is done as following:

- a) Fix for each characteristics j , based on its meaning, an interval Δx_j , whose value allows us to consider the parameters x'_{ij} and $x'_{ij} \pm \Delta x_j$ as equivalents.

- b) Normalize the n elements Δx_j using the following formula:

$$\Delta x'_j = \Delta x_j / s_j$$

- c) Obtain the mean interval $\overline{\Delta x'}$ using the following formula:

$$\overline{\Delta x'} = \frac{1}{n} \sum_{j=1}^n \Delta x'_j$$

- d) Replace the $m \times n$ normalized elements x'_{ij} by the normalized and equalized elements x''_{ij} as defined below:

$$x''_{ij} = x'_{ij} \cdot \overline{\Delta x'} / \Delta x'_j$$

Notice that if Δx_j is equal to s_j , then x'_{ij} will be equal to x'_{ij} .

Step III - Obtaining the Distances Between Samples

The determination of the distances between samples is necessary to measure the affinity between the samples:

It is done as following:

- a) Obtain $\frac{m(m-1)}{2}$ distances a_{fk} between the random samples X_f and X_k , using the expression below:

$$a_{fk} = \sum_{j=1}^n (x'_{fj} - x'_{kj})^2 / 2$$

- b) Define the symmetric matrix of m^2 elements a_{fk} , which indicates the distances a_{fk} between the samples X_f and X_k .

$$\begin{array}{c}
 X_1 \quad X_2 \quad \dots \quad X_m \\
 \begin{array}{c|ccc}
 X_1 & 0 & a_{12} \dots a_{1m} \\
 X_2 & a_{21} & 0 \dots a_{2m} \\
 \vdots & \vdots & \\
 X_m & a_{m1} & a_{m2} \dots 0
 \end{array}
 \end{array}$$

Step IV - Searching Particular Subsets

The searching of particular subsets is done as following:

- a) Choose an integer number N (the bigger the value of N the best is the search).

Define a sequence of values A_r , using the expression:

$$A_r = \frac{r}{N} (a_{\max} - a_{\min}) + a_{\min} \quad r = N, N-1, N-2, \dots, 1$$

Where a_{\max} and a_{\min} correspond to the maximum and minimum values of the distances between samples.

- b) For each value of A_r , define subsets $S_I^r, S_{II}^r, S_{III}^r, \dots$ of order r using the following procedure:

- b-1) Make zero all the elements of the matrix of distances a_{fk} that are greater than A_r .
- b-2) Find for each sample X_k (therefore for each column k) the value of the reduced mean distance α_k^r , using the expression:

$$\alpha_k^r = \frac{1}{(m-u_k^r)} \sum_{f=1}^m a_{fk}$$

where u_k^r is the number of distances that were made zero in (b-1) (because they were greater than A_r) at column k .

- b-3) Call α_I^r the minimum value of α_k^r
- b-4) Define the first subset S_I^r by grouping at column I (associated to α_I^r) the samples corresponding to the distances that were not made zero in (b-1).
- b-5) Exclude in the values of α_k^r , the ones whose columns correspond to samples already included in the first subset S_I^r .
- b-6) Call α_{II}^r the minimum value of the remaining α_k^r , and then define the second subset S_{II}^r as was done before.
- b-7) Define new subsets S_{III}^r, \dots using the values of α_k^r not included in previous subsets.

c) Stop the searching of subsets when, with the decrease of r one of the following situations happens:

- c-1) The number of subsets remains unchanged for various values of r without showing any increase of common elements, which confirms their existence as independent entities.
- c-2) The number of subsets remains unchanged for various values of r , showing a progressive reduction of common elements, which confirms their existence as interdependents entities.
- c-3) The number of subsets increases up to the value m , without verifying either situation (c-1) or (c-2), which does not confirm their existence as real entities.

Samples that are not included on the identified subsets, will be considered as marginals.

The repetition of situations (c-1) during the process of searching identifies

secondary subsets (2nd order subset, 3rd order subset, etc.)

3. Conclusion

After the recognition of the existence of subsets, each subset may be submitted to the conventional statistical treatment.

BIBLIOGRAPHY

1. Miller, R.L.; Kahn, J.S.; Statistical analysis in Geological Sciences, Appendix A; John Wiley and Sons, Inc. New York, 1962.


```

DIMENSION X(120,20),XBAR(20),DELX(20),S(20)
COMMON AFAST(120,120),M
1 FORMAT(3I4)
2 FORMAT(8F10.4)
3 FORMAT(1H1,16HNUM DE AMOSTRAS=,I4,20H NUM DE PARAMETROS=,I4,2X
116HN ARBITRADO COMO, I4, /5X,9HDELTA(J)=,I, /10X,F10.4))
4 FORMAT(1H1,10HMATRIZ X1J)
5 FORMAT(1H ,7HAMOSTRA, I4 /((10X,8F10.4))
6 FORMAT(1H1,18HMATRIZ XLINHA(I,J))
7 FORMAT(1H1,19HMATRIZ X2LINHA(I,J))
10 FORMAT(1H0,10HAFAST MAX=,F10.2,3X,10HAFAST MIN=,F10.2)
C INICIO DAS ETAPAS I E II
1 READ(5,1)M,N,NN
2 DO 444 I=1,M
3 READ(5,2)((X(I,J),J=1,N))
4 DO 444 J=1,N
5 WRITE(6,3)M,N,NN,(DELX(J),J=1,N)
6 WRITE(6,4)
7 DO 333 I=1,M
8 333 WRITE(6,5)I,(X(I,J),J=1,N)
C CALCULO DE XBAR(J) S(J) E DA MATRIZ XLINHA(I,J)
9 DO 100 J=1,N
10 SOMAX=0.0
11 SOMAQX=0.0
12 DO 200 I=1,M
13 SOMAX=SOMAX+X(I,J)
14 200 SOMAQX=SOMAQX+X(I,J)**2
15 XBAR(J)=SOMAX/FLOAT(M)
16 S(J)=SQRT((FLOAT(M)*SOMAQX-SOMAX**2)/(FLOAT(M)*FLOAT(M-1)))
17 DO 210 I=1,M
18 210 X(I,J)=(X(I,J)-XBAR(J))/S(J)
19 IF(DELX(J).EQ.0.) DELX(J)=S(J)
20 DELX(J)=DELX(J)/S(J)
21 100 CONTINUE
22 WRITE(6,6)
23 DO 127 I=1,M
24 127 WRITE(6,5)I,(X(I,J),J=1,N)
C CALCULO DA MATRIZ X2LINHA(I,J)
25 SOMAD=0.0
26 DO 220 J=1,N
27 220 SOMAD=SOMAD+DELX(J)
28 DBARL=SOMAD/FLOAT(N)
29 DO 230 J=1,N
30 DO 230 I=1,M
31 230 X(I,J)=X(I,J)*(DBARL/DELX(J))
32 1 WRITE(6,7)
33 DO 128 I=1,M
34 128 WRITE(6,5)I,(X(I,J),J=1,N)
C FIM DAS ETAPAS I E II
C INICIO DA ETAPA III DETERM AFASTAMEN, E AFAST MAX EMIN
35 AEMIN =32000.
36 AFMAX=0.0
37 DO 240 LF =1,M
38 DO 250 LK =1,M
39 SOMAC=0.0
40 DO 260 J=1,N
41 260 SOMAC=SOMAC + (X(LF,J)-X(LK,J))**2

```

```

AFAST(LF,LK)=SORT(SOMAO)
IF(LF-LK) 888,250,888
888 IF(AFAST(LF,LK)-AFMIN)88,88,89
88 AFMIN=AFAST(LF,LK)
GO TO 250
89 IF(AFAST(LF,LK)-AFMAX)250,188,188
188 AFMAX=AFAST(LF,LK)
250 CONTINUE
240 CONTINUE
WRITE(6,10)AFMAX,AFMIN
C FIM DA ETAPA III INCLUINDO OS AFASTAMENTOS MAXIMO E MINIMO
C INICIO DA ETAPA IV
DELTA=AFMAX-AFMIN
DO 55 IR=1,NN
AR =AFMIN+(FLOAT(NN-IR+1)/FLOAT(NN))*DELTA
WRITE(6,4422)IR
4422 FORMAT(1H1,31HREALIZACAO DA ETAPA IV PARA R =,I3,13H DU SEJA PARA,
1/)
CALL ETAP4(AR)
55 CONTINUE
CALL EXIT
END

```

```

SUBROUTINE ETAP4(AR)
DIMENSION ALFKR(120),ISIR(120)
COMMON AFAST(120,120),M
11 FORMAT(1H0,20HVALORES DE ALFKR(LK),/(5X,20F5.2) )
12 FORMAT(1H ,4HAR =,F10.4)
13 FORMAT(1H0,5X,12HAEAMINIMO =,F10.4 )
14 FORMAT(1H0,5X,56HAMOSTRA PRINCIPAL DO SUBCONJUNTO QUE SE SEGUE E A
1 NUMERO,14/)
DO 2000 LK=1,M
IUKR=0
DO 1000 LF=1,M
IF(AFAST(LF,LK))24,24,301
24 IUKR=IUKR+1
GO TO 1000
301 IF (AFAST(LF,LK)-AR )1000,1000,400
400 AFAST(LF,LK)=0.0
IUKR=IUKR+1
1000 CONTINUE
C SUBTRAIR 1 DE UKR PARA EXCLUIR DISTANCIA ENTRE AMOSTRAS IDENTICAS
IUKR=IUKR-1
C CALCULO DOS ALFA(I,R)
SOMAF=0.0
DO 1010 LF=1,M
1010 SOMAF=SOMAF+ AFAST(LF,LK)
2000 ALFKR(LK)=SOMAF/FLOAT(M-IUKR)
WRITE(6,12)AR
WRITE(6,11)(ALFKR(LK),LK=1,M)

```

1 C. PESQUISA DE ALFAKR(LK) MINIMO

```

4000 ALFMI= 32000.
      DO 1050 LK=1,M
        IF(ALFKR(LK))1544,1050,544

```

544	IF(ALFKR(LK)-ALFMI)44,44,1050								
44	ALFMI=ALFKR(LK)								
	ICOL=LK								
1050	CONTINUE								
	WRITE(6,13)ALFMI								
	C. DETERMINACAO DOS SUBCONJUNTOS S(I,R)								
	WRITE(6,14)ICOL								
	II=0								
	DO 1099 LF=1,M								
	IF(LF-ICOL)99,2099,99								
99	IF(AFAST(LF,ICOL))2099,1099,2099								
2099	ALFKR(LF)=0,0								
	II=II+1								
	ISIR(II)=LF								
1099	CONTINUE								
15	FORMAT(1H,5X,16HAMOSTRAS NUMEROS,7(23X,2014))								
	K=0								
	DO 2355 LK=1,M								
	IF(ALFKR(LK))2354,2355,2354								
2354	K=1								
2355	CONTINUE								
	WRITE(6,11)(ALFKR(LK),LK=1,M)								
	IE(K)5000,5000,4000								
5000	RETURN								
	END								