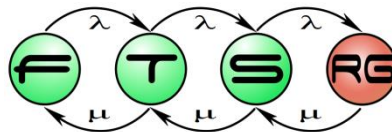


# Outlierdetektálás

Salánki Ágnes

salanki@mit.bme.hu

**Budapest University of Technology and Economics**  
**Fault Tolerant Systems Research Group**



# Definíció

- Kevés van belőlük
- „Gyanús”, hogy **más** a generáló folyamat/forrás

# Definíció

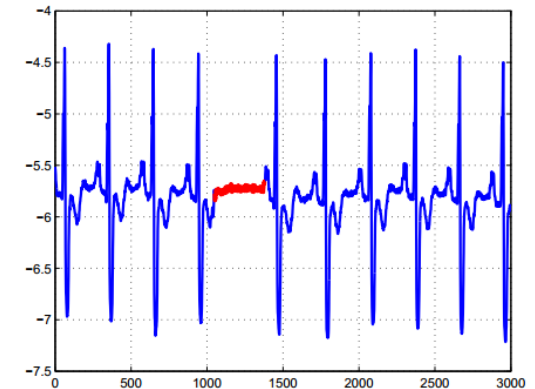
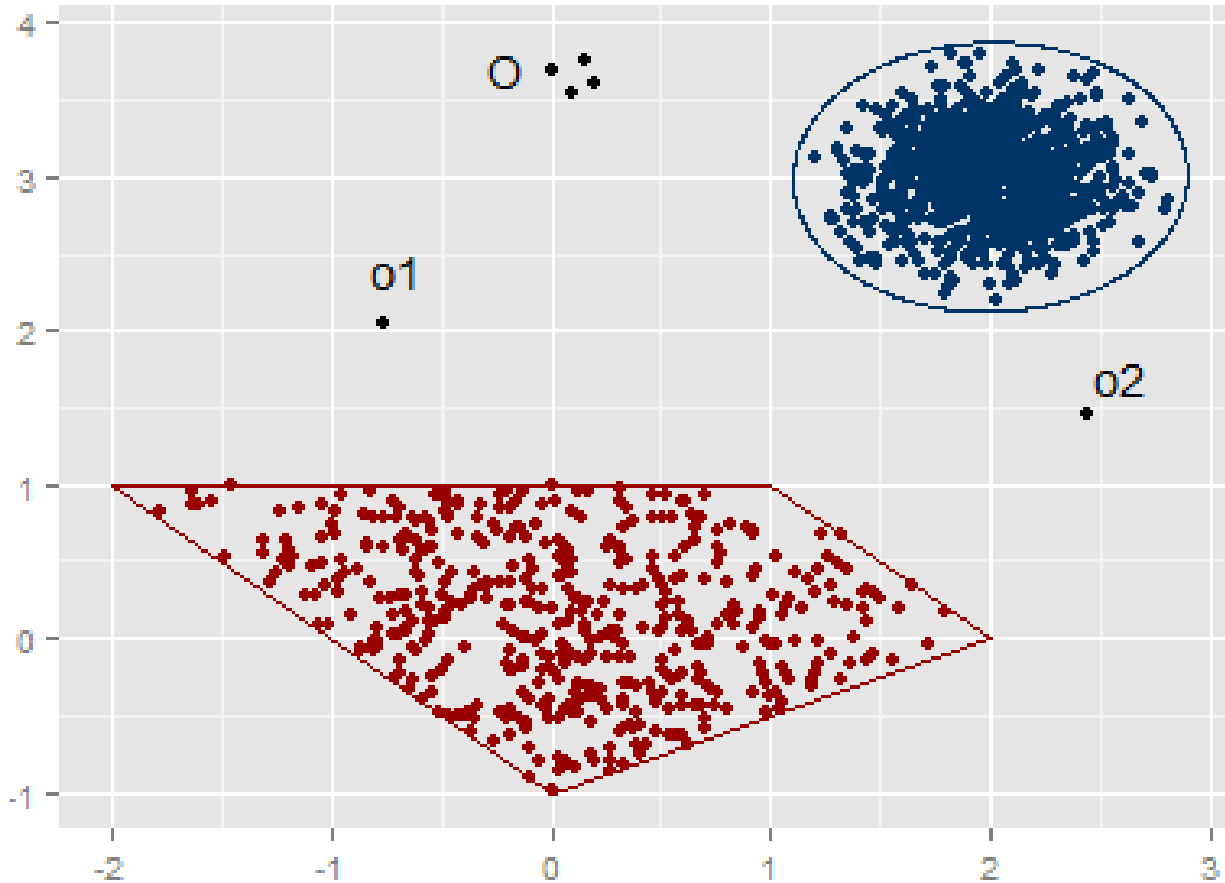
- Kevés van belőlük
- „Gyanús”, hogy **más** a generáló folyamat/forrás
  - Все счастливые семьи похожи друг на друга, каждая несчастливая семья несчастлива **по-своему**.
  - Happy families are all alike; every unhappy family is unhappy **in its own way**.
  - A boldog családok mind hasonlóak egymáshoz, minden boldogtalan család **a maga módján** az.

(Tolsztoj: Anna Karenina)

# Pont- és kollektív anomália

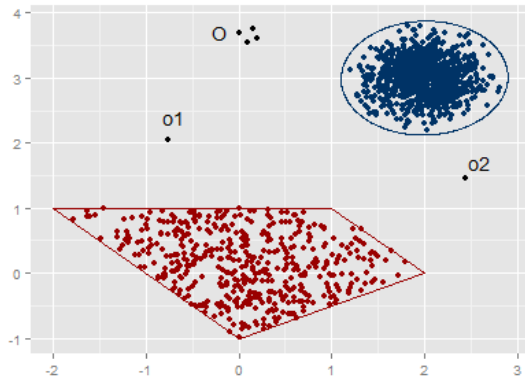
## ■ Pontanomália

## ■ Kollektív anomália

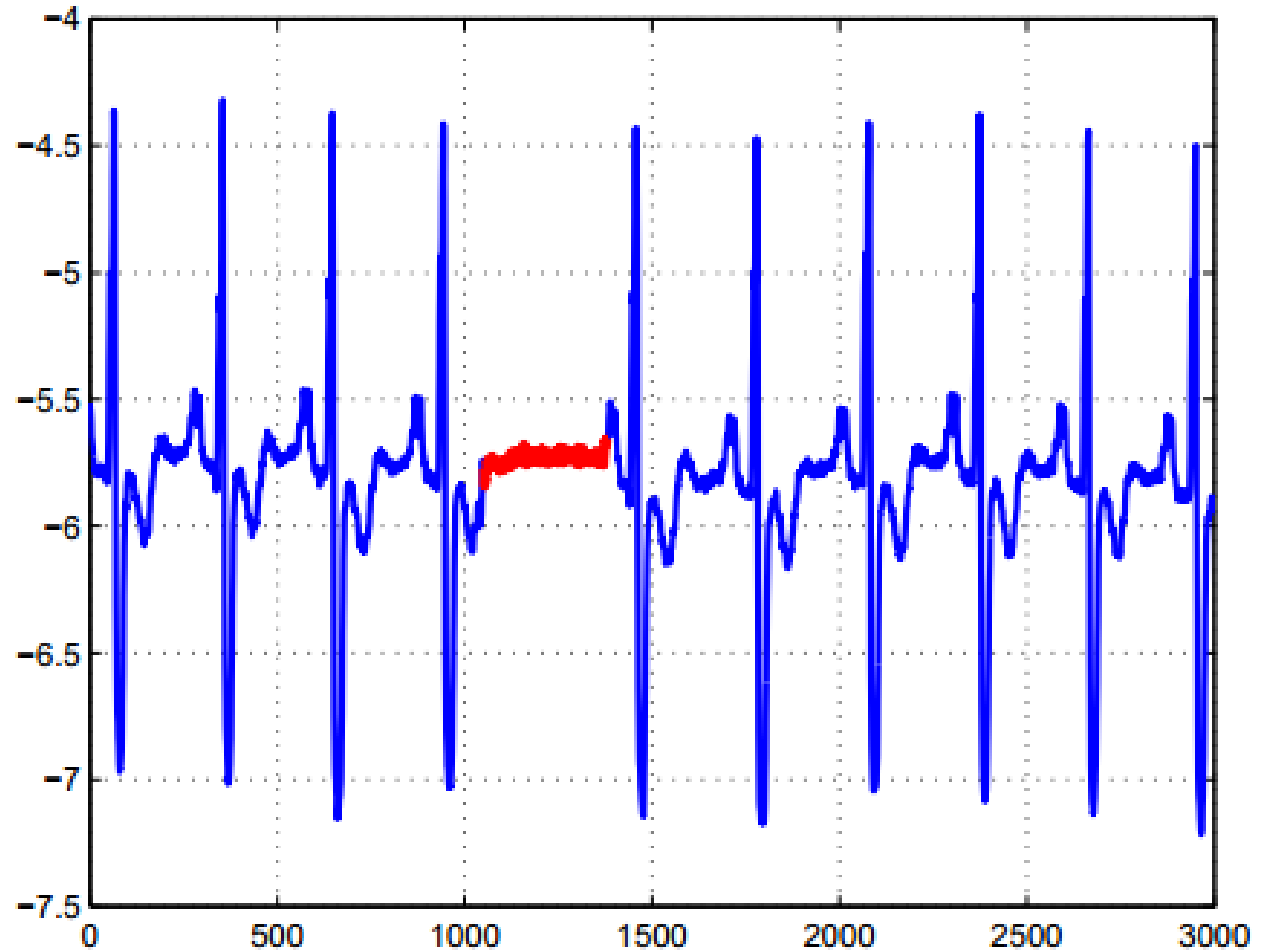


# Pont- és kollektív anomália

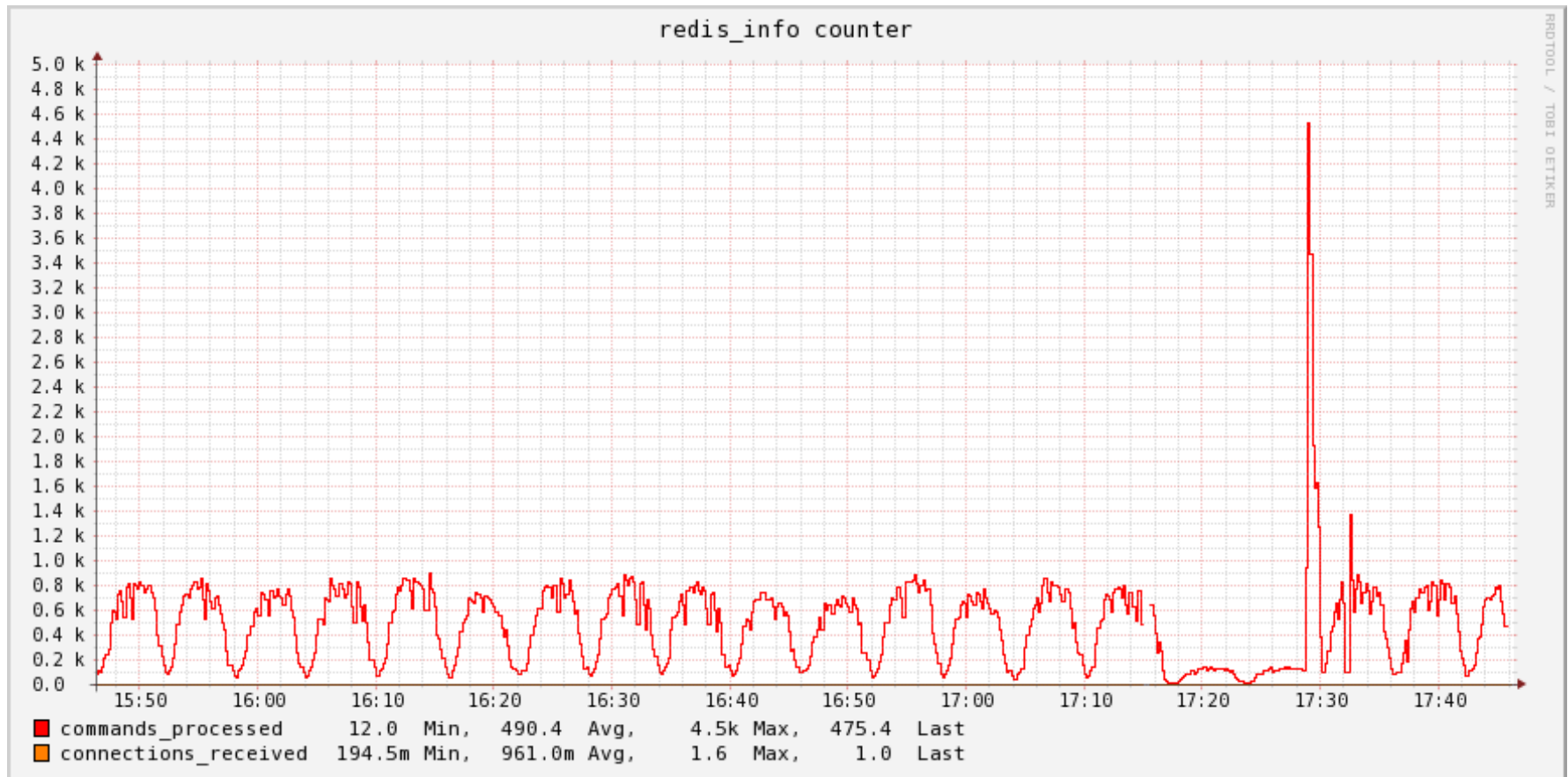
## ■ Pontanomália



## ■ Kollektív anomália



# Pont- és kollektív anomália

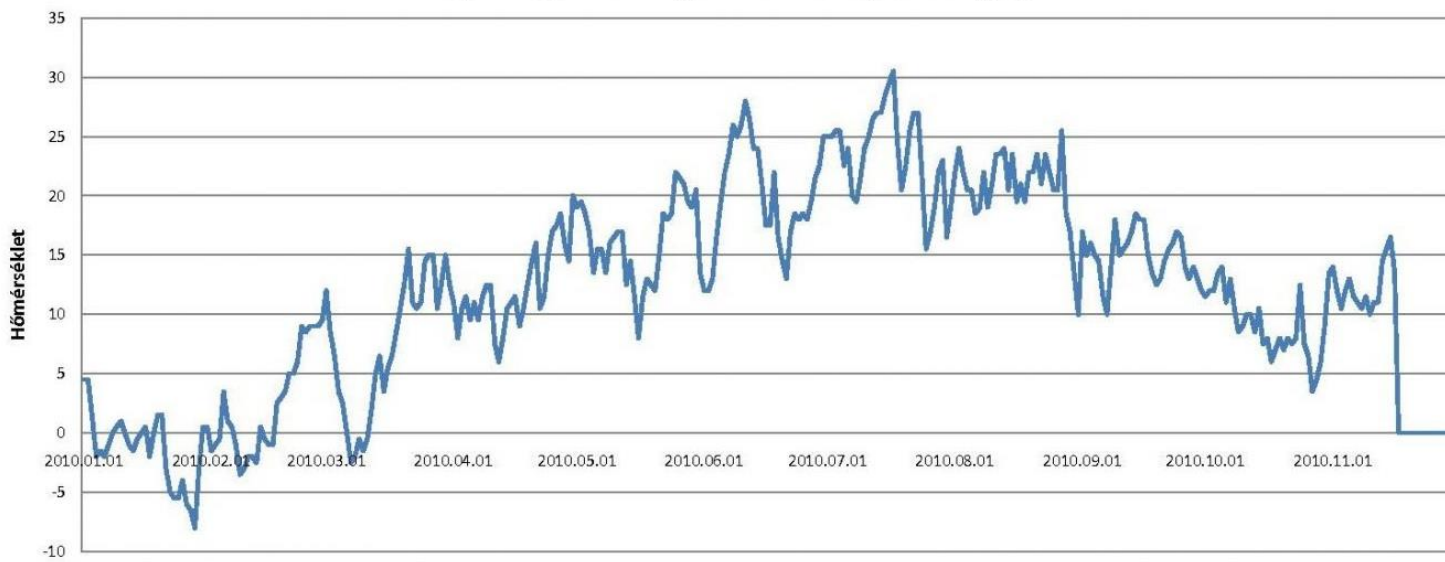


# Viselkedési és kontextusanomália

## ■ Viselkedési

## ■ Kontextus

Napi átlag hőmérséglet 2010-ben (Zala megye)

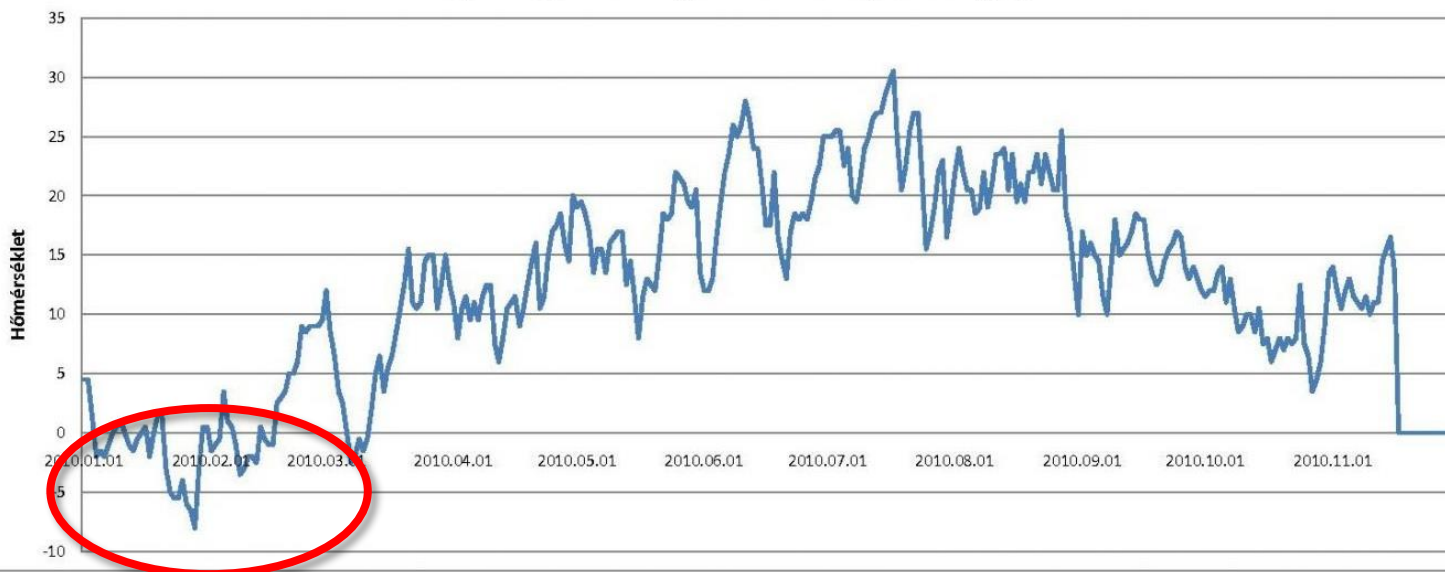


# Viselkedési és kontextusanomália

## ■ Viselkedési

## ■ Kontextus

Napi átlag hőmérséglet 2010-ben (Zala megye)



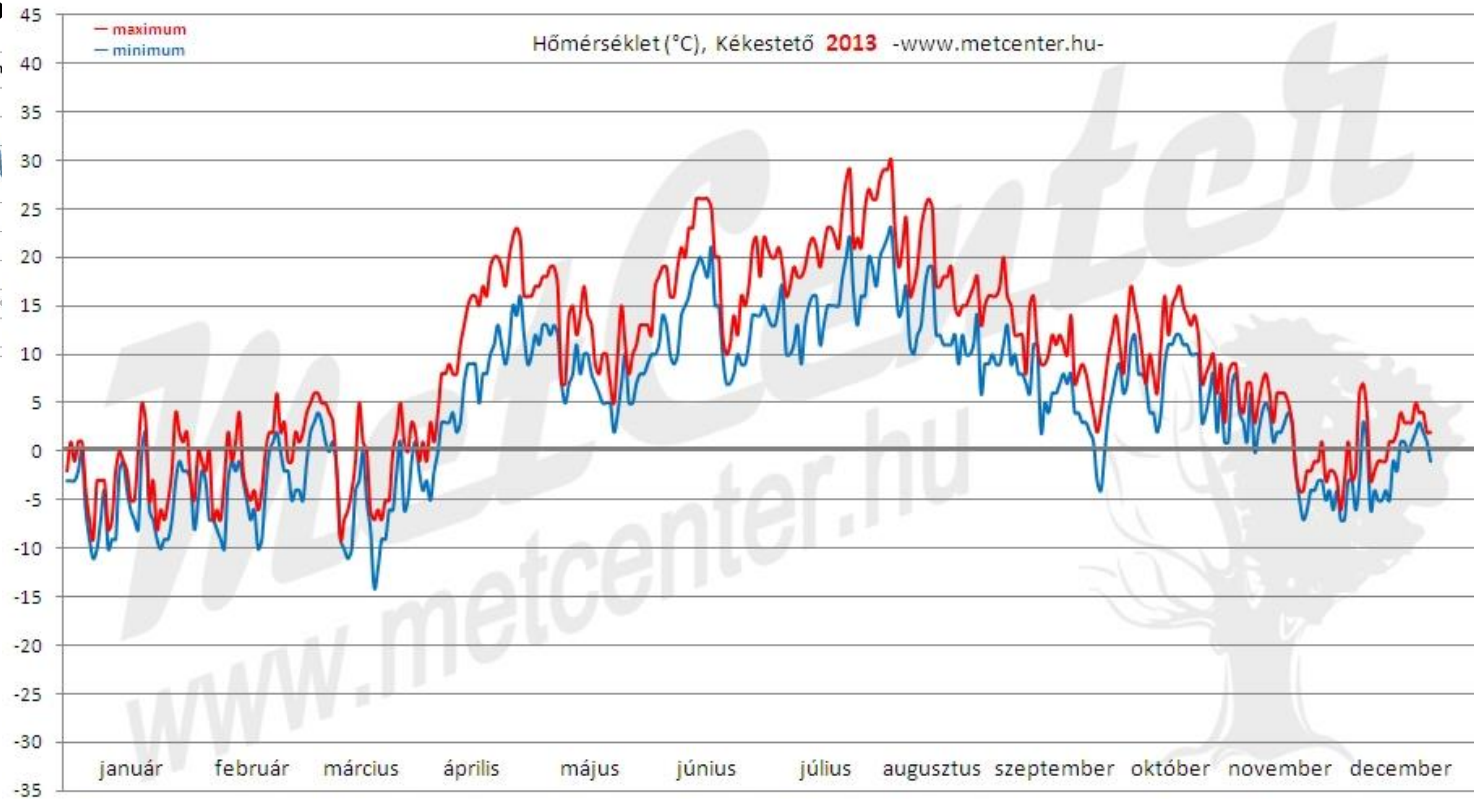


# Viselkedési és kontextusanomália

## ■ Viselkedés:



## ■ Kontextus

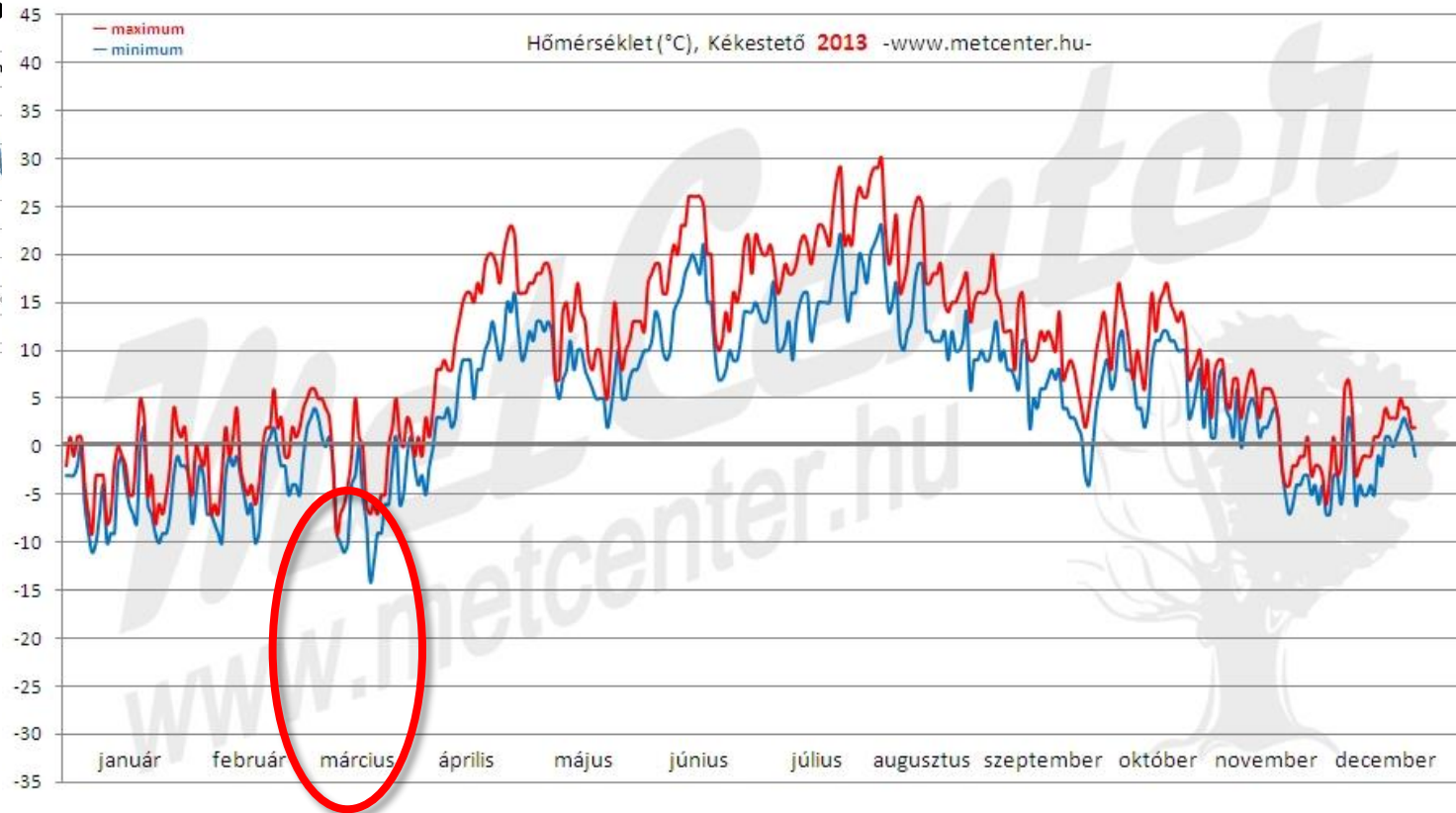


# Viselkedési és kontextusanomália

## ■ Viselkedés:



## ■ Kontextus

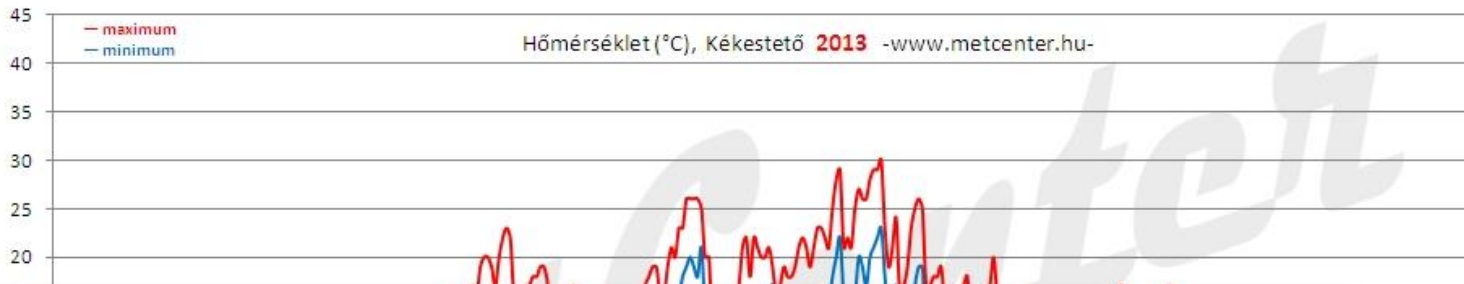


# Viselkedési és kontextusanomália

## ■ Viselkedés:



## ■ Kontextus

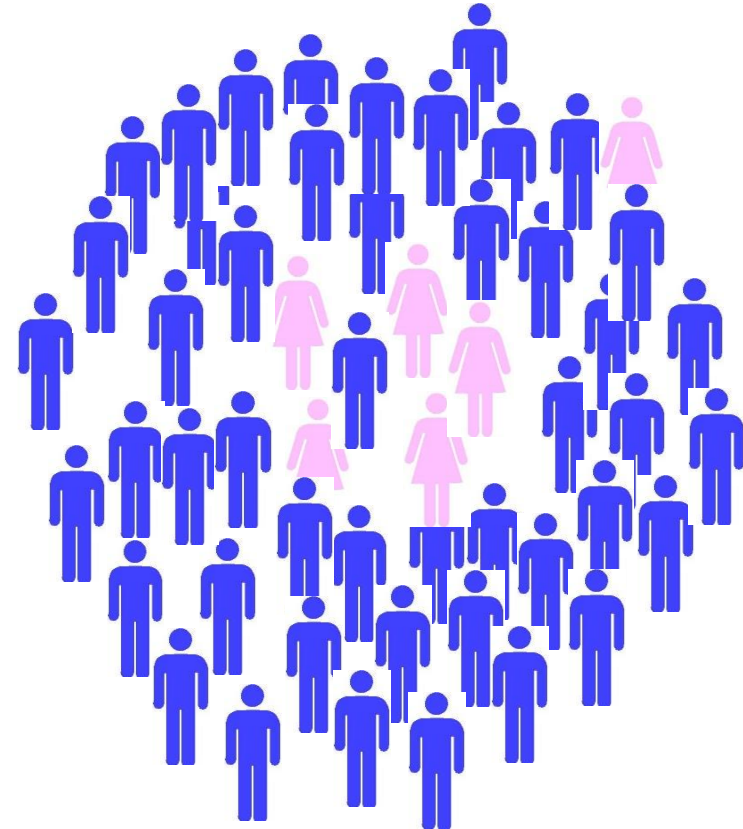
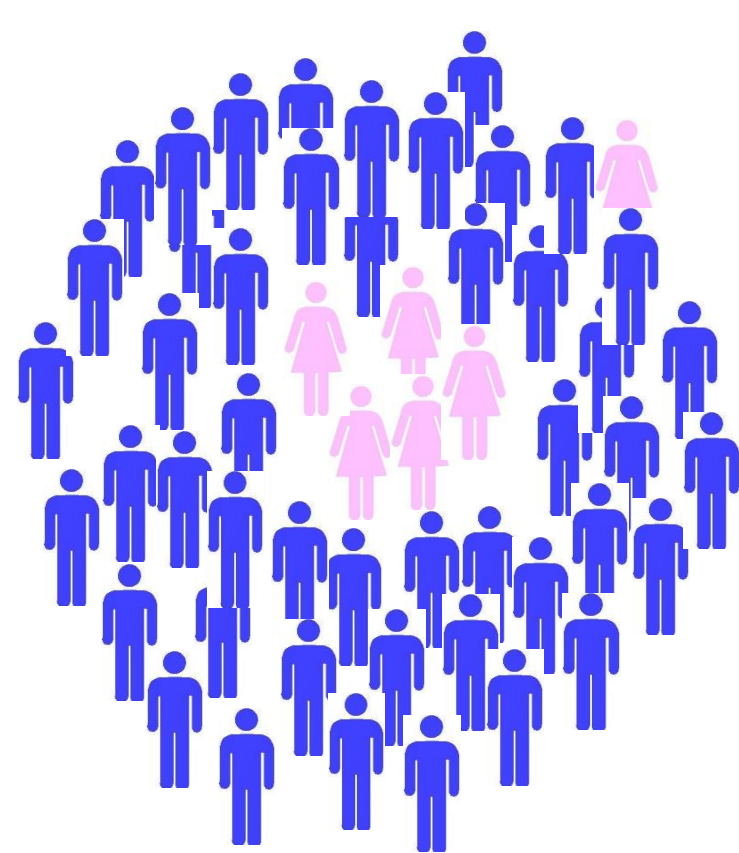


**Itt: viselkedési és pontanomáliák**

# Megközelítések

- Globális outlierek

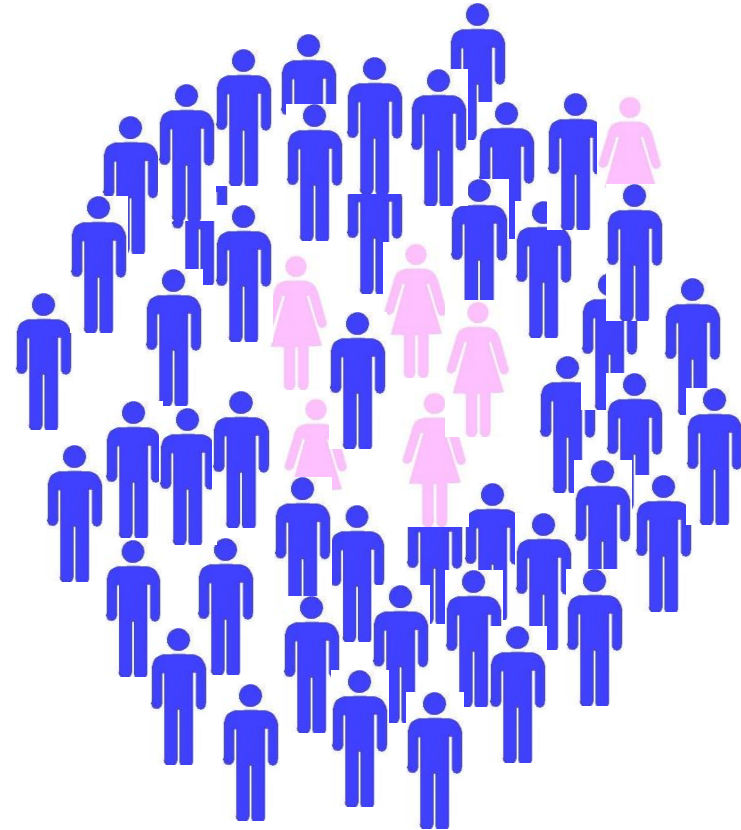
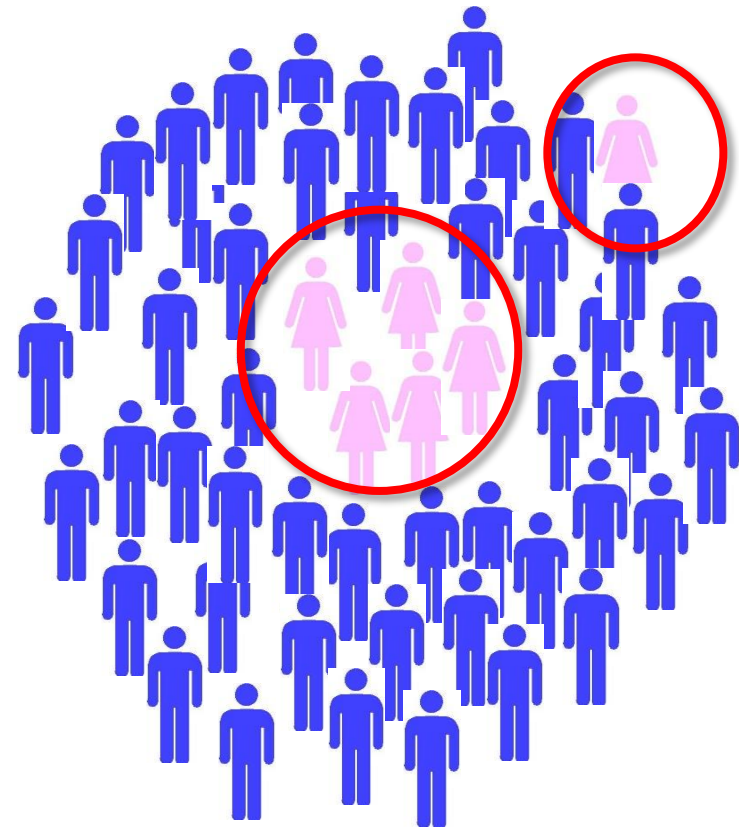
- Lokális outlierek



# Megközelítések

- Globális outlierok

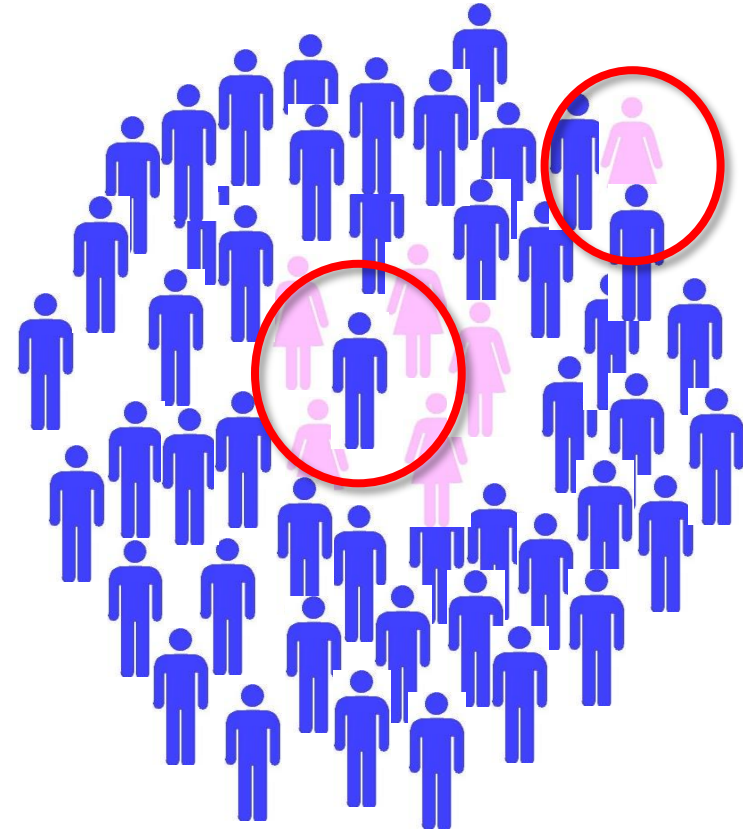
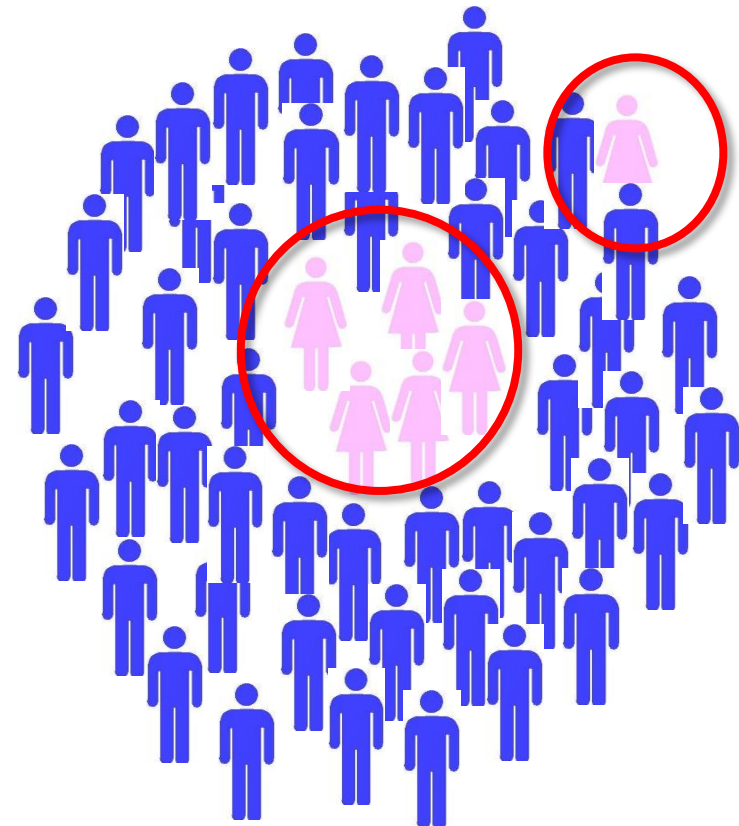
- Lokális outlierok



# Megközelítések

- Globális outlierek

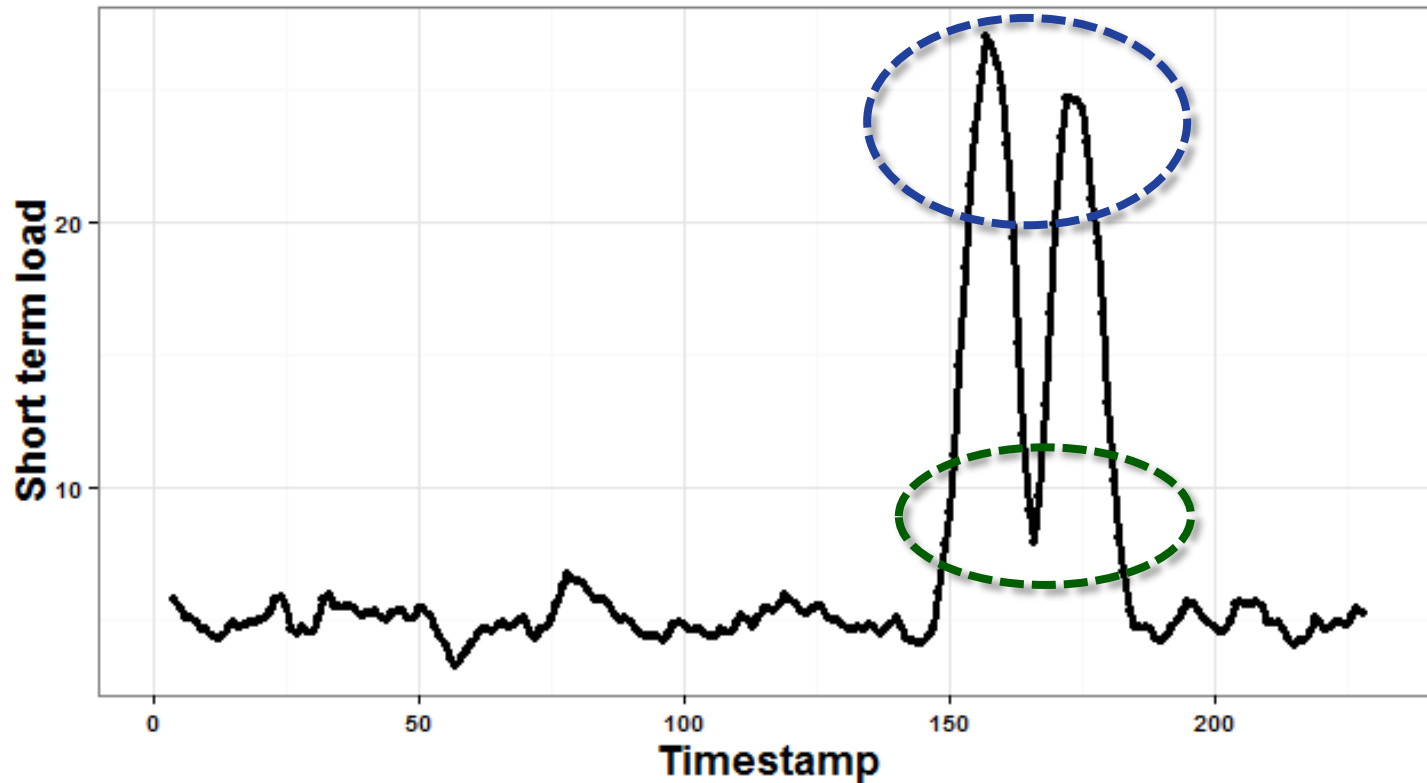
- Lokális outlierek



# Megközelítések

- Globális outlierok

- Lokális outlierok



# Esettanulmány: PISA 2012

- PISA 2012 results

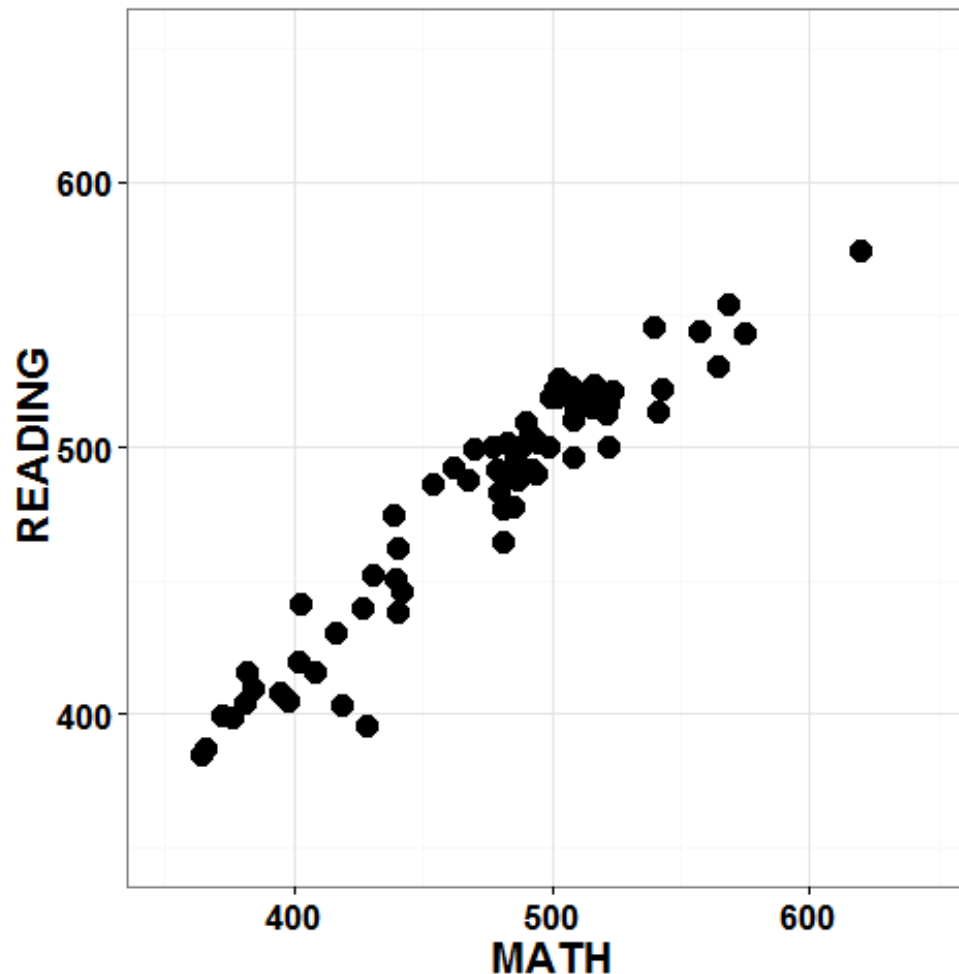
- Most: matematika és értő olvasás eredmények



# Esettanulmány: PISA 2012

## ■ PISA 2012 results

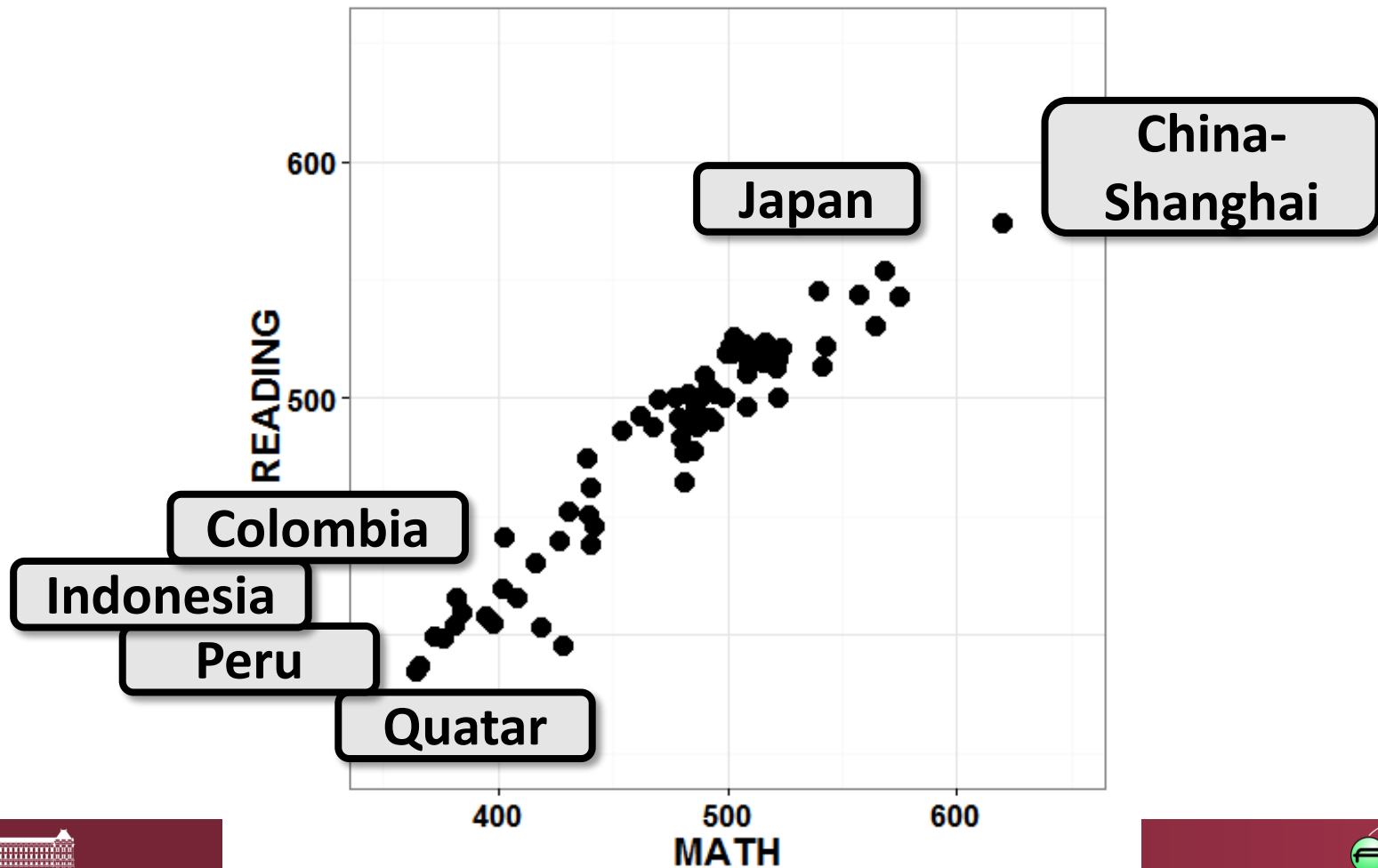
- Most: matematika és értő olvasás eredmények



# Esettanulmány: PISA 2012

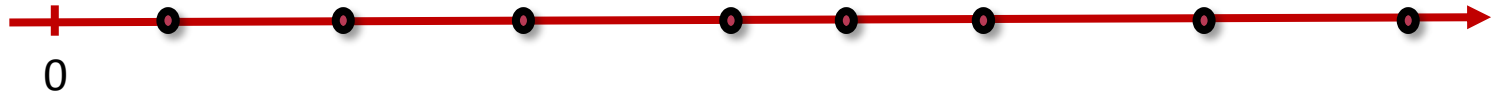
## ■ PISA 2012 results

- Most: matematika és értő olvasás eredmények



# Befoglaló burok

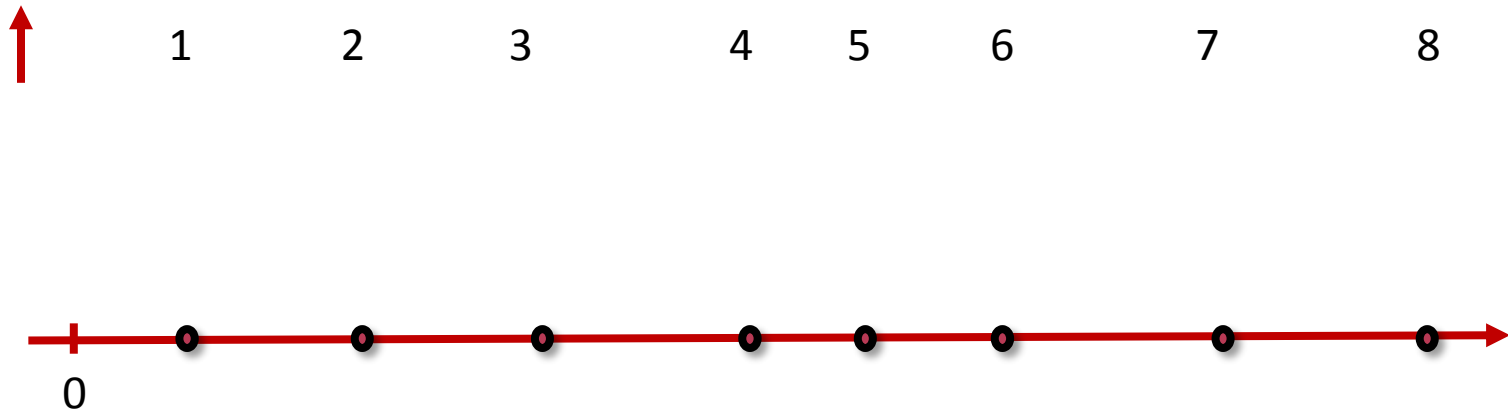
- Féltér-mélység: Tukey, 1974



$$hds(z): \min\{|x_i: x_i \leq z|, |x_j: x_j \geq z|\}$$

# Befoglaló burok

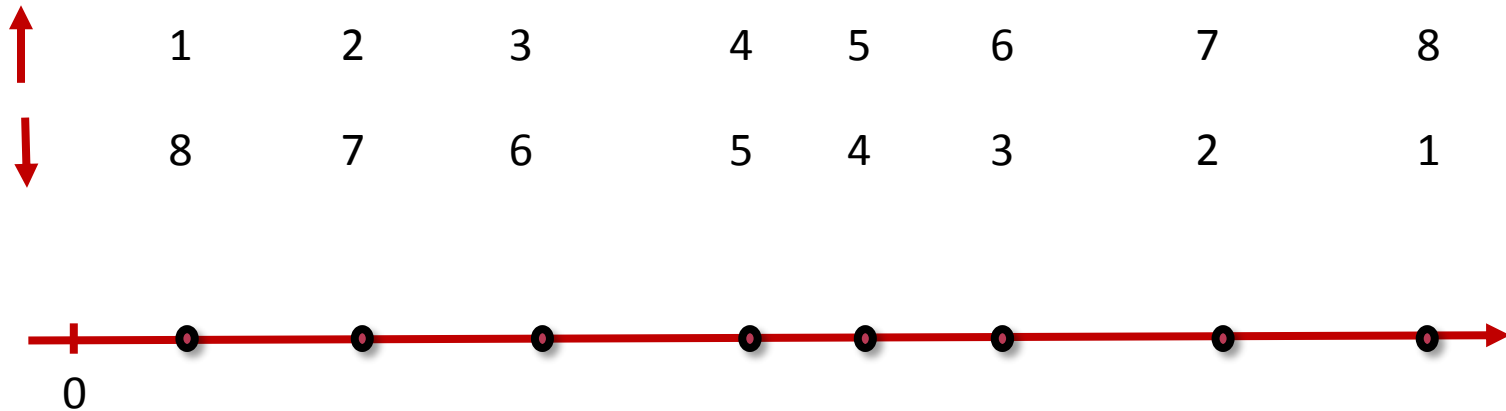
- Féltér-mélység: Tukey, 1974



$$hds(z): \min\{|\{x_i: x_i \leq z\}|, |\{x_j: x_j \geq z\}|\}$$

# Befoglaló burok

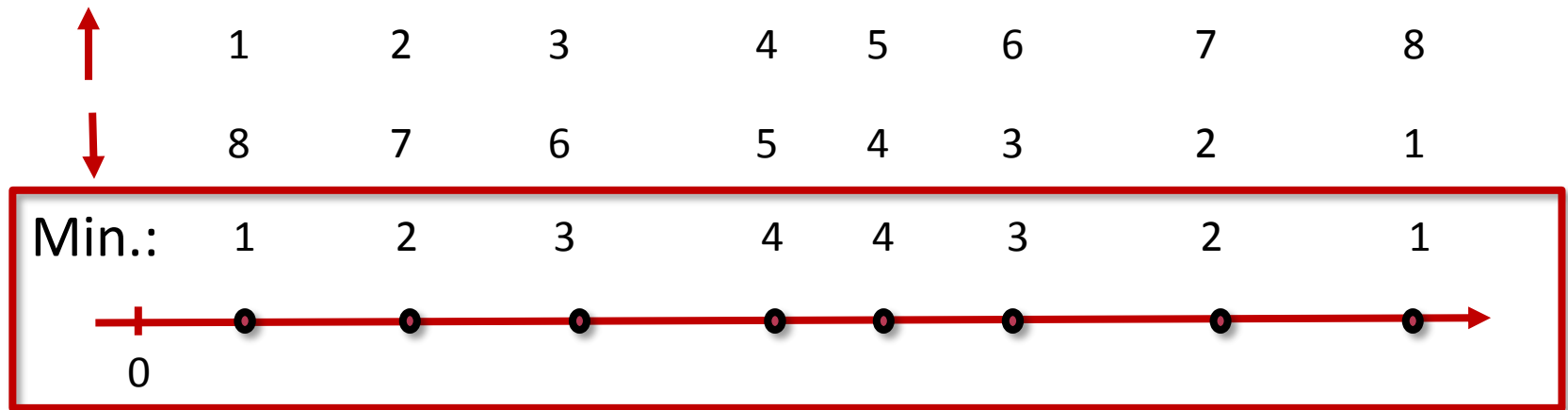
- Féltér-mélység: Tukey, 1974



$$hds(z): \min\{|\{x_i: x_i \leq z\}|, |\{x_j: x_j \geq z\}|\}$$

# Befoglaló burok

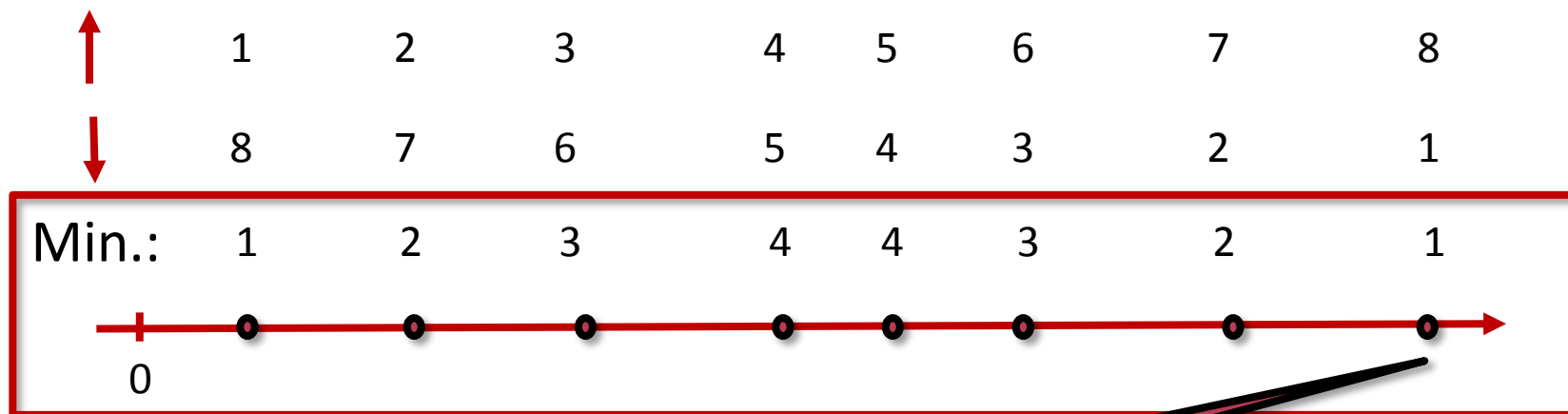
- Féltér-mélység: Tukey, 1974



$$hds(z): \min\{|x_i: x_i \leq z|, |x_j: x_j \geq z|\}$$

# Befoglaló burok

- Féltér-mélység: Tukey, 1974



Extrém  
pontok

$$): \min\{|x_i: x_i \leq z|, |x_j: x_j \geq z|\}$$

# Befoglaló burok

- Féltér-mélység: Tukey, 1974



Extrém  
pontok

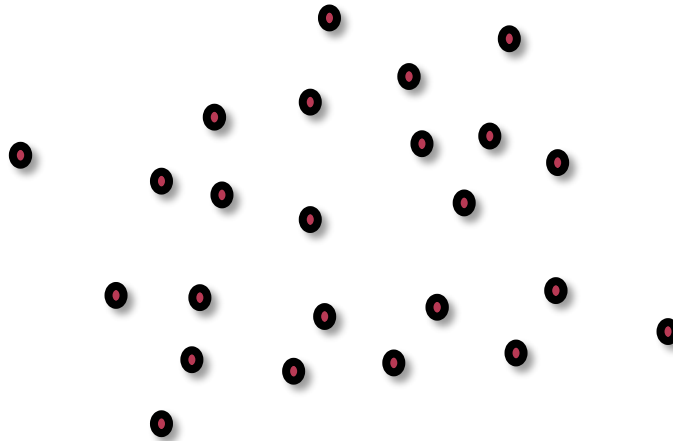
$$): \min\{|x_i: x_i \leq z|$$

Medián: majd a végén...



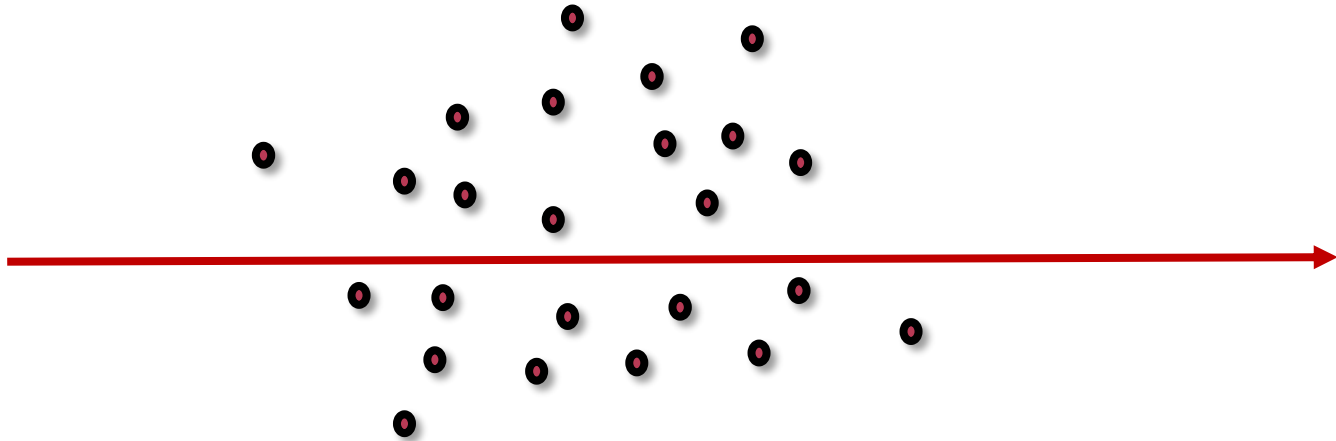
# Befoglaló burok

- Féltér-mélység: Tukey, 1974



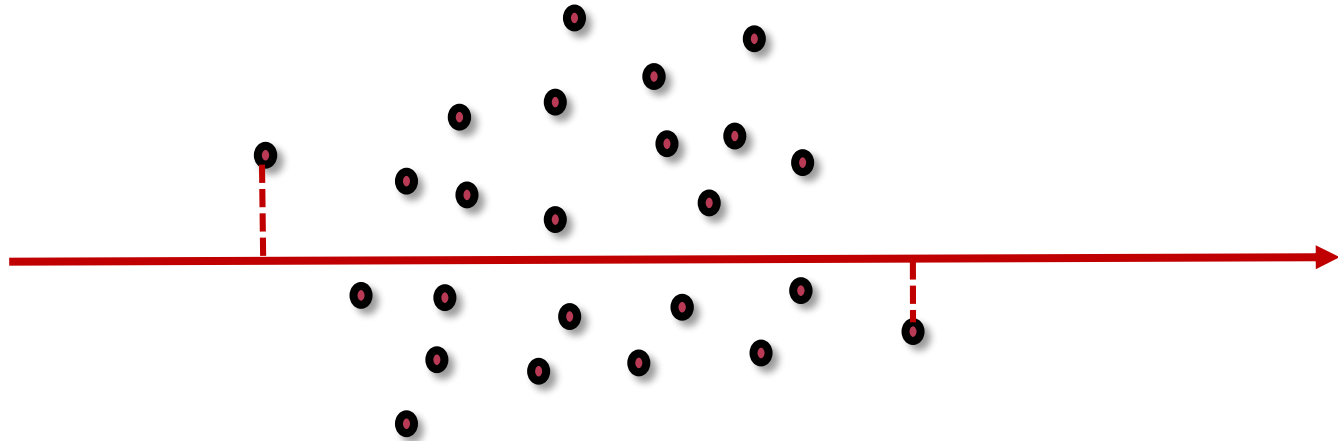
# Befoglaló burok

- Féltér-mélység: Tukey, 1974



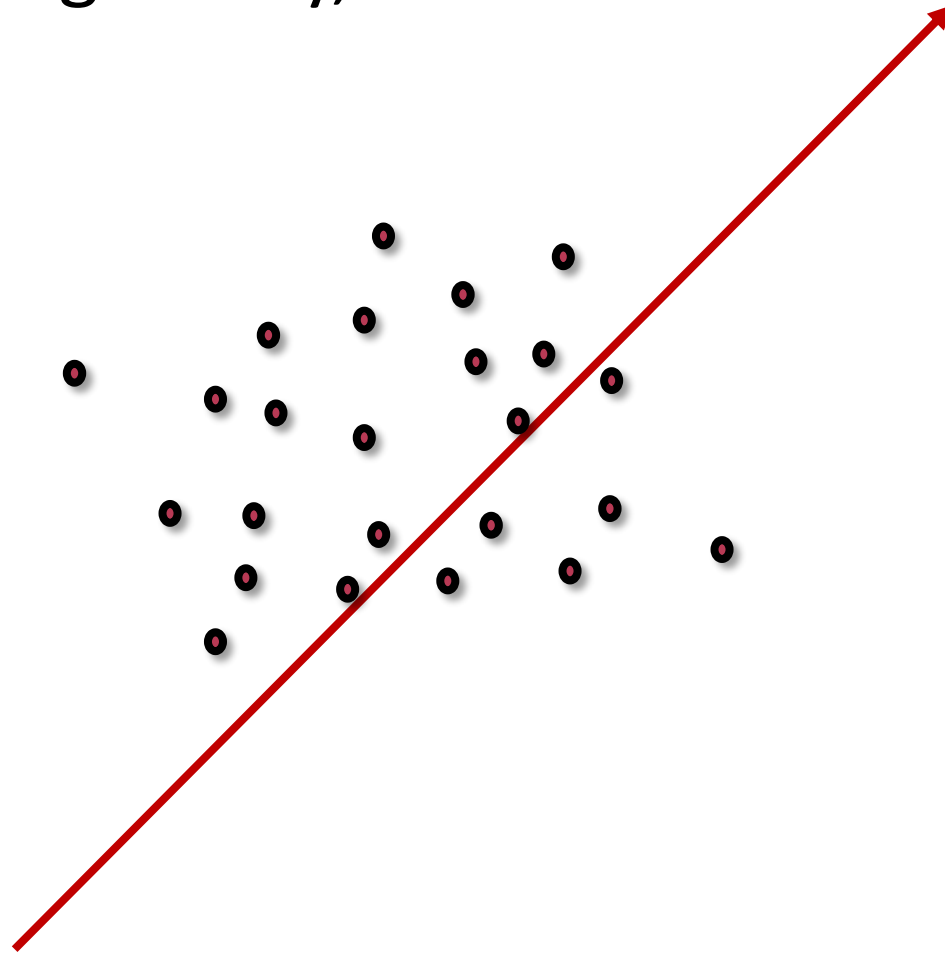
# Befoglaló burok

- Féltér-mélység: Tukey, 1974



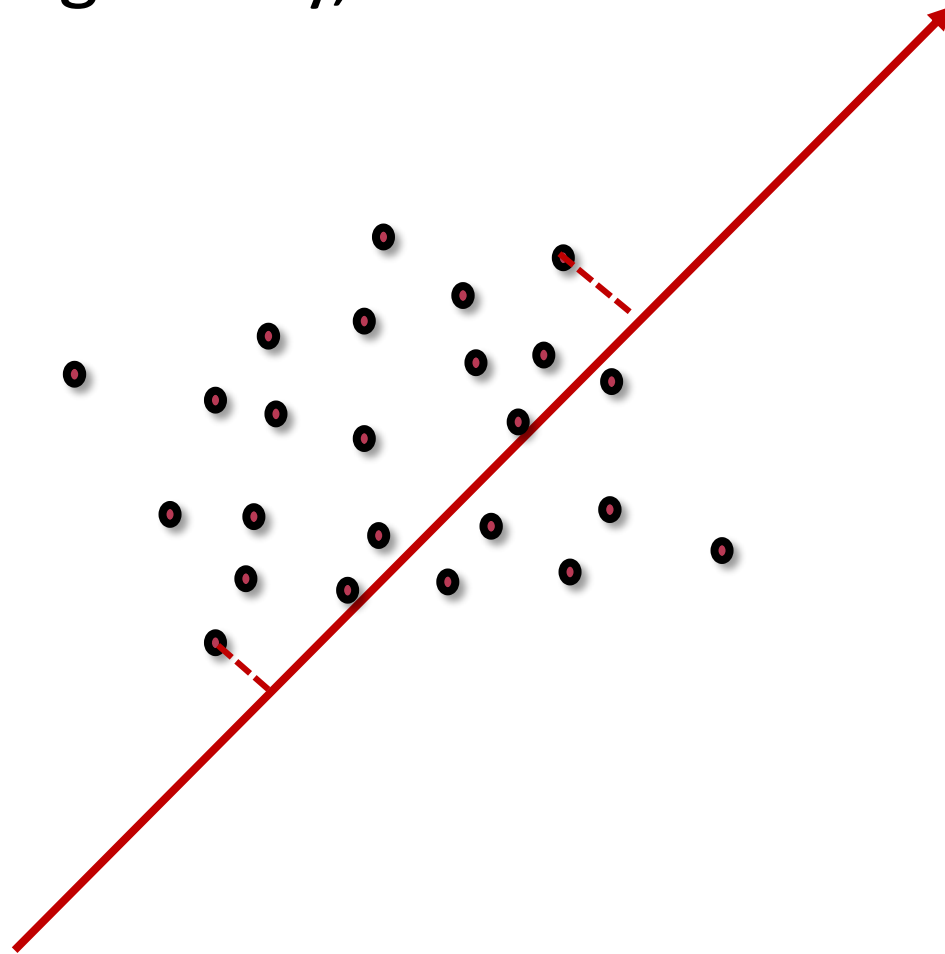
# Befoglaló burok

- Féltér-mélység: Tukey, 1974



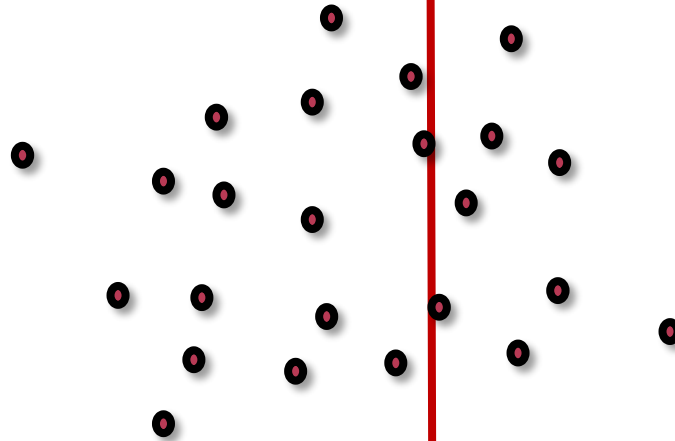
# Befoglaló burok

- Féltér-mélység: Tukey, 1974



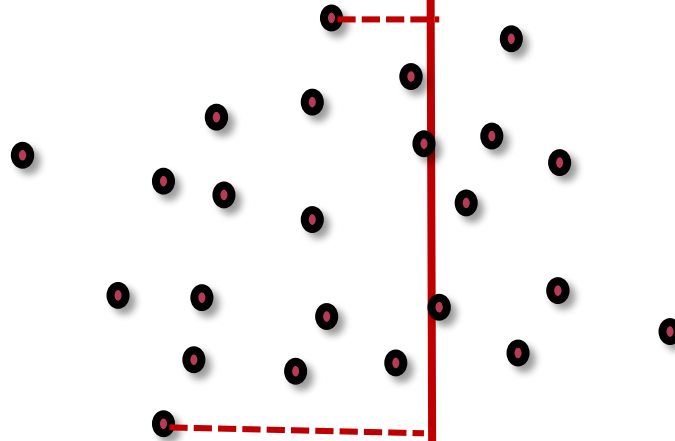
# Befoglaló burok

- Féltér-mélység: Tukey, 1974



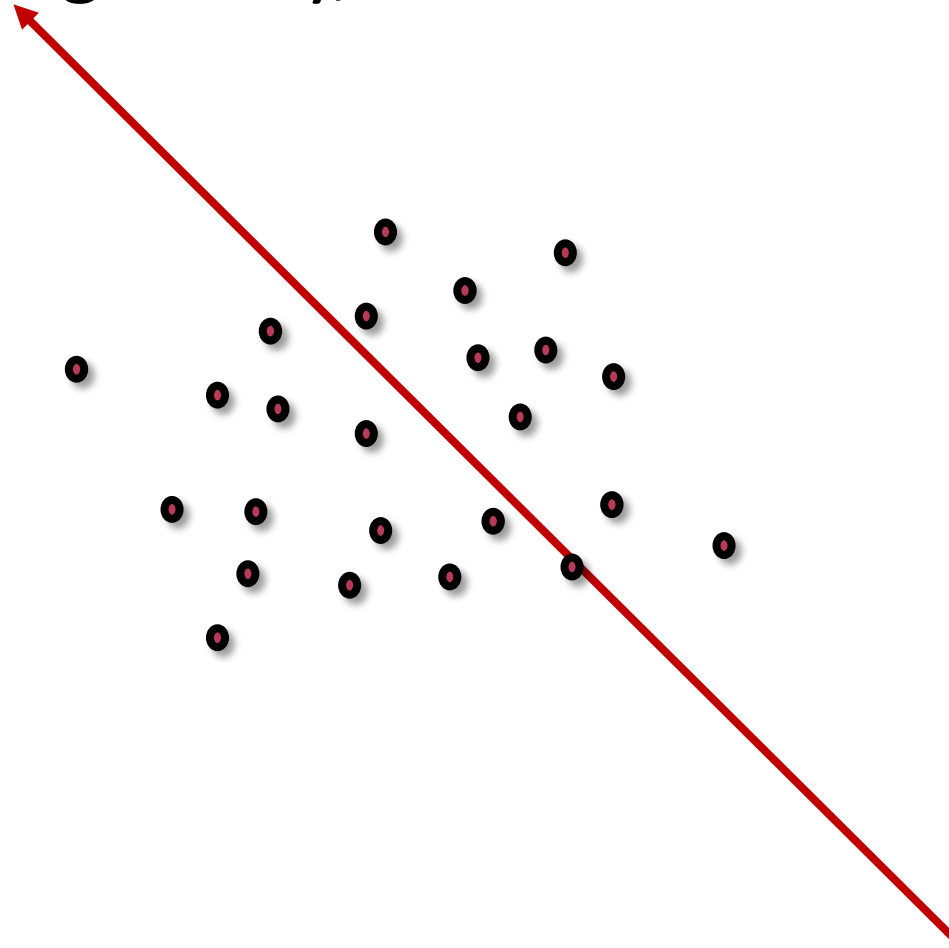
# Befoglaló burok

- Féltér-mélység: Tukey, 1974



# Befoglaló burok

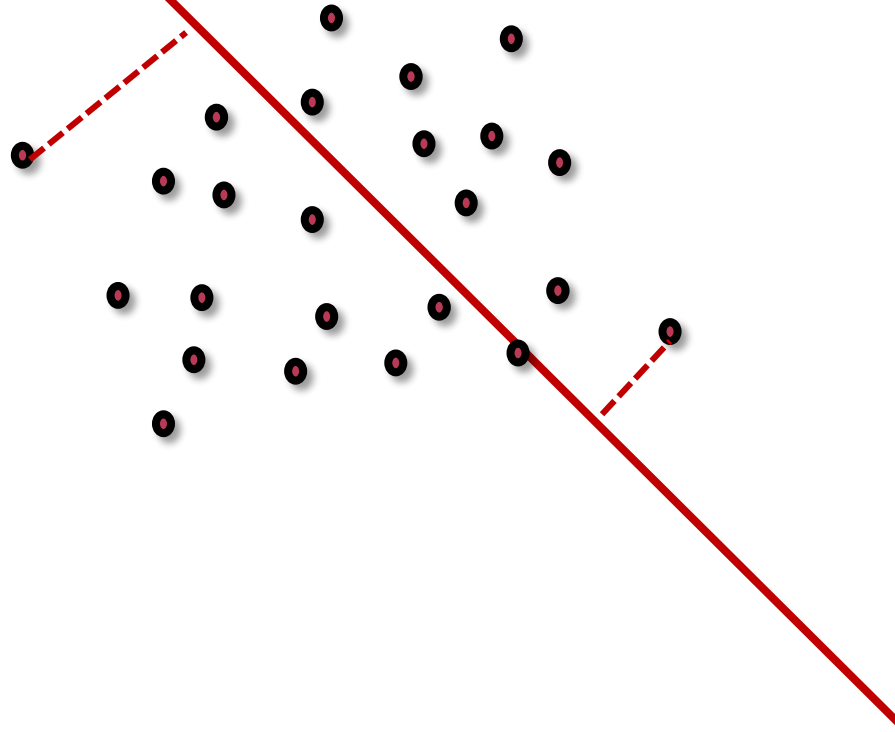
- Féltér-mélység: Tukey, 1974





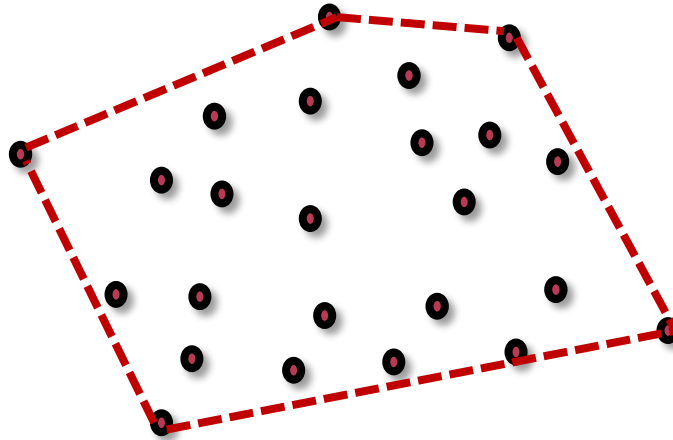
# Befoglaló burok

- Féltér-mélység: Tukey, 1974

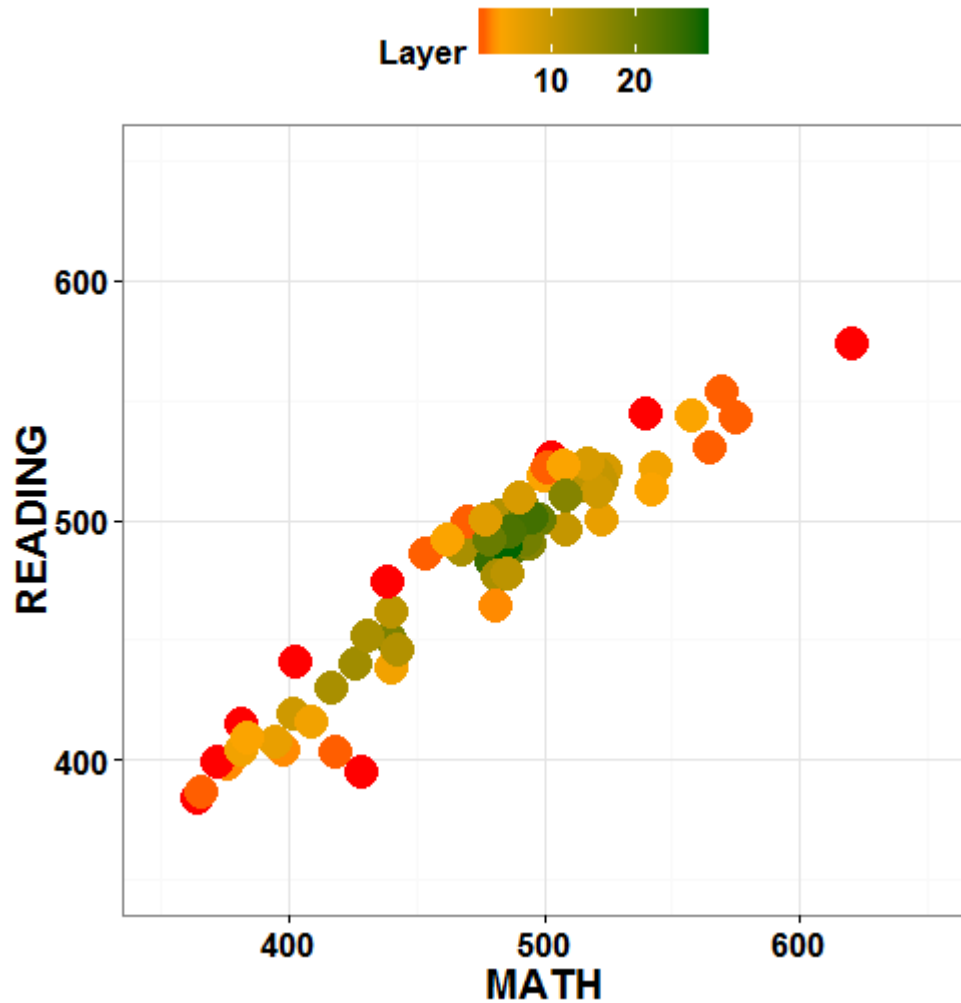


# Befoglaló burok

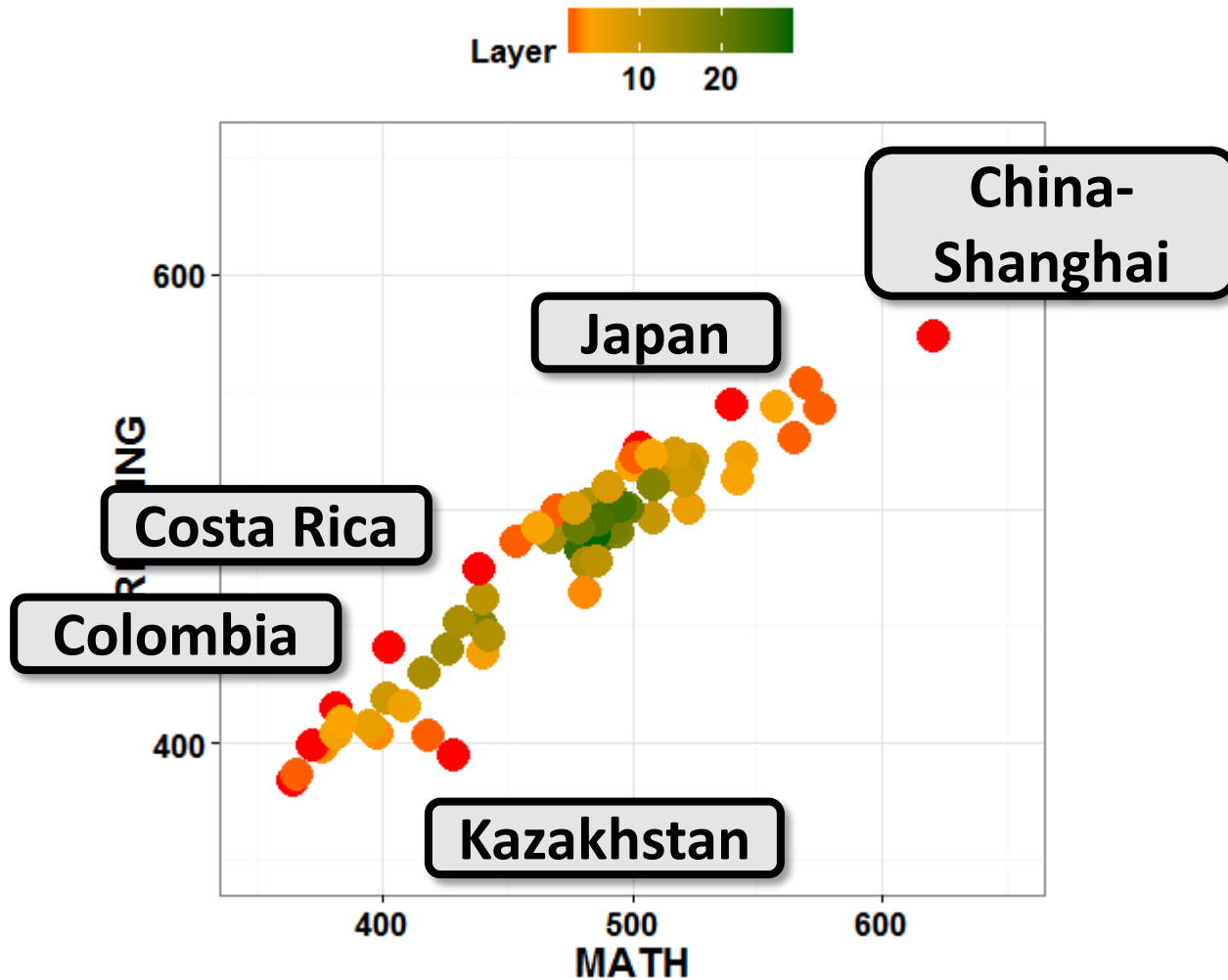
- Féltér-mélység: Tukey, 1974



# Isodepth (Depth-Based)

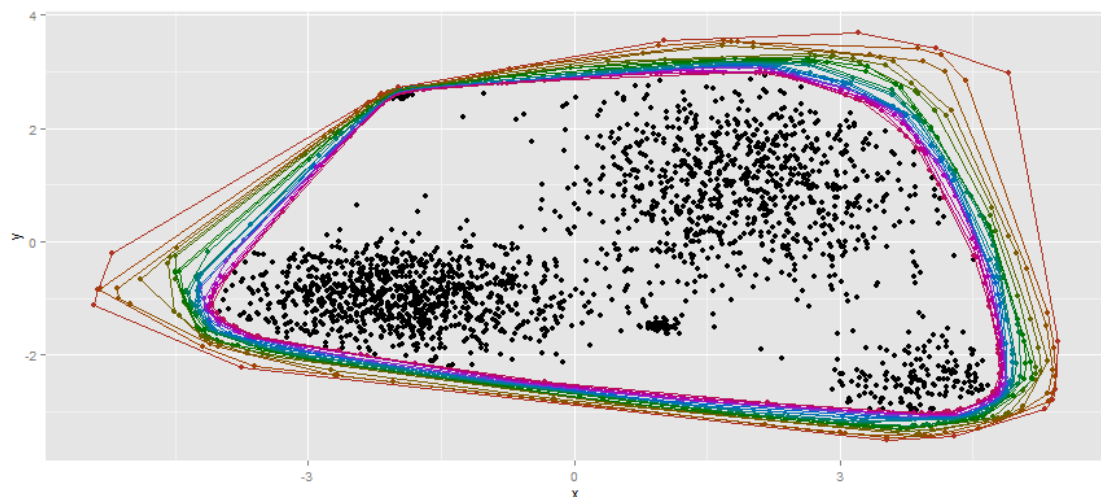
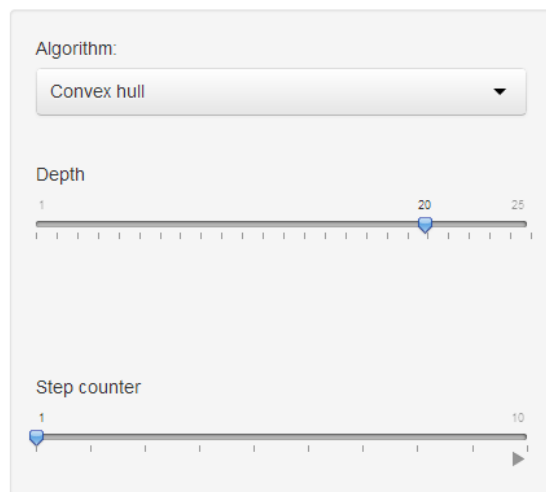


# Isodepth (Depth-Based)



- Csomag: *depth*
- Hasznos függvények: *depth*, *isodepth*
- Paraméterek: *u* pont, *dpth* mélység

### Rare event detection

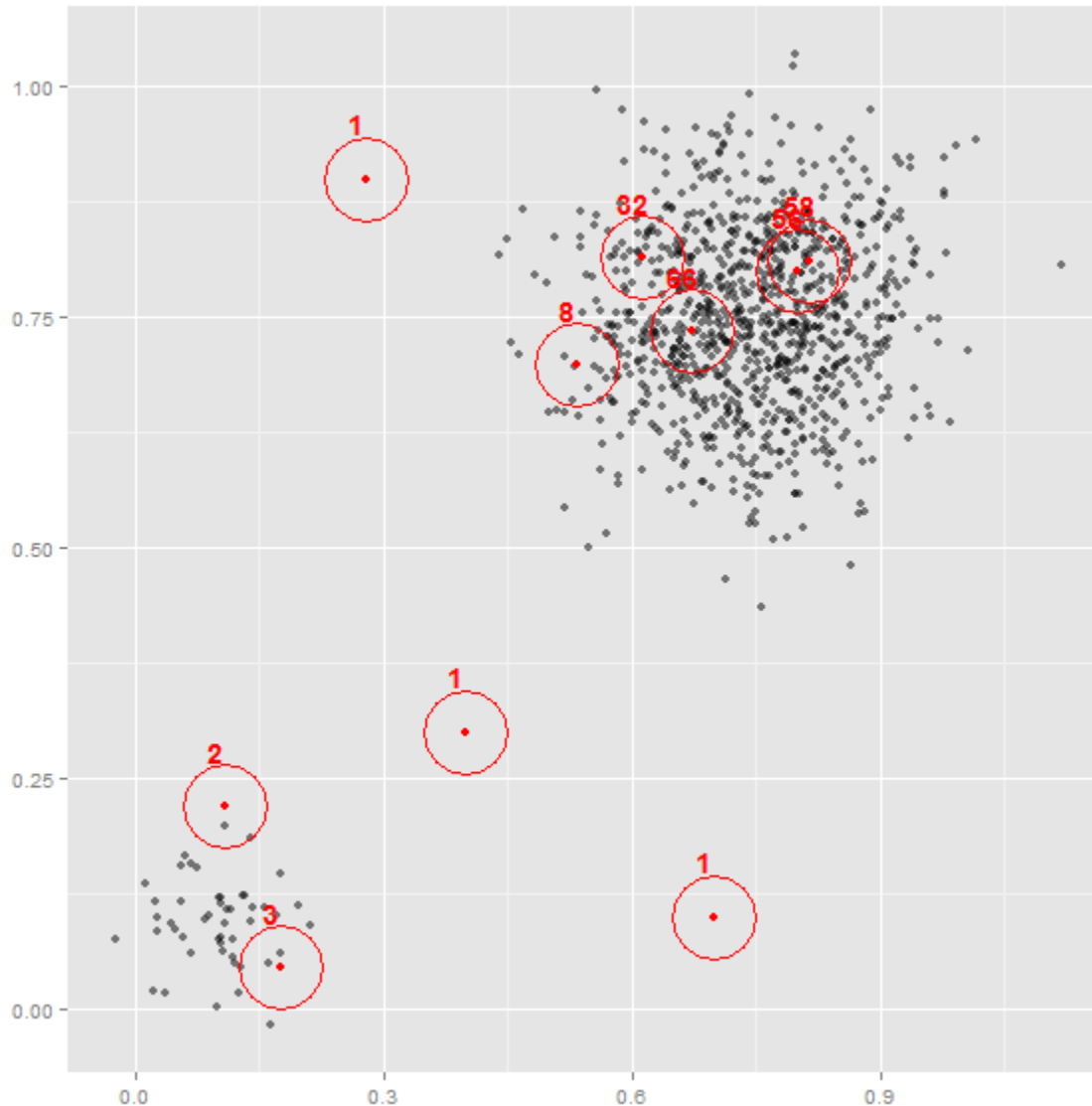


```
[1] "Depth:"  
[1] 20  
[1] "Ratio of chosen elements: "  
[1] 0.4095455
```

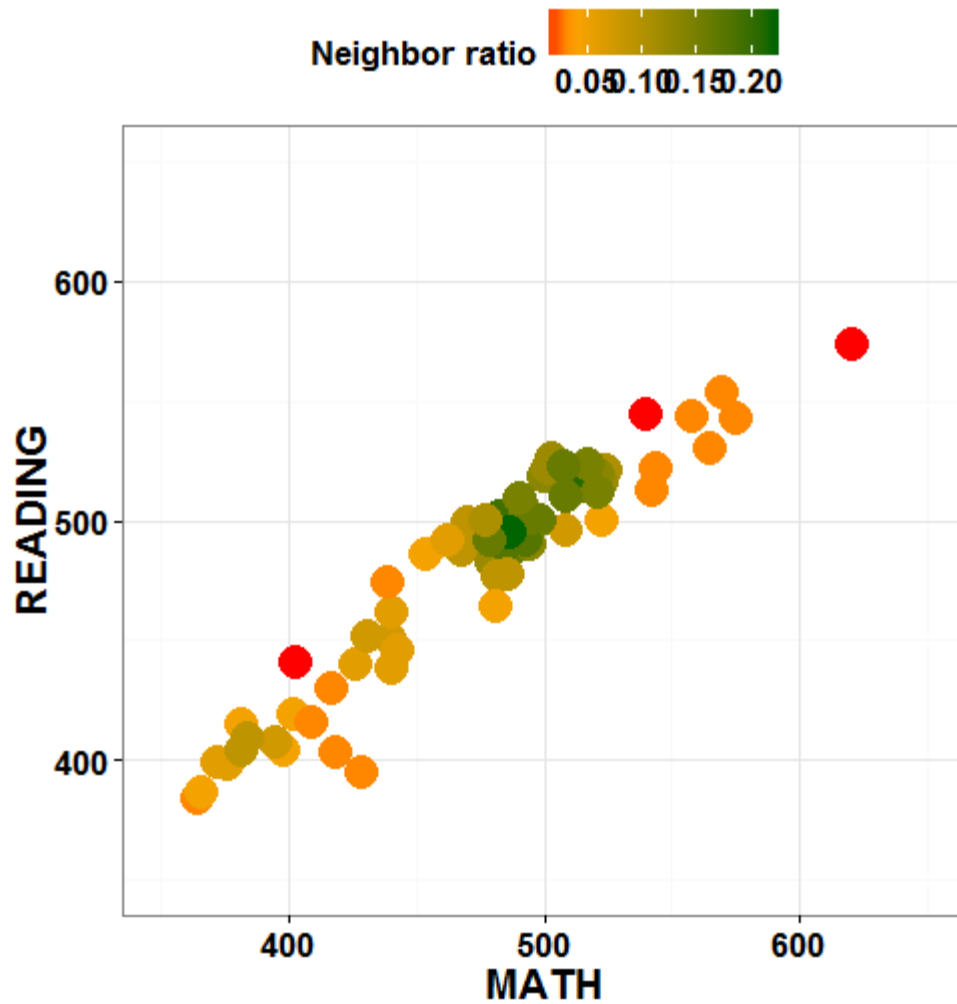
- Distance Based
- Outlier: szomszédok száma alacsony
- Paraméterek
  - $r$  sugarú hipergömb
  - Szomszédok elvárt  $\pi$  aránya

# DB

- Distanc
- Outlier
- Paramé
- $r$  sug
- Szom

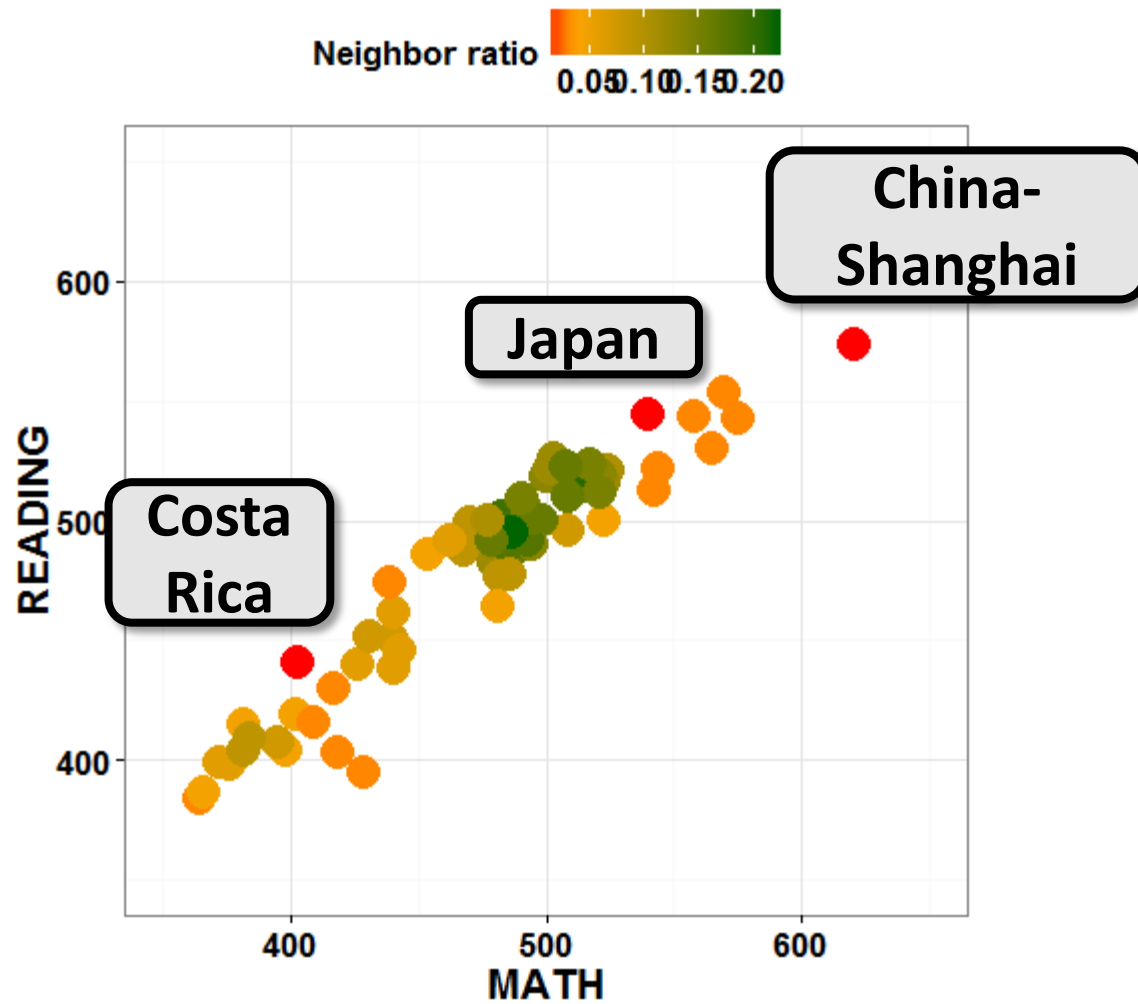


# db – Distance-Based



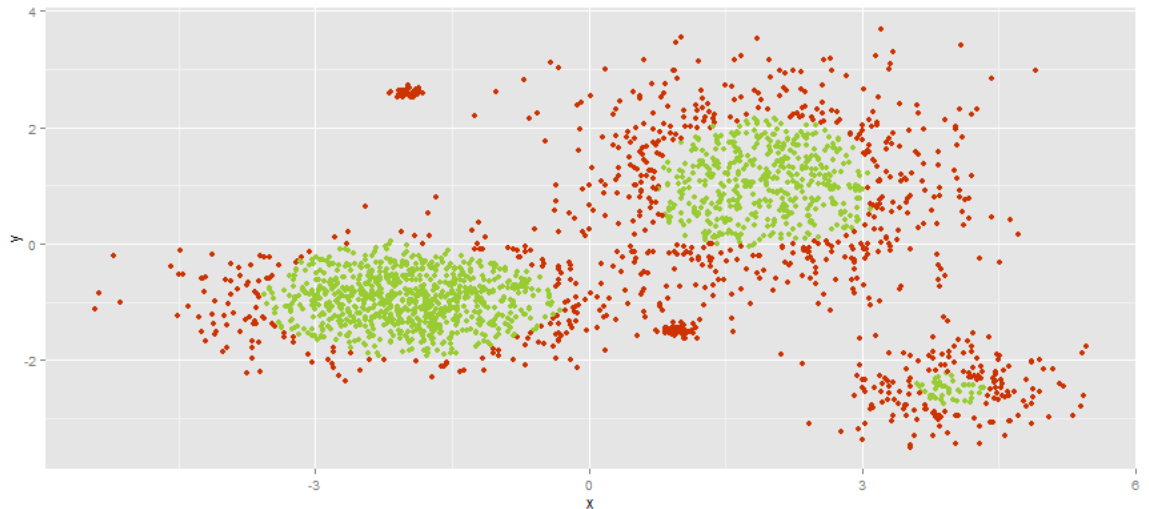
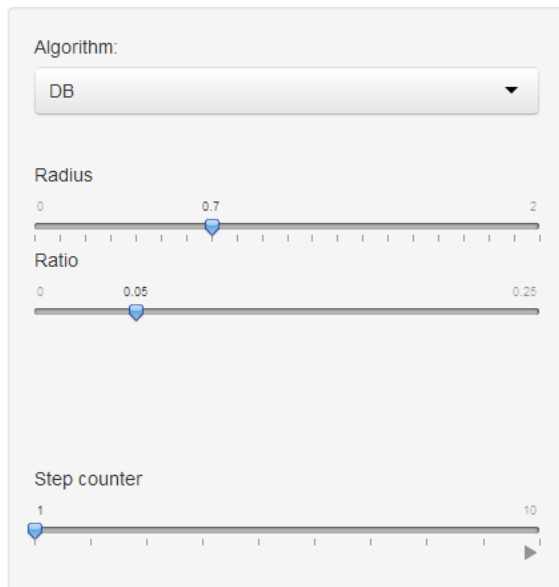


# db – Distance-Based



- Csomag: *fields*
- Függvény: *fields.rdist.near*
- Paraméterek: *delta* sugár

## Rare event detection

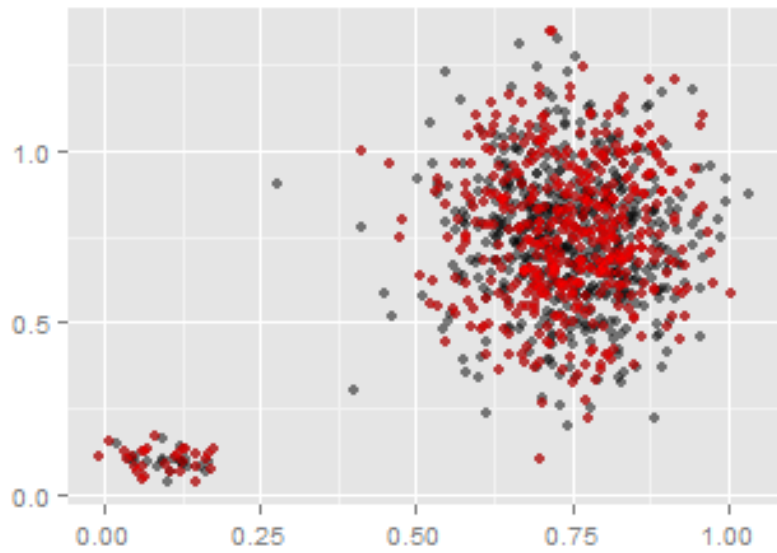


```
[1] "Ratio of chosen elements: "  
[1] 0.4145455
```

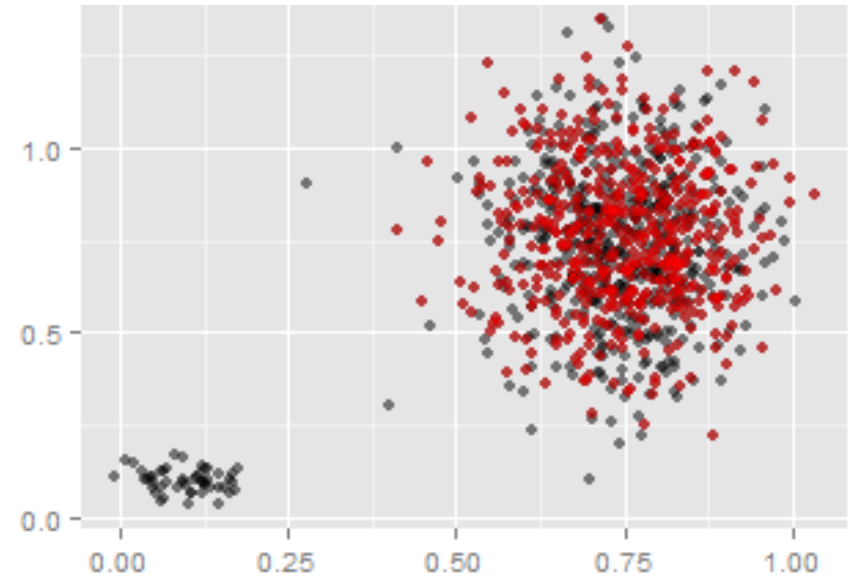
# MCD

- Minimum Covariance Determinant
- Alapötlet
  - Keressük meg a legkompaktabb részhalmazt!

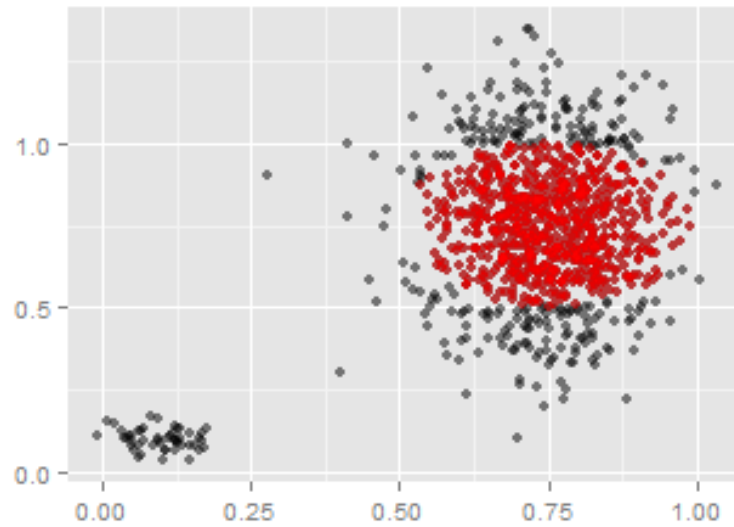
# MCD



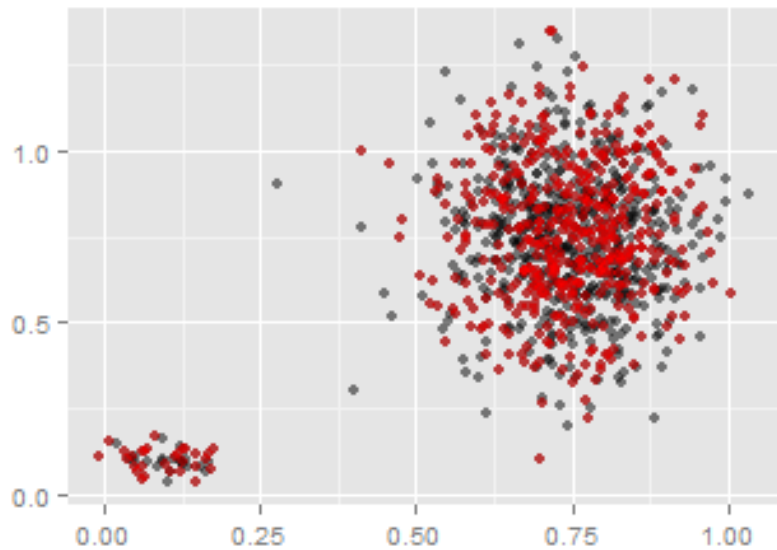
De



ipá



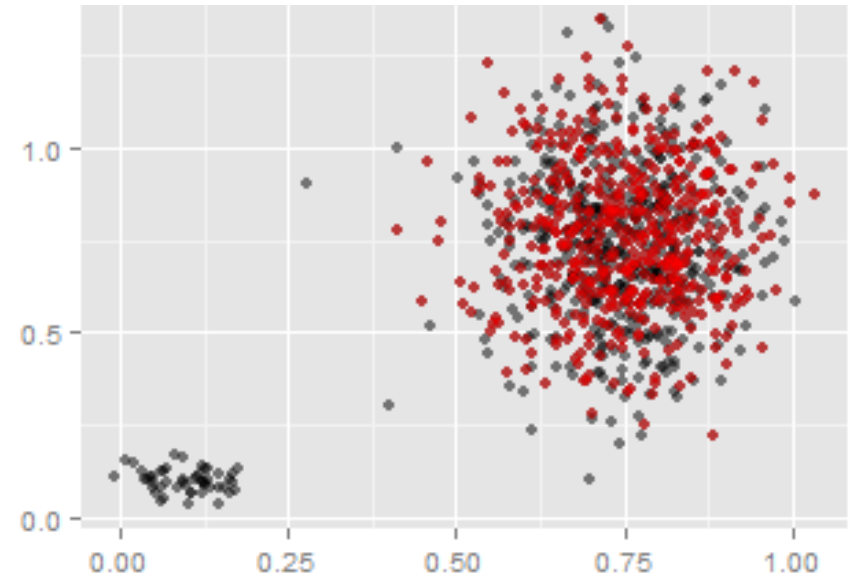
# MCD



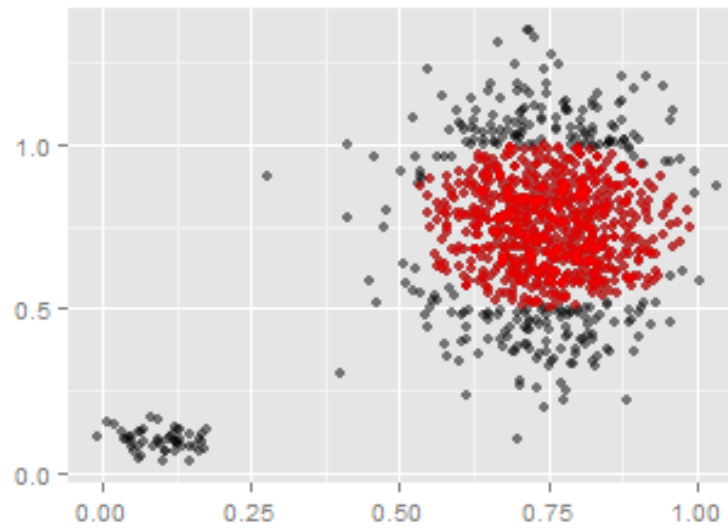
0.0014

$De$

$ipa$

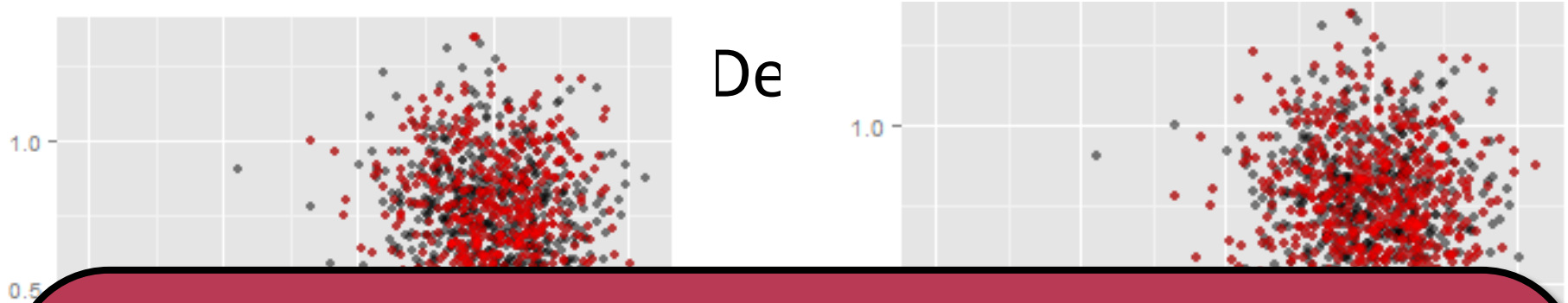


0.00041



0.00011

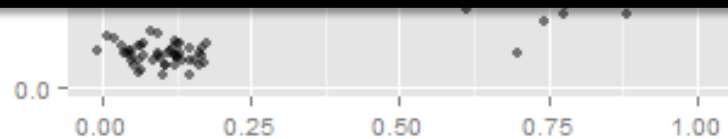
# MCD



Kimerítő keresés?

```
choose (n = 1000, k = 900)
```

```
[1] 6.385051e+139
```



0.00011

# FAST-MCD

- Közelítő algoritmus
- Véletlenszerűen választott kezdőhalmaz
- Iteratív
- Legközelebbi pontok kiválasztása
  - Mahalanobis távolság alapján

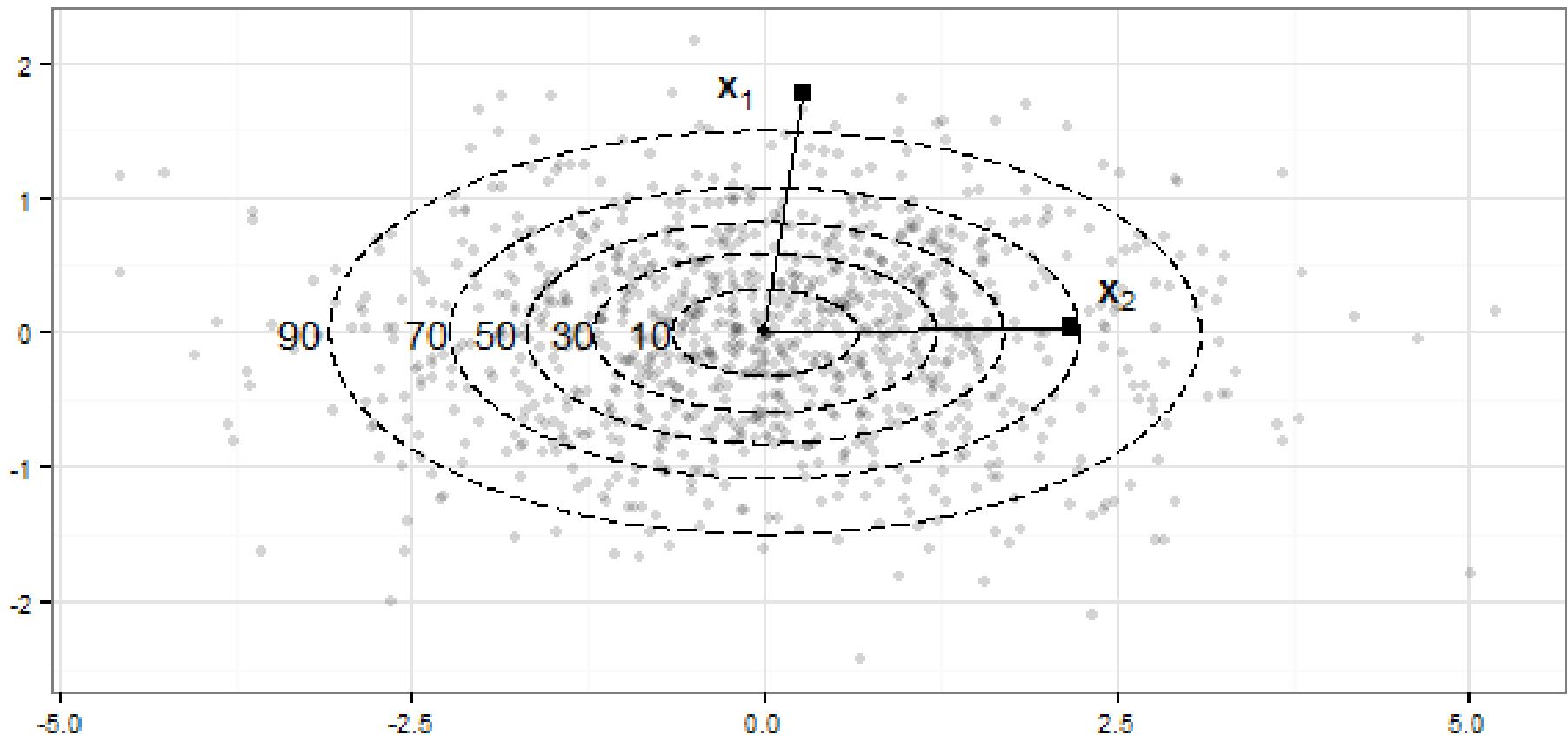
# Mahalanobis távolság

- $D(x, M) = \sqrt{(x - \vartheta)^T S^{-1} (x - \vartheta)}$ 
  - $S$  – kovarianciamátrix
  - $\vartheta$  – súlypont



# Mahalanobis távolság

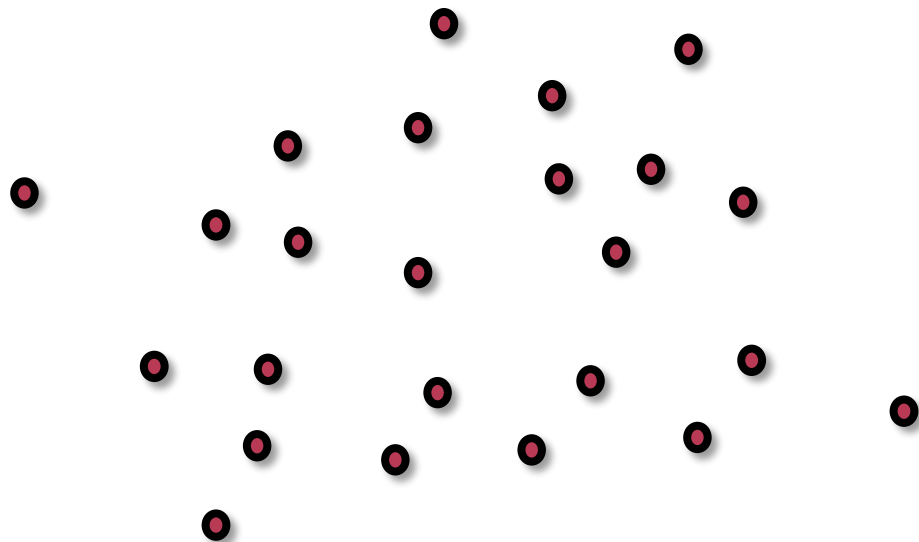
- $D(x, M) = \sqrt{(x - \vartheta)^T S^{-1} (x - \vartheta)}$ 
  - $S$  – kovarianciamátrix



# FAST-MCD

- Közelítő algoritmus
- Véletlenszerűen választott kezdőhalmaz

- Iteratív

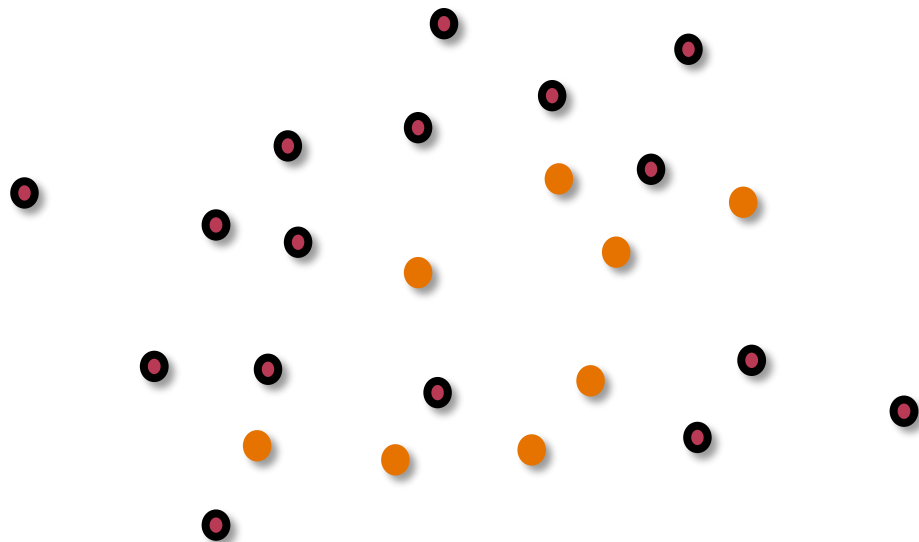


- Legközelebbi pontok kiválasztása
  - Mahalanobis távolság alapján
  - Legközelebbi  $x^0$

# FAST-MCD

- Közelítő algoritmus
- Véletlenszerűen választott kezdőhalmaz

- Iteratív

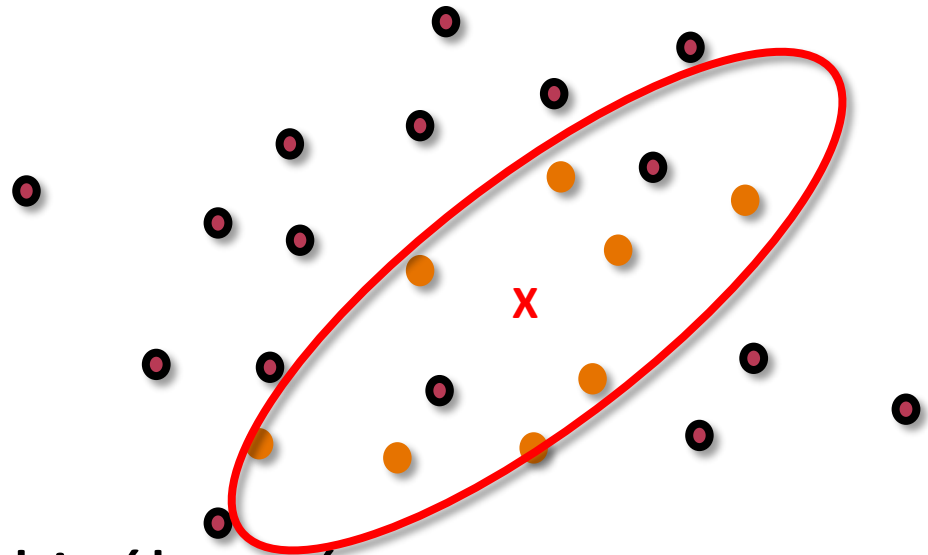


- Legközelebbi pontok kiválasztása
  - Mahalanobis távolság alapján
  - Legközelebbi  $x\%$

# FAST-MCD

- Közelítő algoritmus
- Véletlenszerűen választott kezdőhalmaz

- Iteratív



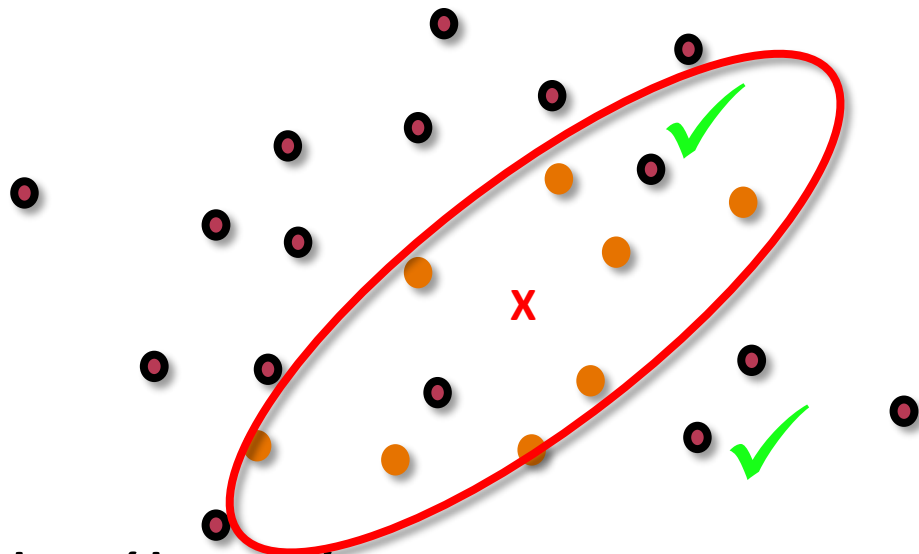
- Legközelebbi pontok kiválasztása
  - Mahalanobis távolság alapján
  - Legközelebbi  $x^0$

# FAST-MCD

- Közelítő algoritmus
- Véletlenszerűen választott kezdőhalmaz

- Iteratív

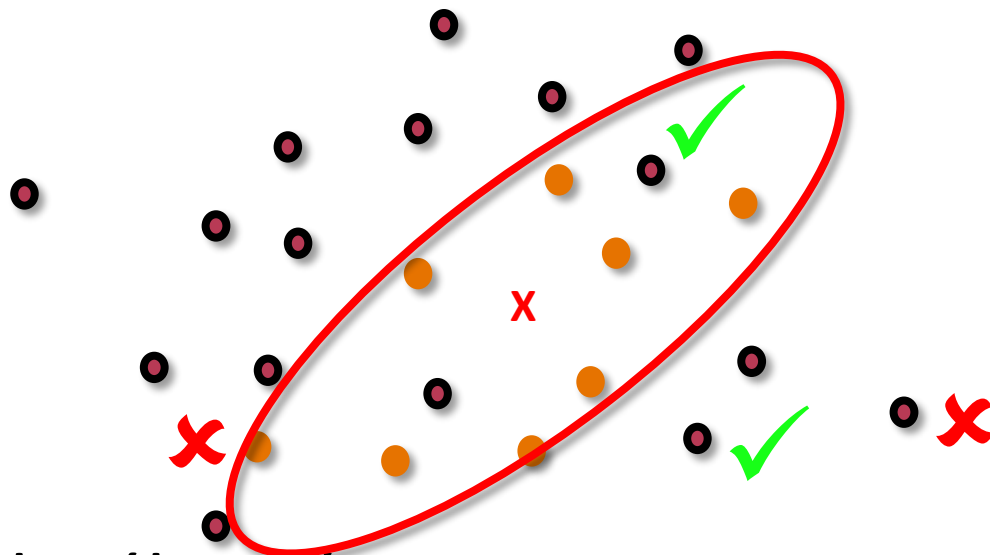
- Legközelebbi pontok kiválasztása
  - Mahalanobis távolság alapján
  - Legközelebbi  $x^0$



# FAST-MCD

- Közelítő algoritmus
- Véletlenszerűen választott kezdőhalmaz

- Iteratív



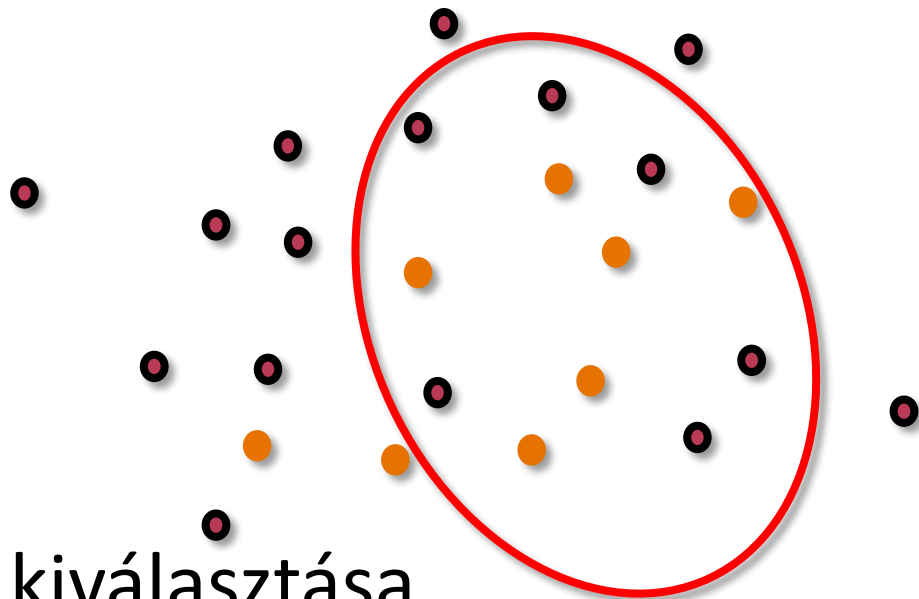
- Legközelebbi pontok kiválasztása
  - Mahalanobis távolság alapján
  - Legközelebbi  $x\%$

# FAST-MCD

- Közelítő algoritmus
- Véletlenszerűen választott kezdőhalmaz

- Iteratív

- Legközelebbi pontok kiválasztása
  - Mahalanobis távolság alapján
  - Legközelebbi  $x\%$

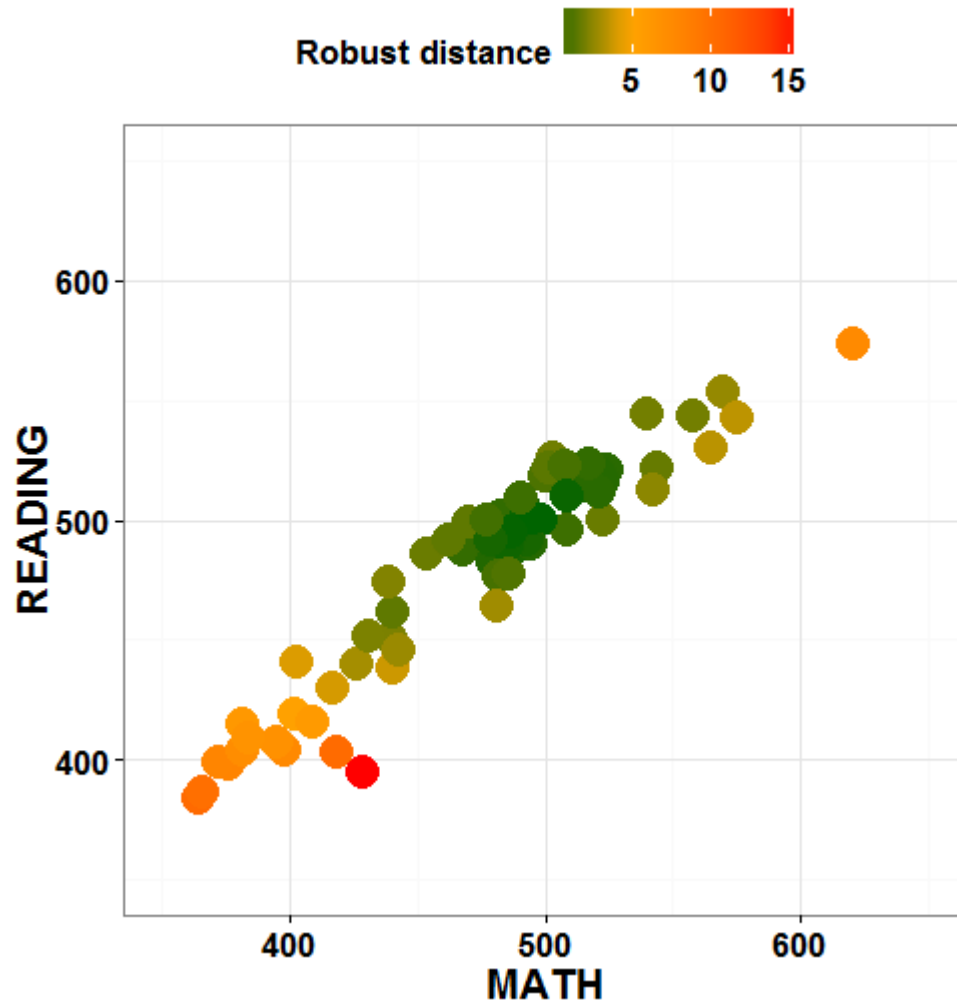


# BACON

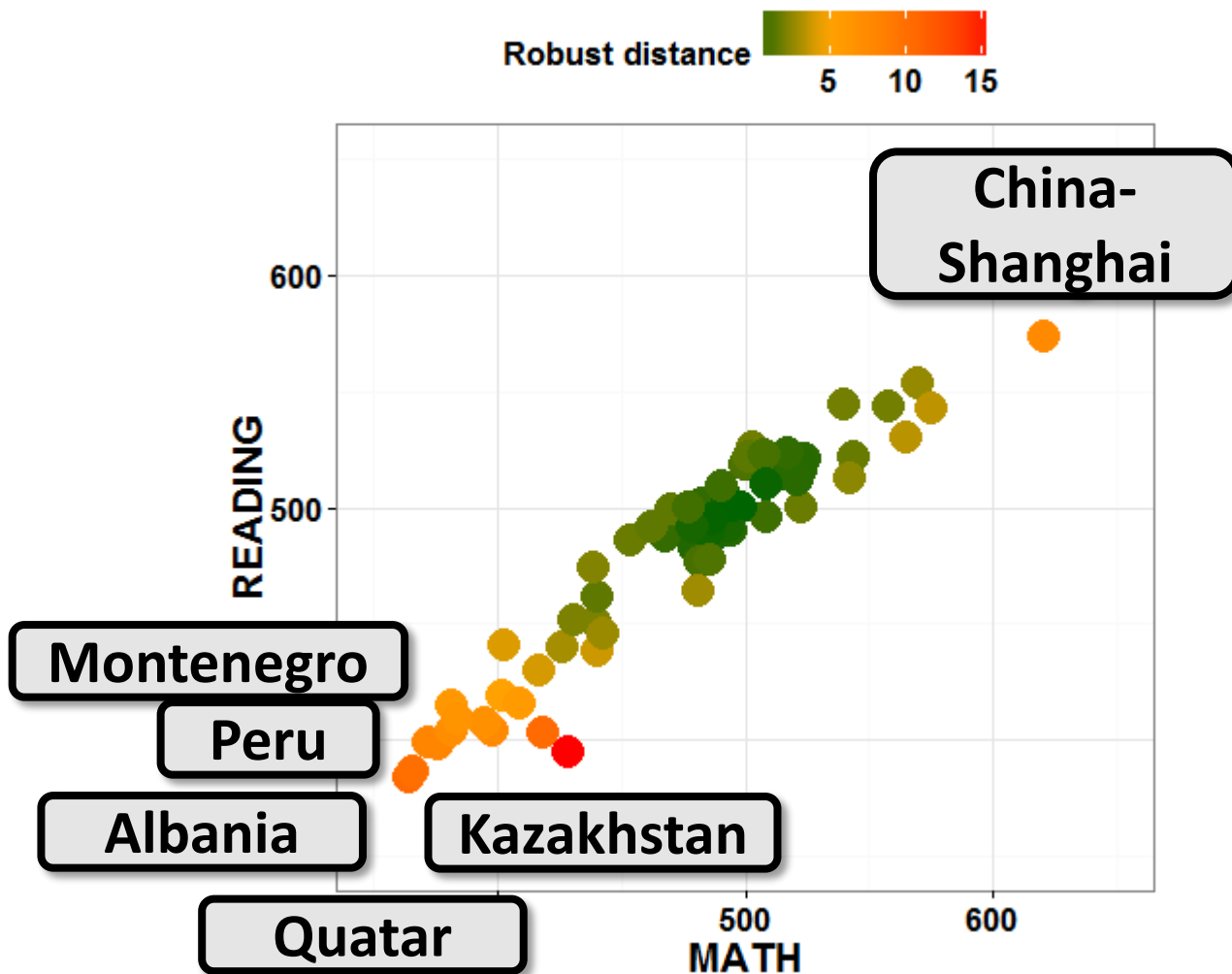
- Blocked Adaptive Computationally Efficient Outlier Nominators
- Kiinduló halmaz félig felügyelt módban is!
- Új halmaz: küszöbérték alapján



# BACON



# BACON



- Csomag: *robustX*
- Függvény: *mvBACON*
- Paraméterek
  - *init.sel* kezdőhalmaz
    - „manual” – *man.sel* kezdőhalmaz
    - „Mahalanobis”, „dUniMedian” – *m* kezdőhalmaz mérete

■ Csomag: *robustX*

## Rare event detection

Algorithm:

BACON

Type:

Manual selection

alpha

0.01

0.95

Initial set:

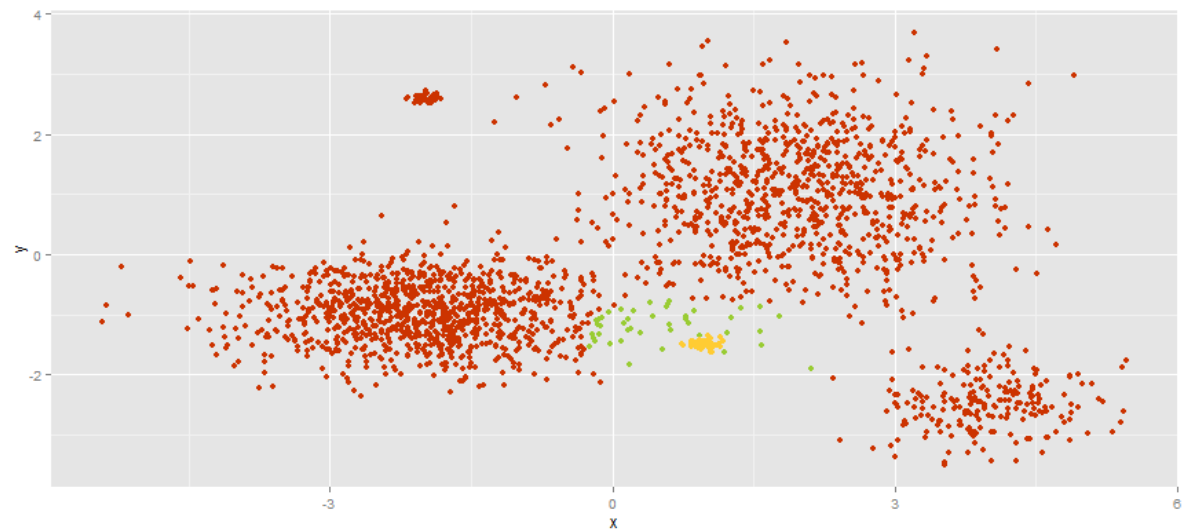
Set 5

Step counter

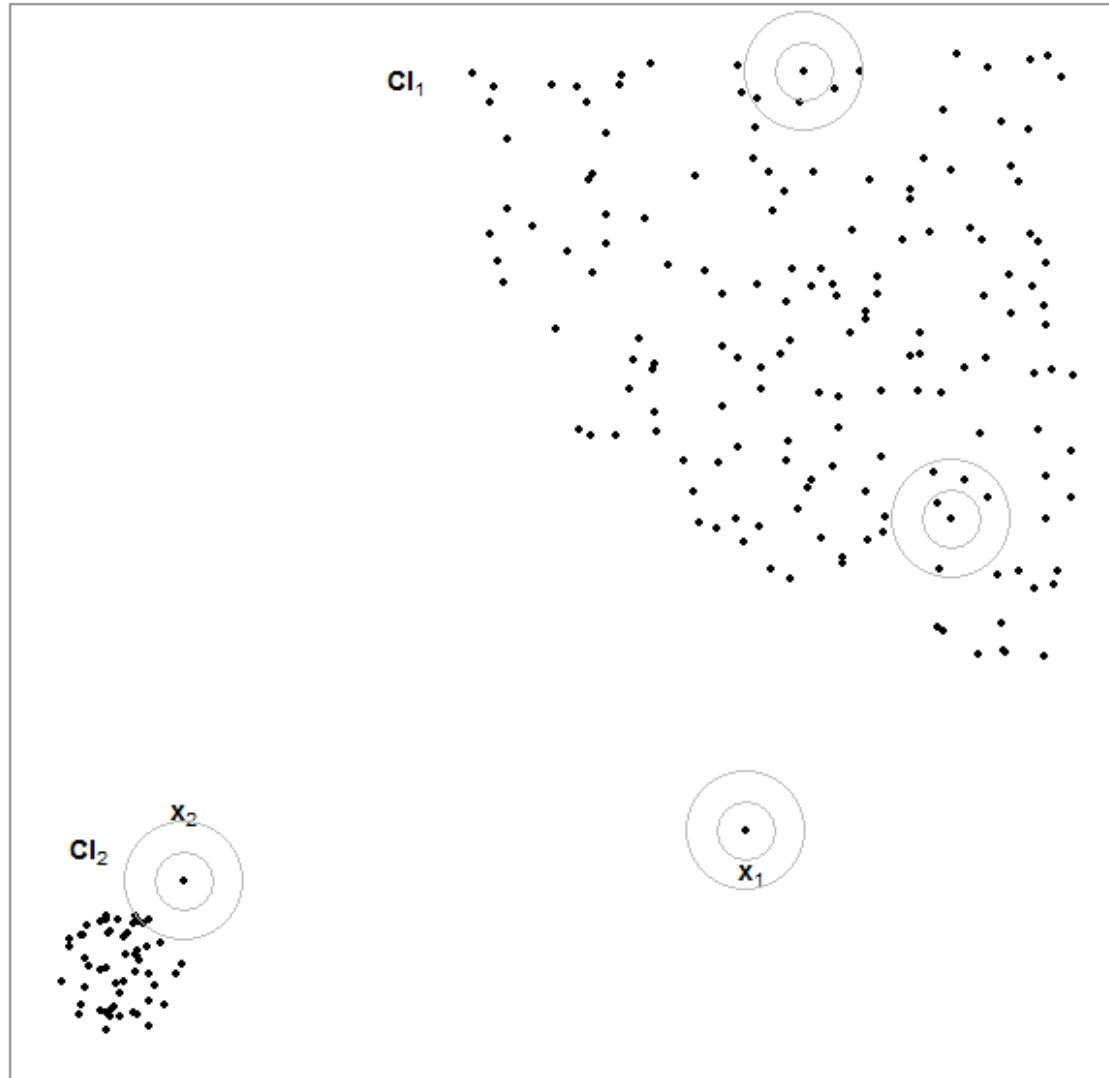
1

2

10



# LOF motiváció



- Local Outlier Factor

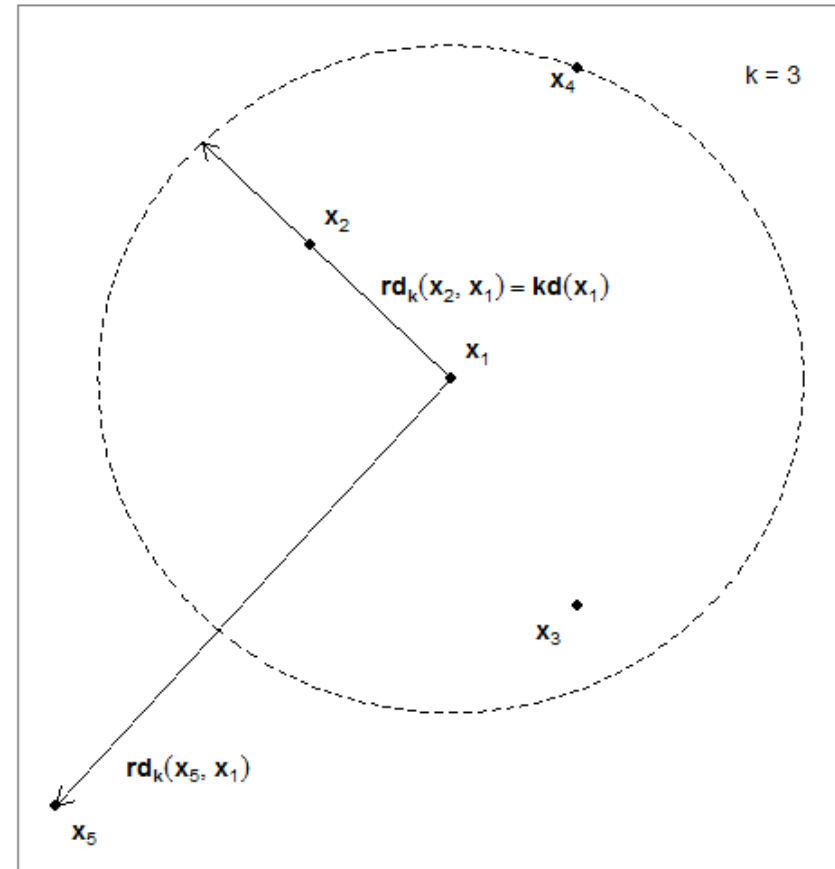
- Alapötlet

- csak a szomszéd számít
- a távolság is módosul
- lokális sűrűség

- Outlier kritérium

- a lokális sűrűség jóval kisebb, mint a szomszédaimnak átlagosan

- Local Outlier Factor
- Alapötlet
  - csak a szomszéd számít
  - a távolság is módosul
  - lokális sűrűség
- Outlier kritérium
  - a lokális sűrűség jóval kisebb átlagosan



- Local Outlier Factor

- Alapötlet

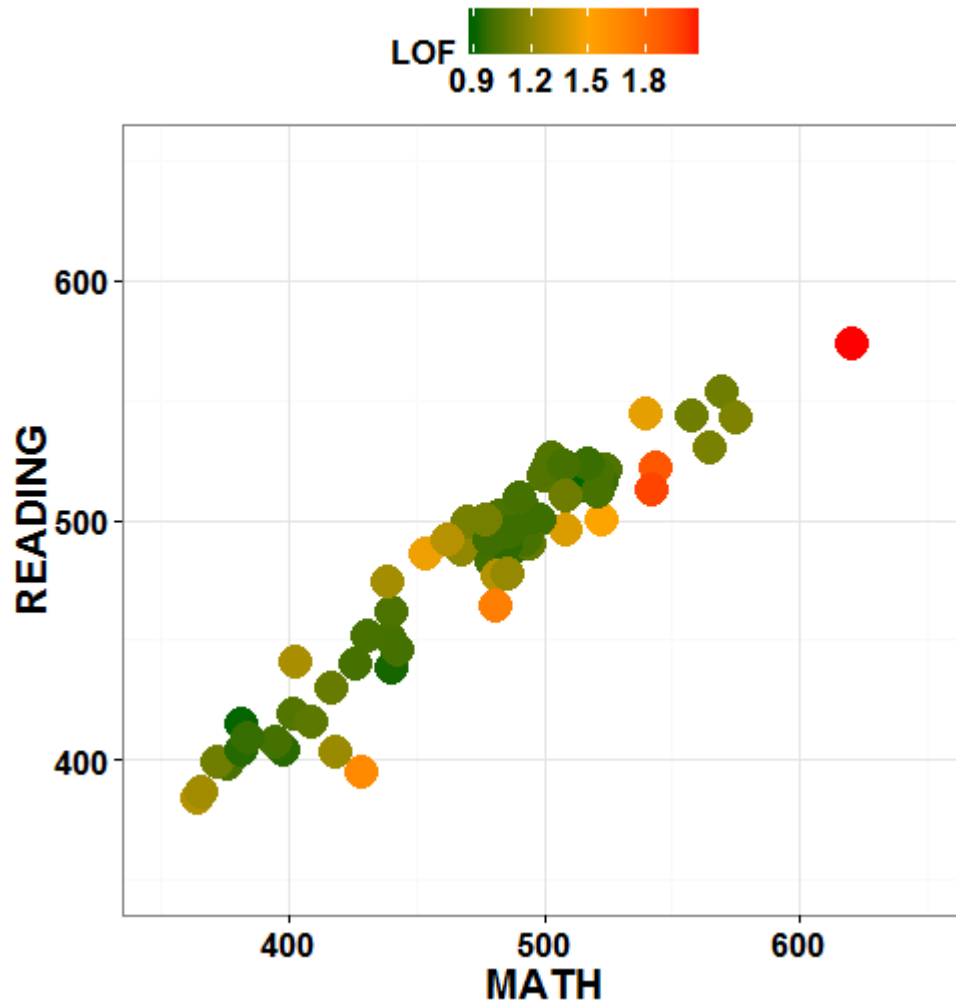
- csak a szomszéd számít
- a távolság is módosul
- lokális sűrűség

- Outlier kritérium

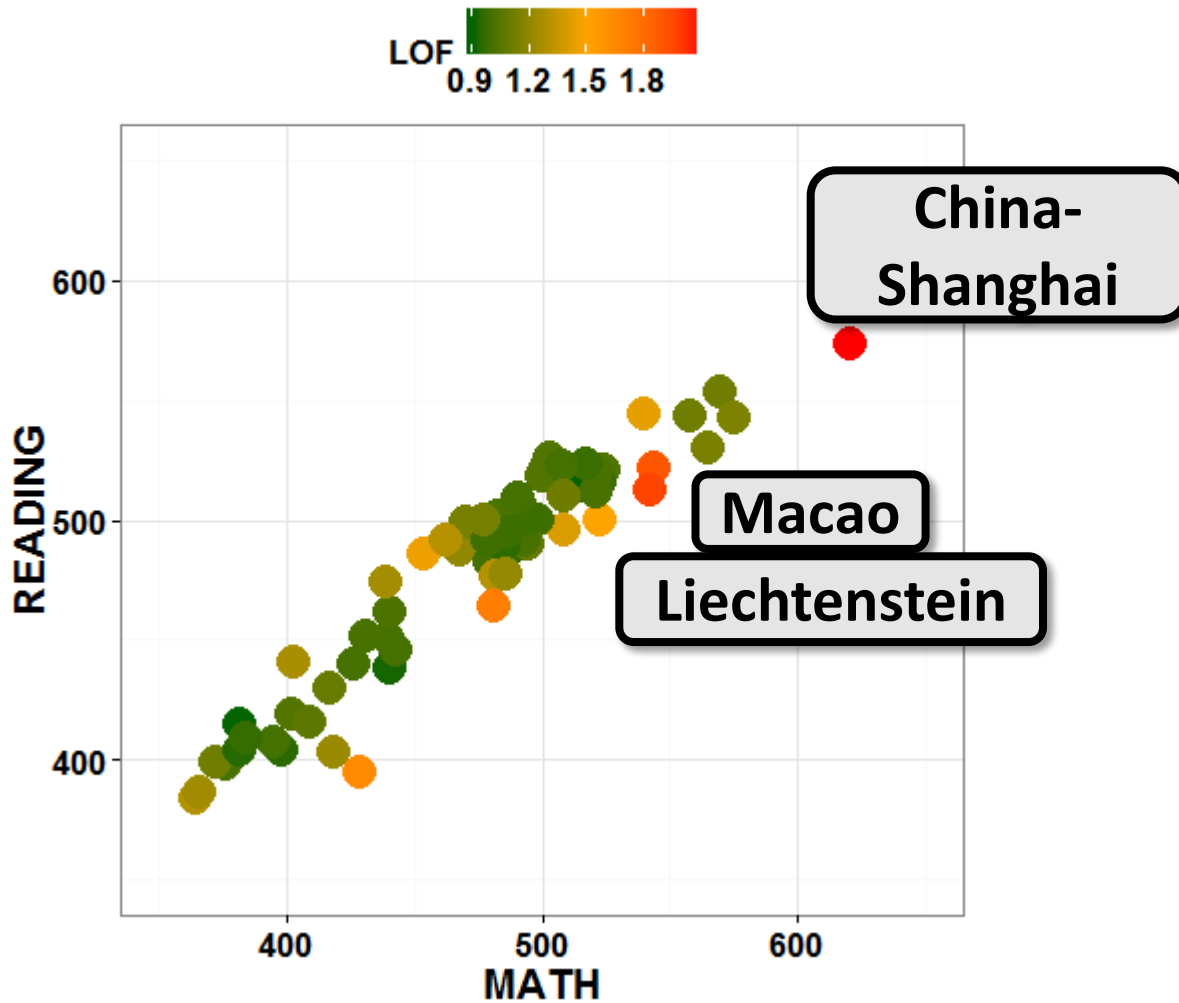
- a lokális sűrűség jóval kisebb, mint a szomszédaimnak átlagosan



# LOF



# LOF



- Csomag: *DMwR* (Data Mining with R)
- Függvény: *lofactor*
- Paraméterek:  $k$  szomszédság mérete

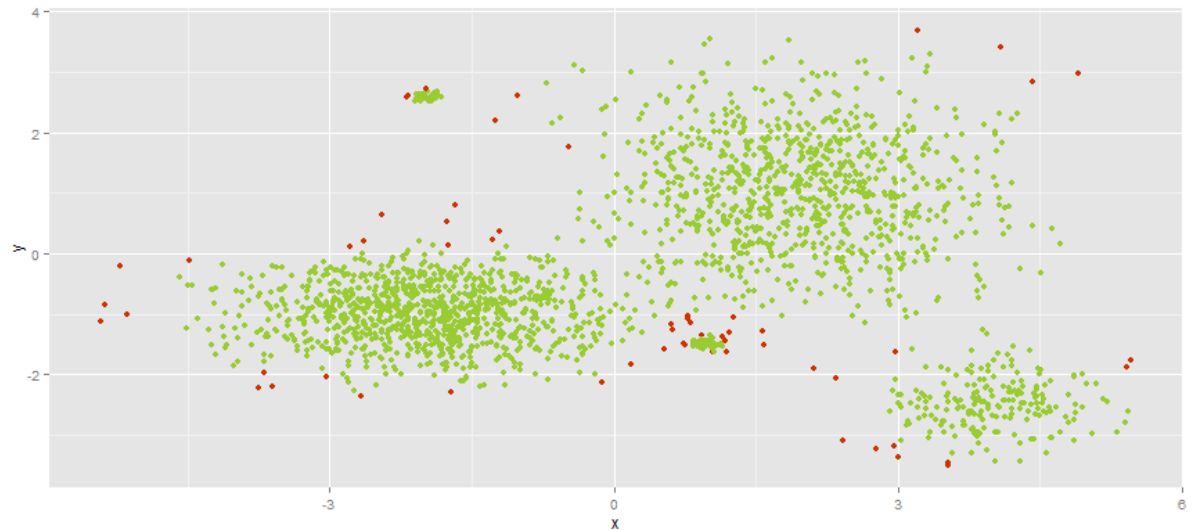
## Rare event detection

Algorithm:  
LOF

k  
1 15 30

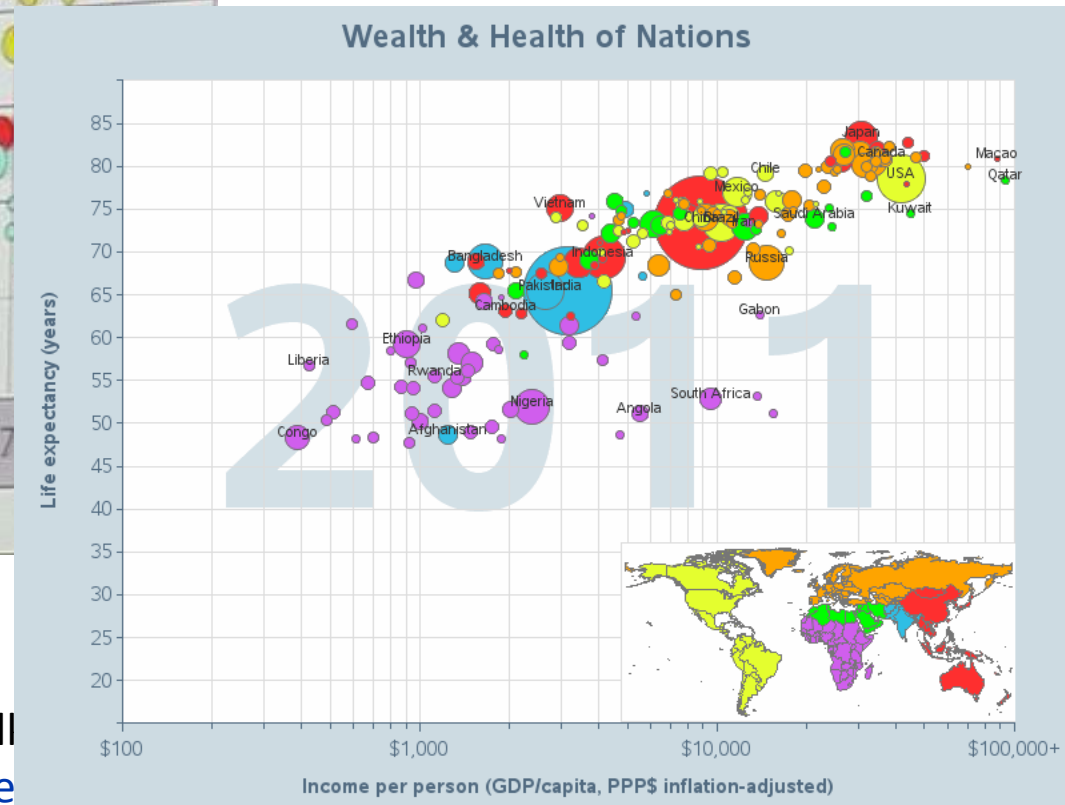
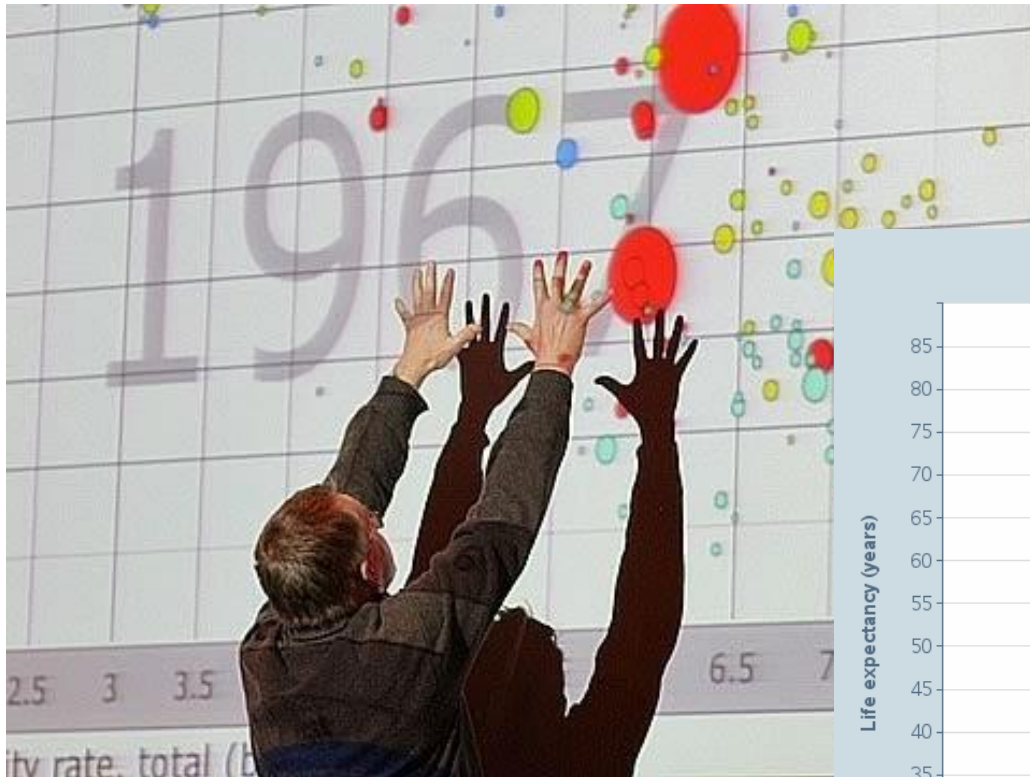
threshold  
1 1.5 3

Step counter  
1 10



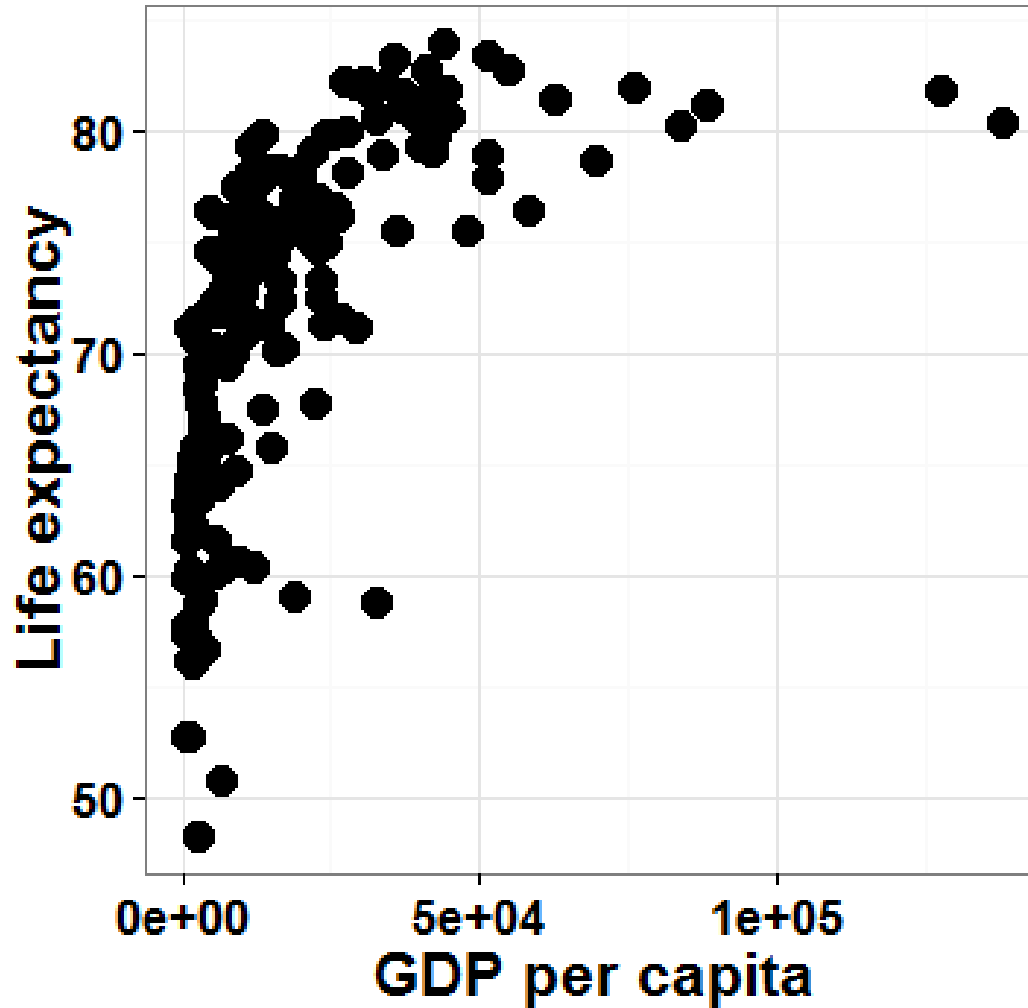
# Esettanulmány: wealth & health of nations

## ■ Hans Rosling 2006-os TED talkja

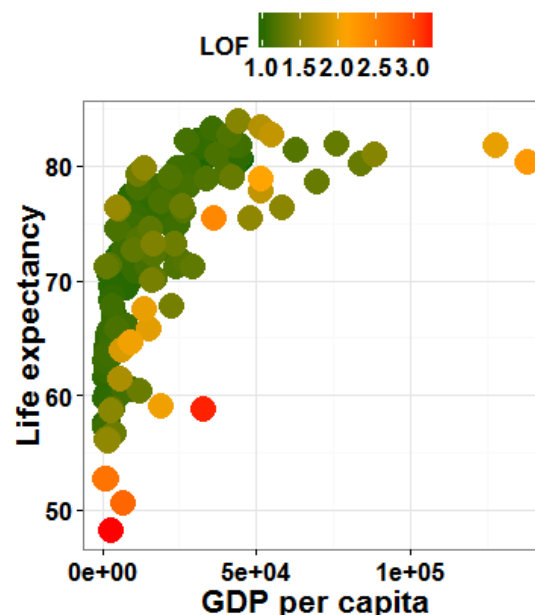
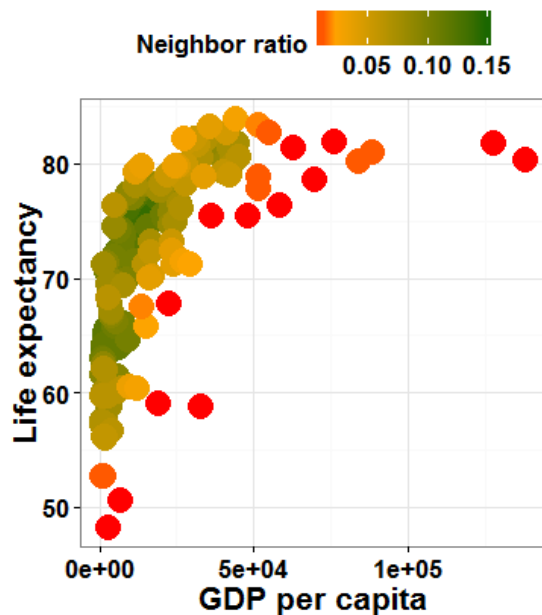
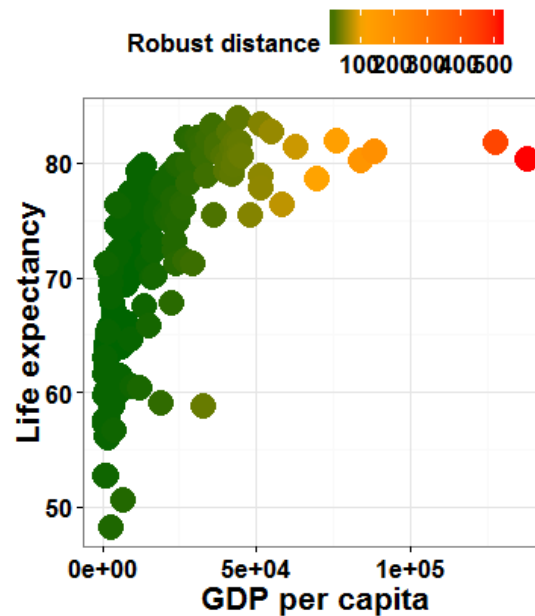
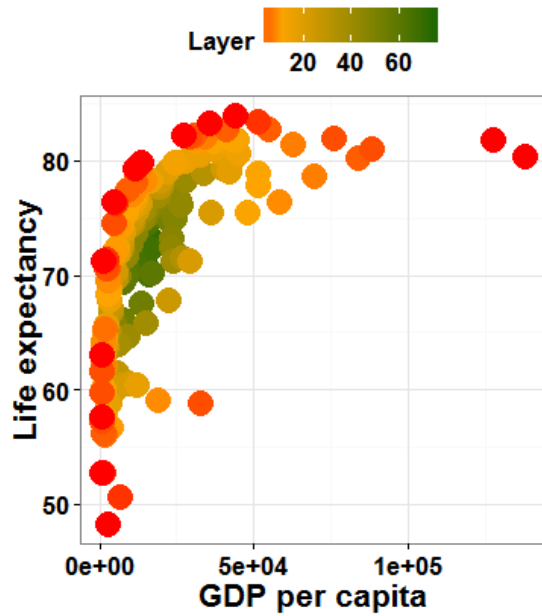


A másik teljesen irreleváns kedvenc TED talk  
[http://www.ted.com/talks/david\\_mccandle](http://www.ted.com/talks/david_mccandle)

# Wealth and health of nations

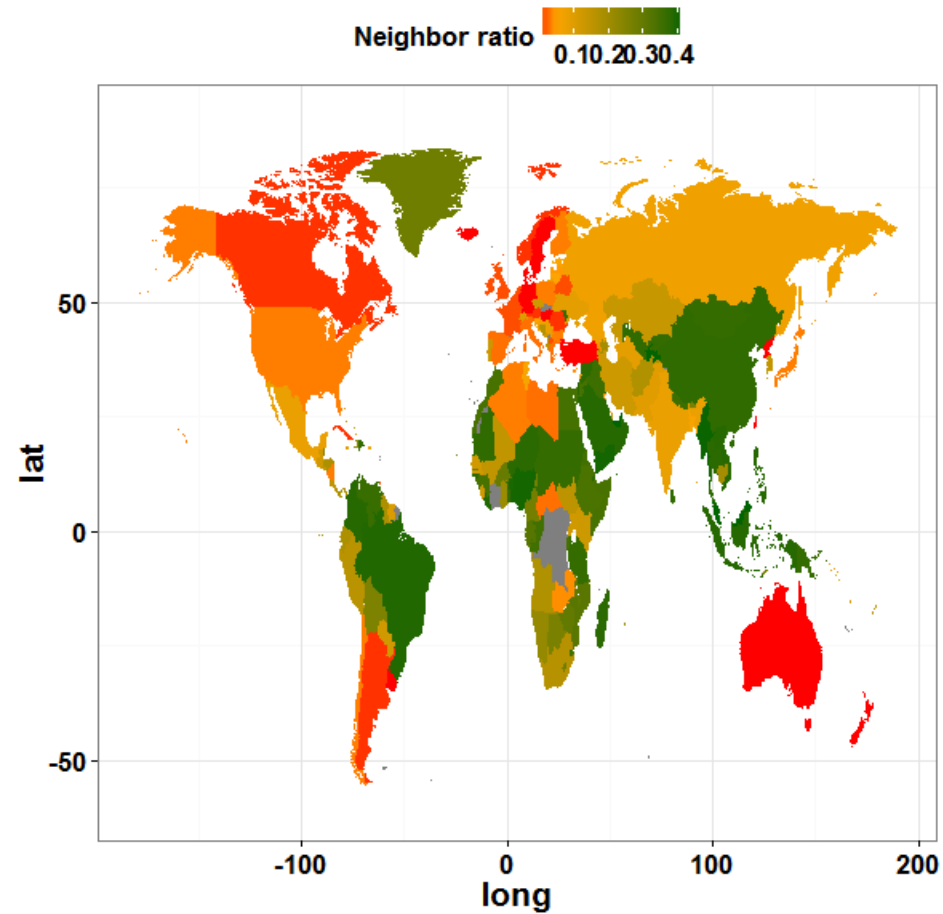
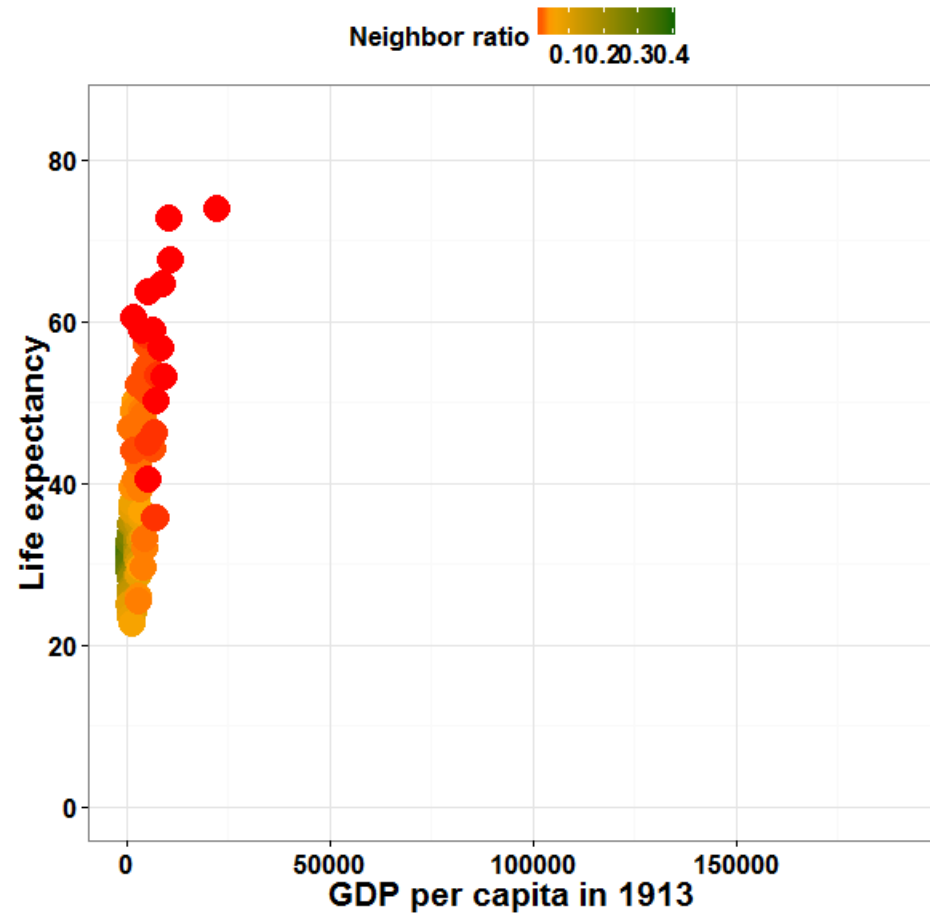


# Wealth and health of nations



Happy families are all alike;  
**every** unhappy family is unhappy **in its own way.**

# Wealth and health of nations időbeli változás



# Hivatkozásjegyzék

- Outlier Detection alapmű
  - Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. ACM Computing Surveys (CSUR), 41(3):15, 2009
- Outlier Detection demo
  - <https://github.com/FTSRG/BigDANTE/tree/master/RareEventDetectionDemo>