

KLASZTERANALÍZIS

Az eseteket homogén csoportokba (ú.n. klaszterekbe) soroljuk. A csoportosítás alapja egy adott metrika szerinti közelség, illetve egy adott hasonlósági mérték szerinti hasonlóság.

DISZKRIMINANCIAALÍZIS

Az esetek egy kategóriaváltozó értékei alapján osztályokba vannak tagolva. A feladat az, hogy a többdimenziós térben az osztályokat szeparáló felületekkel elválasszuk.

OSZTÁLYOZÁS

Ismert kategóriájú esetek segítségével (tananyag) döntésfüggvényt konstruálunk, amivel ismeretlen kategóriájú esetekhez is tudunk osztályokat rendelni.

PÉLDÁK

KLASZTERANALÍZIS

- Milyen csoportok alakíthatók ki az *employee* állományban a fizetési adatok (*salary, salbegin*) alapján?
- Milyen csoportosulások keletkeznek az országok halmazában, ha az egészségügyi helyzetet jellemző változókat tekintjük: *lifeexpf, lifeexpm, babymort, calories, aids_rt, b_to_d*
- Milyen csoportosulások keletkeznek az országok halmazában, ha a gazdasági helyzetet jellemző változókat tekintjük: *gdp_cap, cropgrow, urban*

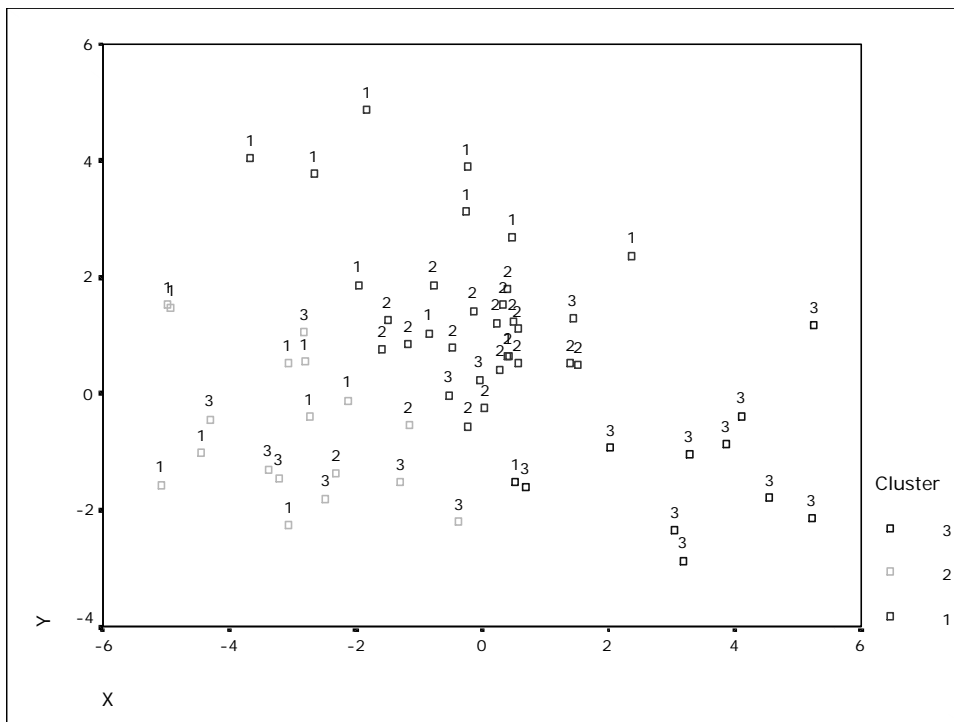
PÉLDÁK

DISZKRIMINANCIAALÍZIS

- A fizetési adatok mennyire választható szét a *jobcat*, *gender* illetve *minority* kategóriaváltozók alapján?
- Mennyire válnak szét az országok a gazdasági tömörülés (*region*) alapján?

OSZTÁLYOZÁS

- Betufelismerés
- Muholdképpontok osztályozása
- Banki rizikóelemzés: kapjon hitelt? Ne kapjon?
- Orvosi diagnosztika: Beteg? Nem beteg?
- Repülésirányítás: Felszálljon? Töröljék?



A *k*-közép módszer (*K-Means Cluster Analysis*)

Olyan dinamikus klaszterező eljárás, amikor előre meg kell adni a klaszterek számát. A klaszter-középpontok térbeli helyzetét iterációban állandóan változtatjuk, amíg egy stabil állapot ki nem alakul. Az esetvektorok a legközelebbi klaszterközépponthez lesznek rendelve.

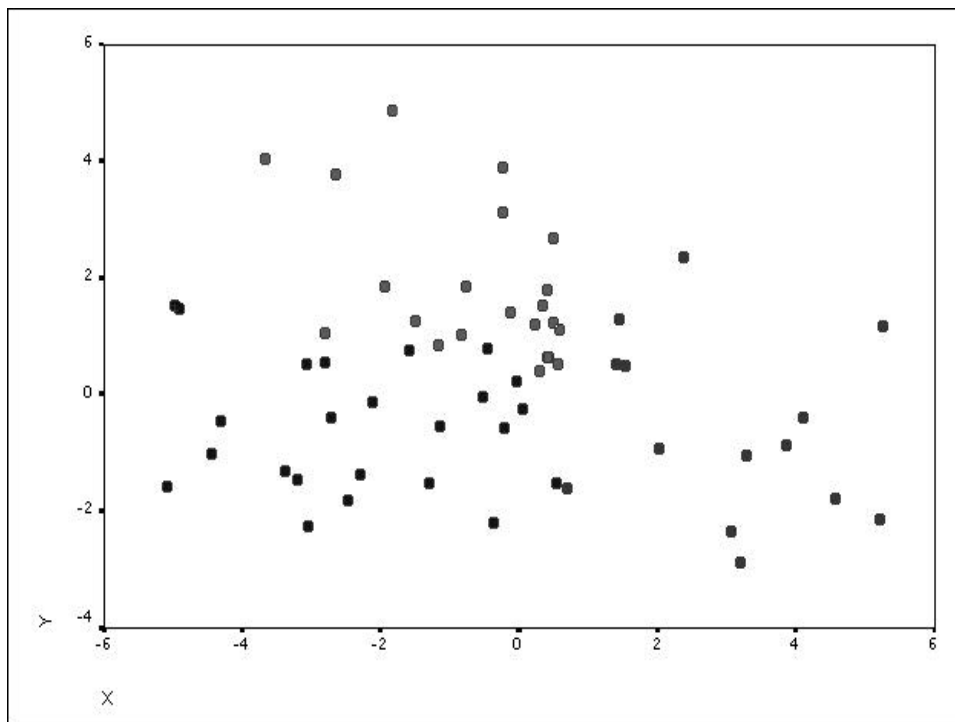
Előnye:

Nagy esetszámú adatmátrix feldolgozható vele.

Hátránya:

A metrika beépített, körülményes a koordinátasúlyozás

Előre meg kell adni a klaszterek számát



A hierarchikus klaszterezés (Hierarchical Cluster Analysis)

Egyelemu klaszterekbol kiindulva, minden lépésben a két legközelebbi fekvő klasztert összevonva csökkentjük a klaszterek számát, amíg minden eset egyetlen klaszterbe nem kerül. A folyamatot regisztráló dendogramot utólag kielemezve, azt a köztes állapotot fogadjuk el, amikor az összevonás erőltetett volt, azaz az összevont klaszterek elég távol vannak egymástól.

Előnye:

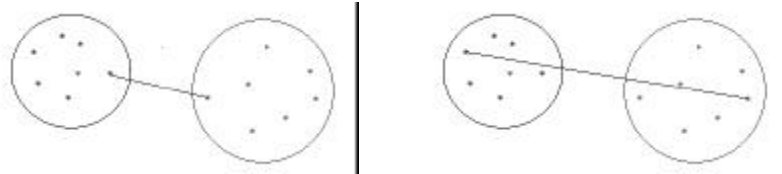
Nem kell előre tudni a klaszterek számát

Változtatható a távolság- és hasonlósági-mérték

Hátránya:

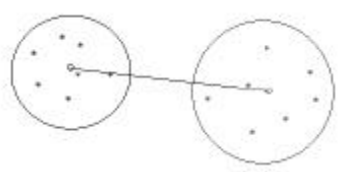
Kis dimenziószám esetén indítható el

KLASZTEREK $d(C_1, C_2)$ TÁVOLSÁGA



A legközelebbi-társ távolság

A legtávolabbi-társ távolság



Klasztercentrumok távolsága

ESETEK $d(\underline{x}, \underline{y})$ TÁVOLSÁGAI

Minkowsky $d(\underline{x}, \underline{y}) = \left(\sum_{i=1}^p |x_i - y_i|^r \right)^{\frac{1}{r}}$

City-blokk $d(\underline{x}, \underline{y}) = \sum_{i=1}^p |x_i - y_i|$

Euklideszi $d(\underline{x}, \underline{y}) = \left(\sum_{i=1}^p |x_i - y_i|^2 \right)^{\frac{1}{2}}$

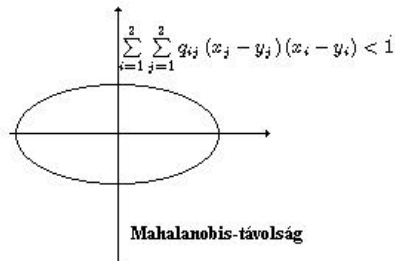
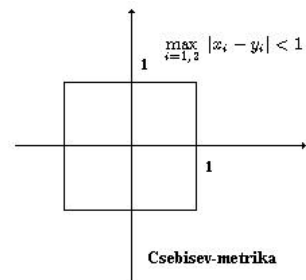
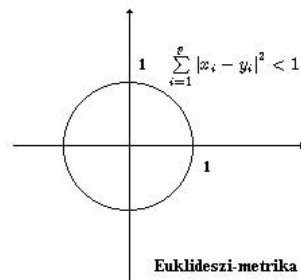
Mahalanobis

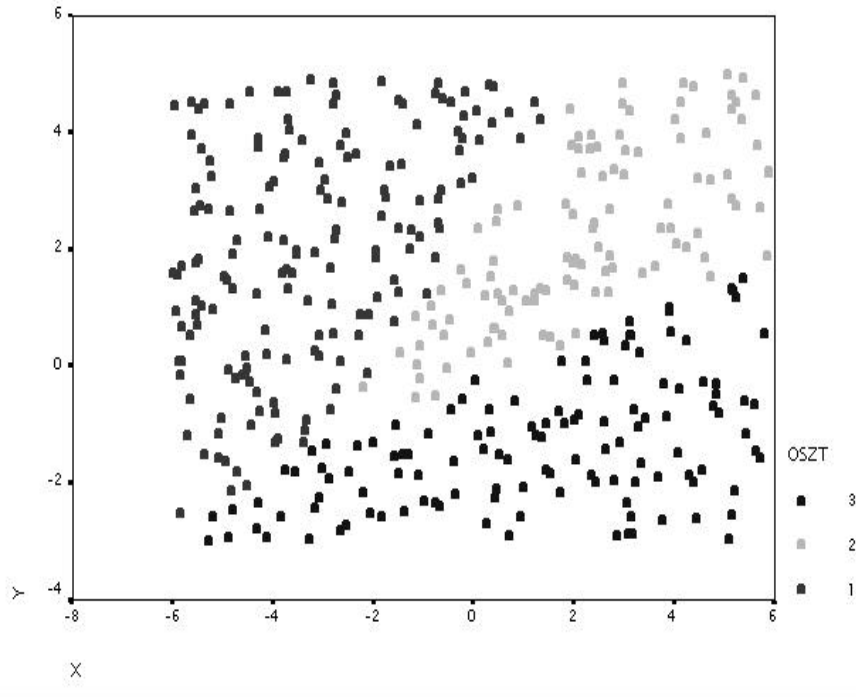
$$d(\underline{x}, \underline{y}) = (\underline{x} - \underline{y})^T \underline{Q} (\underline{x} - \underline{y}) = \sum_{i=1}^p \sum_{j=1}^p q_{ij} (x_j - y_j) (x_i - y_i)$$

Canberra $d(\underline{x}, \underline{y}) = \sum_{i=1}^p \frac{|x_i - y_i|}{|x_i| + |y_i|}$

Csebisev $d(\underline{x}, \underline{y}) = \max_{i=1, \dots, p} |x_i - y_i|$

ESETEK $d(\underline{x}, \underline{y})$ TÁVOLSÁGAI





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 struktúráldása. Jellegzetes tömörülések, irányok fedezhetők fel az elkészült scatter-grafikonon.

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

$\underline{x}^{(1)}$ 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.

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