

dipl. ing. Zoran Radojičić

# LINEARNI STATISTIČKI MODELI

---

Beograd 1999.



---

# SADRŽAJ:

1. Elementi matrice algebre	5
1.1 Vektori	5
1.2 Matrice	6
1.3 Neke važnije osobine matrica	11
2. Vrste podataka i merne skale	12
3. Višedimenzionalne slučajne promenljive	15
3.1 Sredina i kovarijaciona matrica	15
3.2 Podela vektora	16
3.3 Uzoračka sredina i kovarijaciona matrica	17
3.4 Nekoliko važnijih osobina uzoračke sredine i kovarijacione matrice	24
4. Metoda glavnih komponenti	29
4.1 Definicija glavnih komponenti	30
5. Faktorska analiza	38
5.1 Istorijat	39
5.2 Model faktorske analize	40
5.3 Metode izdvajanja faktora (metode ocenjivanja)	46
5.3.1 Metoda glavnih komponenti	46
5.3.2 Metoda glavnih faktora	50
6. Klaster analiza	54
7. Rešeni zadaci	59
8. Literatura	66



# 1. ELEMENTI MATRIČNE ALGEBRE

## 1.1 Vektori

Niz  $\mathbf{x}$  realnih brojeva  $x_1, x_2, \dots, x_n$  se naziva **vektor** i piše

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

ili

$$\mathbf{x}' = [x_1 \quad x_2 \quad | \quad x_n]$$

ili

$$\mathbf{x} = [x_1 \quad x_2 \quad | \quad x_n]'$$

pri čemu se poslednja operacija naziva **transponovanje kolone u vrstu**.

Nad vektorima  $\mathbf{x}$  i  $\mathbf{y}$  su definisane sledeće operacije:

$$c\mathbf{x} = \begin{bmatrix} cx_1 \\ cx_2 \\ \vdots \\ cx_n \end{bmatrix}$$

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

$$\mathbf{x}'\mathbf{y} = x_1y_1 + x_2y_2 + \dots + x_ny_n$$

(unutrašnji proizvod)

$$x'x = x_1^2 + x_2^2 + \dots + x_n^2.$$

Za dva vektora  $\mathbf{x}$  i  $\mathbf{y}$  koji su iste dimenzije se kaže da su **linearno zavisni** ako postoje konstante  $c_1$  i  $c_2$  različite od nule takve da je

$$c_1x + c_2y = 0.$$

U opštem slučaju, za skup vektora  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$  se kaže da su linearno zavisni ako postoje konstante  $c_1, c_2, \dots, c_k$  koje nisu sve jednake nuli tako da je

$$c_1x_1 + c_2x_2 + \dots + c_kx_k = 0.$$

Linearna zavisnost implicira da se barem jedan vektor iz skupa može predstaviti kao linearna kombinacija ostalih.

**Dužina vektora  $\mathbf{x}$** , u oznaci  $L_x$  se definiše kao

$$L_x = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} = \sqrt{x'x}.$$

## 1.2 Matrice

Matrica je pravougaona šema sastavljena od  $p$  redova (vrsta) i  $n$  kolona sa  $p \times n$  elemenata:

$$A_{(p \times n)} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{p1} & a_{p2} & \dots & a_{pn} \end{bmatrix} = [a_{ij}]_{p \times n}.$$

**Transponovana matrica  $A'$**  se dobija uzajamnom razmenom mesta redova i kolona i to onim redom kako se javljaju u matrici  $A$

$$A'_{(n \times p)} = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{p1} \\ a_{12} & a_{22} & \dots & a_{p2} \\ \dots & \dots & \dots & \dots \\ a_{1n} & a_{2n} & \dots & a_{pn} \end{bmatrix} = A^T.$$

Množenje matrice konstantom  $c$ :

$$c A_{(p \times n)} = \begin{bmatrix} ca_{11} & ca_{12} & \dots & ca_{1n} \\ ca_{21} & ca_{22} & \dots & ca_{2n} \\ \dots & \dots & \dots & \dots \\ ca_{p1} & ca_{p2} & \dots & ca_{pn} \end{bmatrix}.$$

Dve matrice  $A$  i  $B$  istih dimenzija se mogu sabrati:

$$A + B = [a_{ij}]_{p \times n} + [b_{ij}]_{p \times n} = [a_{ij} + b_{ij}]_{p \times n}$$

Da bi se dve matrice  $A$  i  $B$  mogle međusobno pomnožiti, neophodno je da budu *saglasne*, tj. da matrica  $A$  ima onoliko kolona koliko matrica  $B$  ima redova:

$$\underset{(p \times k)}{A} \underset{(k \times n)}{B} = \underset{(p \times n)}{C}, \quad \text{tj. } c_{ij} = \sum_{l=1}^k a_{il} b_{lj}.$$

**Kvadratna matrica** je simetrična ako je  $A=A'$  odnosno  $a_{ij}=a_{ji} \forall i,j$ .

**Primer 1.1:**

matrica  $\begin{bmatrix} 3 & 5 \\ 5 & -2 \end{bmatrix}$  je simetrična,

matrica  $\begin{bmatrix} 3 & 6 \\ 4 & -2 \end{bmatrix}$  nije simetrična.

**Nula matrica** predstavlja matricu čiji su svi elementi jednaki nuli i može biti bilo kog tipa:

$$0 = \begin{bmatrix} 0 & 0 & | & 0 \\ 0 & 0 & | & 0 \\ \hline 0 & 0 & | & 0 \end{bmatrix}.$$

**Jedinična matrica**  $I$  predstavlja matricu kod koje su svi članovi na glavnoj dijagonali jednaki jedinici a svi ostali nuli:

$$I = \begin{bmatrix} 1 & 0 & | & 0 \\ 0 & 1 & | & 0 \\ \hline 0 & 0 & | & 1 \end{bmatrix}.$$

**Determinanta** kvadratne matrice ima formu

$$\det A = |A| = \begin{vmatrix} a_{11} & a_{12} & | & a_{1n} \\ a_{21} & a_{22} & | & a_{2n} \\ \hline a_{n1} & a_{n2} & | & a_{nn} \end{vmatrix}.$$

**Minor** elementa  $a_{ij}$  u determinanti  $A$  je determinanta reda  $n-1$  koja se dobija iz  $\det A$  kada se izbrišu njena  $i$ -ta vrsta i  $j$ -ta kolona.

**Rang matrice** je definisan kao veličina najveće determinante koja je različita od nule i koja se može formirati od matrice  $\mathbf{A}$ . Drugačije rečeno, rang matrice je najveći red minora te matrice koji je različit od nule.

Za kvadratnu matricu se kaže da je **singularna** ako je njena determinanta jednaka nuli. Ako je  $\det \mathbf{A} \neq 0$ , kvadratna matrica  $\mathbf{A}$  je **regularna** (nesingularna). Matrica je sigurno singularna u sledećim slučajevima:

- ako su svi elementi jedne vrste ili kolone jednaki nuli,
- ako su vrste ili kolone matrice linearno zavisne.

Singularna matrica nema svoju inverznu matricu.

Kvadratna matrica  $\mathbf{A}^{-1}$  se naziva **inverzna matrica** date kvadratne matrice  $\mathbf{A}$  ako je:

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{A}\mathbf{A}^{-1} = \mathbf{I}.$$

Za svaku **dijagonalnu** matricu

$$D = \left[ \begin{array}{cc|c} d_1 & 0 & 0 \\ 0 & d_2 & 0 \\ \hline 0 & 0 & d_n \end{array} \right]$$

čiji su svi dijagonalni elementi  $d_i \neq 0$ , inverzna će biti:

$$D^{-1} = \left[ \begin{array}{cc|c} \frac{1}{d_1} & 0 & 0 \\ 0 & \frac{1}{d_2} & 0 \\ \hline 0 & 0 & \frac{1}{d_n} \end{array} \right].$$

**Pozitivno definitna matrica** je takva kvadratna matrica  $\mathbf{A}$  koja ima osobinu da je odgovarajuća kvadratna forma

$$\mathbf{X}'\mathbf{A}\mathbf{X} > 0$$

za svaki vektor  $\mathbf{X} \neq 0$ .

Ako važi relacija

$$\mathbf{X}'\mathbf{A}\mathbf{X} \geq 0$$

tada se matrica  $\mathbf{A}$ , kao i odgovarajuća kvadratna forma, naziva **pozitivno semidefinitnom**. Ako je simetrična matrica  $\mathbf{A}$  pozitivno definitna onda je i matrica  $\mathbf{A}^{-1}$  pozitivno definitna. Pozitivno definitna matrica  $\mathbf{A}$  je regularna, tj.  $\det \mathbf{A} \neq 0$ .



Neka je  $A_{(k \times k)}$  kvadratna matrica i  $I_{(k \times k)}$  jedinična matrica. Skalari  $\lambda_1, \lambda_2, \dots, \lambda_k$  koji zadovoljavaju jednačinu

$$|A - \lambda I| = 0 \quad \text{odnosno} \quad \det(A - \lambda I) = 0 \quad (*)$$

nazivaju se **karakteristični koreni** matrice  $A$ . Jednačina (\*) se naziva **karakteristična jednačina**.

**Primer 1.2:**

Za sledeću matricu

$$A = \begin{bmatrix} 1 & 0 \\ 1 & 3 \end{bmatrix}$$

se dobijaju dva karakteristična korena:

$$|A - \lambda I| = \left| \begin{bmatrix} 1 & 0 \\ 1 & 3 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right| = \begin{vmatrix} 1 - \lambda & 0 \\ 1 & 3 - \lambda \end{vmatrix} = (1 - \lambda)(3 - \lambda) = 0$$

$$\lambda_1 = 1, \lambda_2 = 3.$$

**Primer 1.3:**

$$A = \begin{bmatrix} 13 & -4 & 2 \\ -4 & 13 & -2 \\ 2 & -2 & 10 \end{bmatrix}$$

$$|A - \lambda I| = \left| \begin{bmatrix} 13 & -4 & 2 \\ -4 & 13 & -2 \\ 2 & -2 & 10 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right| = \begin{vmatrix} 13 - \lambda & -4 & 2 \\ -4 & 13 - \lambda & -2 \\ 2 & -2 & 10 - \lambda \end{vmatrix} =$$

$$= -\lambda^3 + 36\lambda^2 - 405\lambda + 1458 = 0$$

$$\lambda_1 = 9, \lambda_2 = 9, \lambda_3 = 18.$$

Neka je  $A$  matrica dimenzije  $k \times k$  i  $\lambda$  njen karakteristični koren. Tada, ako važi

$$\underset{(k \times 1)}{x} \neq \underset{(k \times 1)}{0} \quad \text{i} \quad Ax = \lambda x$$

za  $x$  se kaže da predstavlja **karakteristični vektor** matrice  $A$  pridružen karakterističnom korenu  $\lambda$ .

Ekvivalentan uslov je:

$$(A - \lambda I)x = 0.$$

**Primer 1.4:**

Odrediti karakteristične vektore matrice

$$A = \begin{bmatrix} 1 & 0 \\ 1 & 3 \end{bmatrix}$$

Iz prethodnog primera smo dobili da je  $\lambda_1=1$  i  $\lambda_2=3$ .

$$Ax = \lambda_1 x \Leftrightarrow \begin{bmatrix} 1 & 0 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \lambda_1 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Leftrightarrow \begin{bmatrix} 1 & 0 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 1 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (*)$$

$$x_1 = x_1$$

$$x_1 + 3x_2 = x_2 \Rightarrow x_1 = -2x_2$$

$$Ax = \lambda_2 x \Leftrightarrow \begin{bmatrix} 1 & 0 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \lambda_2 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Leftrightarrow \begin{bmatrix} 1 & 0 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 3 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (**)$$

$$x_1 = 3x_1$$

$$x_1 + 3x_2 = 3x_2 \Rightarrow x_1 = 0$$

Iz (\*) sledi da ima beskonačno mnogo rešenja za  $x_1$  i  $x_2$ . Tako, ako uzmemo da je  $x_2=1$  dobijamo da je  $x_1=-2$ , tj.

$$x = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

karakteristični vektor koji odgovara karakterističnom korenu  $\lambda_1=1$ .

Iz izraza (\*\*) dobijamo da je  $x_1=0$  a za  $x_2$  uzimamo proizvoljnu vrednost  $1$  i dobijamo

$$x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

karakterističan vektor koji odgovara karakterističnom korenu  $\lambda_2=3$ .

### 1.3 Neke važnije osobine matrica

Neka su  $A$  i  $B$   $k \times k$  kvadratne matrice. Tada važi:

- a)  $|A| = |A'|$   
 b) Ako su svi elementi jedne vrste (kolone) matrice  $A$  jednaki nuli, tada je:

$$|A| = 0$$

- c) Ako su bilo koje dve vrste (kolone) matrice  $A$  identične, tada je:

$$|A| = 0$$

- d) Ako je  $A$  nesingularna matrica, tada je:

$$|A| = \frac{1}{|A^{-1}|}, \quad \text{tj. } |A||A^{-1}| = 1$$

- e)  $|AB| = |A||B|$

- f)  $|cA| = c^k |A|$ , gde je  $c$  skalar.

Neka je  $A$   $k \times k$  kvadratna matrica. **Trag matrice  $A$** , u oznaci  $tr(A)$  je jednak zbiru elemenata na glavnoj dijagonali:

$$tr(A) = \sum_{i=1}^k a_{ii}.$$

Neka su  $A$  i  $B$   $k \times k$  kvadratne matrice i  $c$  skalar. Tada važi:

- a)  $tr(cA) = c tr(A)$   
 b)  $tr(A \pm B) = tr(A) \pm tr(B)$   
 c)  $tr(AB) = tr(BA)$   
 d)  $tr(B^{-1}AB) = tr(A)$   
 e)  $tr(AA') = \sum_{i=1}^k \sum_{j=1}^k a_{ij}^2$

## 2. VRSTE PODATAKA I MERNE SKALE

Statistička obeležja mogu biti **kvantitativna** i **kvalitativna** (nominalna). Kvantitativna obeležja se dele na **neprekidna** (jedinica mere se može beskonačno deliti) i **prekidna**.

$$X = (X_1 \quad X_2 \quad \dots \quad X_p)^T, \quad p \geq 1 \quad \text{vektor obeležja (promenljivih),}$$

$$E = \{e_1, e_2, \dots, e_N\}, \quad N \geq p \quad \text{skup jedinica posmatranja (entiteta).}$$

### Primer 2.1:

Iz skupa objekata (skup ljudi) izabrana su 3 objekta (3 muškarca) stara između 30 i 40 godina i nazvana  $e_j$ ,  $e_k$  i  $e_q$ . Pretpostavimo da nas o objektima  $e_j$ ,  $e_k$  i  $e_q$  zanimaju sledeće karakteristike:

- $X_1$  - broj dece (dobijen brojanjem)
- $X_2$  - broj članova uže porodice (dobijen brojanjem)
- $X_3$  - telesna masa (dobijena merenjem u gramima)
- $X_4$  - inteligencija (kao broj rešenih zadataka od ukupno 60)
- $X_5$  - obrazovanje (stepen stručne spreme: 1 magistrat, 2 fakultet, 3 viša škola)
- $X_6$  - sportska igra koju najviše vole.

Dobijeni su sledeći rezultati:

	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$
$e_j$	2	4	64000	20	3	fudbal
$e_k$	0	1	86520	40	1	košarka
$e_q$	4	6	128000	0	5	fudbal

**Tabela 2.1**

Skup objekata predstavlja skup jedinica posmatranja (entiteta): {Pera, Mika, Žika}. Karakteristike predstavljaju obeležja. Tabelu 2.1 nazivamo *matrica podataka* i ona je oblika

$$X = [x_{ij}]_{N \times p}$$

odnosno, u našem slučaju

$$X = [x_{ij}]_{3 \times 6}$$

gde je  $x_{ij}$  **realizacija**  $x_i$  na  $e_j$ .

Kako za svaku od ispitanih karakteristika entiteti mogu imati različite rezultate to ćemo rezultate dobijene utvrđivanjem neke karakteristike nazivati **promenljivim** (varijablama).

Osobine definisanih varijabli:

- $X_1$
- broj dece (prekidna kvantitativna varijabla)
  - uzima vrednost iz skupa pozitivnih celobrojnih vrednosti uključujući i nulu
  - rezultati se ne mogu na smislen način podvrgnuti nikakvoj transformaciji već moraju ostati takvi kakvi jesu
  - kako su merne jedinice definisane onim što se meri reč je o **apsolutnim skalama**
- $X_2$
- broj članova uže porodice (prekidna kvantitativna)
  - varijabla je tako definisana da i sam entitet ulazi u broj članova uže porodice (npr.  $e_k$  je očigledno samac bez dece)
  - iz ovoga bi se moglo zaključiti da se u zaključivanju ništa nebi promenilo ako od svakog rezultata oduzmemo jedinicu pa prema tome  $X_2$  dobija vrednosti:  
3  
0  
5
- prema tome, veličina varijabli ovog tipa je određena do neke proizvoljne konstante koju možemo dodati ili oduzeti pa se zato ove skale nazivaju **aditivne skale**
- $X_3$
- telesna masa (neprekidna kvantitativna)
  - ovde je reč o **skali odnosa** koja ima sledeće osobine:
    - a) količnik ma koje dve vrednosti ima smislenu interpretaciju  
npr.  $e_q$  je dva puta teži od  $e_j$
    - b) rastojanje između dva objekta mereno na ma kom delu ove skale je jednako  
npr.  $e_j$  ima 22520 grama manje od entiteta  $e_k$ , a  $e_q$  64000 grama više od  $e_j$
    - c) vrednostima pozicioniranim na skali mogu se dodeliti rangovi od višeg ka nižem  
npr.  $e_q$  ima viši rang od  $e_k$

- $X_4$
- inteligencija (neprekidna kvantitativna)
  - nema realnu nultu tačku (nema ljudskog stvora sa inteligencijom nula)
  - merne jedinice su proizvoljne (bodovi na testu)
  - poseduje osobine  $b$  i  $c$  (za drugi test druge vrednosti)
  - ove skale se nazivaju **intervalnim** jer je jedino interval (razmak) između rezultata definisan  
npr. razlika u inteligenciji između entiteta  $e_j$  i  $e_q$  je jednaka kao i između entiteta  $e_j$  i  $e_k$
- $X_5$
- obrazovanje (kvalitativna varijabla koja u ovom slučaju uslovno može da prođe kao kvantitativna, kao stepen stručne spreme)
  - vrednosti varijable ne označavaju količinu već rang, pa se zato nazivaju **ordinarnim varijablama** ili varijablama koje leže na nekoj ordinarnoj skali
  - poseduju samo osobinu  $c$
- $X_6$
- sportska igra koju najviše vole (kvalitativna varijabla)
  - rezultat ne govori ništa ni o količini ni o redosledu već samo o tome koje svojstvo ima neki entitet iz grupe međusobno isključivih svojstva
  - kako oznaka pripadanja ne označava nikakvu veličinu, varijable ovog tipa se nazivaju **kvalitativnim** ili nominalnim varijablama

Klasifikacija metoda multivarijacione analize u zavisnosti od vrste promenljivih:

kvantitativne promenljive:

- metoda glavnih komponenti
- faktorska analiza
- klaster analiza
- višestruka korelacija
- višestruka regresija
- *man*-ova metoda (multidimenziona analiza varijanse)
- kanonična korelaciona analiza

kvalitativne promenljive:

- loglinearni modeli
- diskriminaciona analiza (samo zavisna promenljiva mora da bude kvalitativna)

### 3. VIŠEDIMENZIONALNE SLUČAJNE PROMENLJIVE

#### 3.1 Sredina i kovarijaciona matrica

Označimo  $p$  jednodimenzionalnih slučajnih promenljivih sa  $X_1, X_2, \dots, X_p$ . Skup ovih slučajnih promenljivih nazivamo **slučajni vektor**  $X$  (vektor obeležja)

$$X = (X_1 \quad X_2 \quad \vdots \quad X_p)^T, \quad p \geq 1.$$

Realizaciju  $X$  na  $e_i$  označavamo sa

$$x_{(i)} = (x_{i1} \quad x_{i2} \quad \vdots \quad x_{ip})^T.$$

Za svaku jednodimenzionalnu slučajnu promenljivu možemo odrediti **sredinu**  $\mu_i = E(X_i)$  i **varijansu**  $\sigma_i^2 = E(X_i - \mu_i)^2$ ,  $i = 1, 2, \dots, p$ .

Sredina slučajnog vektora  $X$  je

$$\mu = E(X) = (E(X_1) \quad E(X_2) \quad \vdots \quad E(X_p))^T = (\mu_1 \quad \mu_2 \quad \vdots \quad \mu_p)^T.$$

Za ma koji par slučajnih promenljivih  $X_i$  i  $X_j$  definišemo **kovarijansu**

$$\sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)] = Cov(X_i, X_j)$$

$$\sigma_{ii} = \sigma_i^2 = Var(X_i)$$

$$\sigma_{ij} = \sigma_{ji}.$$

Za slučajni vektor  $\mathbf{X}$  se definiše  $p \times p$  simetrična matrica koja se naziva **kovarijaciona matrica** od  $\mathbf{X}$  sa znakom  $\Sigma$ :

$$\Sigma = [\sigma_{ij}]_{p \times p} = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \vdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \vdots & \sigma_{2p} \\ \hline \sigma_{p1} & \sigma_{p2} & \vdots & \sigma_{pp} \end{bmatrix}.$$

Kovarijaciona matrica se može iskazati i kao očekivana vrednost **slučajne matrice**

$$\Sigma = \begin{bmatrix} (X_1 - \mu_1)^2 & (X_1 - \mu_1)(X_2 - \mu_2) & \vdots & (X_1 - \mu_1)(X_p - \mu_p) \\ (X_2 - \mu_2)(X_1 - \mu_1) & (X_2 - \mu_2)^2 & \vdots & (X_2 - \mu_2)(X_p - \mu_p) \\ \hline (X_p - \mu_p)(X_1 - \mu_1) & (X_p - \mu_p)(X_2 - \mu_2) & \vdots & (X_p - \mu_p)^2 \end{bmatrix}.$$

Na kraju, kovarijacionu matricu možemo iskazati kao očekivanu vrednost sledeće **matrice odstupanja**:

$$W = [(X_i - \mu_i)(X_j - \mu_j)]_{p \times p}$$

$$\Sigma = E(W) = E[(X - \mu)(X - \mu)^T].$$

### 3.2 Podela vektora

Često se tokom istraživanja javlja potreba da  $p$  promenljivih posmatramo razdvojeno:

$$\mathbf{X} = (X_1 \quad X_2 \quad \vdots \quad X_q \quad | \quad X_{q+1} \quad X_{q+2} \quad \vdots \quad X_p)^T.$$

U tom slučaju vršimo podelu sredine  $\mu$  i kovarijacione matrice  $\Sigma$

$$\mu = (\mu_1 \quad \mu_2)^T$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix}$$

gde su

$$\Sigma_{11} = [\sigma_{ij}^{(1)}]_{q \times q}$$

$$\Sigma_{22} = [\sigma_{ij}^{(2)}]_{(p-q) \times (p-q)}$$

$$\Sigma_{12} = [\sigma_{ij}^{(3)}]_{(p-q) \times q}$$



### 3.3 Uzoračka sredina i kovarijaciona matrica

#### Primer 3.1:

Za matricu podataka

	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$
$e_j$	2	4	64000	20	3	fudbal
$e_k$	0	1	86520	40	1	košarka
$e_q$	4	6	128000	0	5	fudbal

odrediti realizovanu vrednost sredine vektora  $X$  za prve dve varijable ( $X_1$  - broj dece,  $X_2$  - broj članova uže porodice).

Slučajni vektor (vektor obeležja)

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = (X_1 \quad X_2 \quad \vdots \quad X_p)^T$$

u našem slučaju je

$$X = (X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5 \quad X_6)^T, \quad p = 6.$$

Skup jedinica posmatranja (entiteta)

$$E = \{e_1, e_2, \dots, e_N\}$$

u našem slučaju je

$$E = \{e_j, e_k, e_q\}, \quad N = 3.$$

Matrica podataka

$$X = [x_{ij}]_{N \times p}$$

u našem slučaju je

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} \end{bmatrix}$$

tj.

$$X = \begin{bmatrix} 2 & 4 & 64000 & 20 & 3 & fudbal \\ 0 & 1 & 86520 & 40 & 1 & kosarka \\ 4 & 6 & 128000 & 0 & 5 & fudbal \end{bmatrix}.$$

U daljem radu se nećemo sretati sa nominalnim varijablama u matrici podataka dok ćemo matricu podataka zadavati uglavnom na ovaj način.

U ovom primeru je reč o uzorku iz **višedimenzionalnog rasporeda**. Višedimenzionalni raspored se opisuje nekom **merom centralne tendencije** i nekom **merom disperzije** (rasipanja).

Za meru centralne tendencije najčešće se uzima **sredina slučajnog vektora**  $X$  koja se označava sa  $\mu$

$$\mu = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} = (\mu_1 \quad \mu_2 \quad \vdots \quad \mu_p)^T.$$

Za meru varijabiliteta najčešće se uzima **kovarijaciona matrica** od  $X$  koja se označava sa  $Cov(X)$  ili  $\Sigma$

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \vdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \vdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \vdots & \sigma_{pp} \end{bmatrix}$$

gde je  $\sigma_{ii} = \sigma_i^2$ ,  $i=1, 2, \dots, p$ .

Ako definišemo simetričnu matricu kvadrata

$$\begin{bmatrix} (X_1 - \mu_1)^2 & (X_1 - \mu_1)(X_2 - \mu_2) & \vdots & (X_1 - \mu_1)(X_p - \mu_p) \\ (X_2 - \mu_2)(X_1 - \mu_1) & (X_2 - \mu_2)^2 & \vdots & (X_2 - \mu_2)(X_p - \mu_p) \\ \vdots & \vdots & \ddots & \vdots \\ (X_p - \mu_p)(X_1 - \mu_1) & (X_p - \mu_p)(X_2 - \mu_2) & \vdots & (X_p - \mu_p)^2 \end{bmatrix}$$

koja je proizvod slučajnih vektora odstupanja od sredine, tj.

$$(X - \mu)(X - \mu)'$$

onda je njena očekivana vrednost jednaka kovarijacionoj matrici:

$$E[(X - \mu)(X - \mu)'] = \Sigma.$$

Iz višedimenzionalne raspodele se može izračunati i **korelaciona matrica**  $\rho$  i to preko kovarijacione matrice. U matričnoj notaciji to izgleda ovako:

$$\rho = \left( D^{\frac{1}{2}} \right)^{-1} \Sigma \left( D^{\frac{1}{2}} \right)^{-1}$$

gde je  $D$  dijagonalna matrica koja sadrži elemente na glavnoj dijagonali kovarijacione matrice (**matrica standardnih devijacija**), tj.

$$\rho = \begin{bmatrix} 1 & 0 & 0 \\ \sqrt{\sigma_{11}} & & \\ 0 & 1 & 0 \\ & \sqrt{\sigma_{22}} & \\ \hline 0 & 0 & 1 \\ & & \sqrt{\sigma_{pp}} \end{bmatrix} \begin{bmatrix} \sigma_{11} & \sigma_{12} & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \sigma_{2p} \\ \sigma_{p1} & \sigma_{p2} & \sigma_{pp} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ \sqrt{\sigma_{11}} & & \\ 0 & 1 & 0 \\ & \sqrt{\sigma_{22}} & \\ \hline 0 & 0 & 1 \\ & & \sqrt{\sigma_{pp}} \end{bmatrix} =$$

$$= \begin{bmatrix} 1 & \rho_{12} & \rho_{1p} \\ \rho_{21} & 1 & \rho_{2p} \\ \rho_{p1} & \rho_{p2} & 1 \end{bmatrix}$$

gde je  $\rho_{ij}$  koeficijent korelacije između  $X_i$  i  $X_j$ :

$$\rho_{ij} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}} \sqrt{\sigma_{jj}}}$$

Na osnovu uspostavljene relacije može se napisati da je

$$\Sigma = D^{\frac{1}{2}} \rho D^{\frac{1}{2}}$$

Na osnovu matrice podataka smo u mogućnosti da izračunamo pokazatelje uzorka: sredinu uzorka, kovarijacionu matricu uzorka i korelacionu matricu uzorka.

**Uzoračka sredina** je data sa:

$$\bar{X} = (\bar{X}_1 \quad \bar{X}_2 \quad \vdots \quad \bar{X}_p)^T = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \\ \vdots \\ \bar{X}_p \end{bmatrix}$$

gde je

$$\bar{X}_j = \frac{1}{N} \sum_{i=1}^N X_{ij}, \quad j = 1, 2, \dots, p.$$

**Uzoračku kovarijacionu matricu** definišemo preko matrice **centriranih podataka**  $X^*$  i matrice **uzajamnih proizvoda centriranih podataka**  $X^{*'} X^*$

$$X^* = \begin{bmatrix} X_{11} - \bar{X}_1 & X_{12} - \bar{X}_2 & \dots & X_{1p} - \bar{X}_p \\ X_{21} - \bar{X}_1 & X_{22} - \bar{X}_2 & \dots & X_{2p} - \bar{X}_p \\ \dots & \dots & \dots & \dots \\ X_{N1} - \bar{X}_1 & X_{N2} - \bar{X}_2 & \dots & X_{Np} - \bar{X}_p \end{bmatrix}$$

odnosno

$$X^{*'} X^* = \begin{bmatrix} \sum_{i=1}^N (X_{i1} - \bar{X}_1)^2 & \sum_{i=1}^N (X_{i1} - \bar{X}_1)(X_{i2} - \bar{X}_2) & \dots & \sum_{i=1}^N (X_{i1} - \bar{X}_1)(X_{ip} - \bar{X}_p) \\ \sum_{i=1}^N (X_{i2} - \bar{X}_2)(X_{i1} - \bar{X}_1) & \sum_{i=1}^N (X_{i2} - \bar{X}_2)^2 & \dots & \sum_{i=1}^N (X_{i2} - \bar{X}_2)(X_{ip} - \bar{X}_p) \\ \dots & \dots & \dots & \dots \\ \sum_{i=1}^N (X_{ip} - \bar{X}_p)(X_{i1} - \bar{X}_1) & \sum_{i=1}^N (X_{ip} - \bar{X}_p)(X_{i2} - \bar{X}_2) & \dots & \sum_{i=1}^N (X_{ip} - \bar{X}_p)^2 \end{bmatrix}_{p \times p}$$

pa je uzoračka kovarijaciona matrica  $S_{(p \times p)}$

$$S = \frac{1}{N-1} X^{*'} X^* = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})(X_i - \bar{X})'.$$

**Uzoračka korelaciona matrica** je

$$R = \begin{bmatrix} 1 & r_{12} & r_{13} & \dots & r_{1p} \\ r_{12} & 1 & r_{23} & \dots & r_{2p} \\ r_{13} & r_{23} & 1 & \dots & r_{3p} \\ \dots & \dots & \dots & \dots & \dots \\ r_{1p} & r_{2p} & r_{3p} & \dots & 1 \end{bmatrix}$$

gde su  $r_{ij}$  **uzorački koeficijenti korelacije**

$$r_{ij} = \frac{S_{ij}}{\sqrt{S_{ii}} \sqrt{S_{jj}}}, \quad i, j = 1, 2, \dots, p$$

a  $S_{ij}$  **elementi uzoračke kovarijacione matrice**

$$S_{ij} = \frac{1}{N-1} \sum_{k=1}^N (X_{ki} - \bar{X}_i)(X_{kj} - \bar{X}_j), \quad i, j = 1, 2, \dots, p.$$

Do uzoračke kovarijacione matrice smo mogli doći i koristeći relaciju koju smo uspostavili između kovarijacione i korelacione matrice populacije:

$$R = \left( D^{\frac{1}{2}} \right)^{-1} S \left( D^{\frac{1}{2}} \right)^{-1}.$$

Vratimo se primeru. Interesuje nas matrica podataka

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{bmatrix} = \begin{bmatrix} 2 & 4 \\ 0 & 1 \\ 4 & 6 \end{bmatrix} \quad i = 1, 2, 3 \Rightarrow N = 3, \quad j = 1, 2 \Rightarrow p = 2.$$

Treba odrediti  $\bar{X}$ :

$$\bar{X} = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \\ \bar{X}_p \end{bmatrix} = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{3}(2+0+4) \\ \frac{1}{3}(4+1+6) \end{bmatrix} = \begin{bmatrix} 2 \\ 3.67 \end{bmatrix}$$

što govori da je prosečan broj dece  $\bar{X}_1 = 2$  a prosečan broj članova uže porodice  $\bar{X}_2 = 3.67$ .

**Primer 3.2:**

Realizovane vrednosti slučajnog uzorka su date matricom podataka

$$X = \begin{bmatrix} -2 & 3 \\ 4 & 2 \\ 1 & 1 \end{bmatrix}.$$

Odrediti realizovanu vrednost uzoračke kovarijacione matrice.

Prvo ćemo izračunati uzoračku sredinu  $\bar{X}$ :

$$\bar{X} = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{3}(-2+4+1) \\ \frac{1}{3}(3+2+1) \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad N = 3, p = 2.$$

Zatim računamo matricu centriranih podataka  $X^*$ :

$$X^* = \begin{bmatrix} X_{11} - \bar{X}_1 & X_{12} - \bar{X}_2 & \vdots & X_{1p} - \bar{X}_p \\ X_{21} - \bar{X}_1 & X_{22} - \bar{X}_2 & \vdots & X_{2p} - \bar{X}_p \\ \vdots & \vdots & \ddots & \vdots \\ X_{N1} - \bar{X}_1 & X_{N2} - \bar{X}_2 & \vdots & X_{Np} - \bar{X}_p \end{bmatrix}.$$

Kako je  $N=3$  i  $p=2$

$$X^* = \begin{bmatrix} X_{11} - \bar{X}_1 & X_{12} - \bar{X}_2 \\ X_{21} - \bar{X}_1 & X_{22} - \bar{X}_2 \\ X_{31} - \bar{X}_1 & X_{31} - \bar{X}_2 \end{bmatrix} = \begin{bmatrix} -2-1 & 3-2 \\ 4-1 & 2-2 \\ 1-1 & 1-2 \end{bmatrix} = \begin{bmatrix} -3 & 1 \\ 3 & 0 \\ 0 & -1 \end{bmatrix}.$$

Nakon toga računamo matricu uzajamnih proizvoda centriranih podataka  $X^*X^*$

$$X^{*'}X^* = \begin{bmatrix} -3 & 3 & 0 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} -3 & 1 \\ 3 & 0 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} 18 & -3 \\ -3 & 2 \end{bmatrix}$$

pa je uzoračka kovarijaciona matrica

$$S = \frac{1}{N-1} X^{*'}X^* = \frac{1}{3-1} \begin{bmatrix} 18 & -3 \\ -3 & 2 \end{bmatrix} = \begin{bmatrix} 9 & -\frac{3}{2} \\ -\frac{3}{2} & 1 \end{bmatrix}.$$

Oдавде je  $S_1^2=9$ ,  $S_2^2=1$  i  $S_{12}=S_{21}=-\frac{3}{2}$  (negativna linearna veza).

Za domaću vežbu naći kovarijacionu matricu uzorka za podatke date u primeru 3.1.

### **Primer 3.3:**

Za podatke iz prethodnog primera odrediti realizovanu vrednost uzoračke korelacione matrice koristeći relaciju uspostavljenu između nje i kovarijacione matrice.

Na osnovu izračunate kovarijacione matrice

$$S = \begin{bmatrix} 9 & -\frac{3}{2} \\ -\frac{3}{2} & 1 \end{bmatrix}$$

formiramo dijagonalnu matricu

$$D = \begin{bmatrix} 9 & 0 \\ 0 & 1 \end{bmatrix}$$

a zatim izračunavamo korelacionu matricu

$$R = \left( D^{\frac{1}{2}} \right)^{-1} S \left( D^{\frac{1}{2}} \right)^{-1}$$

$$R = \begin{bmatrix} \frac{1}{3} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 9 & -\frac{3}{2} \\ -\frac{3}{2} & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{3} & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}.$$

Često se tokom istraživanja javlja potreba da  $p$  promenljivih posmatramo razdvojeno, odnosno da ih tretiramo kao da pripadaju dvema grupama promenljivih

$$X = (X_1 \quad X_2 \quad \vdots \quad X_q \quad | \quad X_{q+1} \quad X_{q+2} \quad \vdots \quad X_p)^T$$

tako da dobijamo

$$X_1 = (X_1 \quad X_2 \quad \vdots \quad X_q)^T \quad X_1 = [x_{ij}]_{N \times q},$$

$$X_2 = (X_{q+1} \quad X_{q+2} \quad \vdots \quad X_p)^T \quad X_2 = [x_{ij}]_{N \times s}.$$

Tako npr. ako ispitujemo uticaj demografskih faktora na potrošnju proizvoda široke potrošnje tada ćemo vektor  $X$  dimenzije  $p \times 1$  podeliti na vektor  $X_1$  u koji će ući sledeći podskup promenljivih: pol, starost i stručna sprema, i na vektor  $X_2$  u koji će ući promenljive potrošnje hleba, mesa, mleka, povrća i voća.

$X_1$	- pol	→ podela vektora
$X_2$	- starost	
$X_3$	- stručna sprema	
$X_4$	- potrošnja hleba	
$X_5$	- potrošnja mesa	
$X_6$	- potrošnja mleka	
$X_7$	- potrošnja povrća	
$X_8$	- potrošnja voća	

$$X = (X_1 \quad X_2 \quad X_3 \quad | \quad X_4 \quad X_5 \quad X_6 \quad X_7 \quad X_8)^T$$

$$X_1 = (X_1 \quad X_2 \quad X_3)^T, \quad q = 3$$

$$X_2 = (X_4 \quad X_5 \quad X_6 \quad X_7 \quad X_8)^T, \quad s = 5$$

U skladu sa podelom slučajnog vektora  $\mathbf{X}$  vršimo podelu sredine  $\mu$  i kovarijacione matrice  $\Sigma$ :

$$\mu = (\mu_1 \quad \mu_2)^T$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

gde su:

$$\mu_{1(q \times 1)} \text{ i } \mu_{2(q \times 2)} \quad \text{sredine slučajnih vektora } \mathbf{X}_1 \text{ i } \mathbf{X}_2$$

$$\Sigma_{11(q \times q)} \text{ i } \Sigma_{22(s \times s)} \quad \text{kovarijacione matrice slučajnih vektora } \mathbf{X}_1 \text{ i } \mathbf{X}_2$$

dok su elementi matrice  $\Sigma_{12(q \times s)}$  kovarijanse između elemenata slučajnih vektora  $\mathbf{X}_1$  i  $\mathbf{X}_2$  sa osobinom  $\Sigma_{21} = \Sigma_{12}^T$ .

### 3.4 Nekoliko važnijih osobina uzoračke sredine i kovarijacione matrice

Neka su slučajne promenljive  $X_j$  i  $X_k$  linearno transformisane tj. neka su definisane nove slučajne promenljive

$$cX_j + a \text{ i } dX_k + b, \quad a, b, c \text{ i } d \text{ su realne konstante.}$$

Tada, na osnovu definicije kovarijanse, sledi

$$\text{Cov}(cX_j + a, dX_k + b) = cd \text{Cov}(X_j, X_k).$$

Uopštimo ovaj slučaj na linearnu kombinaciju  $p$  slučajnih promenljivih iz slučajnog vektora  $\mathbf{X}$  sa sredinom  $\mu$  i kovarijacionom matricom  $\Sigma$ . Linearnom kombinacijom

$$Y = a_1 X_1 + a_2 X_2 + \dots + a_p X_p = a^T X$$

za dati vektor linearne kombinacije

$$a^T = (a_1 \quad a_2 \quad \vdots \quad a_p)$$

definišemo novu slučajnu promenljivu  $Y$  čija je očekivana vrednost

$$\mu_Y = E(Y) = E(a^T X) = a^T \mu$$





**Primer 3.4:**

Izračunati sredine, varijanse i kovarijanse sledeće dve linearne kombinacije na osnovu podataka iz primera 3.2:

- a) direktno, na osnovu opservacija za  $Y_1$  i  $Y_2$   
 b) indirektno na osnovu izračunate uzoračke sredine  $\bar{X}$  i uzoračke kovarijacione matrice  $S$ .

$$A = \begin{bmatrix} -1 & 2 \\ 3 & -2 \end{bmatrix} \quad X = \begin{bmatrix} -2 & 3 \\ 4 & 2 \\ 1 & 1 \end{bmatrix}.$$

Treba da dobijemo

$$Y = \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ y_{31} & y_{32} \end{bmatrix}.$$

Imamo dve linearne kombinacije:

$$Y_1 = a_{11}X_1 + a_{12}X_2 = a_1^T X = \begin{bmatrix} -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = -x_1 + 2x_2$$

$$Y_2 = a_{21}X_1 + a_{22}X_2 = a_2^T X = \begin{bmatrix} 3 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 3x_1 - 2x_2$$

- a) Za prvu linearnu kombinaciju

$$Y_1 = a_1^T X = \begin{bmatrix} -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

imamo realizovane vrednosti opservacija  $y_{11}$ ,  $y_{21}$  i  $y_{31}$  tj.

$$y_{11} = (-1)(-2) + 2 \cdot 3 = 8$$

$$y_{21} = (-1) \cdot 4 + 2 \cdot 2 = 0$$

$$y_{31} = (-1) \cdot 1 + 2 \cdot 1 = 1$$

**Napomena:**

$$X_1 = (x_1 \quad x_2 \quad x_3)^T$$

$$\bar{y}_1 = \frac{1}{3}(8 + 0 + 1) = 3$$

$$S_{Y_1}^2 = \frac{1}{3-1}[(8-3)^2 + (0-3)^2 + (1-3)^2] = 19.$$

Za drugu linearnu kombinaciju je

$$Y_2 = a_2^T X = [3 \quad -2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$y_{12} = 3 \cdot (-2) + (-2) \cdot 3 = -12$$

$$y_{22} = 3 \cdot 4 + (-2) \cdot 2 = 8$$

$$y_{32} = 3 \cdot 1 + (-2) \cdot 1 = 1$$

pa je

$$\bar{y}_2 = \frac{1}{3}(-12 + 8 + 1) = -1$$

$$S_{Y_2}^2 = \frac{1}{3-1}[(-12+1)^2 + (8+1)^2 + (1+1)^2] = 103$$

$$\begin{aligned} \text{Cov}(Y_1, Y_2) &= S_{12} \frac{\sum_{k=1}^3 (y_{k1} - \bar{y}_1)(y_{k2} - \bar{y}_2)}{3-1} = \\ &= \frac{(8-3)(-12+1) + (0-3)(8+1) + (1-3)(1+1)}{2} = -43 \end{aligned}$$

$$S_Y = \begin{bmatrix} 19 & -43 \\ -43 & 103 \end{bmatrix}.$$

b) Indirektno:

$$Y = AX \Leftrightarrow \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\bar{y} = A\bar{x} \Leftrightarrow \begin{bmatrix} \bar{y}_1 \\ \bar{y}_2 \end{bmatrix} = \begin{bmatrix} -1 & 2 \\ 3 & -2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 3 \\ -1 \end{bmatrix}.$$

Za uzoračku kovarijacionu matricu imamo iste relacije kao i za populaciju, tj.

$$\Sigma_Y = A\Sigma_X A'$$

$$S_Y = AS_X A'$$

$$S_Y = \begin{bmatrix} -1 & 2 \\ 3 & -2 \end{bmatrix} \begin{bmatrix} 9 & -\frac{3}{2} \\ -\frac{3}{2} & 1 \end{bmatrix} \begin{bmatrix} -1 & 2 \\ 3 & -2 \end{bmatrix} = \begin{bmatrix} 19 & -43 \\ -43 & 103 \end{bmatrix}.$$

### Primer 3.5:

Zavod za statistiku u svojoj redovnoj anketi o potrošnji domaćinstava prikuplja podatke o njihovim prihodima i rashodima po različitim kategorijama za poljoprivredna, mešovita i nepoljoprivredna domaćinstva. Izabran je slučajni uzorak od šest domaćinstava i izmereni su sledeći podaci:

	$X_1$	$X_2$	$X_3$	$X_4$
	6	9.4	6.2	0.20
	7	12.1	7.2	0.10
	4	10.7	3.1	0.22
	5	9.3	5.4	0.12
	2	11.2	2.1	0.26
	3	13.3	1.8	0.30

- $X_1$  - broj članova domaćinstva,  
 $X_2$  - godišnja raspoloživa sredstva,  
 $X_3$  - izdaci za ishranu,  
 $X_4$  - izdaci za obrazovanje i razonodu.

Izračunati kovarijacionu i korelacionu matricu uzorka.

Rešenje:

$$S = \begin{bmatrix} 3.50 & -0.86 & 4.10 & -0.12 \\ -0.86 & 2.416 & -1.51 & 0.0524 \\ 4.10 & -1.51 & 5.152 & -0.1568 \\ -0.12 & 0.0524 & -0.1568 & 0.00608 \end{bmatrix}$$

$$R = \begin{bmatrix} 1 & -0.2957 & 0.9655 & -0.8226 \\ & 1 & -0.4280 & 0.4323 \\ & & 1 & -0.8859 \\ & & & 1 \end{bmatrix}.$$

## 4. METODA GLAVNIH KOMPONENTI

Ideja metode glavnih komponenti je da se originalne promenljive transformišu u nove promenljive, odnosno, linearne kombinacije.

Polazimo od vektora obeležja

$$X = (X_1 \quad X_2 \quad \dots \quad X_p)^T, \quad p \geq 1.$$

Ako je, na primer,  $p=20$  (20 obeležja) njihova međusobna povezanost iskazana je pomoću 190 koeficijenata korelacije što predstavlja problem kako za razumevanje, tako i za interpretaciju. Zato ćemo ovih  $p$  promenljivih zameniti manjim brojem linearnih kombinacija tipa

$$L_i' X, \quad i = 1, 2, \dots$$

ali tako da su te linearne kombinacije nezavisne među sobom i da sadrže što je moguće više *inicijalnog varijabiliteta*.

Prva glavna komponenta je konstruisana tako da obuhvata najveći deo varijanse originalnog skupa podataka. Druga glavna komponenta obuhvata najveći deo preostale varijanse originalnog skupa podataka i tako redom.

Analiza glavnih komponenti je bazirana na pretpostavci da će relativno mali broj glavnih komponenti dobro aproksimirati kovarijacionu strukturu skupa originalnih promenljivih. Ovim metoda glavnih komponenti postiže dva cilja:

1. vrši redukciju originalnog skupa podataka
2. olakšava interpretaciju kovarijacione strukture originalnih promenljivih na bazi manjeg broja međusobno nekorelisanih glavnih komponenti

Tako ćemo umesto  $p=20$  promenljivih u daljoj analizi koristiti 4 linearne kombinacije.

## 4.1 Definicija glavnih komponenti

Neka je

$$X = (X_1 \quad X_2 \quad \dots \quad X_p)^T, \quad p \geq 1$$

slučajni vektor sa kovarijacionom matricom  $\Sigma$ . Prvu glavnu komponentu određujemo kao

$$Y_1 = l_{11}X_1 + l_{12}X_2 + \dots + l_{1p}X_p = L_1'X$$

linearnu kombinaciju elemenata slučajnog vektora  $X$ , gde su  $l_{11}, l_{12}, \dots, l_{1p}$  koeficijenti linearne kombinacije.

Od ranije znamo da je

$$\text{Var}(Y_1) = \text{Var}(L_1'X) = L_1'\Sigma L_1.$$

Naš zadatak je da odredimo vektor koeficijenata  $L_1$  tako da se maksimizira varijansa od  $Y_1$ . Ovde se uvodi uslov ortogonalnosti same transformacije jer se to pokazalo jako korisno u praksi i omogućuje se optimizacija, tj. uvodimo uslov

$$L_1'L_1 = 1$$

odnosno hoćemo da maksimiziramo  $L_1'\Sigma L_1$  uz ograničenje  $L_1'L_1 = 1$ .

To se radi uz pomoć **Lagranžovih množitelja**. Diferenciranjem Lagranžove funkcije po koeficijentima  $L_1$  a zatim izjednačavanjem dobijenog izraza sa nulom dobijamo

$$\Sigma L_1 - \lambda L_1 = 0$$

ili

$$(\Sigma - \lambda I)L_1 = 0$$

gde je  $I_{(p \times p)}$  jedinična matrica,  $\lambda$  jedan od karakterističnih korena kovarijacione matrice  $\Sigma$  i to najveći  $\lambda = \lambda_1$ , a  $L_1$  je karakteristični vektor.

Da bi se dobilo netrivialno rešenje za  $L_1$  treba da je ispunjeno

$$|\Sigma - \lambda I| = 0.$$

Svaku sledeću linearnu kombinaciju određujemo slično uz dodatni uslov da se radi o nekorelisanim veličinama (da je kovarijansa između prve i druge glavne komponente jednaka nuli)

$$\Sigma L_2 - \lambda L_2 = 0$$

ili

$$(\Sigma - \lambda I)L_2 = 0$$

gde za  $\lambda$  biramo što je moguće veću vrednost  $\lambda = \lambda_2$ ,  $L_2$  je karakteristični vektor a linearna kombinacija

$$Y_2 = L_2' X$$

predstavlja drugu glavnu komponentu.

**Rezime:**

1. Ako su svi karakteristični koreni matrice  $\Sigma$  međusobno različiti (uredili smo ih u opadajući niz  $\lambda_1 > \lambda_2 > \dots > \lambda_p \geq 0$ ) tada postoji  $p$  glavnih komponenti

$$Y_1, Y_2, \dots, Y_p$$

tj.

$$Y_j = L_j' X, \quad j = 1, 2, \dots, p$$

2. Vektori koeficijenata  $L_1, L_2, \dots, L_p$  predstavljaju karakteristične vektore matrice  $\Sigma$  koji su pridruženi karakterističnim korenima  $\lambda_j$ .
- 3.

$$E(Y_j) = 0$$

$$\text{Var}(Y_j) = L_j' \Sigma L_j = \lambda_j$$

$$\text{Cov}(Y_i, Y_j) = 0, \quad i \neq j$$

$$\Lambda = \left[ \begin{array}{cc|c} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ \hline 0 & 0 & \lambda_p \end{array} \right]$$

**Primer 4.1:**

Za datu kovarijacionu matricu

$$\Sigma = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & -2 \\ 0 & -2 & 4 \end{bmatrix}$$

odrediti sve glavne komponente.

Glavne komponente su:

$$Y_1 = L_1' X$$

$$Y_2 = L_2' X$$

$$Y_3 = L_3' X$$

Prvu glavnu komponentu računamo iz uslova

$$(\Sigma - \lambda_1 I)L_1 = 0$$

odnosno češće se umesto  $L_1$  piše  $P_1$  (*principal components*)

$$(\Sigma - \lambda_1 I)P_1 = 0$$

odnosno

$$\Sigma P_1 = \lambda_1 P_1$$

pa sledi da je potrebno izračunati karakteristične korene kovarijacione matrice  $\Sigma$

$$\Sigma = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & -2 \\ 0 & -2 & 4 \end{bmatrix}$$

$$|\Sigma - \lambda_1 I| = \begin{vmatrix} 4 - \lambda & 0 & 0 \\ 0 & 2 - \lambda & -2 \\ 0 & -2 & 4 - \lambda \end{vmatrix} = (4 - \lambda)(\lambda^2 - 6\lambda + 4) = 0$$

$$\lambda_1 = 4, \lambda_2 = 0.769, \lambda_3 = 5.2361.$$

Šta ne valja?

Po definiciji glavnih komponenti,  $\lambda_1$  je najveći karakteristični koren,  $\lambda_2$  drugi po veličini, itd. pa je

$$\lambda_1 = 5.2361, \lambda_2 = 4, \lambda_3 = 0.769$$

pa će biti

$$\Sigma P_1 = \lambda_1 P_1$$



$$\begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & -2 \\ 0 & -2 & 4 \end{bmatrix} \begin{bmatrix} p_{11} \\ p_{21} \\ p_{31} \end{bmatrix} = 5.2361 \begin{bmatrix} p_{11} \\ p_{21} \\ p_{31} \end{bmatrix}$$

$$4p_{11} + 0p_{21} + 0p_{31} = 5.2361p_{11}$$

$$0p_{11} + 2p_{21} - 2p_{31} = 5.2361p_{21}$$

$$0p_{11} - 2p_{21} + 4p_{31} = 5.2361p_{31}$$

$$-1.2361p_{11} = 0 \Rightarrow p_{11} = 0$$

$$-3.2361p_{21} - 2p_{31} = 0$$

$$-2p_{21} - 1.2361p_{31} = 0$$

Rešavanjem ovog sistema dobija se

$$P_1' = (0 \quad -0.53 \quad 0.85).$$

Sada računamo drugu glavnu komponentu:

$$(\Sigma - \lambda_2 I)P_2 = 0$$

$$\Sigma P_2 = \lambda_2 P_2$$

$$\begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & -2 \\ 0 & -2 & 4 \end{bmatrix} \begin{bmatrix} p_{12} \\ p_{22} \\ p_{32} \end{bmatrix} = 4 \begin{bmatrix} p_{12} \\ p_{22} \\ p_{32} \end{bmatrix}$$

$$4p_{12} + 0p_{22} + 0p_{32} = 5.2361p_{12}$$

$$0p_{12} + 2p_{22} - 2p_{32} = 5.2361p_{22}$$

$$0p_{12} - 2p_{22} + 4p_{32} = 5.2361p_{32}$$

$$P_2' = (1 \quad 0 \quad 0).$$

Treća komponenta:

$$(\Sigma - \lambda_3 I)P_3 = 0$$

$$\Sigma P_3 = \lambda_3 P_3$$

$$\begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & -2 \\ 0 & -2 & 4 \end{bmatrix} \begin{bmatrix} p_{13} \\ p_{23} \\ p_{33} \end{bmatrix} = 0.769 \begin{bmatrix} p_{13} \\ p_{23} \\ p_{33} \end{bmatrix}$$

$$4p_{13} + 0p_{23} + 0p_{33} = 5.2361p_{13}$$

$$0p_{13} + 2p_{23} - 2p_{33} = 5.2361p_{23}$$

$$0p_{13} - 2p_{23} + 4p_{33} = 5.2361p_{33}$$

$$P_3' = (0 \quad -0.8507 \quad -0.5257).$$

Na kraju, glavne komponente su:

$$Y_1 = L_1' X = P_1' X = -0.53X_2 + 0.85X_3$$

$$Y_2 = L_2' X = P_2' X = X_1$$

$$Y_3 = L_3' X = P_3' X = -0.8507X_2 - 0.5257X_3$$

Zadovoljeni su uslovi normiranosti  $L'L=I$ :

$$\text{za } P_1' = (0 \quad -0.53 \quad 0.85) \text{ imamo } 0^2 + (-0.5257)^2 + 0.8507^2 = 1$$

$$\text{za } P_2' = (1 \quad 0 \quad 0) \text{ imamo } 1^2 = 1$$

$$\text{za } P_3' = (0 \quad -0.8507 \quad -0.5257) \text{ važi isto kao i za } P_1':$$

#### **Primer 4.2:**

Na prethodnom primeru pokazati da je varijansa dobijenih glavnih komponenti jednaka odgovarajućim karakterističnim korenima a da je kovarijansa između svakog para glavnih komponenti jednaka nuli.

Treba dokazati da je  $Var(Y_j) = \lambda_j$  za  $j=1, 2, \dots, p$ , tj.

$$Var(Y_1) = Var(-0.5257X_2 + 0.8507X_3) = \lambda_1 = 5.2361$$

itd.

Provera korelisanosti prve i druge glavne komponente:

$$Cov(Y_1, Y_2) = Cov(-0.5257X_2 + 0.8507X_3; X_1) = 0.$$

**Napomena:**

**Generalizovana varijansa** se određuje po dve definicije.

Po prvoj definiciji generalizovana varijansa je jednaka determinanti kovarijacione matrice. **Kovarijaciona matrica** glavnih komponenti je dijagonalna matrica čiji su elementi karakteristični koreni matrice  $\Sigma$  i označava se sa  $\Lambda$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & \vdots & 0 \\ 0 & \lambda_2 & \vdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \vdots & \lambda_p \end{bmatrix}.$$

Prema tome generalizovana varijansa vektora  $Y$  je jednaka determinanti

$$|\Lambda|.$$

Po drugoj definiciji generalizovana varijansa je jednaka tragu kovarijacione matrice. Trag kovarijacione matrice je jednak zbiru karakterističnih korena  $\lambda_j$ .

Korišćenjem ortogonalne matrice  $P_{(p \times p)}$  čiji su redovi karakteristični vektori kovarijacione matrice  $\Sigma$ , možemo izvršiti ortogonalnu dekompoziciju matrice  $\Sigma$  čiji su koreni različiti.

Imamo da je

$$\Sigma = P' \Lambda P$$

ili u razvijenom obliku

$$\Sigma = \begin{bmatrix} L_1 & L_2 & \vdots & L_p \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \vdots & 0 \\ 0 & \lambda_2 & \vdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \vdots & \lambda_p \end{bmatrix} \begin{bmatrix} L_1' \\ L_2' \\ \vdots \\ L_p' \end{bmatrix} = \sum_{j=1}^p \lambda_j L_j L_j'.$$

Takođe je

$$\Lambda = P \Sigma P'.$$

Uz korišćenje osobine traga matrice  $tr(BC)=tr(CB)$  imamo da je

$$tr(\Lambda) = tr(P \Sigma P') = tr(P' P \Sigma) = tr(\Sigma)$$

pošto je  $P'P=I$ .

**Primer 4.3:**

Neka su za vektor  $\mathbf{X}^T=(X_1, X_2, X_3)$  poznati korelacioni koeficijenti  $r_{12}=0.5$ ,  $r_{13}=0.3$  i  $r_{23}=0.2$  i neka su date prve dve glavne komponente kao

$$Y_1 = 0.84X_1 + 0.79X_2 + 0.6X_3 \quad i \quad Y_2 = 0.17X_1 + 0.42X_2 - 0.79X_3.$$

Izračunati:

- Koliko treba uzeti glavnih komponenti da bi se obuhvatilo barem 85% varijabiliteta i koliko procenata ukupnog varijabiliteta nosi svaka od glavnih komponenti  $Y_1$ ,  $Y_2$  i  $Y_3$ ? Izračunati generalizovanu varijansu vektora  $\mathbf{X}$  korišćenjem varijanse glavnih komponentata i pokazati zašto je to moguće.
- Aproksimirati korelacionu matricu na osnovu projekcije trodimenzionalnog vektora  $\mathbf{X}$  u prostor prve glavne komponente.
- Korelaciju između prve glavne komponente  $Y_1$  i varijable  $X_1$  kao i između druge glavne komponente  $Y_2$  i varijable  $X_2$ .

$$\mathbf{X}^T = (X_1 \quad X_2 \quad X_3)$$

$$Y_1 = 0.84X_1 + 0.79X_2 + 0.6X_3$$

$$Y_2 = 0.17X_1 + 0.42X_2 - 0.79X_3$$

$$R = \begin{bmatrix} 1 & 0.5 & 0.3 \\ & 1 & 0.2 \\ & & 1 \end{bmatrix}$$

Da bi važio uslov  $\alpha^T\alpha=1$ , inicijalne glavne komponente treba podeliti pa se dobija:

$$P_1^T(0.646 \quad 0.608 \quad 0.461)$$

$$P_2^T(0.187 \quad 0.462 \quad -0.868)$$

a)

$$\begin{bmatrix} 1 & 0.5 & 0.3 \\ 0.5 & 1 & 0.2 \\ 0.3 & 0.2 & 1 \end{bmatrix} \begin{bmatrix} 0.646 \\ 0.608 \\ 0.461 \end{bmatrix} = \lambda_1 \begin{bmatrix} 0.646 \\ 0.608 \\ 0.461 \end{bmatrix} \Rightarrow \lambda_1 = 1.685$$

$$\begin{bmatrix} 1 & 0.5 & 0.3 \\ 0.5 & 1 & 0.2 \\ 0.3 & 0.2 & 1 \end{bmatrix} \begin{bmatrix} 0.187 \\ 0.462 \\ -0.868 \end{bmatrix} = \lambda_2 \begin{bmatrix} 0.187 \\ 0.462 \\ -0.868 \end{bmatrix} \Rightarrow \lambda_2 = 0.830$$

$$\lambda_1 + \lambda_2 + \lambda_3 = 3 \Rightarrow \lambda_3 = 0.485$$

Prva glavna komponenta nosi 56.17% varijabiliteta, druga 27.67% a treća 16.17% varijabiliteta. Da bi se obuhvatilo 85% varijabiliteta potrebne su sve tri glavne komponente.

b)

$$\lambda_1 \alpha_1 \alpha_1^T = 1.685 \begin{bmatrix} 0.646 \\ 0.608 \\ 0.461 \end{bmatrix} \begin{bmatrix} 0.646 & 0.608 & 0.461 \end{bmatrix} = \begin{bmatrix} 0.703 & 0.662 & 0.502 \\ 0.662 & 0.623 & 0.472 \\ 0.502 & 0.472 & 0.358 \end{bmatrix}$$

c)

$$\rho_{Y_1 X_1} = \alpha_{11} \sqrt{\lambda_1} = 0.646 \sqrt{1.685} = 0.839$$

$$\rho_{Y_2 X_2} = \alpha_{22} \sqrt{\lambda_2} = 0.462 \sqrt{0.83} = 0.421$$

## 5. FAKTORSKA ANALIZA

Faktorska analiza je metod multivarijacione analize koji se koristi za opis međusobne zavisnosti velikog broja promenljivih korišćenjem manjeg broja osnovnih ali *neopažljivih* slučajnih promenljivih (faktora). Slično kao metoda glavnih komponenti i faktorska analiza se koristi za redukciju skupa podataka aproksimirajući kovarijacionu ili korelacionu matricu originalnih promenljivih.

Između metode glavnih komponenti i faktorske analize postoje sledeće razlike:

1. Analiza glavnih komponenti se bazira na varijansama - interesuju nas dijagonalni elementi kovarijacione matrice, dok nas kod faktorske analize interesuju vandijagonalni elementi.
2. Za razliku od analize glavnih komponenti, faktorska analiza pretpostavlja postojanje teorijskog modela kojim se uspostavlja relacija između opservacija  $p$ -dimenzionalne promenljive i manjeg broja zajedničkih faktora.
3. Često se analiza glavnih komponenti tretira kao specifičan slučaj ili prva faza u faktorskoj analizi (*SPSS*).
4. Analiza glavnih komponenti izučava ukupan varijabilitet skupa podataka. Nasuprot njoj, faktorska analiza polazi od razlaganja varijabiliteta na dva dela: zajednički deo i specifični deo. Zajednički deo je onaj deo varijacija promenljive koji ona deli sa ostalim promenljivama. Sa druge strane, specifični deo je onaj deo varijacija promenljive koji je poseban (specifičan) za tu promenljivu. Faktorska analiza izučava zajednički deo varijabiliteta za sve promenljive
5. U obe metode se javljaju dve vrste promenljivih: *opažljive* promenljive koje čini originalni skup podataka i *neopažljive* (latentne) promenljive. Međutim, u analizi glavnih komponenti na osnovu linearne kombinacije opažljivih promenljivih se formiraju glavne komponente kao neopažljive komponente, dok se u faktorskoj analizi na osnovu faktora kao neopažljivih komponenti izražavaju originalne promenljive

## 5.1 Istorijat

Krajem XIX veka *Galton* i *Spirman* postavljaju problem inteligencije:

*Da li je inteligencija bazirana na jednom bazičnom tj. opštem faktoru ili se bazira na nekoliko zajedničkih faktora kao što su npr. verbalna sposobnost, matematička sposobnost, memorija?*

### Primer 5.1:

*Čarls Spirman* je 1904. godine izučavao pripremljenost dece za školu i dobio je korelacionu matricu rezultata testova iz klasike ( $X_1$ ), francuskog ( $X_2$ ), engleskog ( $X_3$ ) i matematike ( $X_4$ ) u obliku

$$\begin{bmatrix} 1 & 0.83 & 0.78 & 0.70 \\ 0.83 & 1 & 0.67 & 0.67 \\ 0.78 & 0.67 & 1 & 0.64 \\ 0.70 & 0.67 & 0.64 & 1 \end{bmatrix}$$

*Spirman* je uočio proporcionalnost ma koja dva reda ili kolone u ovoj matrici (ako se zanemare elementi na glavnoj dijagonali). Tako za elemente drugog i trećeg reda imamo

$$\frac{0.83}{0.67} \approx \frac{0.78}{0.64} \approx 1.2.$$

Zato *Spirman* predlaže redukciju problema sa  $p=4$  na  $p=1$  tako što rezultate svih testova ( $X_i$ ,  $i=1, 2, 3$  i  $4$ ) iskazuje preko modela

$$X_i = \beta_i F + \varepsilon_i, \quad i = 1, 2, 3, 4$$

gde su

- $F$  - zajednički faktor,
- $\beta_i$  - koeficijenti (faktorska opterećenja),
- $\varepsilon_i$  - slučajne greške (specifični faktori).

Dakle, *Spirman* je formulisao dvofaktorsku teoriju testova inteligencije prema kojoj se rezultat svakog testa može dekomponovati na dva dela: prvi, koji je zajednički za sve testove ( $F$ ) i koji se može interpretirati kao "opšta sposobnost" ili "inteligencija", i drugi, koji je specifičan za svaki test ( $\varepsilon_i$ ).

### Napomena:

Kasnijim istraživanjima prvobitni model je proširen tako što je uvedeno nekoliko zajedničkih faktora (pretpostavlja se da postoje posebne vrste inteligencija - sposobnosti), a i specifičan faktor je razložen na dva dela od kojih prvi pokazuje u kom stepenu se nečija individualna sposobnost u npr. matematici razlikuje od opšte sposobnosti i drugi koji

predstavlja činjenicu da su rezultati testiranja samo aproksimacija sposobnosti individue u konkretnoj oblasti (ipak samo test - papir).

## 5.2 Model faktorske analize

Neka je  $X$   $p$ -dimenzioni vektor opažljivih promenljivih sa sredinom  $\mu$  i kovarijacionom matricom  $\Sigma$ . Pretpostavimo da se vektor opažljivih promenljivih  $X$  može izraziti preko skupa od  $m$  neopažljivih promenljivih koje nazivamo zajednički faktori u oznaci  $F_1, F_2, \dots, F_m$  ( $m \ll p$ ) i preko skupa  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$  specifičnih ali neopažljivih faktora.

Tada model u razvijenom obliku izgleda

$$\begin{aligned}(X_1 - \mu_1) &= \beta_{11}F_1 + \beta_{12}F_2 + \dots + \beta_{1m}F_m + \varepsilon_1 \\(X_2 - \mu_2) &= \beta_{21}F_1 + \beta_{22}F_2 + \dots + \beta_{2m}F_m + \varepsilon_2 \\ \hline (X_p - \mu_p) &= \beta_{p1}F_1 + \beta_{p2}F_2 + \dots + \beta_{pm}F_m + \varepsilon_p\end{aligned}$$

ili ekvivalentno u matricnoj notaciji

$$X - \mu = B F + \varepsilon$$

$(p \times 1) \quad (p \times m) (m \times 1) \quad (p \times 1)$

gde je

$$X - \mu = \begin{bmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \\ \hline X_p - \mu_p \end{bmatrix} \quad F = \begin{bmatrix} F_1 \\ F_2 \\ \hline F_m \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \hline \varepsilon_p \end{bmatrix}$$

$$B = \begin{bmatrix} \beta_{11} & \beta_{12} & \vdots & \beta_{1m} \\ \beta_{21} & \beta_{22} & \vdots & \beta_{2m} \\ \hline \beta_{p1} & \beta_{p2} & \vdots & \beta_{pm} \end{bmatrix}$$

Poslednja matrica predstavlja matricu **faktorskih opterećenja** a  $\beta_{ij}$  faktorsko opterećenje  $i$ -te promenljive na  $j$ -ti faktor.

### **Napomena:**

Na prvi pogled model faktorske analize liči na model višestruke regresije, međutim ovde  $p$  odstupanja  $(X_1 - \mu_1), \dots, (X_p - \mu_p)$  izražavamo preko  $m+p$  slučajnih promenljivih  $F_1, F_2, \dots, F_m$  i  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$  koje su neopažljive za razliku od regresionog modela gde su nezavisne promenljive opažljive.

U model se uvode sledeća dodatna ograničenja:



$$E(F) = 0$$

$$\text{Cov}(F) = E(FF') = \Phi$$

i dalja izlaganja se baziraju na specifičnom slučaju faktorske analize kada je

$$\Phi = I$$

odnosno, radi se o ortogonalnom modelu kod koga su faktori  $F_1, F_2, \dots, F_m$  međusobno nezavisni.

Što se tiče specifičnih faktora važi

$$E(\varepsilon) = 0$$

$$\text{Cov}(\varepsilon) = E(\varepsilon\varepsilon') = \Psi = \begin{bmatrix} \Psi_1 & 0 & | & 0 \\ 0 & \Psi_2 & | & 0 \\ \hline 0 & 0 & | & \Psi_p \end{bmatrix}$$

Takođe pretpostavljamo da su zajednički faktori nezavisni od specifičnih tj. da je

$$\text{Cov}(\varepsilon, F) = E(\varepsilon F') = 0.$$

Vežu između odstupanja opažljivih promenljivih od njihove sredine i neopažljivih faktora

$$X - \mu = BF + \varepsilon$$

odnosno

$$\underset{(p \times 1)}{X} = \underset{(p \times 1)}{\mu} + \underset{(p \times m)}{B} \underset{(m \times 1)}{F} + \underset{(p \times 1)}{\varepsilon}$$

zajedno sa navedenim pretpostavkama nazivamo **ortogonalni model** faktorske analize sa  $m$  zajedničkih faktora. Ovaj model omogućuje razlaganje kovarijacione matrice  $\Sigma$

$$\Sigma = \text{Cov}(X) = BB' + \Psi.$$

**Napomena:**

U opštem slučaju, kad ne važi uslov ortogonalnosti tj. kada je  $\Phi \neq I$ , razlaganje kovarijacione matrice  $\Sigma$  je

$$\Sigma = B\Phi B' + \Psi.$$

Kako je

$$\text{Cov}(X, F) = E(BF + \varepsilon)F' = BE(FF') + E(\varepsilon F') = B$$

znači da su elementi matrice faktorskih opterećenja kovarijanse između originalnih promenljivih i faktora.

Korelacionu matricu promenljivih  $X$  i faktora  $F$  nazivamo matrica **faktorske strukture**.

**Napomena:**

U slučaju ortogonalnog modela faktorske analize važi da je

$$\text{Cov}(X, F) = \text{Cor}(X, F).$$

U opštem slučaju se ove dve vrednosti međusobno razlikuju.

Na osnovu razlaganja kovarijacione matrice kod ortogonalnog modela faktorske analize

$$\Sigma = BB' + \Psi$$

imamo da je varijansa  $i$ -te promenljive

$$\sigma_{ii} = \text{Var}(X_i) = \beta_{i1}^2 + \beta_{i2}^2 + \dots + \beta_{im}^2 + \Psi_i$$

tj. varijansa  $i$ -te originalne promenljive je podeljena na dva dela:

- prvi deo je

$$h_i^2 = \beta_{i1}^2 + \beta_{i2}^2 + \dots + \beta_{im}^2$$

i predstavlja varijansu objašnjenu zajedničkim faktorima i nazivamo ga **zajednička varijansa** ili **komunalitet**

- drugi deo nazivamo specifična varijansa

**Primer 5.2:**

Data je kovarijaciona matrica

$$\Sigma = \begin{bmatrix} 15 & 8 & 3 & 16 \\ 8 & 10 & 1 & 10 \\ 3 & 1 & 16 & 2 \\ 16 & 10 & 2 & 24 \end{bmatrix}$$

i jednakost

$$\begin{bmatrix} 15 & 8 & 3 & 16 \\ 8 & 10 & 1 & 10 \\ 3 & 1 & 16 & 2 \\ 16 & 10 & 2 & 24 \end{bmatrix} = \begin{bmatrix} 3 & 2 \\ 2 & 1 \\ -1 & 3 \\ 4 & 2 \end{bmatrix} \begin{bmatrix} 3 & 2 & -1 & 4 \\ 2 & 1 & 3 & 2 \end{bmatrix} + \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix}.$$

Dokazati da važi

$$\Sigma = BB' + \Psi.$$

Iz jednakosti sledi postojanje dva zajednička faktora ( $m=2$ ) pa je

$$B = \begin{bmatrix} 3 & 2 \\ 2 & 1 \\ -1 & 3 \\ 4 & 2 \end{bmatrix} \quad \Psi = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix}.$$

Potrebno je da na osnovu ortogonalnog modela faktorske analize odredimo komunalitete sve četiri originalne promenljive i razložimo njihove varijanse.

Komunalitet prve promenljive:

$$h_1^2 = \beta_{11}^2 + \beta_{12}^2 = 3^2 + 2^2 = 13$$

jer je

$$B = \begin{bmatrix} \beta_{11} & \beta_{12} \\ \beta_{21} & \beta_{22} \\ \beta_{31} & \beta_{32} \\ \beta_{41} & \beta_{42} \end{bmatrix}$$

pa je razlaganje varijanse  $\sigma_{11}$  na komunalitet i specifičnu varijansu

$$\sigma_{11} = h_1^2 + \Psi_1$$

$$15 = 13 + 2$$

jer je

$$\Psi = \begin{bmatrix} \Psi_1 & 0 & 0 & 0 \\ 0 & \Psi_2 & 0 & 0 \\ 0 & 0 & \Psi_3 & 0 \\ 0 & 0 & 0 & \Psi_4 \end{bmatrix}.$$

Slično se dobija da je

$$\sigma_{22} = h_2^2 + \Psi_2 \Leftrightarrow 10 = 5 + 5$$

$$\sigma_{33} = h_3^2 + \Psi_3 \Leftrightarrow 16 = 10 + 6$$

$$\sigma_{44} = h_4^2 + \Psi_4 \Leftrightarrow 24 = 20 + 4$$

odnosno

$$h_2^2 = \beta_{21}^2 + \beta_{22}^2 = 2^2 + 1^2 = 5$$

$$h_3^2 = \beta_{31}^2 + \beta_{32}^2 = (-1)^2 + 3^2 = 10$$

$$h_4^2 = \beta_{41}^2 + \beta_{42}^2 = 4^2 + 2^2 = 20$$

***Napomena:***

Korišćenjem korelacione umesto ovarijacione matrice dobijamo da se jedinična varijansa standardizovane promenljive sastoji iz dva dela:

$$\text{Var}(X_i) = 1 = h_i^2 + \Psi_i.$$

***Generalizovana varijansa*** (ukupna varijansa) originalnog skupa promenljivih je oblika

$$\text{tr}(\Sigma) = \sum_{i=1}^p \sigma_{ii} = \sum_{i=1}^p \sum_{j=1}^m \beta_{ij}^2 + \sum_{i=1}^p \Psi_i.$$

Ako sa  $h$  označimo ukupan komunalitet od  $\mathbf{X}$

$$h = \sum_{i=1}^p h_i^2 = \sum_{i=1}^p \sum_{j=1}^m \beta_{ij}^2$$

tada se može pisati

$$\text{tr}(\Sigma) = h + \text{tr}(\Psi).$$

Znači da je generalizovana varijansa od  $\mathbf{X}$  jednaka zbiru dve komponente: ukupnog ***komunaliteta*** i ukupne ***varijanse*** specifičnih faktora. Formiranjem količnika

$$\frac{\sum_{i=1}^p \beta_{ij}^2}{h}, \quad j = 1, 2, \dots, m$$

dobija se proporcija ukupnog varijabiliteta koja se može pripisati  $j$ -tom zajedničkom faktoru.

***Primer 5.3:***

Na osnovu ortogonalnog modela faktorske analize iz prethodnog primera odrediti ukupnu varijansu, ukupan komunalitet, ukupnu specifičnu varijansu i doprinose svakog faktora ukupnom komunalitetu.

$$\Sigma = \begin{bmatrix} 15 & 8 & 3 & 16 \\ 8 & 10 & 1 & 10 \\ 3 & 1 & 16 & 2 \\ 16 & 10 & 2 & 24 \end{bmatrix} \quad B = \begin{bmatrix} 3 & 2 \\ 2 & 1 \\ -1 & 3 \\ 4 & 2 \end{bmatrix} \quad \Psi = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

Ukupna varijansa:

$$tr(\Sigma) = \sigma_{11} + \sigma_{22} + \sigma_{33} + \sigma_{44} = 15 + 10 + 16 + 24 = 65.$$

Ukupan komunalitet:

$$h = h_1^2 + h_2^2 + h_3^2 + h_4^2 = 13 + 5 + 10 + 20 = 48.$$

Ukupna specifična varijansa:

$$tr(\Psi) = \Psi_1 + \Psi_2 + \Psi_3 + \Psi_4 = 2 + 5 + 6 + 4 = 17.$$

Doprinos prvog zajedničkog faktora ukupnom komunalitetu,  $j=1$ :

$$\frac{\beta_{11}^2 + \beta_{21}^2 + \beta_{31}^2 + \beta_{41}^2}{h} = \frac{3^2 + 2^2 + (-1)^2 + 4^2}{48} = \frac{30}{48} = 62.5\% .$$

Doprinos drugog zajedničkog faktora ukupnom komunalitetu,  $j=2$ :

$$\frac{\beta_{12}^2 + \beta_{22}^2 + \beta_{32}^2 + \beta_{42}^2}{h} = \frac{2^2 + 1^2 + 3^2 + 2^2}{48} = \frac{18}{48} = 37.5\% .$$

**Napomena:**

Faktorska analiza se najčešće obavlja na *standardizovanim* promenljivama, odnosno vršimo razlaganje korelacione matrice.

### 5.3 Metode izdvajanja faktora (metode ocenjivanja)

Na osnovu uzorka uzetog iz  $p$ -dimenzione populacije ocenjuje se model faktorske analize, odnosno proveravamo da li važi razlaganje kovarijacione (ili korelacione) matrice koje je implicirano tim modelom.

Postoje dve osnovne metode ocenjivanja modela faktorske analize:

- metoda glavnih komponenti (glavnih faktora)
- metoda maksimalne verodostojnosti

**Metoda glavnih komponenti** se javlja u dve varijante:

1. direktno korišćenje metode glavnih komponenti na kovarijacionu ili korelacionu matricu radi istovremenog ocenjivanja komunaliteta i matrice faktorskih opterećenja (metoda glavnih komponenti)
2. korelaciona matrica se modifikuje u duhu modela faktorske analize a zatim se iterativno i odvojeno ocenjuju komunaliteti i faktorska opterećenja primenom metode glavnih komponenti na tu modifikovanu korelacionu matricu (metoda glavnih faktora)

### 5.3.1 Metoda glavnih komponenti

U uzoračkom slučaju postupak primene metode glavnih komponenti u ocenjivanju faktorskog modela zahteva određivanje karakterističnih korena  $\tilde{\lambda}_j$  i karakterističnih vektora  $\tilde{\alpha}_j$ .

Na osnovu prvih  $m$  faktora formiramo matricu ocenjenih faktorskih opterećenja prema sledećem izrazu:

$$\tilde{B} = \left[ \begin{array}{c|c|c} \sqrt{\tilde{\lambda}_1} \tilde{\alpha}_1 & \sqrt{\tilde{\lambda}_2} \tilde{\alpha}_2 & \dots & \sqrt{\tilde{\lambda}_p} \tilde{\alpha}_p \end{array} \right].$$

Ocenjene specifične varijanse dobijamo kao dijagonalne elemente matrice

$$S - \tilde{B}\tilde{B}'$$

pa je

$$\tilde{\Psi}_i = s_{ii} - \sum_{j=1}^m \tilde{\beta}_{ij}^2, \quad i = 1, 2, \dots, p$$

a ocenjeni komunaliteti su

$$\tilde{h}_i^2 = \tilde{\beta}_{i1}^2 + \tilde{\beta}_{i2}^2 + \dots + \tilde{\beta}_{im}^2, \quad i = 1, 2, \dots, p.$$

Broj zadržanih faktora možemo odrediti *a priori* ili posmatranjem **matrice reziduala**

$$S - (\tilde{B}\tilde{B}' + \Psi)$$

koja je dobijena kao rezultat aproksimacije uzoračke kovarijacione matrice sa prvih  $m$  faktora. Broj faktora koje smo uključili u aproksimaciju povećavamo sve dok ne procenimo da su elementi reziduala dovoljno mali.

**Primer 5.4:**

Na osnovu uzorka iz  $p$ -dimenzione populacije ( $p=5$ ) dobijena je uzoračka korelaciona matrica

$$R = \begin{bmatrix} 1 & 0 & 0.8 & 0.3 & 0 \\ 0 & 1 & 0.1 & 0.6 & 0.8 \\ 0.8 & 0.1 & 1 & 0.4 & 0.1 \\ 0.3 & 0.6 & 0.4 & 1 & 0.7 \\ 0 & 0.8 & 0.1 & 0.7 & 1 \end{bmatrix}.$$

Karakteristični vektori i koreni matrice  $R$  su dati u tabeli:

<b>Karakteristični vektori <math>\tilde{\alpha}_j</math></b>				
prvi	drugi	treći	četvrti	peti
0.2775	0.6365	-0.1965	0.6831	-0.1123
0.4893	-0.3508	-0.6102	-0.1294	-0.4985
0.3426	0.5871	-0.1201	-0.6802	0.2468
0.5510	-0.0451	0.7568	-0.0213	-0.3481
0.5124	0.3538	-0.0444	0.2314	0.7462
<b>Karakteristični koreni <math>\tilde{\lambda}_j</math></b>				
2.5835	1.7166	0.3338	0.1941	0.1720

Korišćenjem metode glavnih komponenti oceniti model faktorske analize.

Potrebno je prvo da odredimo broj glavnih komponenti koje ćemo da zadržimo u analizi. Kada se koristi korelaciona matrica uobičajen kriterijum je da se zadrže one glavne komponente kod kojih je varijansa (karakteristični koren) veća od jedinice. Ovaj kriterijum se naziva "**kriterijum jediničnog korena - Kaiser-ov kriterijum**".

Prema ovom kriterijumu ćemo zadržati prve dve glavne komponente jer je

$$\tilde{\lambda}_1 = 2.5835 \text{ i } \tilde{\lambda}_2 = 1.7166.$$

Oni vuku

$$\frac{\tilde{\lambda}_1 + \tilde{\lambda}_2}{p} = \frac{2.5835 + 1.7166}{5} = 0.86$$

tj. 86% uzoračke varijanse.

Sada je potrebno da uporedimo korelacionu matricu na osnovu modela faktorske analize sa dva zajednička faktora

$$\tilde{B}\tilde{B}' + \Psi$$

sa uzoračkom korelacionom matricom  $R$ , i ako je dobijena razlika zanemariva, konstatujemo da ova dva zajednička faktora uspešno reprodukuju korelacionu strukturu pet originalnih promenljivih.

**Napomena:**

Za ortogonalni model važi da je  $Cov(X, F) = Cor(X, F)$ .

Daklem, treba da odredimo

$$\tilde{B}\tilde{B}' + \tilde{\Psi}$$

$$\tilde{B} = \begin{bmatrix} \sqrt{\tilde{\lambda}_1} \tilde{\alpha}_1 & \sqrt{\tilde{\lambda}_2} \tilde{\alpha}_2 \end{bmatrix}$$

$$\tilde{B} = \begin{bmatrix} \sqrt{2.5835} \cdot 0.2775 & \sqrt{1.7166} \cdot 0.6365 \\ \sqrt{2.5835} \cdot 0.4893 & \sqrt{1.7166} \cdot (-0.3508) \\ \sqrt{2.5835} \cdot 0.3426 & \sqrt{1.7166} \cdot 0.5871 \\ \sqrt{2.5835} \cdot 0.5510 & \sqrt{1.7166} \cdot (-0.0451) \\ \sqrt{2.5835} \cdot 0.5124 & \sqrt{1.7166} \cdot (-0.3538) \end{bmatrix} = \begin{bmatrix} 0.4460 & 0.8339 \\ 0.7865 & -0.4596 \\ 0.5507 & 0.7692 \\ 0.8857 & -0.0591 \\ 0.8236 & -0.4635 \end{bmatrix}$$

Da bi odredili  $\tilde{\Psi}$  potrebno je da odredimo komunalitet

$B_1$	$B_2$	$h_i^2 = B_1^2 + B_2^2$	$\tilde{\Psi}_i = 1 - h_i^2$
0.4460	0.8339	$0.4460^2 + 0.8339^2 = 0.8944$	0.1056
0.7865	-0.4596	0.8297	0.1703
0.5507	0.7692	0.8949	0.1051
0.8857	-0.0591	0.7879	0.2121
0.8236	-0.4635	0.8932	0.1068

pa je



$$\begin{aligned}
\tilde{B}\tilde{B}' + \Psi &= \begin{bmatrix} 0.4460 & 0.8339 \\ 0.7865 & -0.4596 \\ 0.5507 & 0.7692 \\ 0.8857 & -0.0591 \\ 0.8236 & -0.4635 \end{bmatrix} \begin{bmatrix} 0.4460 & 0.7865 & 0.5507 & 0.8837 & 0.8236 \\ 0.8339 & -0.4596 & 0.7692 & -0.0591 & -0.4635 \end{bmatrix} + \\
&+ \begin{bmatrix} 0.1056 & 0 & 0 & 0 & 0 \\ 0 & 0.1703 & 0 & 0 & 0 \\ 0 & 0 & 0.1051 & 0 & 0 \\ 0 & 0 & 0 & 0.2121 & 0 \\ 0 & 0 & 0 & 0 & 0.1068 \end{bmatrix} = \\
&= \begin{bmatrix} 1 & -0.0325 & 0.8871 & 0.3451 & -0.0192 \\ & 1 & 0.0796 & 0.7237 & 0.8608 \\ & & 1 & 0.4423 & 0.0971 \\ & & & 1 & 0.7569 \\ & & & & 1 \end{bmatrix}.
\end{aligned}$$

Na osnovu poređenja ove matrice sa matricom  $R$  možemo zaključiti da zajednički faktori relativno uspešno reprodukuju korelacionu strukturu originalnih promenljivih. Kvalitet aproksimacije smo takođe mogli sagledati i na osnovu matrice reziduala:

$$R - (\tilde{B}\tilde{B}' + \Psi) = \begin{bmatrix} 0 & 0.0325 & -0.0871 & -0.0457 & 0.0192 \\ & 0 & 0.0204 & -0.1237 & -0.0608 \\ & & 0 & -0.0423 & 0.0029 \\ & & & 0 & -0.0569 \\ & & & & 0 \end{bmatrix}.$$

### 5.3.2 Metoda glavnih faktora

Ako analizu zasnivamo na korelacionoj matrici, nju razlažemo prema faktorskom modelu

$$\rho = BB' + \Psi$$

gde su dijagonalni elementi

$$\rho_{ii} = 1 = h_i^2 + \Psi_i, \quad i = 1, 2, \dots, p.$$

Zajednički faktori objašnjavaju vandijagonalne elemente matrice  $\rho$  i deo dijagonalnih elemenata  $h_i^2$ . Zato formiramo novu matricu

$$\rho - \Psi = BB' = \begin{bmatrix} h_1^2 & \rho_{12} & \vdots & \rho_{1p} \\ \rho_{21} & h_2^2 & \vdots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \vdots & h_p^2 \end{bmatrix}$$

koja predstavlja korelacionu matricu zajedničkih faktora i naziva se **redukovana korelaciona matrica**.

Do ocene korelacione matrice zajedničkih faktora dolazimo na osnovu razlike ocenjene korelacione matrice i matrice specifičnih varijansi.

Kako je

$$\Psi_i = 1 - h_i^2$$

to je potrebno oceniti komunalitet pre nego što se primeni metod glavnih faktora za ocenu modela faktorske metode. Postoje dve osnovne grupe metoda za ocenjivanje komunaliteta:

1. Za ocenu  $i$ -tog komunaliteta uzima se maksimalna vrednost koeficijenta korelacije  $i$ -te promenljive i ma koje od preostalih  $p-1$  promenljivih. Apsolutna vrednost tog koeficijenta korelacije zamenjuje zatim  $i$ -tu jedinicu na glavnoj dijagonali korelacione matrice (dobra metoda za veliko  $p$ ).
2. Koristi se celokupna korelaciona matrica i ocena za  $h_i^2$  je kvadrat višestrukog koeficijenta korelacije promenljive  $X_i$  i svih preostalih  $p-1$  promenljivih. Ova metoda daje donju granicu ocene komunaliteta i bazira se na matrici  $R^{-1}$ , tako da je ocena komunaliteta

$$\tilde{h}_i^2 = 1 - \tilde{\Psi}_i = 1 - \frac{1}{r^{ii}}, \quad i = 1, 2, \dots, p$$

gde je  $r^{ii}$   $i$ -ti dijagonalni element matrice  $R^{-1}$ .

Nakon ocene komunaliteta i formiranja redukovane korelacione matrice vrši se izbor broja faktora istim postupcima kao kod metode glavnih komponenti. Potom se određuju karakteristični vektori redukovane korelacione matrice koji odgovaraju prvim  $m$  pozitivnim karakterističnim korenima te matrice.

Prema tome, ocena matrice faktorskih opterećenja je

$$B^* = \left[ \begin{array}{ccc|c} \sqrt{\tilde{\lambda}_1^*} \tilde{\alpha}_1^* & & & \\ & \sqrt{\tilde{\lambda}_2^*} \tilde{\alpha}_2^* & & \\ & & \ddots & \\ & & & \sqrt{\tilde{\lambda}_m^*} \tilde{\alpha}_m^* \end{array} \right]$$

gde su

$$(\tilde{\lambda}_j^*, \tilde{\alpha}_j^*), \quad j = 1, 2, \dots, m$$

parovi karakterističnih korena i pridruženih karakterističnih vektora redukovane korelacione matrice.

Ocena specifične varijanse je

$$\tilde{\Psi}_i^* = 1 - \sum_{j=1}^m \tilde{\beta}_{ij}^{*2}, \quad i = 1, 2, \dots, p$$

gde su  $\tilde{\beta}_{ij}^*$  ocenjena faktorska opterećenja.

Na osnovu ovako ocenjene specifične varijanse možemo ponovo oceniti komunalitete

$$\tilde{h}_i^{*2} = \sum_{j=1}^m \tilde{\beta}_{ij}^{*2}, \quad i = 1, 2, \dots, p$$

a zatim ih u narednoj iteraciji koristiti za formiranje redukovane korelacione matrice i određivanje njenih karakterističnih vektora. Ovaj iterativni proces ocenjivanja modela faktorske analize nastavljamo do momenta kada promene u sukcesivnim ocenama komunaliteta ne budu zanemarive.

### **Primer 5.5:**

Neka je poznata kovarijaciona matrica  $\Sigma$  vektora  $X^T = (X_1, X_2, X_3, X_4)$

$$\Sigma = \begin{bmatrix} 19 & 30 & 2 & 12 \\ & 57 & 5 & 23 \\ & & 38 & 47 \\ & & & 68 \end{bmatrix}$$

i neka se znaju faktorska opterećenja za prvi faktor  $F_1$

$$a_{11} = 4, \quad a_{12} = 7, \quad a_{13} = -1 \quad i \quad a_{14} = 1$$

kao i doprinos specifičnog faktora ukupnoj varijansi

$$\rho_1^2 = 2, \rho_2^2 = 4, \rho_3^2 = 1 \text{ i } \rho_4^2 = 3.$$

Izračunati:

- Faktorska opterećenja za faktor  $F_2$  ako se pretpostavi da je model sa dva faktora i datim uticajem specifičnog faktora u potpunosti sposoban da objasni kovarijacionu matricu.
- Proporcije varijabiliteta varijabli  $X_1, X_2, X_3$  i  $X_4$  objašnjene faktorima  $F_1$  i  $F_2$ .
- Promene u reprezentovanju kovarijacione matrice vektora  $\mathbf{X}$ , ako se inicijalno ocenjeni faktori rotiraju matricom  $\mathbf{G}$  za koju važi da je  $\mathbf{G}^T = \mathbf{G}^{-1}$ . Skicirati izgled rotacije od  $30^\circ$  i promene pozicije faktorskih opterećenja.

Rešenje:

a)

$$h_i^2 = \sigma_{ii}^2 - \rho_i, \quad h_i^2 = \beta_{1i}^2 + \beta_{2i}^2 = a_{1i}^2 + a_{2i}^2$$

$$\begin{aligned} h_1^2 &= 17 & h_2^2 &= 53 & h_3^2 &= 37 & h_4^2 &= 65 \\ a_{21}^2 &= 1 & a_{22}^2 &= 4 & a_{23}^2 &= 36 & a_{24}^2 &= 64 \\ a_{21} &= 1 & a_{22} &= 2 & a_{23} &= 6 & a_{24} &= 8 \end{aligned}$$

Faktorska opterećenja za faktor  $F_2$ :

$$F_2(a_{21} = 1 \quad a_{22} = 2 \quad a_{23} = 6 \quad a_{24} = 8)$$

b) Faktorima  $F_1$  i  $F_2$  je objašnjeno:

za	$X_1$	0.8947	varijabiliteta,
za	$X_2$	0.9298	varijabiliteta,
za	$X_3$	0.9736	varijabiliteta,
za	$X_4$	0.9559	varijabiliteta.

c) Ako se inicijalno ocenjeni faktori rotiraju ortogonalnom matricom  $\mathbf{G}$  za koju važi da je  $\mathbf{G}^T = \mathbf{G}^{-1}$ , kovarijaciona matrica se neće promeniti. Rotacija za  $30^\circ$ :

$$M = \begin{bmatrix} \cos 30^\circ & \sin 30^\circ \\ -\sin 30^\circ & \cos 30^\circ \end{bmatrix} = \begin{bmatrix} 0.866 & 0.5 \\ -0.5 & 0.866 \end{bmatrix}$$

$$\tilde{\Gamma} = \tilde{B}M$$

$$\tilde{B} = \begin{bmatrix} 4 & 1 \\ 7 & 2 \\ -1 & 6 \\ 1 & 8 \end{bmatrix}, \quad \tilde{\Gamma} = \begin{bmatrix} 4 & 1 \\ 7 & 2 \\ -1 & 6 \\ 1 & 8 \end{bmatrix} \begin{bmatrix} 0.866 & 0.5 \\ -0.5 & 0.866 \end{bmatrix} = \begin{bmatrix} 2.964 & 2.868 \\ 5.062 & 5.236 \\ -3.866 & 4.708 \\ -3.134 & 7.444 \end{bmatrix}.$$

Rotacijom za  $30^\circ$  se dobija matrica

$$\tilde{\Gamma} = \begin{bmatrix} 2.964 & 2.868 \\ 5.062 & 5.236 \\ -3.866 & 4.708 \\ -3.134 & 7.444 \end{bmatrix}.$$

## 6. KLASITER ANALIZA

### Primer 6.1:

U tabeli su dati podaci o proizvodnji belog i crnog vina za pet proizvođača (u hiljadama litara):

redni broj	proizvođač	belo vino	crno vino
1	NAVIP	5	4
2	SMEDEREVKA	2	4
3	VINO ŽUPA	4	3
4	KRAJINA VINO	7	6
5	KOSOVO VINO	10	9

Izračunati apsolutno i kvadratno odstojanje između proizvođača.

$$E = \{NAVIP, SMEDEREVKA, VINO ŽUPA, KRAJINA VINO, KOSOVO VINO\}$$

$$E = \{e_1, e_2, e_3, e_4, e_5\}, \quad N = 5$$

$$X = (\text{belo vino} \quad \text{crno vino})$$

$$X = (X_1 \quad X_2), \quad p = 2$$

- apsolutno odstojanje:

$$d_{ij} = \sum_{k=1}^p |x_{ik} - x_{jk}|, \quad k = 1, 2$$

$$d_{11} = 0$$

$$i = 1, j = 2:$$

$$d_{12} = |x_{11} - x_{21}| + |x_{12} - x_{22}| = |5 - 2| + |4 - 4| = 3$$

$$i = 1, j = 3:$$

$$d_{13} = |x_{11} - x_{31}| + |x_{12} - x_{32}| = |5 - 4| + |4 - 3| = 2$$

$$d_{14} = |5 - 7| + |4 - 6| = 2 + 2 = 4$$

$$d_{15} = |5 - 10| + |4 - 9| = 5 + 5 = 10$$

$$d_{21} = d_{12}$$

$$d_{22} = 0$$

$$i = 2, j = 3:$$

$$d_{23} = |x_{21} - x_{31}| + |x_{22} - x_{32}| = |2 - 4| + |4 - 3| = 2 + 1 = 3$$

$$d_{24} = |2 - 7| + |4 - 6| = 5 + 2 = 7$$

$$d_{25} = |2 - 10| + |4 - 9| = 8 + 5 = 13$$

$$d_{34} = |4 - 7| + |3 - 6| = 3 + 3 = 6$$

$$d_{35} = |4 - 10| + |3 - 9| = 6 + 6 = 12$$

$$d_{45} = |7 - 10| + |6 - 9| = 3 + 3 = 6$$

$$D = \begin{bmatrix} 0 & 3 & 2 & 4 & 10 \\ & 0 & 3 & 7 & 13 \\ & & 0 & 6 & 12 \\ & & & 0 & 6 \\ & & & & 0 \end{bmatrix}$$

- kvadratno euklidsko odstojanje:

$$d_{12}^2 = (x_{11} - x_{21})^2 + (x_{12} - x_{22})^2 = (5 - 2)^2 + (4 - 4)^2$$

$$d_{13}^2 = (5 - 4)^2 + (4 - 3)^2 = 1^2 + 1^2 = 2$$

$$d_{14}^2 = (5 - 7)^2 + (4 - 6)^2 = 4 + 4 = 8$$

...

$$D = \begin{bmatrix} 0 & 9 & 2 & 8 & 50 \\ & 0 & 5 & 29 & 89 \\ & & 0 & 18 & 72 \\ & & & 0 & 18 \\ & & & & 0 \end{bmatrix}$$

**Primer 6.2:**

Na osnovu podataka iz prethodnog primera formirati hijerarhijsku klasifikaciju korišćenjem metode pojedinačnog povezivanja (*single linkage*) a zatim nacrtati odgovarajući *dendogram*.

$$\begin{array}{l} 1 \text{ NAVIP} \\ 2 \text{ SMEDEREVKA} \\ 3 \text{ VINO ŽUPA} \\ 4 \text{ KRAJINA VINO} \\ 5 \text{ KOSOVO VINO} \end{array} \begin{bmatrix} 0 & 9 & \textcircled{2} & 8 & 50 \\ & 0 & 5 & 29 & 89 \\ & & 0 & 18 & 72 \\ & & & 0 & 18 \\ & & & & 0 \end{bmatrix} = D^{(1)}$$

Prvo se povezuju ona dva elementa koja imaju najmanje odstojanje

$$\min d_{ij} = 2 = d_{13}$$

Znači, 1 NAVIP i 3 VINO ŽUPA čine jedan klaster. Sada primenjujemo algoritam klasifikovanja

$$D(e_i, e_j) = \min d_{sl}$$

$$\begin{array}{l} 1 \text{ U } 3 \text{ NAVIP U VINO ŽUPA} \\ 2 \text{ SMEDEREVKA} \\ 4 \text{ KRAJINA VINO} \\ 5 \text{ KOSOVO VINO} \end{array} \begin{bmatrix} 0 & \textcircled{5} & 8 & 50 \\ & 0 & 29 & 89 \\ & & 0 & 18 \\ & & & 0 \end{bmatrix} = D^{(2)}$$

$$\min d_{ij} = 5$$

Klasteru NAVIP U VINO ŽUPA se pridružuje klaster 2 SMEDEREVKA.

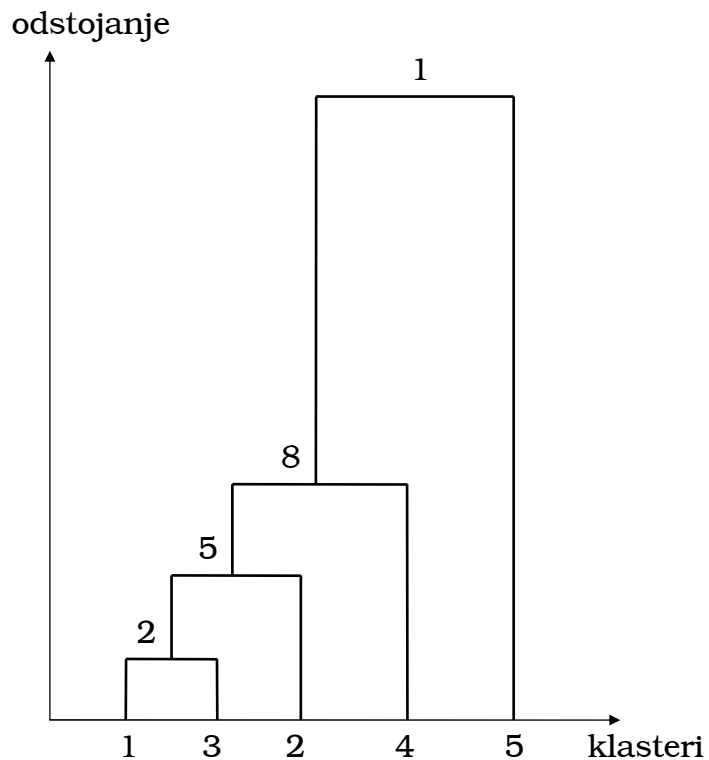
$$\begin{array}{l} 1 \text{ U } 3 \text{ U } 2 \\ 4 \text{ KRAJINA VINO} \\ 5 \text{ KOSOVO VINO} \end{array} \begin{bmatrix} 0 & \textcircled{8} & 50 \\ & 0 & 18 \\ & & 0 \end{bmatrix} = D^{(3)}$$

Klasteru 1 U 3 U 2 se pridružuje klaster 4 KRAJINA VINO.



$$\begin{array}{l} 1 \text{ U } 3 \text{ U } 2 \text{ U } 4 \\ 5 \text{ KOSOVO VINO} \end{array} \begin{bmatrix} 0 & 18 \\ & 0 \end{bmatrix}$$

Dendogram:



$$D^{(h)} = \begin{bmatrix} 0 & 5 & 2 & 8 & 18 \\ & 0 & 5 & 8 & 18 \\ & & 0 & 8 & 18 \\ & & & 0 & 18 \\ & & & & 0 \end{bmatrix}$$

**Primer 6.3:**

Neka su dati podaci o indeksu rezultata i snage sedam desetbojaca:

takmičar	rang snage	rang rezultata
<i>A</i>	10	8
<i>B</i>	2	2
<i>V</i>	2	5
<i>G</i>	4	2
<i>D</i>	2	3
<i>E</i>	8	9
<i>Z</i>	8	7

**K-mean** algoritmom podeliti skup desetobojaca na dve grupe sa početnim rešenjima za prvi centroid  $A$ , a za drugi centroid  $D$  i izračunati međusobno odstojanje konačno formiranih klasa.

$$A^0 = (10,8), B^0 = (2,3)$$

$$\begin{array}{lll} d(e_1, A^0) = 14 & d(e_1, B^0) = 1 & e_1 \rightarrow B' \\ d(e_2, A^0) = 11 & d(e_2, B^0) = 2 & e_2 \rightarrow B' \\ d(e_3, A^0) = 12 & d(e_3, B^0) = 3 & e_3 \rightarrow B' \\ d(e_4, A^0) = 3 & d(e_4, B^0) = 12 & e_4 \rightarrow A' \\ d(e_5, A^0) = 3 & d(e_5, B^0) = 10 & e_5 \rightarrow A' \end{array}$$

$$A' = (8.67, 8), B' = (2.5, 3)$$

$$\begin{array}{lll} d(e_1, A') = 12.67 & d(e_1, B') = 1.5 & e_1 \rightarrow B'' \\ d(e_2, A') = 9.67 & d(e_2, B') = 2.5 & e_2 \rightarrow B'' \\ d(e_3, A') = 10.67 & d(e_3, B') = 2.5 & e_3 \rightarrow B'' \\ d(e_4, A') = 1.67 & d(e_4, B') = 11.5 & e_4 \rightarrow A'' \\ d(e_5, A') = 1.67 & d(e_5, B') = 9.5 & e_5 \rightarrow A'' \end{array}$$

Konačno rešenje:

$$A'' \{A, E, Z\} \quad B'' \{B, V, G, D\}.$$

Međusobno odstojanje klasa:

$$d_{I,II} = 11.17.$$

## 7. REŠENI ZADACI

### Zadatak 1.

Neka je data kovarijaciona matrica

$$\Sigma = \begin{bmatrix} 13 & -4 & 2 \\ & 13 & -2 \\ & & 10 \end{bmatrix}$$

i neka su poznata dva korena jednačine  $|\Sigma - \lambda I| = 0$

$$\lambda_1 = 9 \text{ i } \lambda_2 = 9.$$

- Odrediti varijabilitet prve glavne komponente.
- Odrediti analitički oblik druge glavne komponente.
- Koliko treba uzeti glavnih komponentata da bi bilo obuhvaćeno barem 70% ukupnog varijabiliteta vektora  $\mathbf{X}$ ?

a)

$$\begin{aligned} \text{tr}(\Sigma) = \text{tr}(\Lambda) &\Rightarrow 13 + 13 + 10 = \lambda_1 + \lambda_2 + \lambda_3 \\ 36 &= 18 + \lambda_3 \Rightarrow \lambda_3 = 36 - 18 = 18 \end{aligned}$$

Varijabilitet prve glavne komponente je

$$\max \lambda_i = 18.$$

b)

$$\Sigma P_2 = \lambda_2 P_2 \Leftrightarrow \begin{bmatrix} 13 & -4 & 2 \\ -4 & 13 & -2 \\ 2 & -2 & 10 \end{bmatrix} \begin{bmatrix} p_{12} \\ p_{22} \\ p_{32} \end{bmatrix} = 9 \begin{bmatrix} p_{12} \\ p_{22} \\ p_{32} \end{bmatrix}$$

$$\begin{aligned}
13p_{12} - 4p_{22} + 2p_{32} &= 9p_{12} & 4p_{12} - 4p_{22} + 2p_{32} &= 0 \\
-4p_{12} + 13p_{22} - 2p_{32} &= 9p_{22} \Leftrightarrow -4p_{12} + 4p_{22} - 2p_{32} &= 0 \\
2p_{12} - 2p_{22} + 10p_{32} &= 9p_{32} & 2p_{12} - 2p_{22} + p_{32} &= 0
\end{aligned}$$

Dobili smo tri jednačine sa tri nepoznate ali su prve dve redundantne. Za  $P_{21}$  i  $P_{22}$  odgovara bilo koje rešenje, npr  $P_{21}=1$  i  $P_{22}=1$ , Iz treće jednačine dobijamo da je  $P_{32}=0$ . Vektor

$$(1 \ 1 \ 0)^T$$

zadovoljava sistem jednačina ali ne i uslov normiranosti karakterističnih vektora. Da bi se to postiglo, rešenja se dele sa korenom zbira kvadrata, tj. sa

$$\sqrt{1^2 + 1^2 + 0^2} = \sqrt{2}$$

pa je konačno

$$P_2 = \left( \frac{1}{\sqrt{2}} \quad \frac{1}{\sqrt{2}} \quad 0 \right)^T.$$

c)

$$\frac{\lambda_1}{tr(\Lambda)} = 0.5, \quad \frac{\lambda_2}{tr(\Lambda)} = 0.25, \quad \frac{\lambda_3}{tr(\Lambda)} = 0.25, \quad tr(\Sigma) = tr(\Lambda) = 36$$

$$V(P_1^T X) < 0.7 \cdot 36$$

$$V(P_1^T X) + V(P_2^T X) > 0.7 \cdot 36$$

Oдавde zaključujemo da treba uzeti prve dve glavne komponente da bi se obuhvatilo najmanje 70% ukupnog varijabiliteta.

**Zadatak 2.**

Ako su poznati varijabiliteti prve dve glavne komponente  $V(Y_1)=18$  i  $V(Y_2)=9$  i ako se zna ukupni varijabilitet vektora  $\mathbf{X}$

$$\text{tr}(\Sigma) = 36$$

pri čemu su glavne komponente date sa

$$\begin{aligned} Y_1 &= \frac{2}{3}X_1 - \frac{2}{3}X_2 + \frac{1}{3}X_3 \\ Y_2 &= \frac{1}{\sqrt{2}}X_1 + \frac{1}{\sqrt{2}}X_2 \\ Y_3 &= \frac{1}{\sqrt{18}}X_1 - \frac{1}{\sqrt{18}}X_2 - \frac{4}{\sqrt{18}}X_3 \end{aligned}$$

- a) odrediti procenat varijabiliteta koji nosi treća komponenta,  
b) odrediti kovarijacionu matricu vektora  $\mathbf{X}$ .

a)

$$V(Y_3) = \lambda_3 = \text{tr}(\lambda) - (\lambda_1 + \lambda_2) = \text{tr}(\lambda) - \lambda_1 - \lambda_2 = 36 - 27 = 9$$

$$\frac{\lambda_3}{\text{tr}(\Lambda)} \cdot 100 = \frac{9}{36} \cdot 100 = 25\% \text{ varijabiliteta.}$$

b)

$$V(\mathbf{X}) = \Sigma = P\Lambda P^T = P \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix} P^T = \lambda_1 P_1 P_1^T + \lambda_2 P_2 P_2^T + \lambda_3 P_3 P_3^T =$$

$$= 18 \cdot \begin{bmatrix} \frac{2}{3} \\ -\frac{2}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} \begin{bmatrix} \frac{2}{3} & -\frac{2}{3} & \frac{1}{3} \end{bmatrix} + 9 \cdot \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \end{bmatrix} + 9 \cdot \begin{bmatrix} \frac{1}{\sqrt{18}} \\ -\frac{1}{\sqrt{18}} \\ -\frac{4}{\sqrt{18}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{18}} & -\frac{1}{\sqrt{18}} & -\frac{4}{\sqrt{18}} \end{bmatrix} =$$

$$= 18 \cdot \begin{bmatrix} \frac{4}{9} & -\frac{4}{9} & \frac{2}{9} \\ -\frac{4}{9} & \frac{4}{9} & -\frac{2}{9} \\ \frac{2}{9} & -\frac{2}{9} & \frac{1}{9} \end{bmatrix} + 9 \cdot \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix} + 9 \cdot \begin{bmatrix} \frac{1}{18} & -\frac{1}{18} & -\frac{4}{18} \\ -\frac{1}{18} & \frac{1}{18} & \frac{4}{18} \\ -\frac{4}{18} & \frac{4}{18} & \frac{16}{18} \end{bmatrix} =$$

$$= \begin{bmatrix} 13 & -4 & 2 \\ -4 & 13 & -2 \\ 2 & -2 & 10 \end{bmatrix} = \Sigma.$$

**Zadatak 3.**

Neka je dat vektor

$$X = (X_1 \quad X_2 \quad X_3)^T$$

i njegova kovarijaciona matrica

$$\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

sa karakterističnim vektorima

$$P_1^T = (0.383 \quad -0.924 \quad 0)$$

$$P_2^T = (0 \quad 0 \quad 1)$$

$$P_3^T = (0.924 \quad 0.383 \quad 0)$$

- a) Izračunati glavne komponente  $Y_1$ ,  $Y_2$  i  $Y_3$  i njihove varijanse.  
 b) Izračunati kovarijansu i korelaciju prve glavne komponente  $Y_1$  sa varijablom  $X_1$ .

a)

$$Y_1 = P_1^T X = 0.383X_1 - 0.924X_2$$

$$Y_2 = P_2^T X = X_3$$

$$Y_3 = P_3^T X = 0.924X_1 + 0.383X_2$$

$$\begin{aligned} V(Y_1) &= V(0.383X_1 - 0.924X_2) = 0.383^2V(X_1) + 0.924^2V(X_2) + \\ &\quad + 2 \cdot 0.383 \cdot (-0.924)Cov(X_1, X_2) = 0.147 \cdot 1 + 0.854 \cdot 5 - 0.708 \cdot (-2) = \\ &= 5.83 = \lambda_1 \end{aligned}$$

$$V(Y_2) = V(X_3) = 2 = \lambda_2$$

$$V(Y_3) = tr(\Sigma) - (\lambda_1 + \lambda_2) = 8 - 7.83 = 0.17$$

b)  $X_1$  se može napisati kao

$$(1 \ 0 \ 0) \cdot \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = L_1^T X$$

$$\begin{aligned} \text{Cov}(X_1, Y_1) &= \text{Cov}(L_1^T X, P_1^T X) = L_1^T \text{Cov}(X, X) P_1 = \\ &= L_1^T V(X) P_1 = L_1^T \Sigma P_1 = L_1^T \lambda_1 P_1 \end{aligned}$$

jer je  $\Sigma P_1 = \lambda_1 P_1$ . Stoga je

$$\text{Cov}(X_1, Y_1) = (1 \ 0 \ 0) \cdot 5.83 \cdot \begin{bmatrix} 0.383 \\ -0.924 \\ 0 \end{bmatrix} = 5.83 \cdot 0.383 = 2.23$$

$$\text{Cor}(X_1, Y_1) = \frac{\text{Cov}(X_1, Y_1)}{\sqrt{V(X_1)V(Y_1)}} = \frac{2.23}{\sqrt{1 \cdot 5.83}} = 0.925.$$

#### Zadatak 4.

Neka je data korelaciona matrica

$$R = \begin{bmatrix} 1 & 0.63 & 0.45 \\ & 1 & 0.35 \\ & & 1 \end{bmatrix}$$

a) Pokazati da se  $R$  može generisati modelom faktorske analize sa jednim faktorom tipa

$$X_1^* = 0.9F_1 + \varepsilon_1$$

$$X_2^* = 0.7F_1 + \varepsilon_2$$

$$X_3^* = 0.5F_1 + \varepsilon_3$$

pri čemu je

$$X^* = \frac{X - \mu}{\sigma}$$

$$V(F_1) = 1$$

$$\text{Cov}(F_1, \varepsilon_i) = 0, \quad i = 1, 2, 3$$

$$V(\varepsilon) = \begin{bmatrix} 0.19 & & \\ & 0.51 & \\ & & 0.75 \end{bmatrix}$$

b) Izračunati komunalitete.

a) Model je  $X = AF_1 + \varepsilon$ , a  $R$  se ocenjuje sa

$$\begin{aligned} \tilde{R} &= AA^T + V(\varepsilon) = \begin{bmatrix} 0.9 \\ 0.7 \\ 0.5 \end{bmatrix} \begin{bmatrix} 0.9 & 0.7 & 0.5 \end{bmatrix} + \begin{bmatrix} 0.19 & & \\ & 0.51 & \\ & & 0.75 \end{bmatrix} = \\ &= \begin{bmatrix} 0.81 & 0.63 & 0.45 \\ 0.63 & 0.49 & 0.35 \\ 0.45 & 0.35 & 0.25 \end{bmatrix} + \begin{bmatrix} 0.19 & & \\ & 0.51 & \\ & & 0.75 \end{bmatrix} = R \end{aligned}$$

Kako je  $R - \tilde{R} = 0$  zaključujemo da model generiše  $R$ .

b)

$$h_i^2 = 1 - \varepsilon_i \Rightarrow h_1^2 = 0.81, h_2^2 = 0.49 \text{ i } h_3^2 = 0.25.$$

### Zadatak 5.

Neka je data matrica podataka

	$X_1$	$X_2$
$e_1$	5	3
$e_2$	-1	1
$e_3$	1	-2
$e_4$	-3	-2

a) Podeliti skup  $E = \{e_1, e_2, e_3, e_4\}$  na dva skupa **Quick cluster** algoritmom korišćenjem euklidskog odstojanja  $d_{ij} = \sum_{k=1}^2 (X_{ik} - X_{jk})^2$  uz dato početno rešenje da je



$$A^0 = \{e_1, e_2\} \text{ i } B^0 = \{e_3, e_4\}.$$

b) Izračunati na isti način distancu između konačno formiranih klasa.

a) Centroidi početnog rešenja:

	$X_1$	$X_2$
$A^0$	$\frac{5-1}{2} = 2$	$\frac{3+1}{2} = 2$
$B^0$	$\frac{1+3}{2} = 2$	$\frac{-2-2}{2} = -2$

$$d(e_1, A^0) = (5-2)^2 + (3-2)^2 = 10$$

$$d(e_1, B^0) = (5-1)^2 + (3+2)^2 = 61 \Rightarrow d(e_1, A^0) < d(e_1, B^0) \Rightarrow e_1 \rightarrow A'$$

$$d(e_2, A^0) = 10$$

$$d(e_2, B^0) = 9 \Rightarrow e_2 \rightarrow B'$$

$$d(e_3, A^0) = 17$$

$$d(e_3, B^0) = 1 \Rightarrow e_3 \rightarrow B'$$

$$d(e_4, A^0) = 41$$

$$d(e_4, B^0) = 25 \Rightarrow e_4 \rightarrow B'$$

$$A' = \{e_1\} \text{ i } B' = \{e_2, e_3, e_4\}$$

pa su centroidi

$$A'(5,3) \text{ i } B'(-1,-1)$$

$$d(e_1, A') = 0$$

$$d(e_1, B') = 52 \Rightarrow e_1 \rightarrow A''$$

$$d(e_2, A') = 40$$

$$d(e_2, B') = 0 \Rightarrow e_2 \rightarrow B''$$

$$d(e_3, A') > d(e_3, B') \Rightarrow e_3 \rightarrow B''$$

$$d(e_4, A') > d(e_4, B') \Rightarrow e_4 \rightarrow B''$$

Kako je  $A' \rightarrow A''$  i  $B' \rightarrow B''$  konačna klasifikacija je:

$$A = \{e_1\} \text{ i } B = \{e_2, e_3, e_4\}$$

b)  $d(A, B) = (\bar{X}_{1A} - \bar{X}_{1B})^2 + (\bar{X}_{2A} - \bar{X}_{2B})^2 = 52.$

## 8. LITERATURA

1. Kovačić Zlatko: *"Multivarijaciona analiza"*