

**A
D
A
T
R
E
D
U
K
C
I
Ó**

Olyan statisztikai módszerek tartoznak ide, melyek lehetővé teszik, hogy az adatmátrix méretét csökkentve kisebb költséggel értékelhessük ki a statisztikai sokaságot.

A redukált adatmennyiségből levont statisztikai következtetések érvényesek maradnak az eredeti statisztikai sokaságra is.

A csökkentés vonatkozhat az esetszám csökkentésére és a változók számának a csökkentésére egyaránt.

- **Klaszteranalízis**
- **Ritkítás véletlenszám generálással**
- **Faktoranalízis**
- **Fokkomponens-analízis**
- **Többdimenziós skálázás (MDS)**

**F
A
K
T
O
R
A
N
A
L
Í
Z
I
S**

Nagyszámú, sztochasztikusan erősen összefüggő változónk van. A változók redundáns információt hordoznak.

Ismeretlen, kisszámú faktorváltozót keresünk.

- **Hogyan lehet a változók által közösen magyarázott információt korrelálatlan faktorokkal kifejezni?**
- **A faktorok milyen mértékben magyarázzák az eredeti változókat?**
- **Mely változók vannak ugyanazokkal a faktorokkal kifejezve?**
- **Hogyan lehet ezek alapján a változóinkat csoportosítani?**
- **Mi lehet az egyes faktorok jelentése?**

A VÁLTOZÓK KÖZÖTTI ÖSSZEFÜGGÉS EREJÉNEK MÉRÉSE

$$KMO = \frac{\sum_{i=1}^p \sum_{\substack{j=1 \\ j \neq i}}^p r_{ij}^2}{\sum_{i=1}^p \sum_{\substack{j=1 \\ j \neq i}}^p \rho_{ij}^2 + \sum_{i=1}^p \sum_{\substack{j=1 \\ j \neq i}}^p r_{ij}^2}$$

$$\rho_{ij} = \frac{R_i}{\sqrt{R_i \cdot R_j}}$$

parciális korrelációs együttható

$$r_{ij} = R(X_i, X_j)$$

Kaiser-Meyer-Olkin mérték

korrelációs együttható

- 0.9 ≤ KMO
- 0.8 ≤ KMO < 0.9
- 0.7 ≤ KMO < 0.8
- 0.6 ≤ KMO < 0.7
- 0.5 ≤ KMO < 0.6
- KMO < 0.5

- csodálatos (marvelous)
- dicséretes (meritorious)
- közepes middling)
- mérsékelt (mediocre)
- szánalmas (miserable)
- elfogadhatatlan (unacceptable)

A VÁLTOZÓK KÖZÖTTI ÖSSZEFÜGGÉS EREJÉNEK MÉRÉSE

$$MSA_i = \frac{\sum_{i \neq j} r_{ij}^2}{\sum_{i \neq j} \rho_{ij}^2 + \sum_{i \neq j} r_{ij}^2} \quad (i = 1, 2, \dots, p)$$

measure of sampling adequacy

Az indulási p db változóból azokat érdemes elhagyni, amelyeknél az MSA_i érték a legkisebb.

Elvégezhető még a Bartlett-féle gömb próba.

Itt az a nullhipotézis, hogy a vizsgált változók függetlenek egymástól. Akkor érdemes továbbmenni, ha ez a próba nem szignifikáns!

A k -FAKTOROS MODELL

Adottak az X_1, X_2, \dots, X_p változók, a belőlük alkotott p -dimenziós vektor \underline{X} .

$$\underline{X} = \underline{A} \cdot \underline{F} + \underline{U} + \underline{m}$$

\underline{A} pk -as átviteli mátrix

\underline{F} k -dimenziós közös faktor-vektor

\underline{U} p -dimenziós egyedi faktor-vektor

$E\underline{X} = \underline{m}$ várható érték vektor

A k -FAKTOROS MODELL FELTÉTELEI

F_1, F_2, \dots, F_k páronként korrelálatlanok $R(F_i, F_j) = 0$

$$E(F_i) = 0, \quad D^2(F_i) = 1$$

U_1, U_2, \dots, U_p páronként korrelálatlanok $R(U_i, U_j) = 0$

$$E(U_i) = 0, \quad D^2(U_i) = \Psi_{ii}$$

F_1, F_2, \dots, F_k és U_1, U_2, \dots, U_p páronként korrelálatlanok:

$$R(U_i, U_j) = 0$$

A k -FAKTOROS MODELL FELTÉTELEI

Egy k -faktoros modell pontosan akkor oldható meg, ha

$$\underline{\Sigma} = \underline{A} \cdot \underline{A}^T + \underline{\Psi}$$

$\underline{\Sigma} = \text{cov}(X_i, X_j)$ \underline{X} kovarianciamátrixa

$\underline{\Psi}$ \underline{U} kovarianciamátrixa

Van $p(p+1)/2$ egyenlet, és $p(k+1)$ ismeretlen

$(p+1)/2 > k+1$ esetben az egyenletrendszer túldefiniált

$(p+1)/2 < k+1$ esetben az egyenletrendszer aluldefiniált

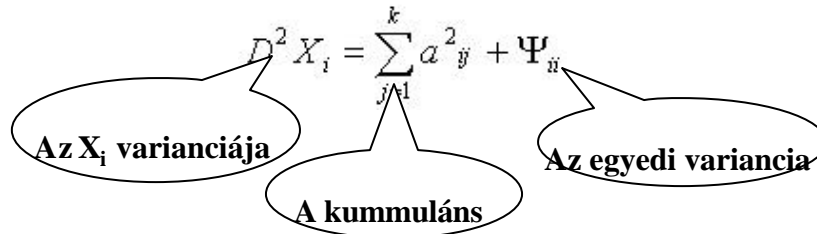
- Maximum likelihood módszer
- Fokkomponens-analízis
- A legkisebb négyzetek módszere
-

A k -FAKTOROS MODELL KOORDINÁNTÁNKÉNT

$$X_i = \sum_{j=1}^k a_{ij} F_j + U_i + m_i$$

\underline{F} koordinátái mindegyik X_i eloállításában szerepelnek

\underline{U} koordinátái közül csak U_i szerepel X_i eloállításában



$\frac{\sum_{j=1}^k a_{ij}^2}{D^2 X_i}$ Ez az arány azt fejezi ki, hány %-ot magyaráznak a közös faktorok.

A FAKTOROK FORGATÁSA (ROTÁCIÓ)

$$\underline{\Sigma} = \underline{A} \cdot \underline{A}^T + \underline{\Psi} = \underline{A} \underline{E} \underline{A}^T + \underline{\Psi} = \underline{A} \underline{G} \underline{G}^T \underline{A}^T + \underline{\Psi} = \underline{A}^* \cdot \underline{A}^{*T} + \underline{\Psi}$$

$$\underline{A}^* = \underline{A} \underline{G}$$

Az új átviteli mátrix

$$\underline{F}^* = \underline{G} \underline{F}$$

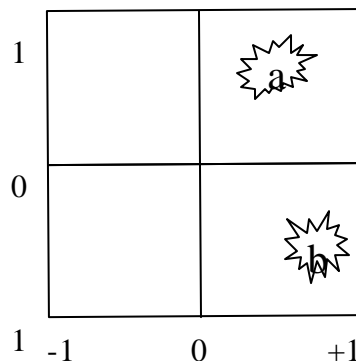
Az új faktorvektor

G ügyes megválasztásával a modell jobban magyarázható lesz!

- Varimax azon változók száma kevés lesz, melyekhez sok faktor szerenel nagy súllyal
- Quartimax a magyarázó faktorok számát minimalizálja
- Equamax a két eljárás keverékét végzi

A FAKTOROK FORGATÁSA (ROTÁCIÓ)

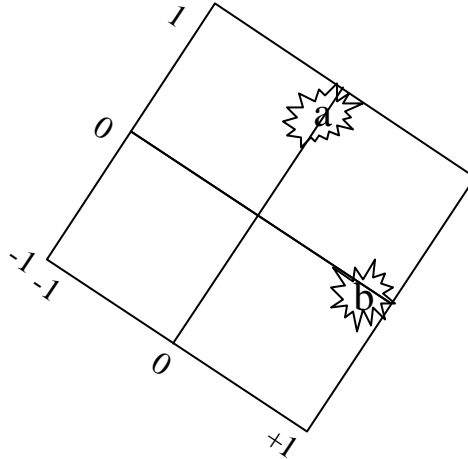
A rotáció szemléltetése egy egyszerű kétdimenziós példán:



Az eredeti változók **a** és **b** csoportja a rotáció nélkül kapott mindkét faktoron jelentős faktorsúllyal rendelkezik.

A FAKTOROK FORGATÁSA (ROTÁCIÓ)

A rotáció szemléltetése egy egyszerű kétdimenziós példán:



Az eredeti változók **a** csoportja csak a rotációval kapott egyik, a **b** csoport pedig csak a másik faktoron rendelkezik jelentős faktorsúllyal.

F O K O M P O N E N S A N A L Í Z I S

A faktoranalízis speciális esete. Dimenziószám csökkentésre használható. Az eredetileg p változóval jellemzett statisztikai sokaságot $k \ll p$ változóval (fokkomponensekkel) jellemezzük. A k -dimenziós statisztikai elemzések következtetései a p -dimenziós sokaságra is érvényesek lesznek. Ezzel jelentős költséget lehet megtakarítani.

Lehetőség van a $p > 3$ dimenziós sokaságot (ha $k < 4$) pontfelhő grafikonon szemléltetni.

A fokkomponensek terében a változók korrelálatlanok lesznek.

A fokkomponens-transzformáció:

$$\underline{F} = \underline{G}^T (\underline{X} - \underline{m}) \quad \text{vagy} \quad \underline{X} = \underline{G} \cdot \underline{F} + \underline{m}$$

$$\underline{F} = (F_1, F_2, \dots, F_p)$$

a fokkomponens-vektor

$$\underline{G} = (\underline{g}_1, \underline{g}_2, \dots, \underline{g}_p)$$

a főirányok mátrixa

A FOKOMPONENS-MODELL TULAJDONSÁGAI

- A fokkomponensek korrelálatlanok: $R(F_i, F_j) = 0$

- A fokkomponensek csökkenő jelentőségűek:

$E(F_1)$ F_1 magyaráz a legtöbbet, F_2 a második legtöbbet, ..., F_p magyaráz a legkevesebbet T -ből.

$$T = \sum_{i=1}^p \lambda_i = \sum_{i=1}^p \lambda_i$$

$\frac{\lambda_i}{T}$ megmutatja, hány %-ot magyaráz F_i

A FOKOMPONENS-MODELL TULAJDONSÁGAI

- A főirányok jelentése: $\sum \underline{g}_i = \lambda_i \underline{g}_i$

\underline{g}_1 ebben az irányban a legnagyobb a variancia

\underline{g}_2 ebben az irányban a legnagyobb a variancia a \underline{g}_1 irányra merőleges irányok között

\vdots

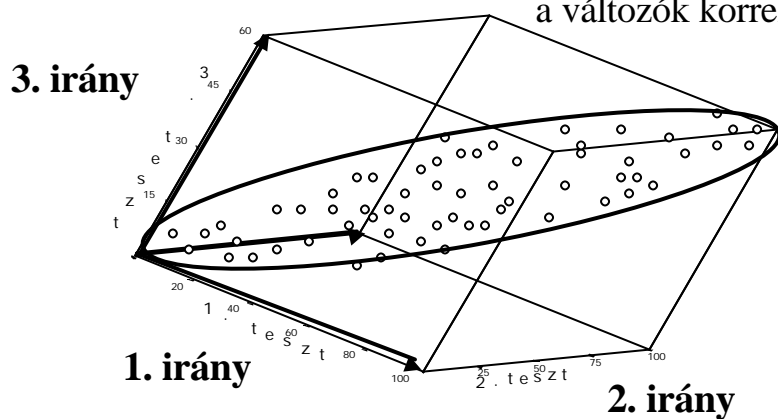
- Dimenziócsökkentés:

Ha \underline{X} helyett az első k főfaktor-alkotta vektorral számolunk,

az elvesztett információ csupán $1 - \frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{T}$

FOKOMPONENSANALÍZIS

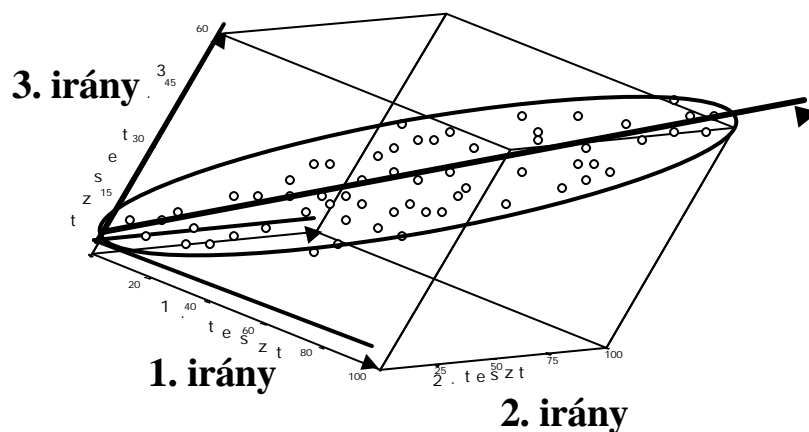
Tengelyek nem derékszögeket zárnak be: a változók korreláltak!



FOKOMPONENSANALÍZIS

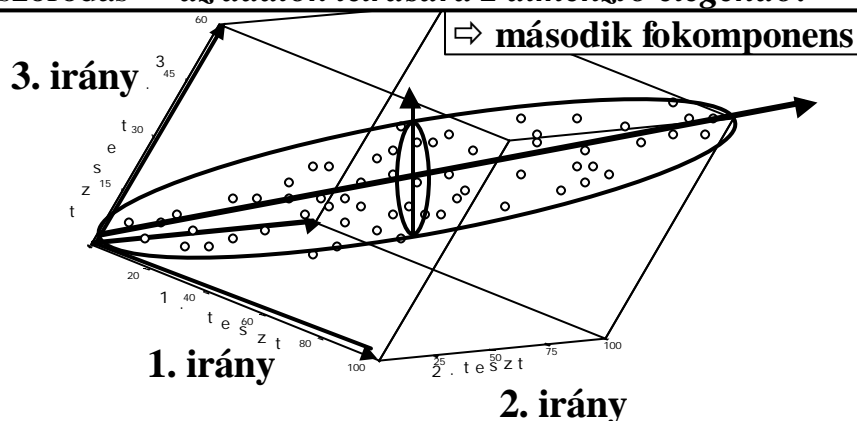
Ebben az irányban tudunk legjobban differenciálni a pontok között. A fokomponensek hosszát (fontosságát) az ún. sajátértékkel (*eigenvalue*) jellemezzük, ami az értelmezett variancia.

első fokomponens



FOKOMPONENSANALÍZIS

Az eljárást folytatni lehetne a harmadik fokkomponens megkeresésével, de ennek a konkrét esetben már nincs értelme, mivel ebben az irányban már jelentéktelen a szóródás \Rightarrow **az adatok leírására 2 dimenzió elegendő!**



**T
Ö
B
B
D
I
M
E
N
Z
I
Ó
S

S
K
A
L
A
Z
Á
S**

Egy p -dimenziós sokaságot lehet egy $k=1,2$ vagy 3 dimenziós Euklideszi pontthalmazzal vizualizálni. A pontthalmaz távolságviszonyai az eredeti sokaság eseteinek távolságviszonyaival nagymértékben egyezik.

A vizualizálás révén tanulmányozható a statisztikai sokaság térbeli struktúrájának. Jellegzetes tömörülések, irányok fedezhetők fel az elkészült scatter-grafikonon.

$$X = \begin{pmatrix} x_1^1 & x_1^2 & \dots & x_1^p \\ x_2^1 & x_2^2 & \dots & x_2^p \\ \vdots & \vdots & \ddots & \vdots \\ x_n^1 & x_n^2 & \dots & x_n^p \end{pmatrix}$$

- $\underline{x}^{(1)}$ 1. eset
- $\underline{x}^{(2)}$ 2. eset
- $\underline{x}^{(n)}$ n. eset

$\underline{\underline{D}} = (d(\underline{x}^{(i)}, \underline{x}^{(j)}))$ Az esetvektorok egymástól vett $n \times n$ -es távolságmátrixa

Megkonstruálhatók olyan $\underline{z}_1, \underline{z}_2, \dots, \underline{z}_n$ $k=1,2$ vagy 3 dimenziós vektorok, melyek $n \times n$ -es Euklideszi távolságmátrixa nagymértékben hasonló $\underline{\underline{D}}$ -hez

$$\sum_{i=1}^n \sum_{j=1}^n \left(d(\underline{x}^{(i)}, \underline{x}^{(j)}) - \rho(\underline{z}_i, \underline{z}_j) \right)^2 \quad \text{„kicsi”}$$