

KONSTANTIN MOMIROVIĆ I ŽIVAN KARAMAN  
Sveučilišni računski centar, Zagreb

Primljeno 13. 4. 1982.

**INDIFF — MODEL, ALGORITAM I PROGRAM ZA  
ANALIZU PROMJENA STANJA NEKOG OBJEKTA  
OPISANOG NAD SKUPOM KVANTITATIVNIH  
VARIJABLI**

*Priloženo  
Karaman E*

**SAŽETAK**

Konstruirani su model i algoritam i napisan je program za komponentnu analizu promjena stanja individualnih objekata. Dekomponiranjem trajektorije promjena na ortogonalne komponente dobiva se uvid u strukturu stohastičkog procesa opisanog kvantitativnim varijablama. Također se analiziraju reziduali kao i relacije među vremenskim točkama.

**0. UVOD**

Model, algoritam i program, koji su ovdje opisani, namijenjeni su komponentnoj analizi promjena stanja nekog objekta opisanog nad skupom kvantitativnih varijabli registriranih kroz neki vremenski period. Potreba za takvom vrstom analize javlja se prilikom istraživanja u kineziologiji, medicini, ekonomiji i drugim znanostima, onda kada je moguće opisati promjene stanja nekog entiteta nekim skupom izmjerivih varijabli. Tako, npr., u kineziologiji to može biti neki sportaš kojemu se u toku treninga mjere fiziološke, motoričke i druge promjene; u medicini to može biti pacijent kojemu se kroz nekoliko dana ili tjedana svakih nekoliko sati registriraju promjene pulsa, krvnog pritiska, temperature i druge medicinski relevantne promjene; u ekonomiji to može biti prodaja nekih proizvoda u nekoj zemlji praćena mjesčno unatrag nekoliko godina, itd.

Tip modela koji je ovdje predložen namijenjen je analizi nestacionarnih stohastičkih procesa; stacionarne stohastičke procese moguće je analizirati pod modelom koji je implementiran u programima serije TALAMBAS. Promjene stanja objekta koje su originalno opisane nekim skupom od  $n$  manifestnih varijabli objašnjavaju se gotovo posve sa  $p < n$  latentnih dimenzija — komponenata promjena. Na taj način dobivamo bitno jednostavniji prikaz promjena stanja proučavanog objekata. Ovakvom analizom dobivamo također uvid u strukturu komponenata promjena, tj. odnose registriranih varijabli i komponenata promjena, kao i jednostavan prikaz relacija među vremenskim točkama. Model je izvorno predložen od L. Tučkera, a opis originalni modela i algoritma nalazi se u Momirović (1972). U ovom radu dan je detaljan opis algoritma koji je implementiran u programu INDIFF, napisanom u meta-jeziku SS (Zakrajšek, Štalec, Momirović, 1974) i implementiranom na računaru UNIVAC 1100 Sveučilišnog računskog centra u Zagrebu.

**1. MODEL**

Neka je nekim eksperimentom dobivena matrica trajektorije promjena  $B$  tipa  $m \times n$ , gdje je  $m$  broj vremenskih točaka, a  $n$  broj varijabli. Pri tome moraju biti zadovoljeni slijedeći uvjeti:

- a)  $m > n$
- b) rang  $(B) = n$  (tj. varijable  $v_1, \dots, v_n$  su linearno nezavisne)
- c) vremenske točke su ekvidistantne
- d) interval između ma koje dvije susjedne vremenske točke je dovoljno kratak da se promjene stanja u svakoj od  $n$  varijabli mogu aproksimirati nekom monotonom funkcijom.

Varijable iz  $B$  standardiziramo i matricu standardiziranih varijabli označimo sa  $Z$ .

Standardiziranu trajektoriju promjena  $Z$  dekomponiramo na  $p < n$  aditivnih ortogonalnih komponenata promjena  $K_1, K_2, \dots, K_p$ , definiranih na skupu vremenskih točaka, takvih da vrijedi

$$K_1 = \sum_{j=1}^n x_{1j} v_j$$

je linearna kombinacija standardiziranih varijabli  $v_1, v_2, \dots, v_n$  čija je varijanca veća od varijance ma koje druge linearne kombinacije čiji koeficijenti zadovoljavaju

$$\text{uvjet } \sum_{j=1}^n x_{1j}^2 = 1$$

$$K_2 = \sum_{j=1}^n x_{2j} v_j$$

je linearna kombinacija standardiziranih varijabli čija je varijanca veća od varijance  $m$  koje druge linearne kombinacije varijabli ortogonalne na  $K_1$ , a čiji koeficijenti

$$\text{zadovoljavaju uvjet } \sum_{j=1}^n x_{2j}^2 = 1.$$

$$K_p = \sum_{j=1}^n x_{pj} v_j$$

je linearna kombinacija standardiziranih varijabli čija je varijanca veća od varijance ma koje druge linearne kombinacije varijabli ortogonalne na  $K_1, K_2, \dots, K_{p-1}$ , a čiji

$$\text{koeficijenti zadovoljavaju uvjet } \sum_{j=1}^n x_{pj}^2 = 1.$$

Broj zadržanih komponenata  $p$  određen je u skladu sa PB kriterijem (Štalec, Momirović, 1971).

Konkateniramo li vektore  $K_i$ ,  $i=1, 1, \dots, p$  u matricu  $K$  (tipa  $m \times p$ ), dobijemo kondenzirani prikaz trajektorije promjena koji sadrži gotovo sve informacije kao i polazna matrica  $Z$ , a mnogo je pregledniji i pogodniji za analizu.

Strukturom komponenata promjena nazivamo korelacije između standardiziranih varijabli ( $Z$ ) i komponenata ( $K$ )

$$H = Z^T K \frac{1}{m}$$

a relacije vremenskih točaka definirane su kao normirani skalarni produkti vektora vremenskih točaka

$$Q = (\text{diag } ZZ^T)^{-1/2} ZZ^T (\text{diag } ZZ^T)^{-1/2}$$

## 2. ALGORITAM

### 2.1. Parametri varijabli

Razmotrimo sada поближе eksperimentom dobijenu matricu  $B = (b_{ij})$ ,  $i=1, \dots, m$ ;  $j=1, \dots, n$ . Svaki redak matrice  $B$  predstavlja jednu vremensku točku  $t_i$ ,  $ET = \{t_i, i=1, 2, \dots, m\}$  u kojoj su registrirane varijable, dok stupci matrice  $B$  odgovaraju pojedinim varijablama  $v_j$ ,  $EV = \{v_j, j=1, \dots, n\}$  praćenim kroz niz vremenskih točaka. Kako ne možemo očekivati da varijable budu u istoj metrcici, potrebno ih je na neki način standardizirati da bismo ih mogli međusobno uspoređivati. To napravimo tako da svaka standardizirana varijabla ima očekivanje 0 i varijancu 1.

$$v_j \rightarrow v_j$$

$$E(v_j) = 0 \quad j = 1, 2, \dots, n$$

$$\text{var}(v_j) = 1.$$

Ako sa  $1^T = (1, 1, \dots, 1)$  označimo  $m$ -dimenzionalni sumacioni vektor, tada je matrica kovarijanci varijabli iz  $V$  dana sa  $C = (B^T B - B^T J B) \frac{1}{m}$ , gdje je  $J = 1^T \frac{1}{m}$ . Matricu varijanci varijabli iz  $V$  označimo sa  $S^2 = \text{diag } C$ , pa je matrica standardiziranih podataka  $Z = (B - J B) S^{-1}$ .

Matrica korelacija varijabli  $v_j$  iz  $V$  biti će

$$R = S^{-1} C S^{-1} = Z^T Z \frac{1}{m}$$

Parcijalne korelacije varijabli iz  $V$  dane su sa

$$P = -UR^{-1}U.$$

gdje je  $U$  dijagonalna matrica definirana sa  $U^2 = (\text{diag } R^{-1})^{-1}$ .

U matrici  $U^2$  nalaze se zapravo univiteti varijabli  $v_j$  iz  $V$ .

Pogledajmo sada detaljnije uvjete navedene u 1.a) — 1.d). Prva dva uvjeta (1.a) i (1.b) nam omogućuju izra-

čunavanje univiteteta i preko toga primjenu PB kriterija za određivanje broja zadržanih glavnih komponenata. Ukoliko ta dva uvjeta nisu ispunjena nemoguće je analizirati promjene na ovdje predložen način. U tom slučaju moguće je primijeniti algoritam i program TALAMBAS I koji univitetete varijabli računa primjenom generaliziranog Moore-ePnrosoovog (Moore, 1935; ePnrose, 1955) inverza  $R^{-}$  matrice korelacija  $R$ . Druga dva uvjeta, (1.c) i (1.d), nisu formalne prirode, te njihovo neispunjavanje neće onemogućiti provođenje analize, ali će tako dobiveni rezultati biti neupotrebljivi budući su uvjeti (1.c) i (1.d) implicitno sadržani u algoritmu. Stoga je važno osigurati da sva četiri uvjeta (1.a) — (1.d) budu zadovoljena.

### 2.2. Komponente promjena

Standardiziranu trajektoriju promjena  $Z$  dekomponirati ćemo na aditivne, ortogonalne komponente. Tražimo takvu linearnu kombinaciju.

$$K^1 = x_{11} v_1 + x_{12} v_2 + \dots + x_{1n} v_n$$

standardiziranih varijabli  $v_j$  iz  $Z$  koja ima maksimalnu varijancu, tj.

$$\text{var}(K^1) = \max$$

uz uvjet da njeni koeficijenti zadovoljavaju dodatni uvjet

$$\sum_{j=1}^n x_{1j}^2 = 1$$

Označimo li sa  $X^T_1 = (x_{11}, x_{12}, \dots, x_{1n})$  prethodne relacije, možemo u matricnom obliku pisati

$$K_i = Z X_i \quad (1)$$

$$K_i^T K_i \frac{1}{m} = \max \quad (2)$$

$$X_i^T X_i = 1 \quad (3)$$

Budući je

$$K_i^T K_i \frac{1}{m} = X_i^T Z^T Z X_i = X_i^T R X_i \quad (3)$$

Iako se pokaže da se problem maksimizacije (2) uz uvjet (3) svodi na rješavanje karakteristične jednadžbe za  $R$

$$(R - \lambda_i I) X_i = 0 \quad (5)$$

Da bismo dobili netrivialno rješenje ( $X_i \neq 0$ ) nužno je da bude

$$|R - \lambda_i I| = 0,$$

a to znači da  $\lambda_i$  mora biti svojstvena vrijednost matrice  $R$ , pa je onda i  $X_i$  pridružen svojstveni vektor. Da bismo odredili koju svojstvenu vrijednost treba uzeti, pomnožimo (5) sa  $X_i^T$ , uvaživši (3), dobijemo

$$X_i^T R X_i - \lambda_i X_i^T X_i = X_i^T R X_i - \lambda_i = 0$$

tj.

$$X_i^T R X_i = K_i^T K_i \frac{1}{m} = \text{var}(K_i) = \lambda_i \quad (6)$$

pa budući želimo da varijanca od  $K_1$  bude maksimalna uzimamo najveću svojstvenu vrijednost matrice korelacija  $R$ , i osim toga iz (6) vidimo da je varijanca prve komponente upravo jednaka najvećoj svojstvenoj vrijednosti matrice  $R$ .

Za drugu glavnu komponentu

$$K_2 = x_{21}v_1 + x_{22}v_2 + \dots + x_{2n}v_n \quad (7)$$

tražimo takvu linearnu kombinaciju standardiziranih varijabli  $v_j$  iz  $Z$ , čija je varijanca veća od varijance ma koje druge linearne kombinacije varijabli, okomite na  $k_1$  i čiji koeficijenti zadovoljavaju uvjet

$$\sum_{j=1}^n x_j^2 = 1.$$

Označimo li sa  $X_2^T = (x_{21}, x_{22}, \dots, x_{2n})$ , pišemo

$$\text{var}(K_2) = \max$$

$$X_1^T X_2 = 0$$

$$X_2^T X_2 = 1$$

Analognom tehnikom kao i za prvu glavnu komponentu dobijemo da su koeficijenti u (7) jednaki elementima svojstvenog vektora pridruženog drugoj po veličini svojstvenoj vrijednosti matrice  $R$ .

Pokazuje se da općenito vrijedi (Morrison, 1967):  $m$ -ta glavna komponenta je linearna kombinacija varijabli

$$K_l = x_{l1}v_1 + x_{l2}v_2 + \dots + x_{ln}v_n,$$

čiji su koeficijenti elementi svojstvenog vektora matrice korelacija varijabli  $R$  pridruženog  $m$ -toj po veličini svojstvenoj vrijednosti  $\lambda_l$ . Ako je  $\lambda_k = \lambda_l$ , koeficijenti  $k$ -te i  $l$ -te komponente su nužno ortogonalni; ako je  $\eta_k = \eta_l$ , elementi se mogu izabrati tako da budu ortogonalni, iako takvih vektora ima beskonačno mnogo. Varijaca  $m$ -te komponente je  $\lambda_l$ , pa je ukupna varijanca

$$\lambda_1 + \lambda_2 + \dots + \lambda_n = \text{trag } R = n$$

Važnost  $m$ -te komponente u ovakvom jednostavnijem prikazu mjeri se sa  $\lambda_l/n$ , pošto vrijedi  $\lambda_1 > \lambda_2 > \dots > \lambda_n$ .

No, mi ne želimo zadržati puni skup komponenata, nego ćemo broj komponenata reducirati na  $p < n$ . Za određivanje broja značajnih komponenata  $p$  postoje mnogi kriteriji, a ovdje se koristi PB kriterij (Štalec, Momirović, 1971).

Označimo li sa  $U^2 = (\text{diag } R^{-1})^{-1}$  dijagonalnu matricu univiteteta, tada je

$$\text{trag } (I - U^2) = \eta$$

zbroj image varijanci varijabli iz  $V$ .

Tada  $p$  mora zadovoljavati nejednakosti

$$\sum_{k=1}^{p-1} \lambda_k < \eta < \sum_{k=1}^p \lambda_k$$

tj. zadržimo toliko glavnih komponenti da zbroj njihovih varijanci bude veći ili jednak zajedničkom varijabilitetu svih varijabli iz  $V$ .

Na taj način dekomponirali smo trajektoriju promjena na  $p$  aditivnih, ortogonalnih komponenata promjena. U matričnom obliku to pišemo

$$K = Z X \quad (8)$$

gdje je

$$K = (K_1, K_2, \dots, K_p) \quad K_i \text{ — glavne komponente}$$

$$X = (X_1, X_2, \dots, X_p) \quad X_i \text{ — svojstveni vektori}$$

pridruženi zadržanim svojstvenim vrijednostima.

Da bismo mogli međusobno uspoređivati  $p$  zadržanih glavnih komponenata, standardiziramo ih na varijancu 1.

$$L = ZX\Lambda^{-1/2}$$

gdje je  $\Lambda$  dijagonalna matrica  $p$  zadržanih svojstvenih vrijednosti matrice  $R$ .

Struktura komponenata promjena pomaže nam pri njihovoj interpretaciji, a definirana je korelacijama varijabli glavnih komponenata

$$H = Z^T L \frac{1}{n} Z^T Z \frac{1}{n} X \Lambda^{-1/2} = R X \Lambda^{-1/2} = X \Lambda^{1/2}$$

### 2.3. Analiza reziduala

Eckart-Youngova (1936) aproksimacija matrice  $Z$  je najbolja aproksimacija (pod kriterijem najmanjih kvadrata) te matrice matricom ranga  $p$

$$Z^* = Z X X^T$$

U toj matrici su sadržane one informacije o promatranim promjenama koje su objašnjive sa  $p$  zadržanih glavnih komponenata. Onaj dio informacija koji se ne može objasniti zadržanim komponentama opisan je rezidualnom matricom standardiziranih podataka

$$Z = Z - Z^* = Z (I - X X^T)$$

Rezidualne korelacije varijabli iz  $V$  definirane su sa

$$R = R - H H^T.$$

a to su upravo oni dijelovi korelacija među varijablama koji se ne mogu objasniti sa  $p$  zadržanih komponenata promjena.

### 2.4. Relacije vremenskih točaka

Relacije vremenskih točaka definiraju se kao normirani skalarni produkti vektora vremenskih točaka u prostoru standardiziranih varijabli.

$$Q^* = Z Z^T$$

$$Q = (\text{diag } Q^*)^{-1/2} Q^* (\text{diag } Q^*)^{-1/2}$$

Teorijske relacije vremenskih točaka definirane su kao skalarni produkti vektora vremenskih točaka u prostoru standardiziranih komponenata promjena, normirani na istu metriku kao i relacije vremenskih točaka  $Q$

$$Q^*_{\tau} = K K^T$$

$$Q_{\tau} = (\text{diag } Q^*)^{-1/2} Q^*_{\tau} (\text{diag } Q^*)^{-1/2},$$

a to je onaj dio relacija vremenskih točaka koji je objašnjiv sa  $p$  zadržanih glavnih komponenata.

Rezidualne relacije vremenskih točaka su razlike relacija vremenskih točaka  $Q$  i teorijskih relacija vremenskih točaka  $Q_T$

$$Q_R = Q - Q_T$$

i to je onaj dio relacija vremenskih točaka koji se ne može objasniti sa  $p$  zadržanih komponenata promjena.

### 3. PROGRAM

Program INDIFF napisan je u meta-jeziku SS i implementiran na računaru UNIVAC 1100/U2 Sveučilišnog računskog centra u Zagrebu. U trenutnoj verziji (5.2/M) sistema SS njime je moguće analizirati promjene stanja objekta opisanog sa do 250 varijabli u ne više od 10.000 distinktnih vremenskih točaka. Podacima je potrebno pridružiti naredbe koje opisuju format podataka u skladu s uputama za sistem SS (jednu SEQUENCE naredbu i onoliko VARIABLE naredbi koliko ima varijabli).

INDIFF ne testira hipoteze o značajnosti komponenata, ali, ako su  $m$  i  $n$  dovoljno veliki brojevi, značajnost neke komponente  $K_i$ ,  $i=1, 2, \dots, p$  može se testirati ovako:

1) neka je  $e_i = m + n + 1 - 2i$

2) neka je  $e = (n-p)(m-p)$  gdje je  $p$  broj komponenata zadržanih po PB kriteriju.

3) neka je

$$f_i = \frac{\lambda_i e_r}{e_i (\sum_{j=p+1}^n \lambda_j)}$$

4) tada je  $f_i$  distribuiran u skladu sa Snedecorovom raspodjelom sa  $e_i$  stupnjeva slobode.

### 4. LITERATURA

1. Eckart, C. T. & Young, G.: The approximation of one matrix by another of lower rank. Psychometrika, 1936, 1, 211—218.
2. Momirović, K.: Metode za transformaciju i kondenzaciju kinezioloških informacija. Institut za kineziologiju Zagreb, 1972, 250—254.
3. Moore, E. H.: General analysis, Part I. Mem. Am. Phil. Soc. 1935, 1, 197.
4. Morrison, D. F.: Multivariate statistical methods. Mc Graw-Hill, New York, 1967.
5. Penrose, R.: A generalized inverse for matrices. Proc. Cambridge Phil. So., 1955, 51, 406—413.
6. Štalec, J., K. Momirović: Ukupna količina valjane varijance kao osnov kriterija za određivanje broja značajnih glavnih komponenata. Kineziologija, 1971, 1, 77—81.
7. Zakrajšek, E., J. Štalec i K. Momirović: SS-Programski sistem za multivarijantnu analizu podataka, I Međunarodni simpozij »Kompjuter na Sveučilištu«, Zagreb 1974

### INDIFF — THE MODEL, ALGORITHM AND PROGRAM FOR ANALYSIS OF CHANGES IN CONDITION OF AN OBJECT DESCRIBED OVER A GROUP OF QUANTITATIVE VARIABLES

The model and algorithm were constructed and the program was written for the component analysis of changes in condition of individual objects. By decomposition of trajectories of changes into orthogonal components we reach an insight into the structure of the stochastic process described by quantitative variables. Also, the residuals and the relation between time points were analysed.

Константин Момирович, Живан Караман

»ИНДИФ« — МОДЕЛЬ, АЛГОРИТМ И ПРОГРАММА ДЛЯ АНАЛИЗА ИЗМЕНЕНИЙ СОСТОЯНИЯ НЕКОТОРОГО ОБЪЕКТА, ОПИСАННОГО НАД МНОЖЕСТВОМ КОЛИЧЕСТВЕННЫХ ПЕРЕМЕННЫХ

Созданы модель и алгоритм и написана программа для компонентного анализа измерения состояний отдельных объектов. При помощи расчленения траектории изменений на ортогональные компоненты выявляется структура стохастического процесса, описанного при помощи количественных переменных. Также проводится анализ результатов и взаимоотношений между различными точками во времени.