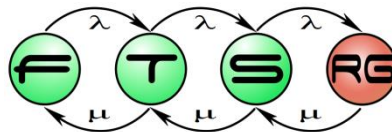


Vizuális adatelemzés

Rendszermodellezés 2019.

Budapest University of Technology and Economics
Fault Tolerant Systems Research Group



Tartalom

Miért vizualizálunk?



Mit vizualizálunk?



Hogyan vizualizálunk?



Mire következtetünk?

Tartalom

Miért vizualizálunk?



Mit vizualizálunk?



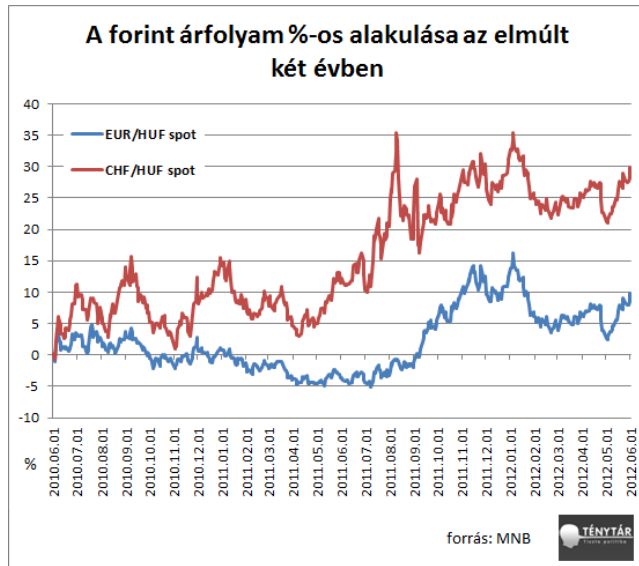
Hogyan vizualizálunk?



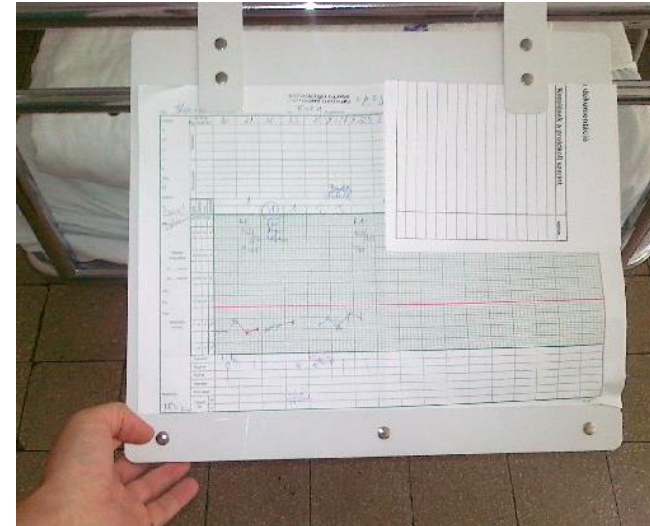
Mire következtetünk?

A vizualizáció alkalmazásai

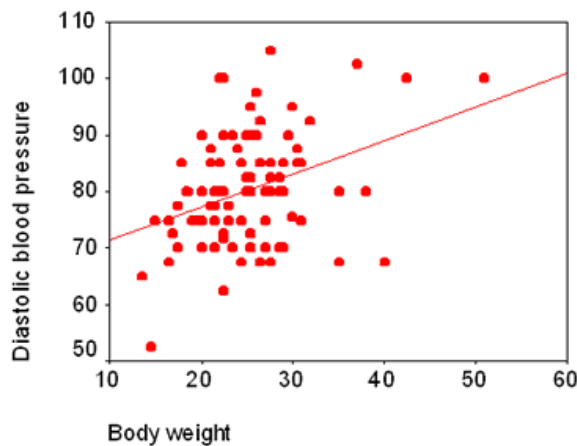
Trend analízis és előrejelzés



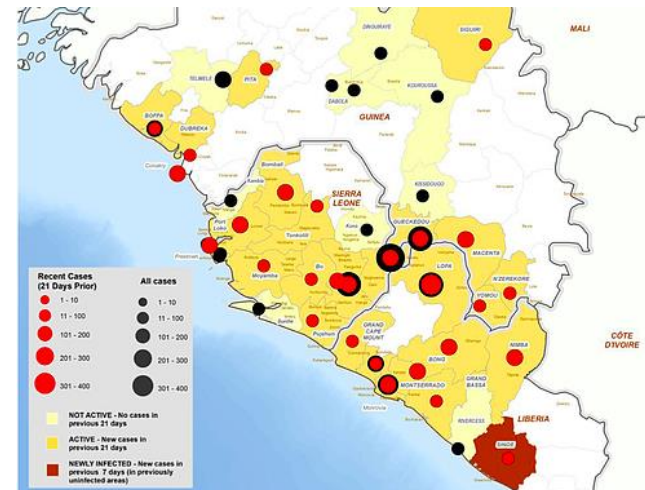
Idősor analízis



Korrelációanalízis

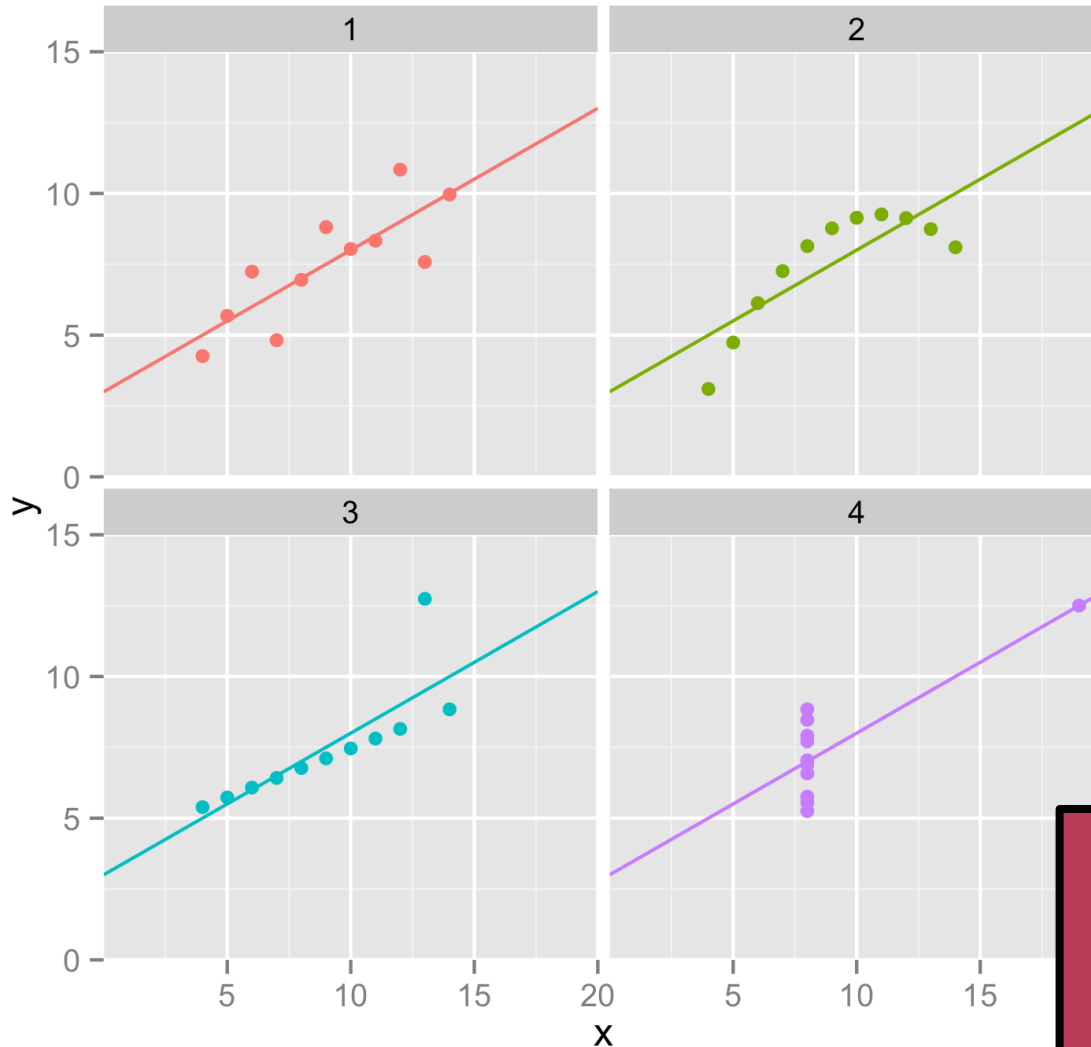


Térbeli analízis



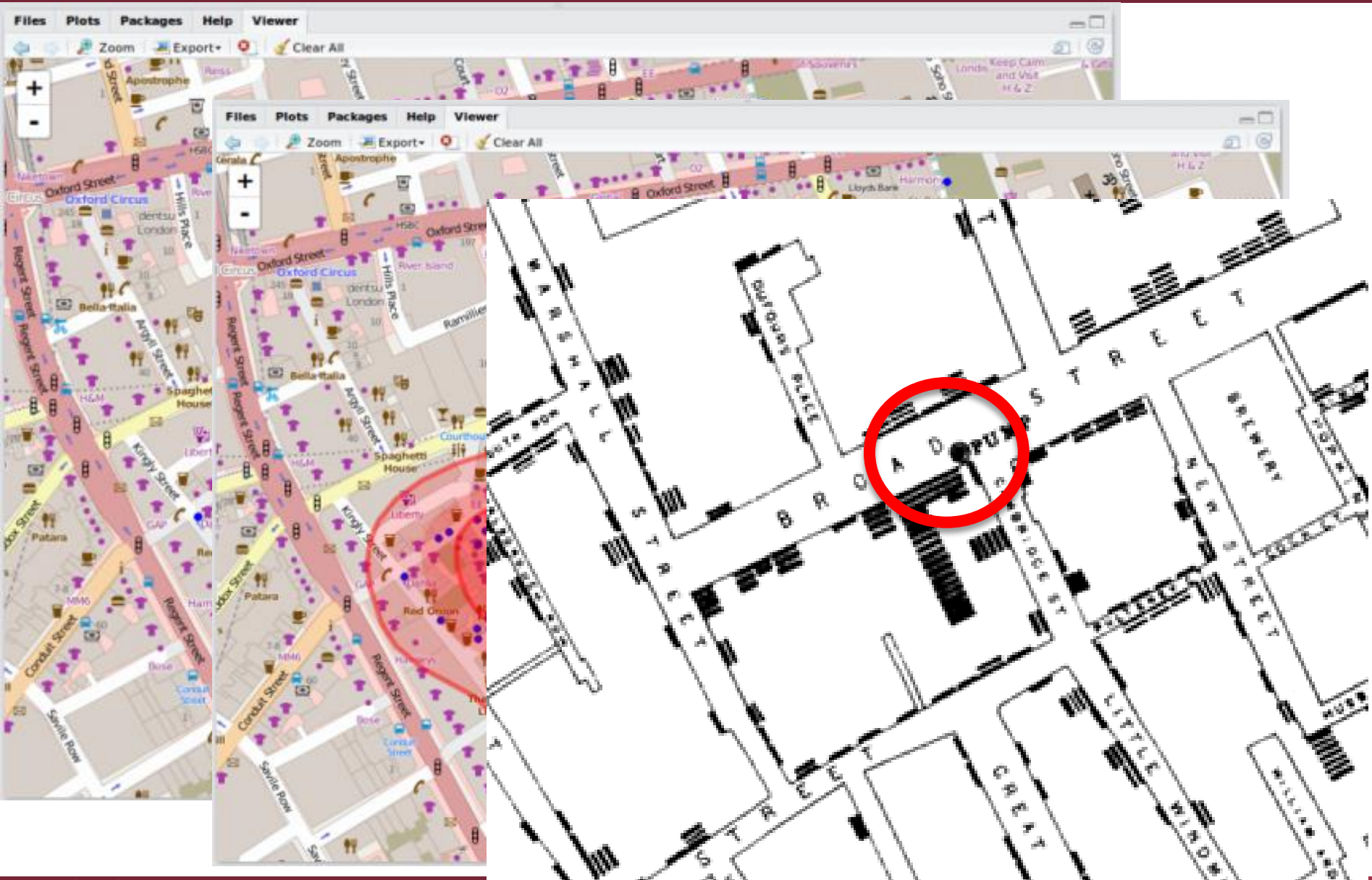
Számítások ellenőrzése

Anscombe's Quartet



Hibás feltételezések
elkerülése... és intuíció

Összefüggések feltárása



Mindent a szemnek!

„Masszív” erőforrások

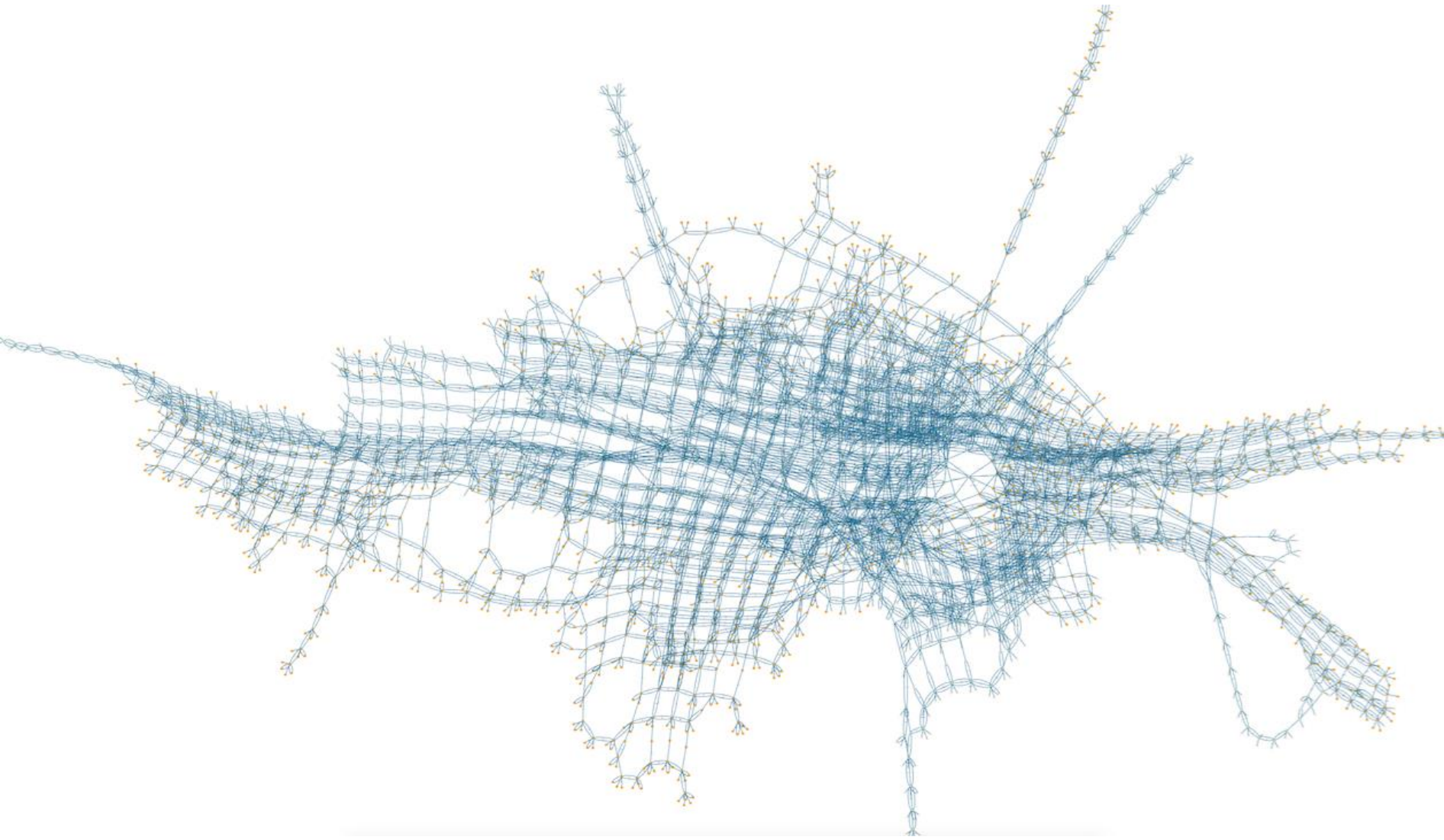
- 120.000.000 szenzor



3. Vizuális kiválasztás és manipuláció

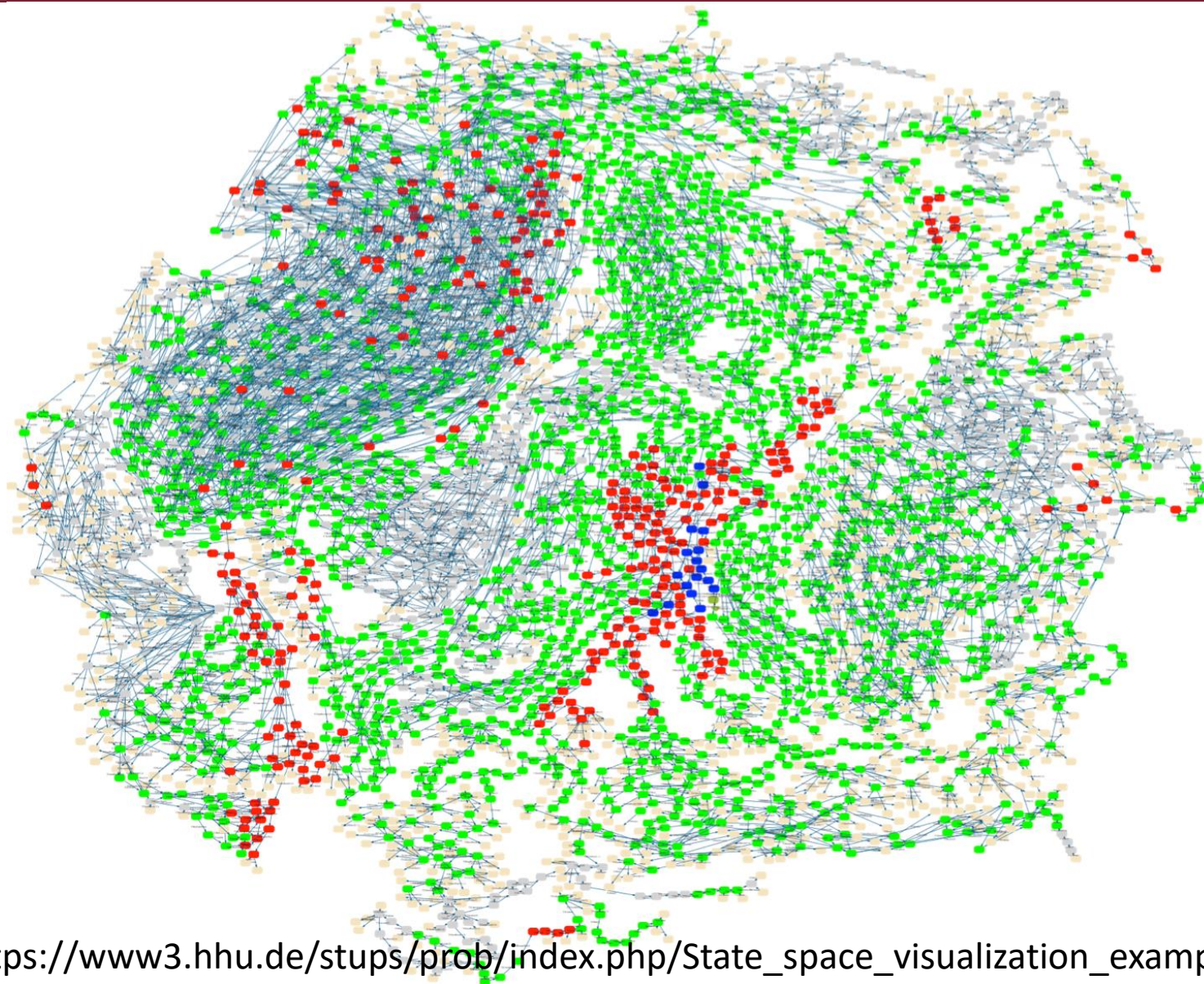
4. Interpretáció, korreláció más modellekkel, kiértékelés

Példa: állapottér vizualizáció (hálózat)



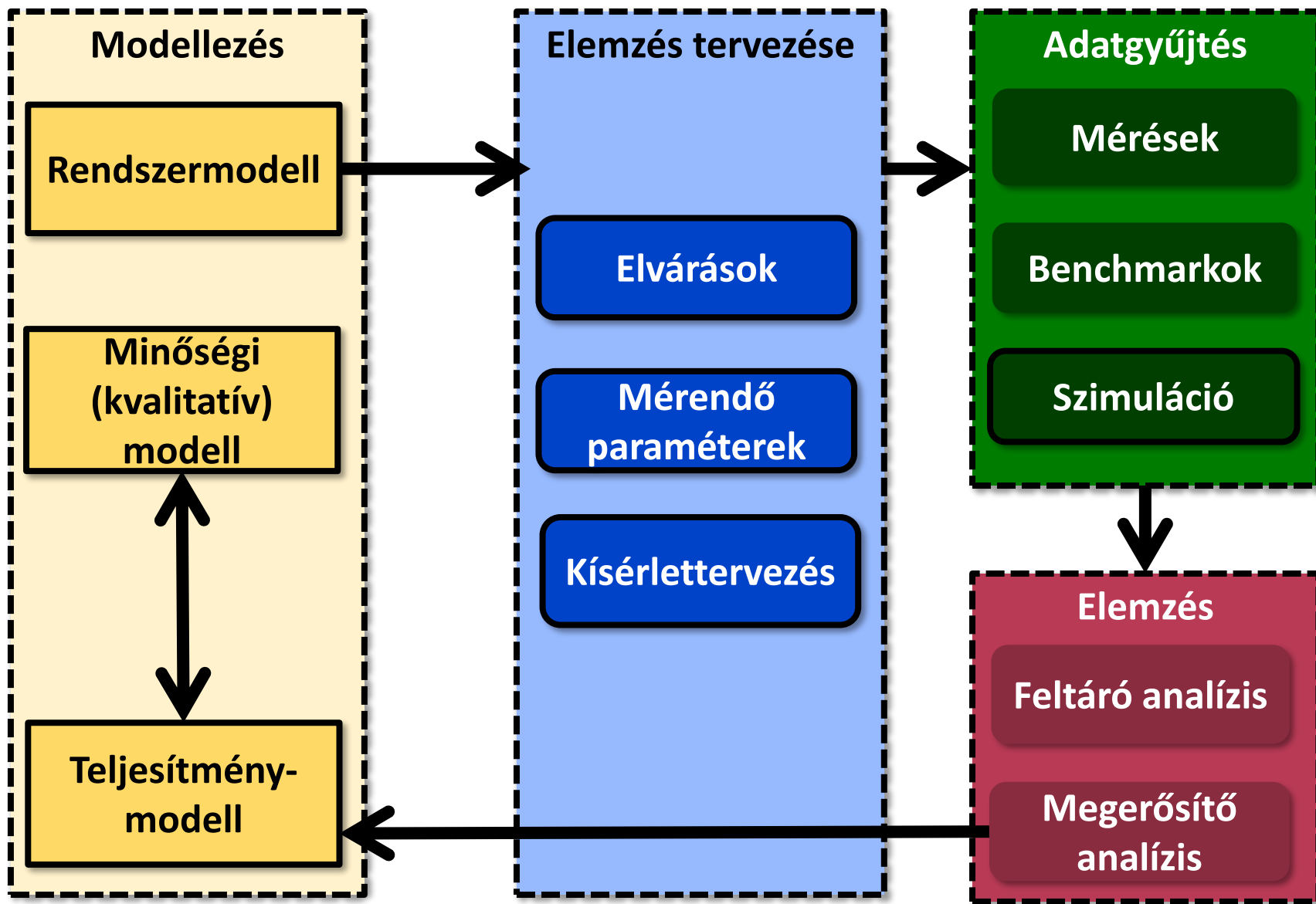
https://www3.hhu.de/stups/prob/index.php/State_space_visualization_examples

Példa: CAN bus állapottér



https://www3.hhu.de/stups/prob/index.php/State_space_visualization_examples

Példa: Rendszermodell → teljesítménymodell



Mi is lesz?

Miért vizualizálunk?



Mit vizualizálunk?



Hogyan vizualizálunk?



Mire következtetünk?

Főbb vizualizációs célok

Reproducible
research

Reporting

Self-Service BI

Dashboard

EDA

Storytelling

Data tour

Data
journalism

Infographics

Emlékeztető: táblázatos ábrázolás

- **Táblázat sora** = modellelem
- **Táblázat oszlopa** = tulajdonság

Név	Típus	Méret (kB)	Utolsó módosítás
Dokumentumok	könyvtár		2016.02.02
szerződés.pdf	fájl	569	2015.11.09
Képek	könyvtár		2016.02.02
logó.png	fájl	92	2015.03.06
alaprajz.jpg	fájl	1226	2016.02.02

- Adatelemzési eszközök (pl. R, Python): **dataframe**
 - Egy sor egy mérés
 - Egyes oszlopoknak **típusai** vannak

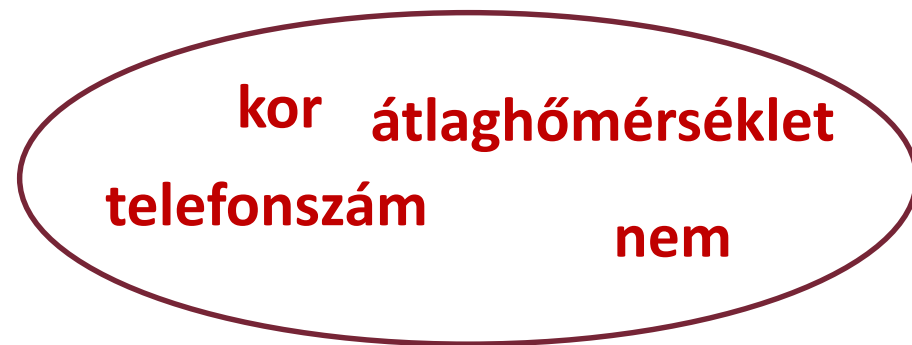
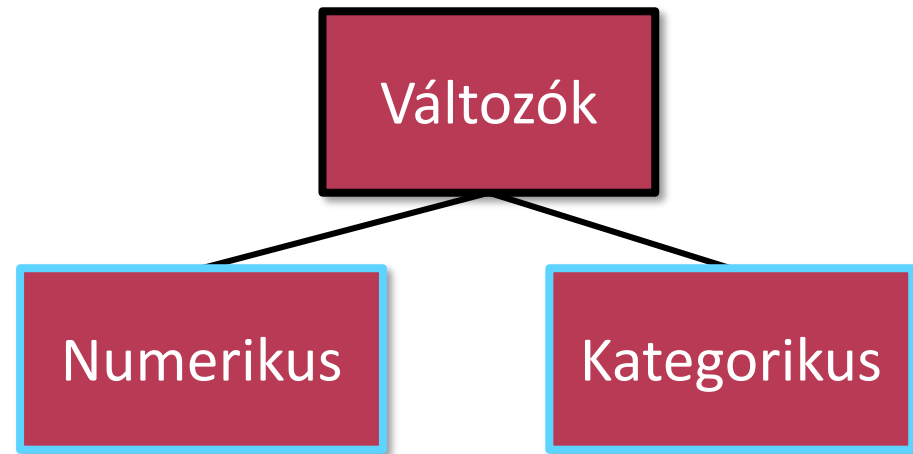
Numerikus és kategorikus változók

- Numerikus (numerical)

- az alapvető aritmetikai műveletek értelmesek

- Kategorikus (categorical)

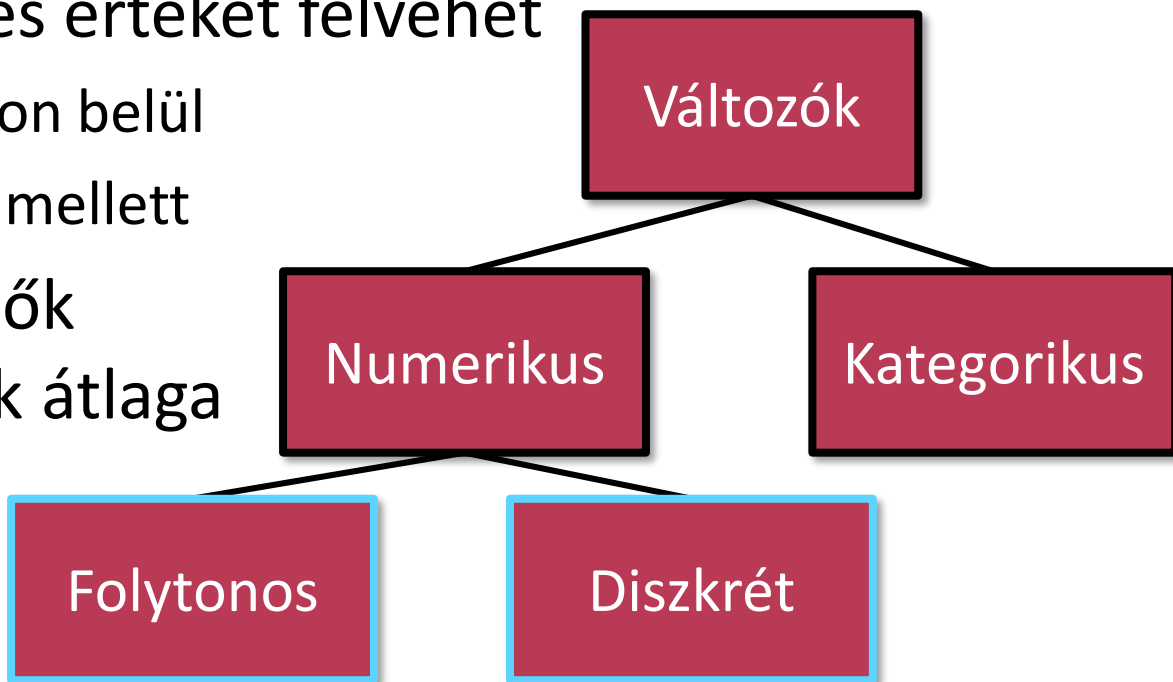
- Matematikai műveletek nem értelmezhetőek rajtuk, legfeljebb sorba rendezés



Numerikus változók

■ Folytonos

- Mért – tetszőleges értéket felvehet
 - adott tartományon belül
 - adott pontosság mellett
- Pl. a teremben ülők ZH pontszámának átlaga



■ Diszkrét

- Számolt – véges sok értéket vehet fel adott tartományban
- Pl. az előadáson ülők száma

Kategorikus változók

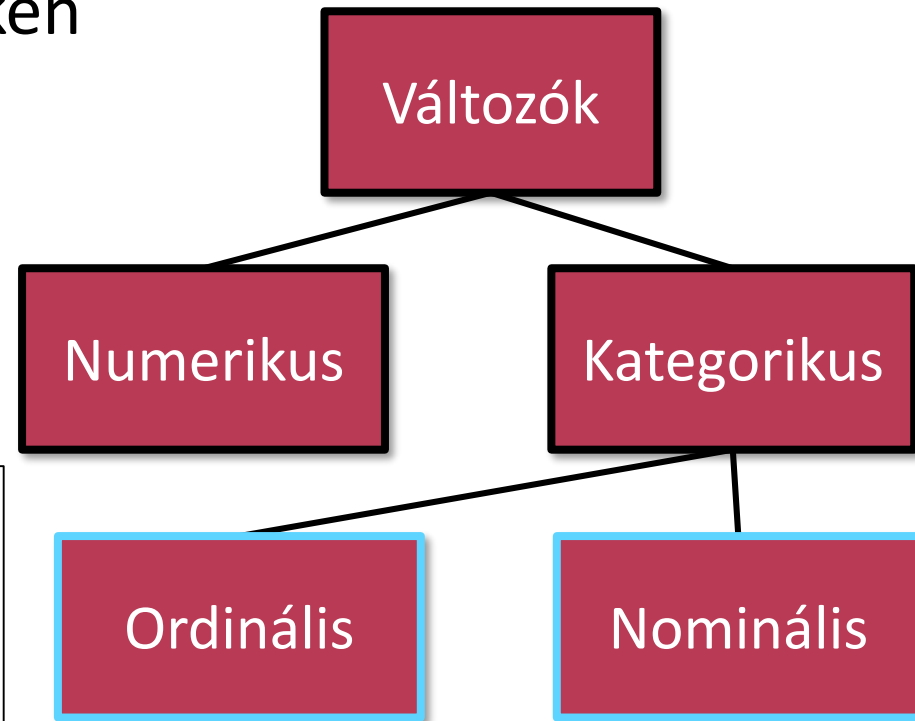
■ Ordinális

- Teljes rendezés az értékeken
- Pl. szállodai csillagok

■ Nominális

9. Ajánlanád-e a tárgyat másoknak?

- Mindenkit rábeszelnék
- Nyugodtan ajánlanám
- Esetleg ajánlanám
- Inkább lebeszelném róla
- Feltétlenül lebeszelném
- Nem kívánok válaszolni



Mi is lesz?

Miért vizualizálunk?



Mit vizualizálunk?



Hogyan vizualizálunk?



Mire következtetünk?

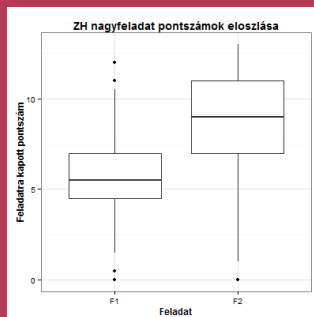
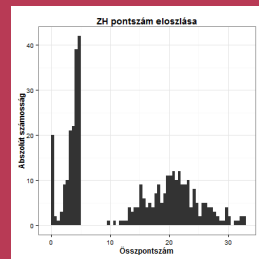
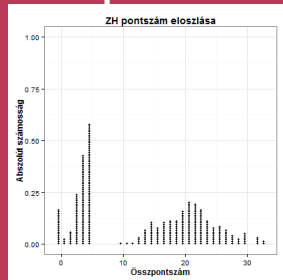
1 változó – eloszlásokra

Változók

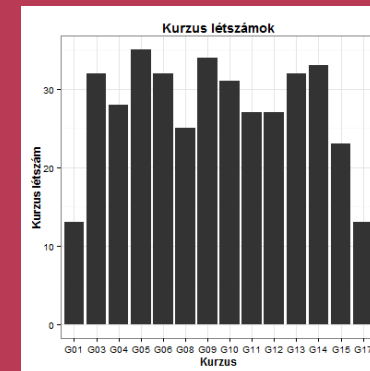
Numerikus

Kategorikus

ZH pontszám: [13, 15, 2, ...]

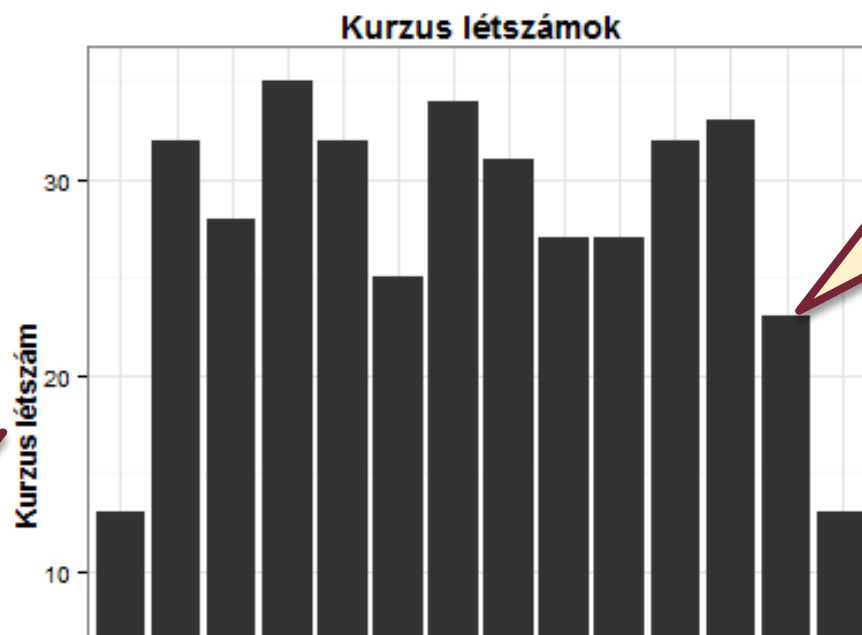


Kurzus: [G01, G03, G15, G17, ...]



Oszlopdiagram

- Bemenő változó: kurzus kód
- Kérdés: az egyes kurzusokra hányan járnak?



abszolút
gyakoriság!

Oszlop-
magasság:
adott érték
gyakorisága

Tervezői döntés: értékkészlet darabolása
Pl.: kedd-csütörtök-péntek

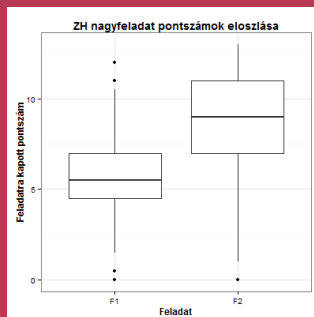
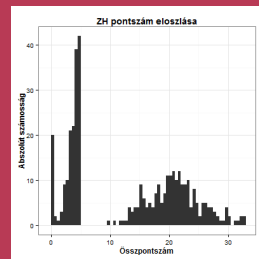
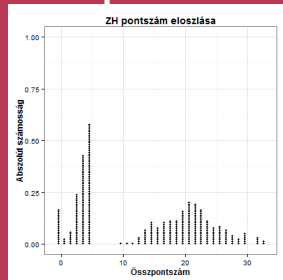
1 változó – eloszlásokra

Változók

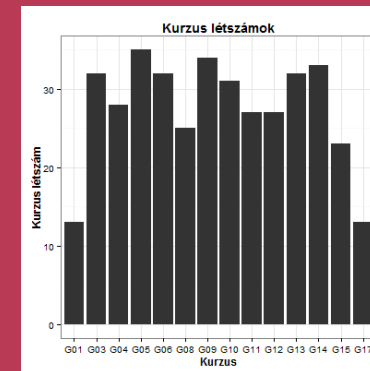
Numerikus

Kategorikus

ZH pontszám: [13, 15, 2, ...]



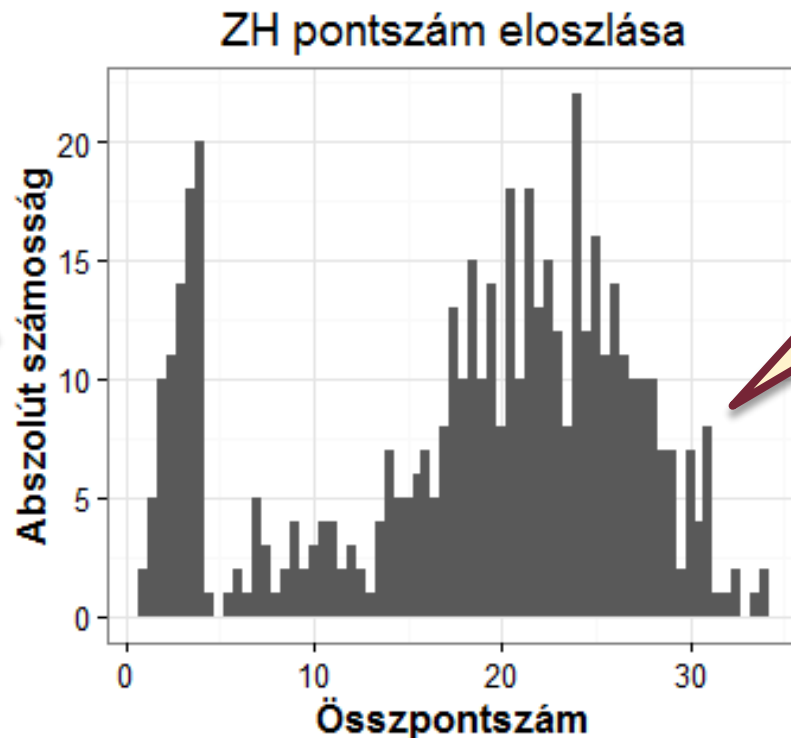
Kurzus: [G01, G03, G15, G17, ...]



Hisztogram

- Bemenő változó: ZH összpontszám
- Kérdés: hogyan alakultak a ZH pontszámok?

abszolút
gyakoriság!

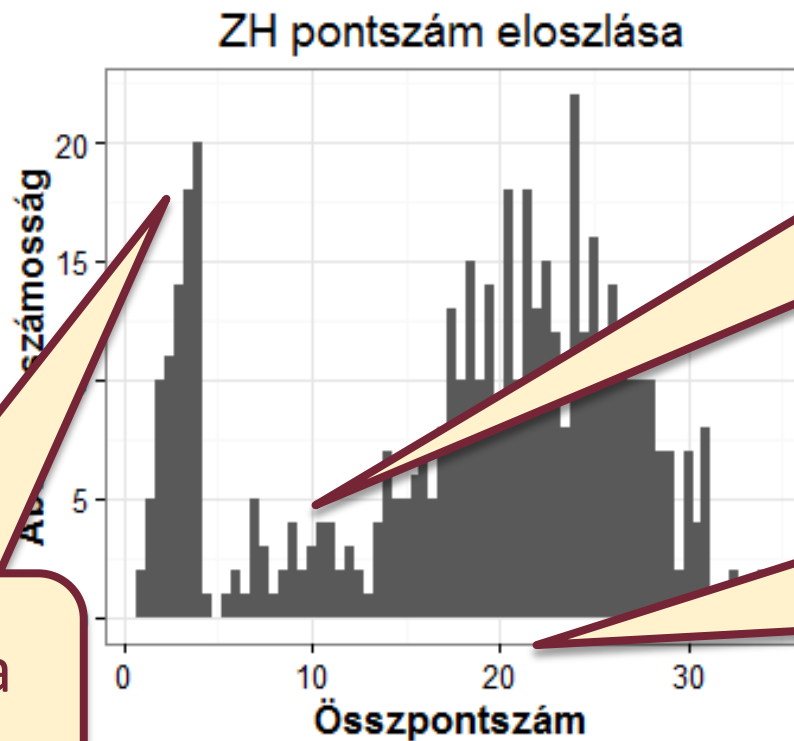


Oszlop-
magasság:
adott
intervallum
számossága

Tervezői döntés: mekkora legyen az intervallum hossza (bin size)?
Pl.: elég 1 pontos felbontással, vagy menjünk fél pontokig?

Hisztogram

- Bemenő változó: ZH összpontszám
- Kérdés: hogyan alakultak a ZH pontszámok?



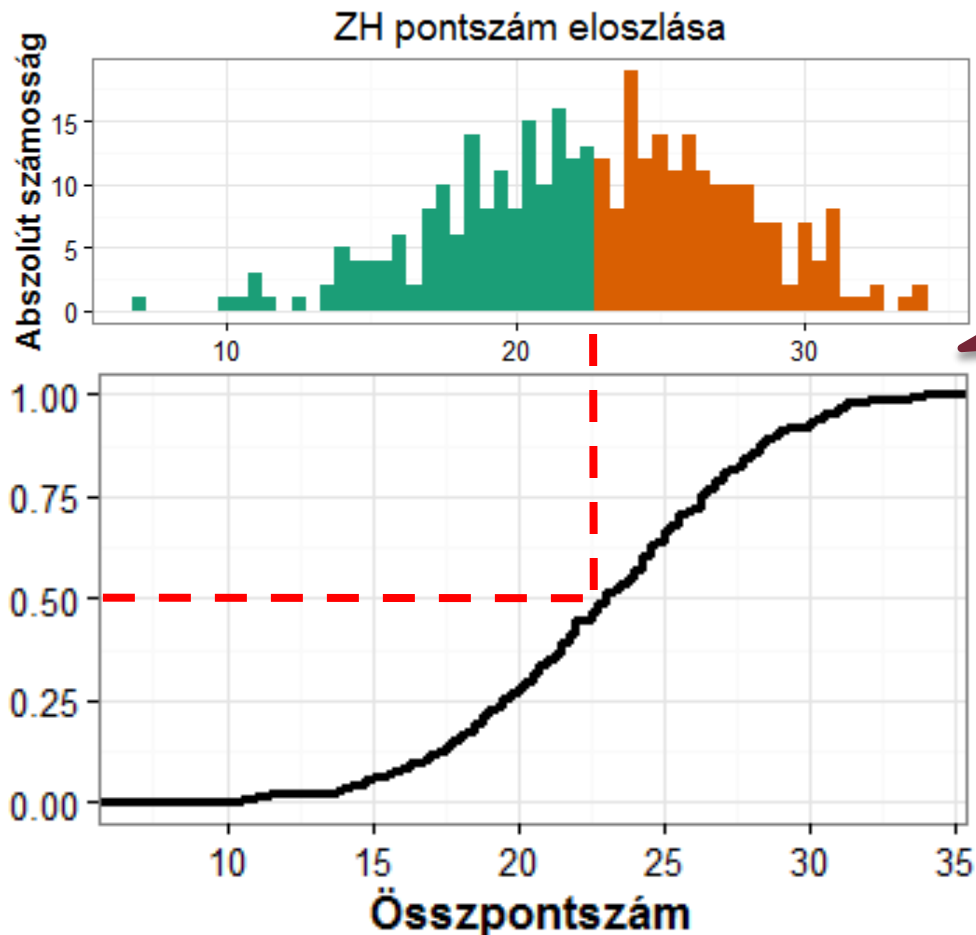
Sokan voltak a határon

Akik átmentek a beugrón, valószínűleg át is mentek

18 pont körül volt az átlag, 20 körül a medián

Egyszerű statisztikai jellemzés

- Hol van az adatok „közepe”?

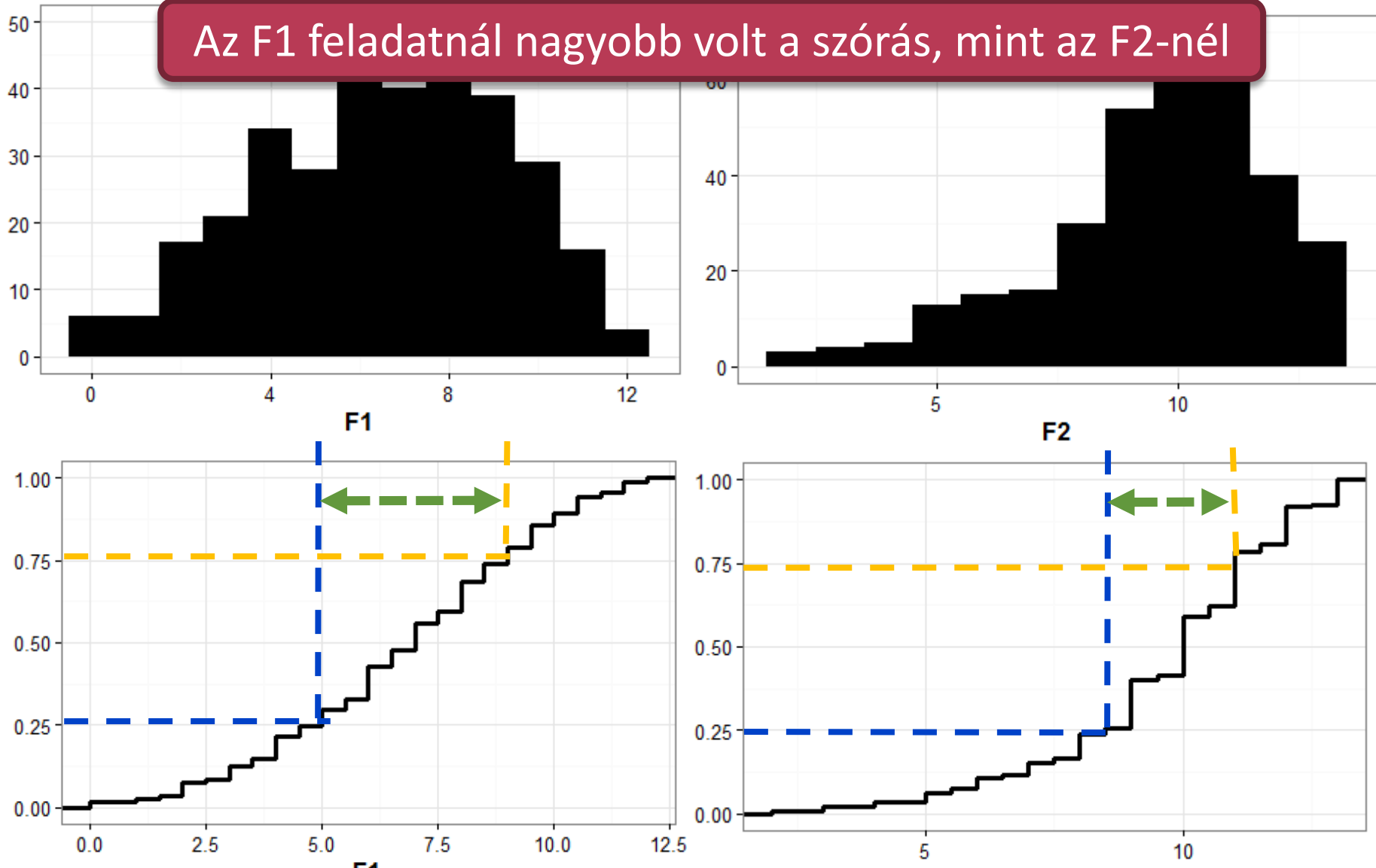


Az átlmentek
összpontszám
mediánja 23

Egyszerű statisztikai jellemzés

■ Mennyire „szórtak” az adatok?

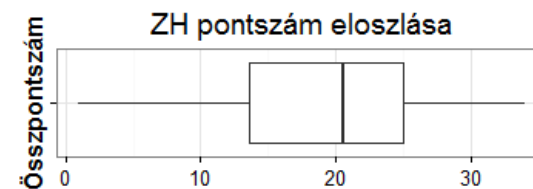
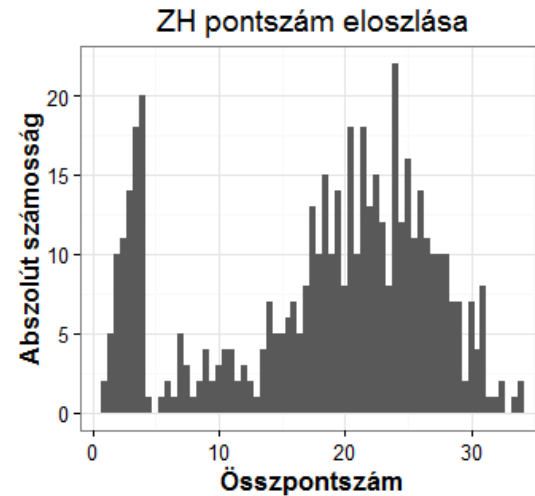
Az F1 feladatnál nagyobb volt a szórás, mint az F2-nél



Boxplot

- Bemenő változó: ZH összpontszám
- Kérdés: hogyan alakultak a ZH pontszámok úgy nagyjából?

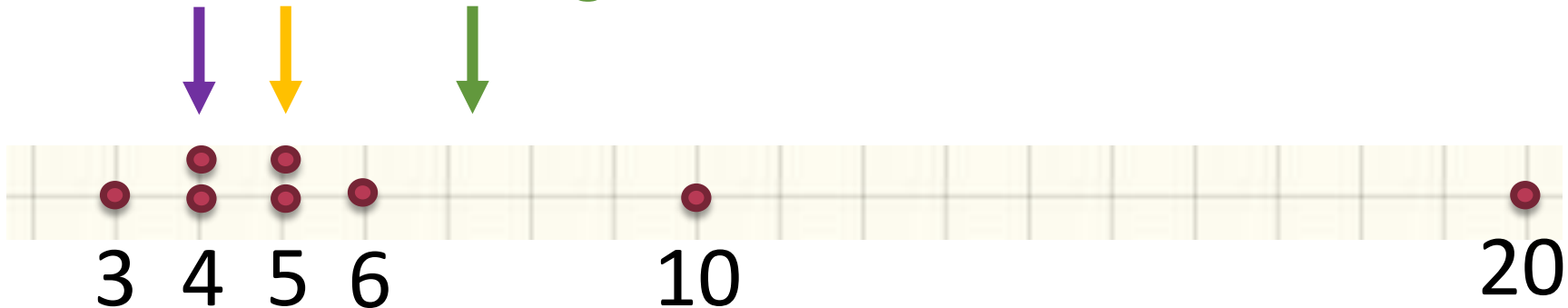
Egyfajta absztrakció itt is:
legyenek intervallumok,
felesleges minden pontot
kirajzolni



(Folytonos) megfigyelések jellemzése

- A „központ” jellemzése
 - Átlag, **medián**, módusz
 - {3, 4, 4, 5, 5, 6, 10, 20}
 - Átlag: ~ 7.125
 - Medián: 5
 - Módusz: 4 és 5

módusz medián átlag



(Folytonos) megfigyelések jellemzése

Ha az értékeket növekvően sorba rendezzük, akkor a középső adat az adathalmaz **mediánja**. Ha nincs középső adat (páros számú érték esetén), akkor a **medián** a két középső érték átlaga (számtani közepe).

A **módusz** az adathalmazban legtöbbször előforduló érték. Ez nem feltétlenül egyértelmű, ilyenkor több módusról beszélünk.

Terjedelem jellemzése: percentilisek

Az n -edik **percentil**nél az értékek $n\%$ -a nem nagyobb.

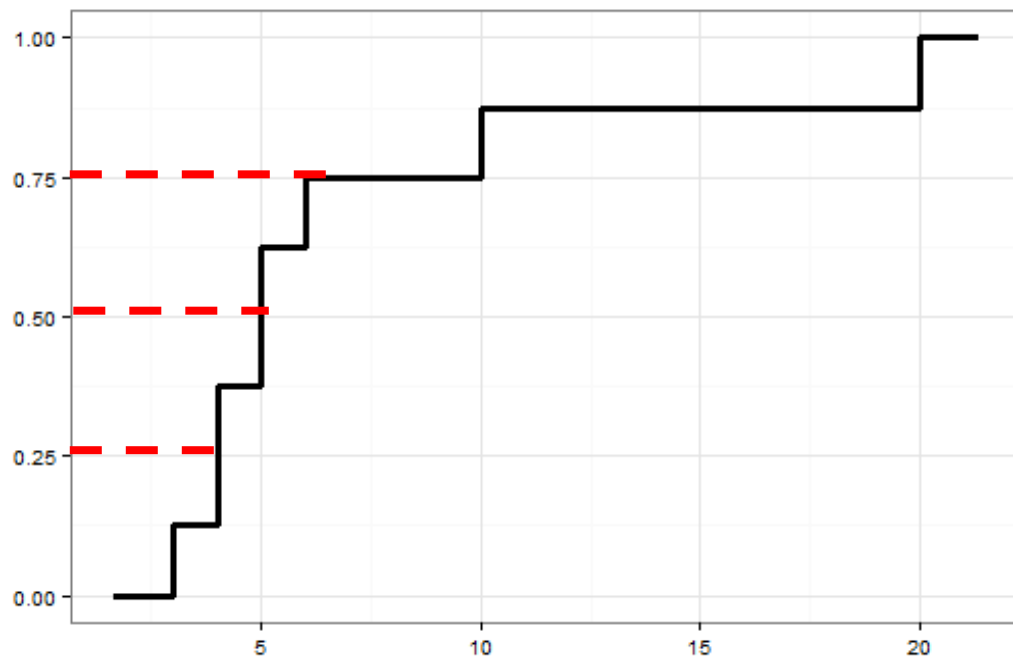
■ Percentilis

○ {3, 4, 4, 5, 5, 6, 10, 20}

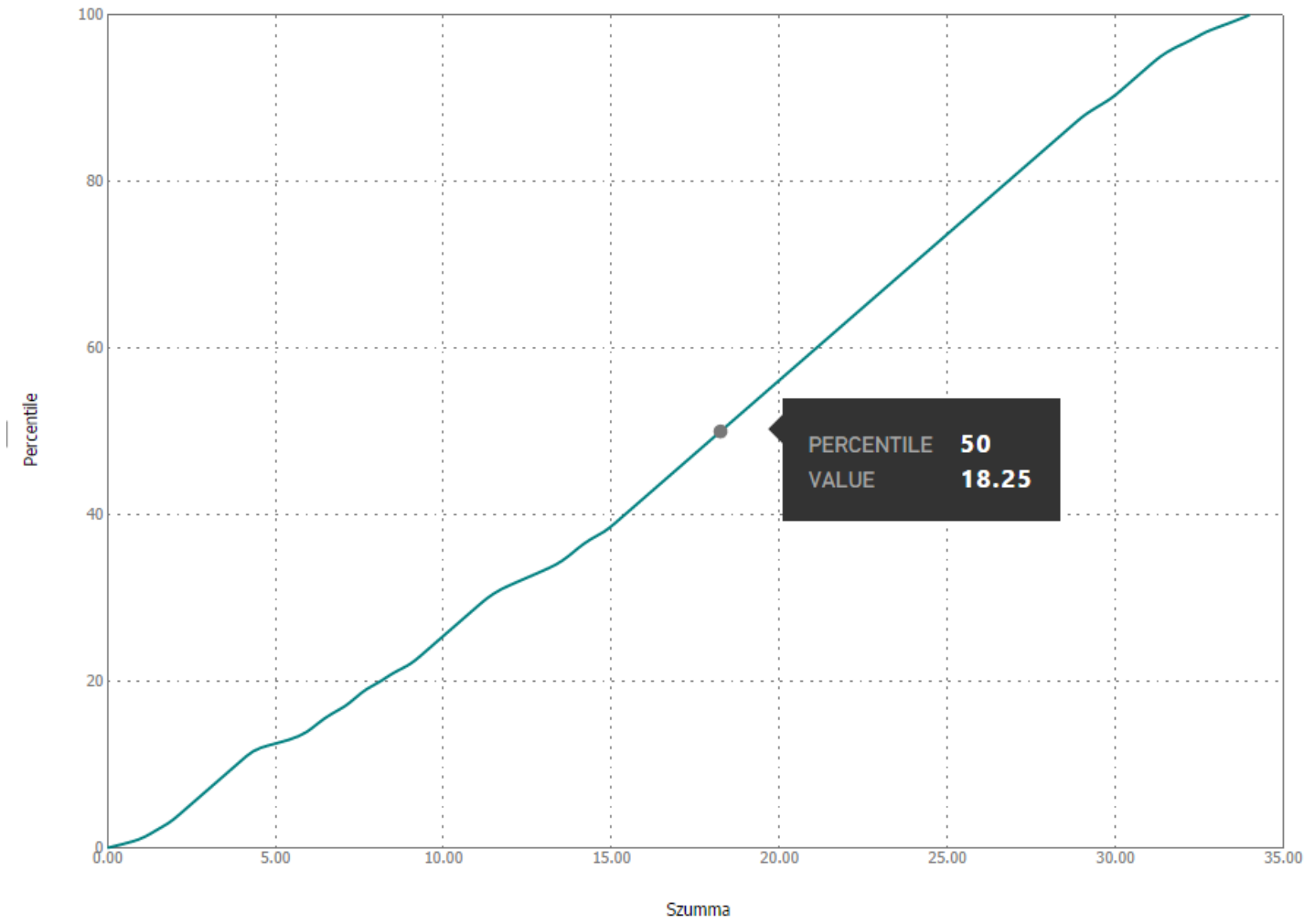
- 50. percentilis: 5
- 25. percentilis: 4
- 75. percentilis: 6

■ Kvartilis

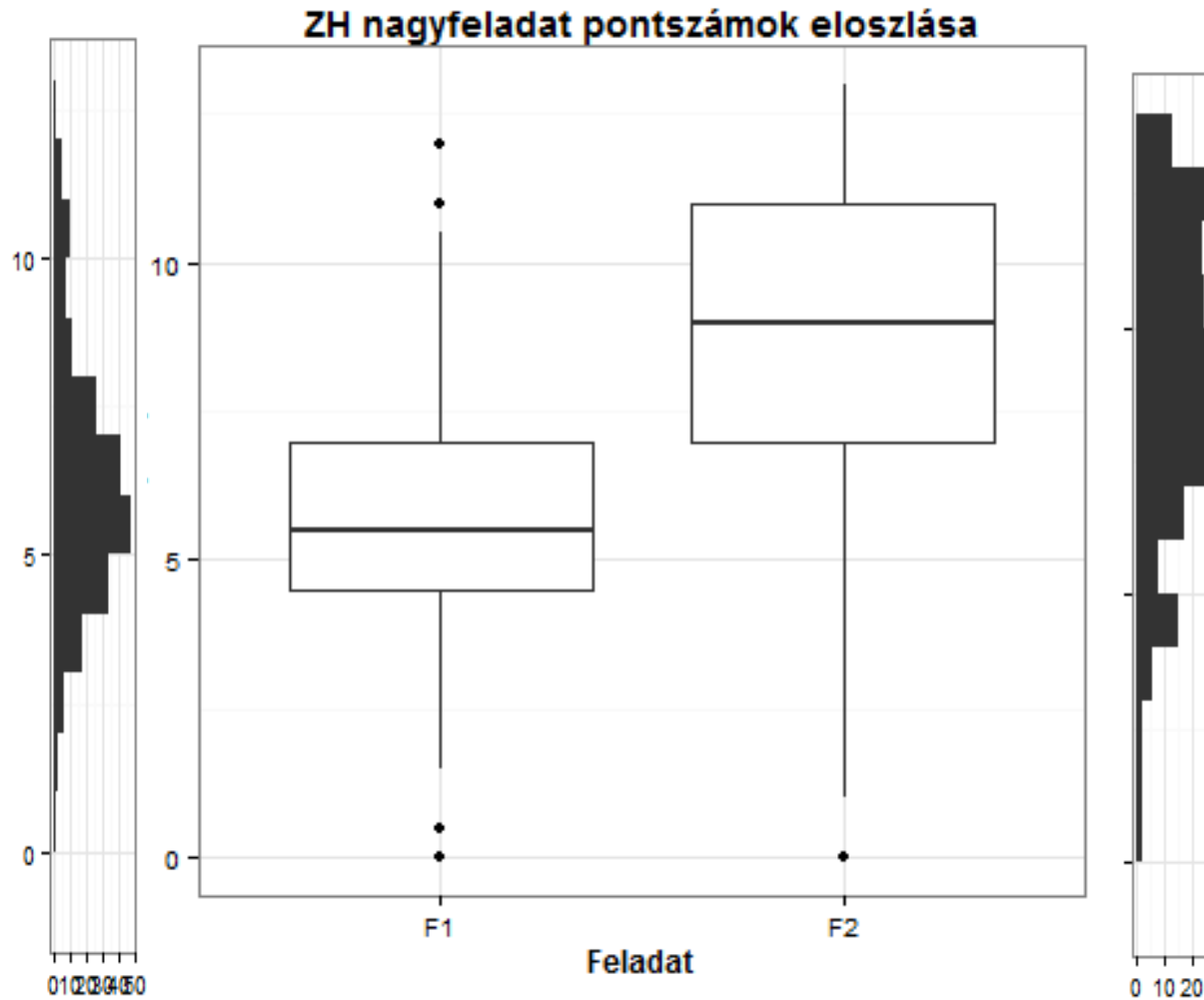
- Q1: 25. percentilis
- Q3: 75. percentilis
- **Q2: medián**



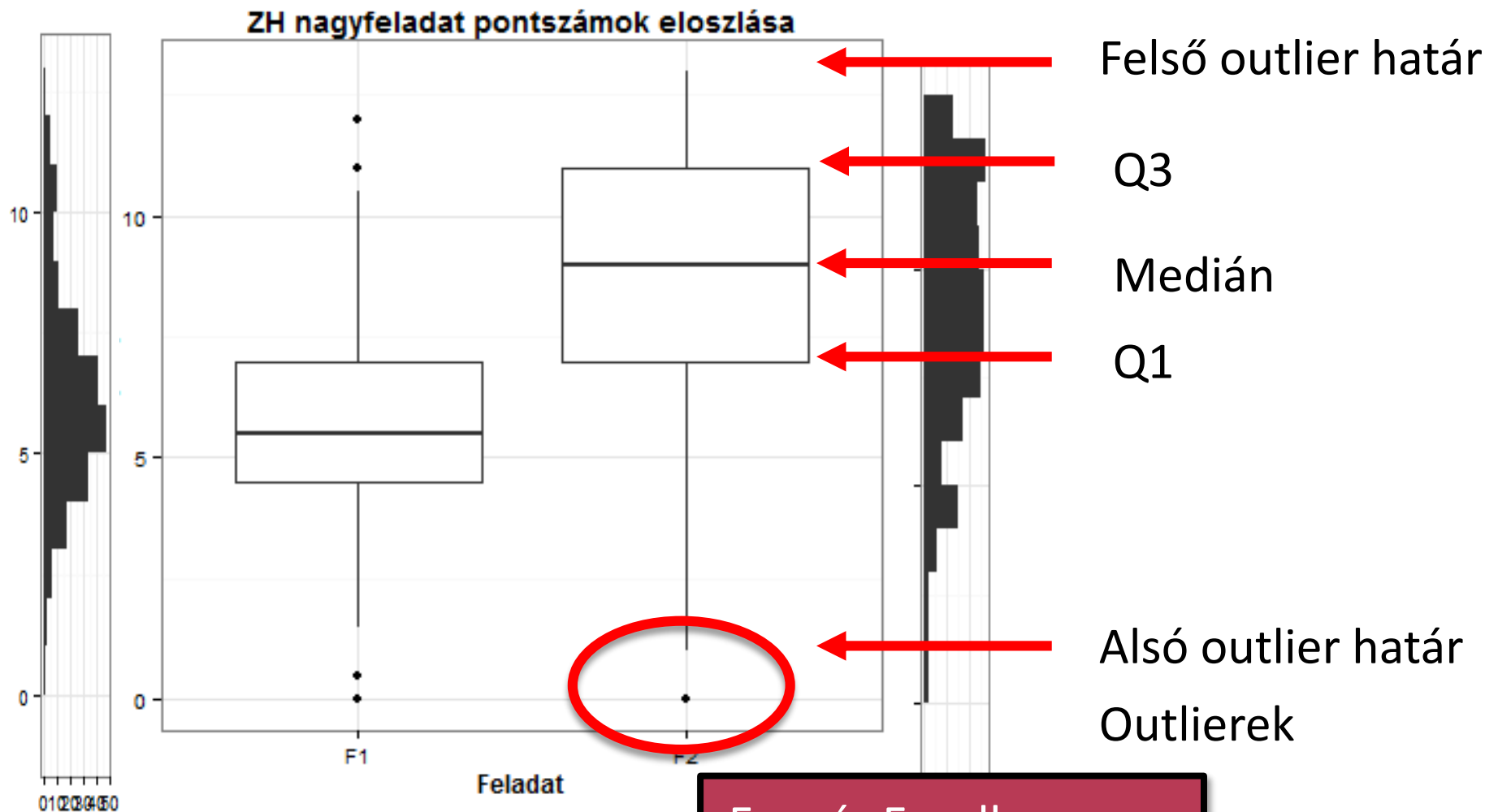
Példa: percentilis ábrázolás



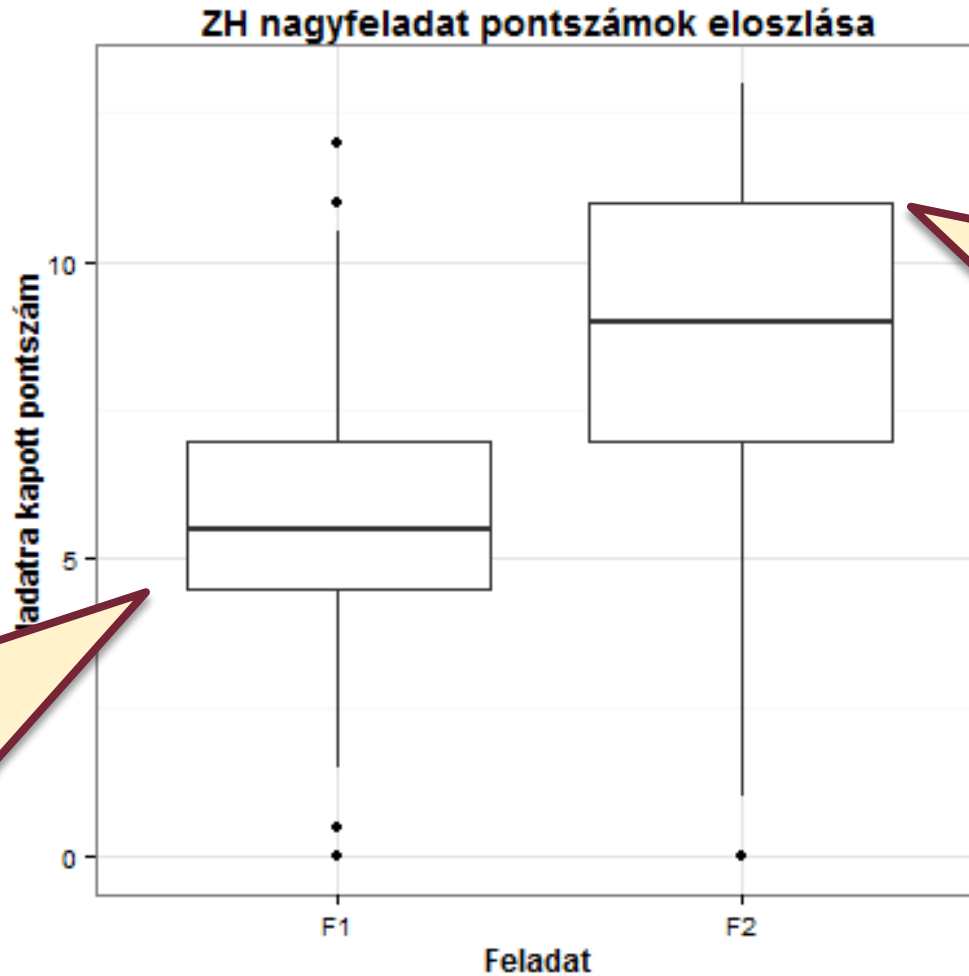
Boxplot (Box and whisker plot)



Boxplot (Box and whisker plot)



Boxplot (Box and whisker plot)

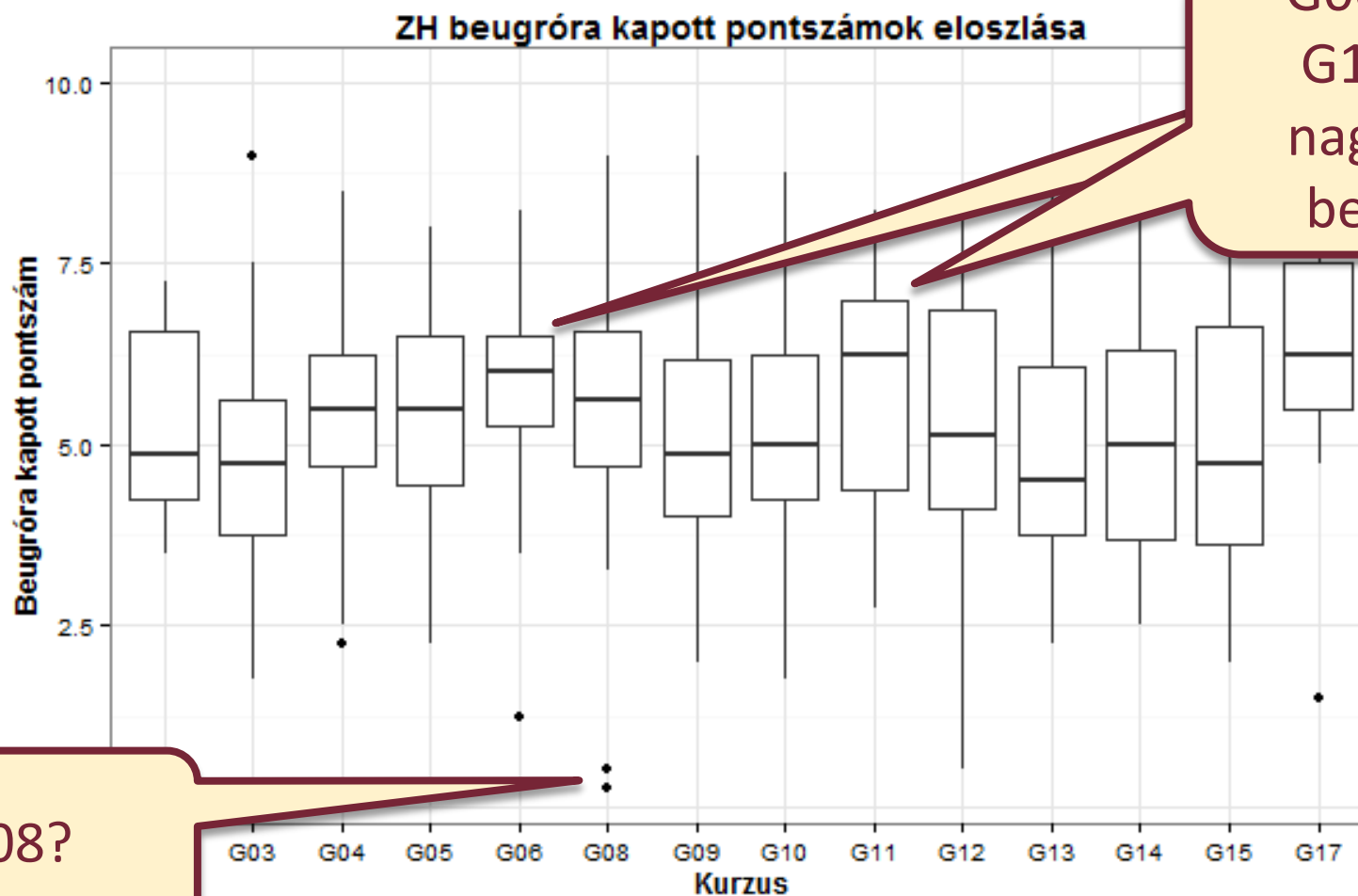


Az F1 pontszámok 50%-a 4.5 és 7.5 között volt

F2-re általában több pontot kaptak, mint F1-re

Boxplot (Box and whisker plot)

- Melyik csoportban hogyan sikerültek a beugrók?

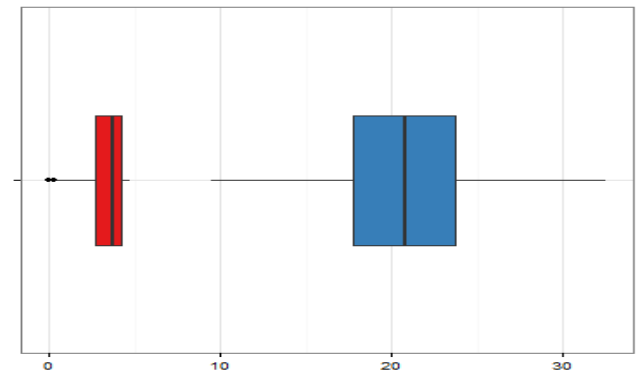
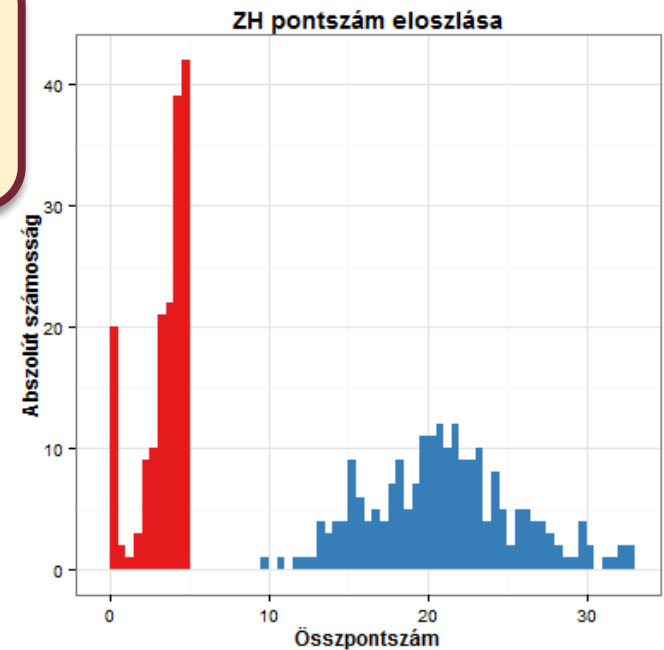
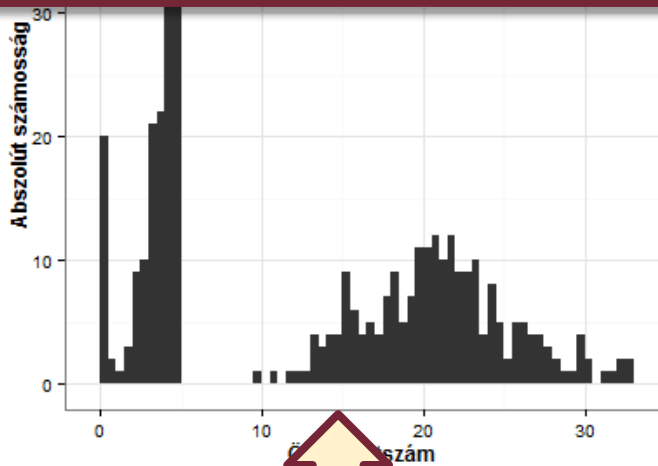


G06, G11,
G17-ben
nagyon jó
beugrók

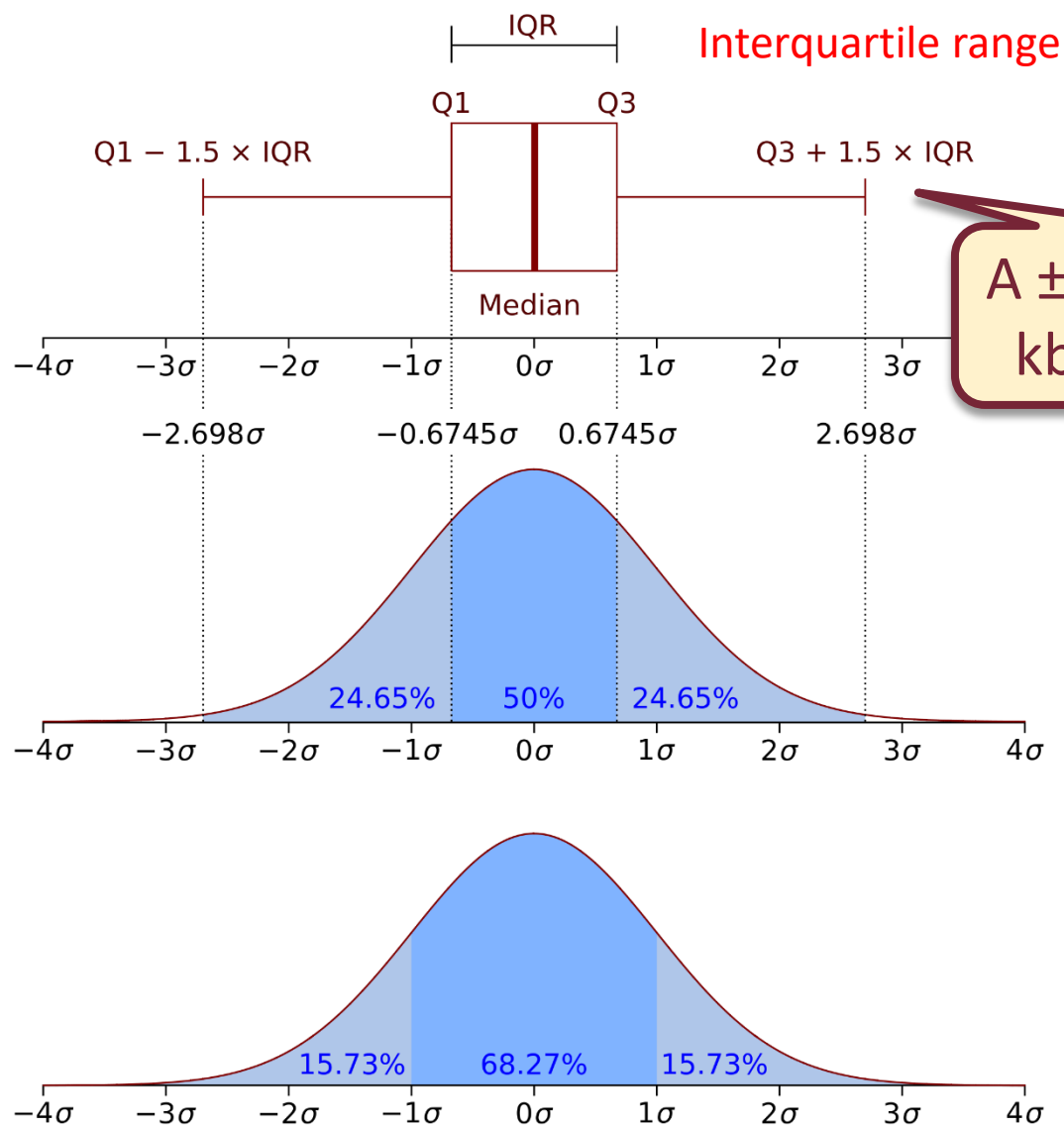
G08?

Boxplot (Box and whisker plot)

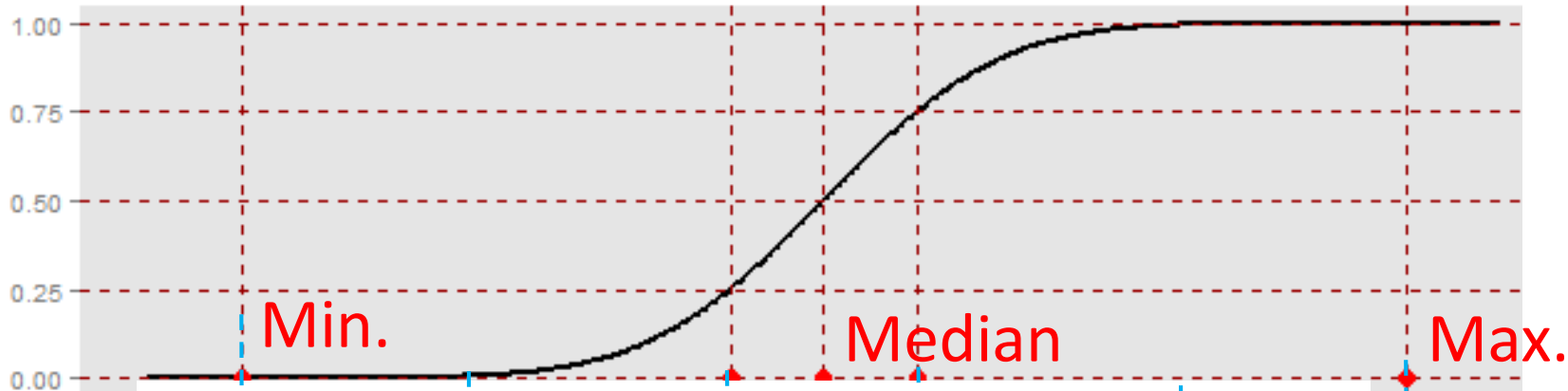
Absztrakció: a boxplottal fontos információt is veszíthetünk!



Boxplot (Box and whisker plot)

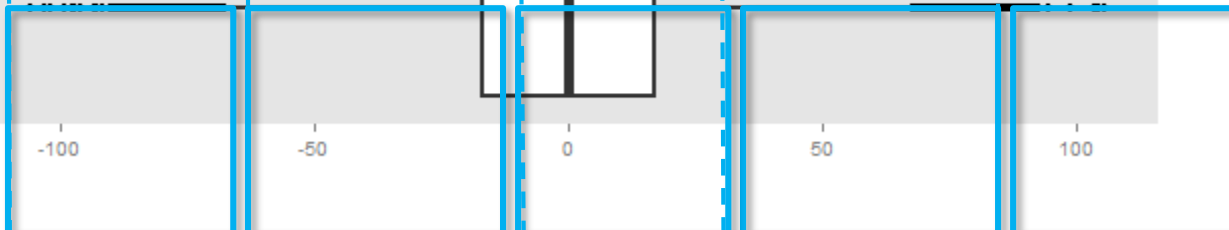


Boxplot: kvalitatív jellemzés



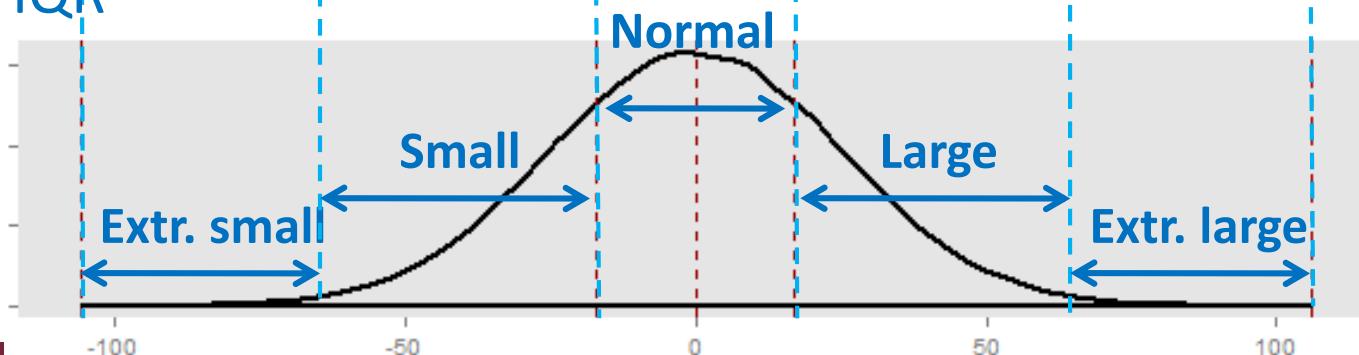
Extr. small Small Normal Large Extr. large

Kvalitatív
tartomány



Q1 – 1.5 IQR

Kvantitatív
tartomány



Miért medián, miért nem átlag?

■ Alaphalmaz

○ 1000 adatpont $\sim U(1, 5)$ egyenletes eloszlás

- *átlag = medián = 3 ms*



3ms \pm 2 ms



Új medián: `sort(resp. times)[501] = 3.02 ms`

Vál. medián



Vál. átlag



Új átlag: $(2 * 10^4 + 3 * 10^3) / 1001 = 25 \text{ ms!}$

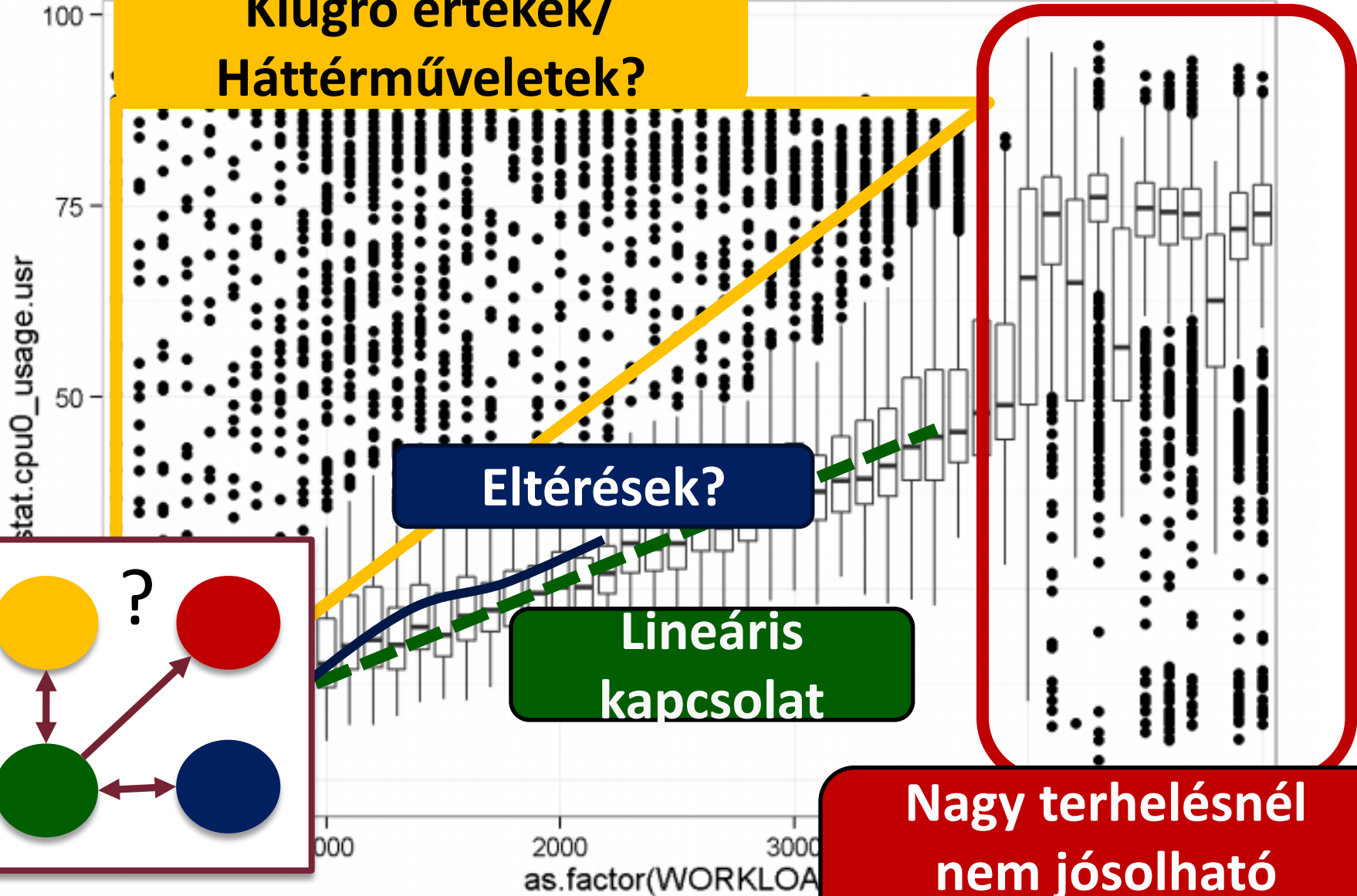
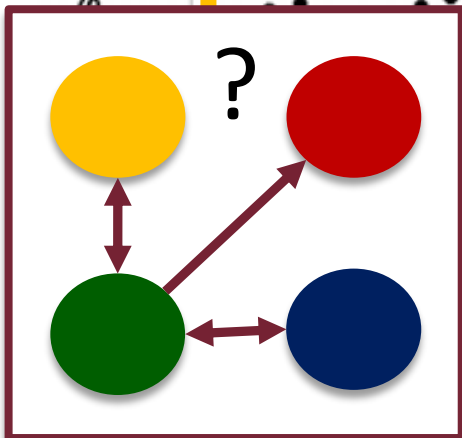
Példa: terhelés vs. kihasználtság

Kiugró értékek/
Háttérműveletek?

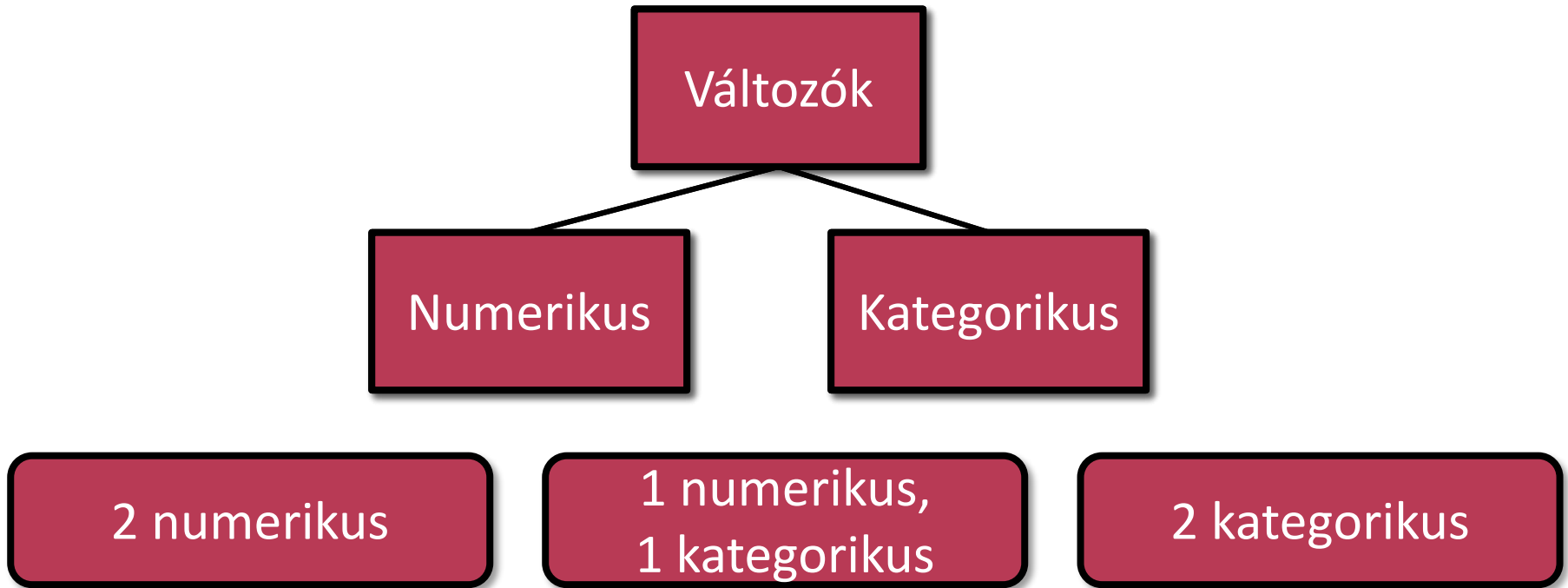
Eltérések?

Lineáris
kapcsolat

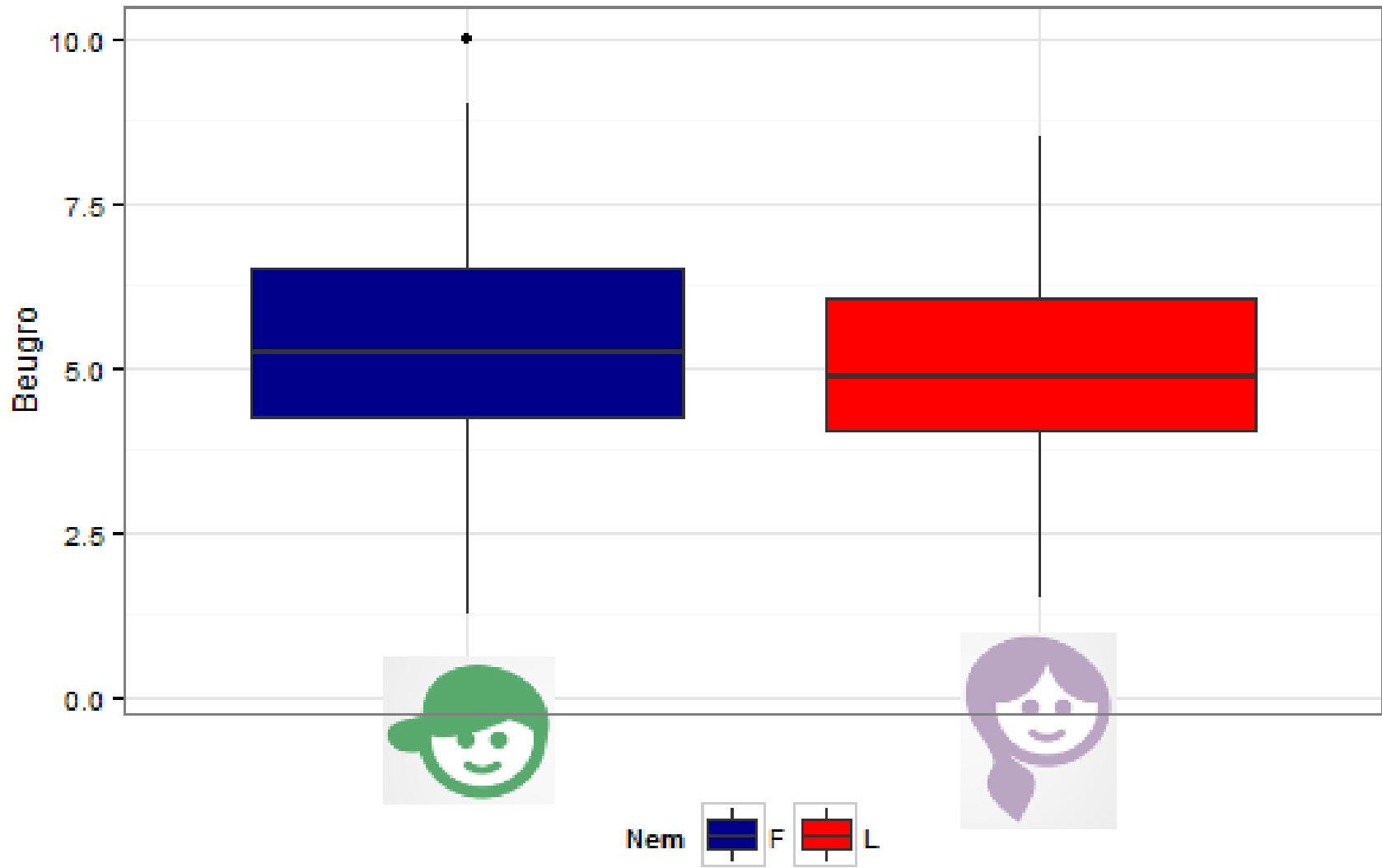
Nagy terhelésnél
nem jósolható
működés



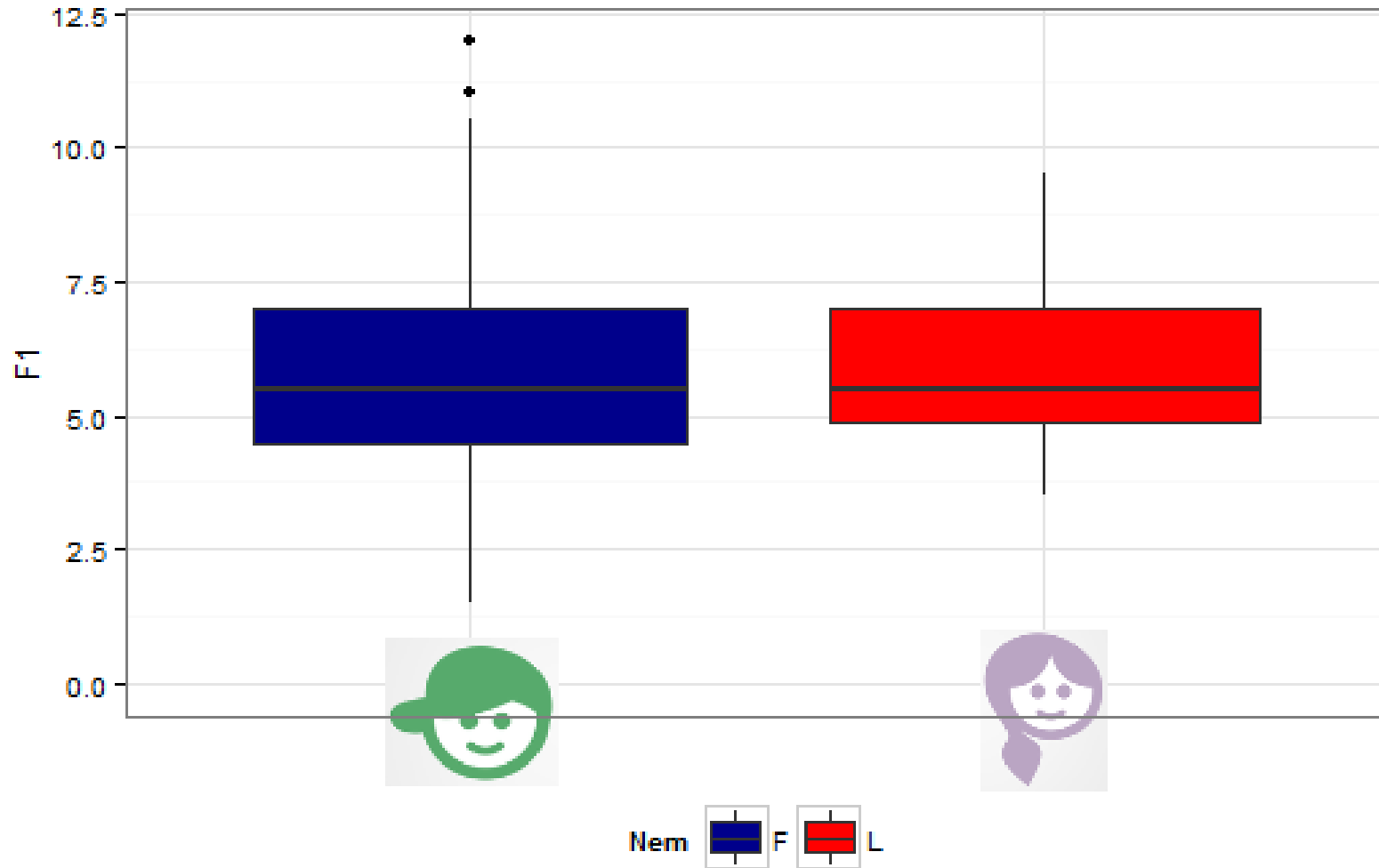
2 változó kapcsolata



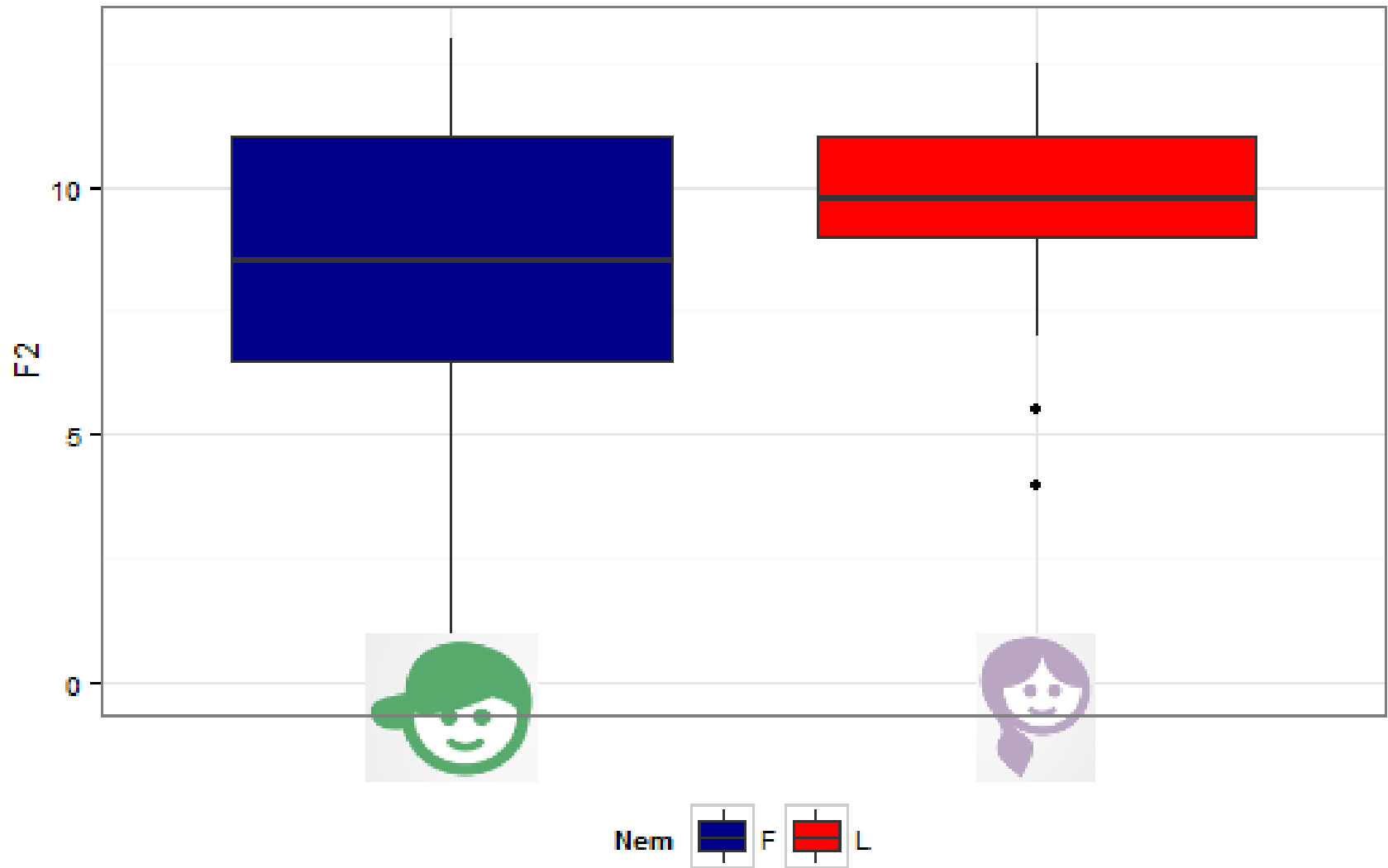
Numerikus kategóriánként



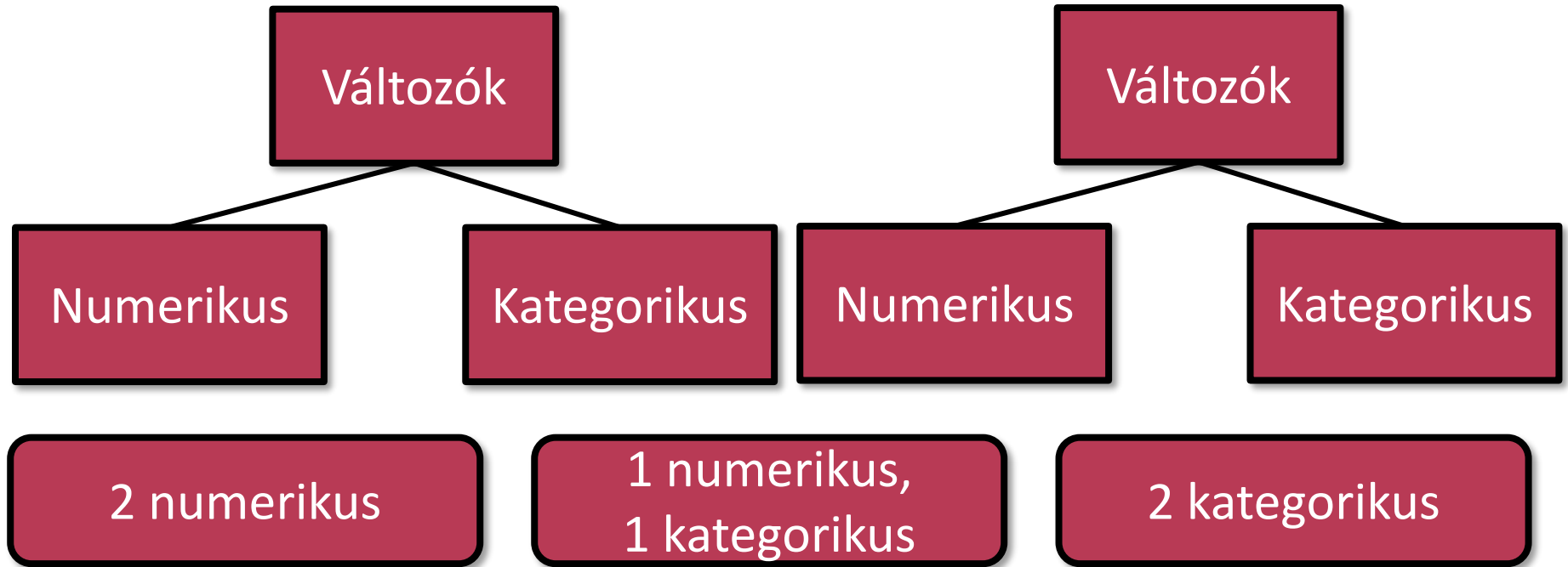
Numerikus kategóriánként



Numerikus kategóriánként



2 változó kapcsolata

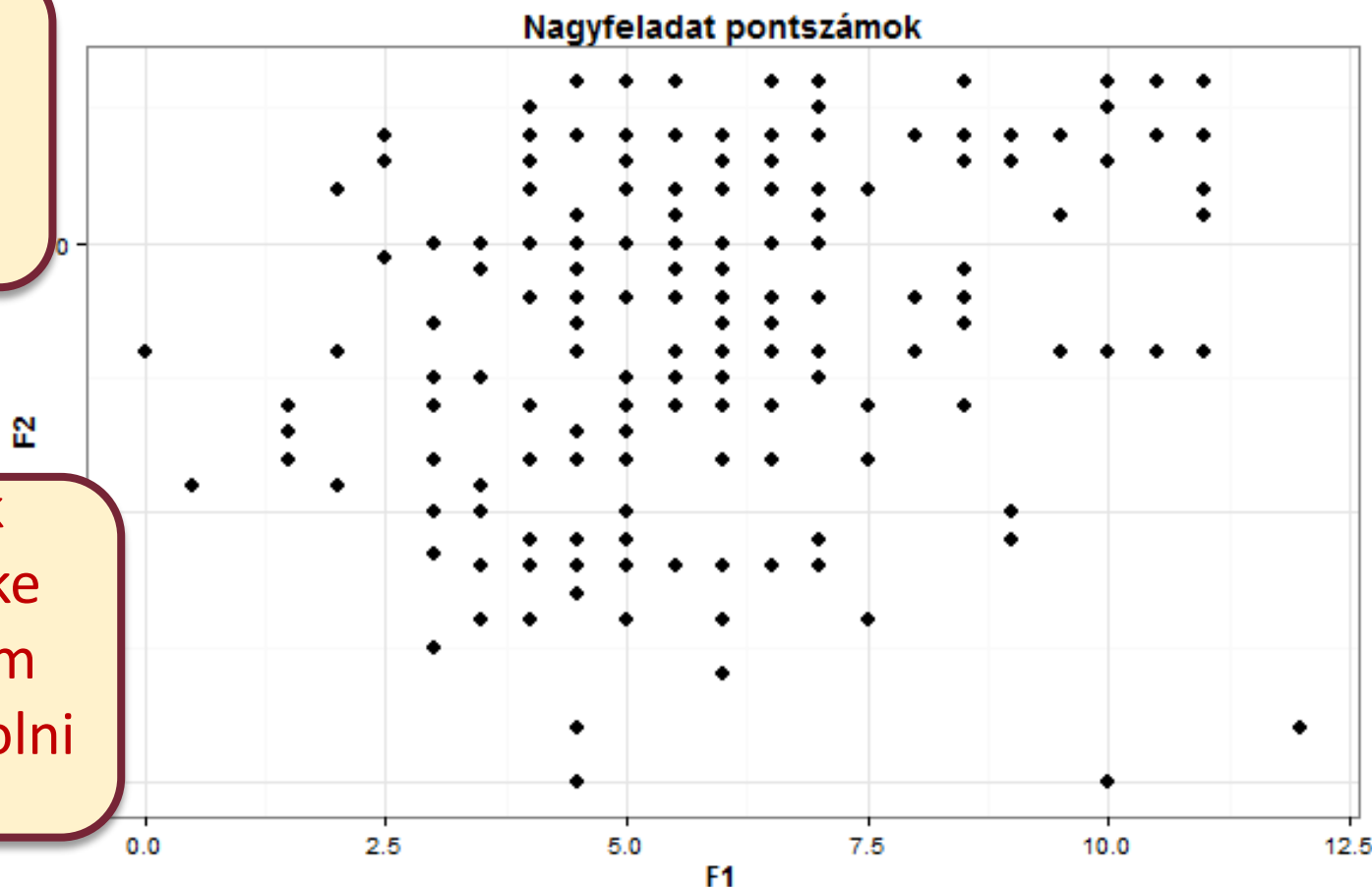


Pont – pont diagram (scatterplot)

- Bemenő változó: nagyfeladatokra kapott pontok
- Kérdés: hogyan viszonyulnak egymáshoz?

Együttesen előforduló pontpárokat vizualizálunk

Ha az egyik változó értéke hiányzik, nem tudjuk felrajzolni



Pont – pont diagram (scatterplot)

- Bemenő változó: nagyfeladatokra kapott pontok
- Kérdés: hogyan viszonyulnak egymáshoz?

Nem biztos,
hogy akinek
megy az F1,
megy az F2 is

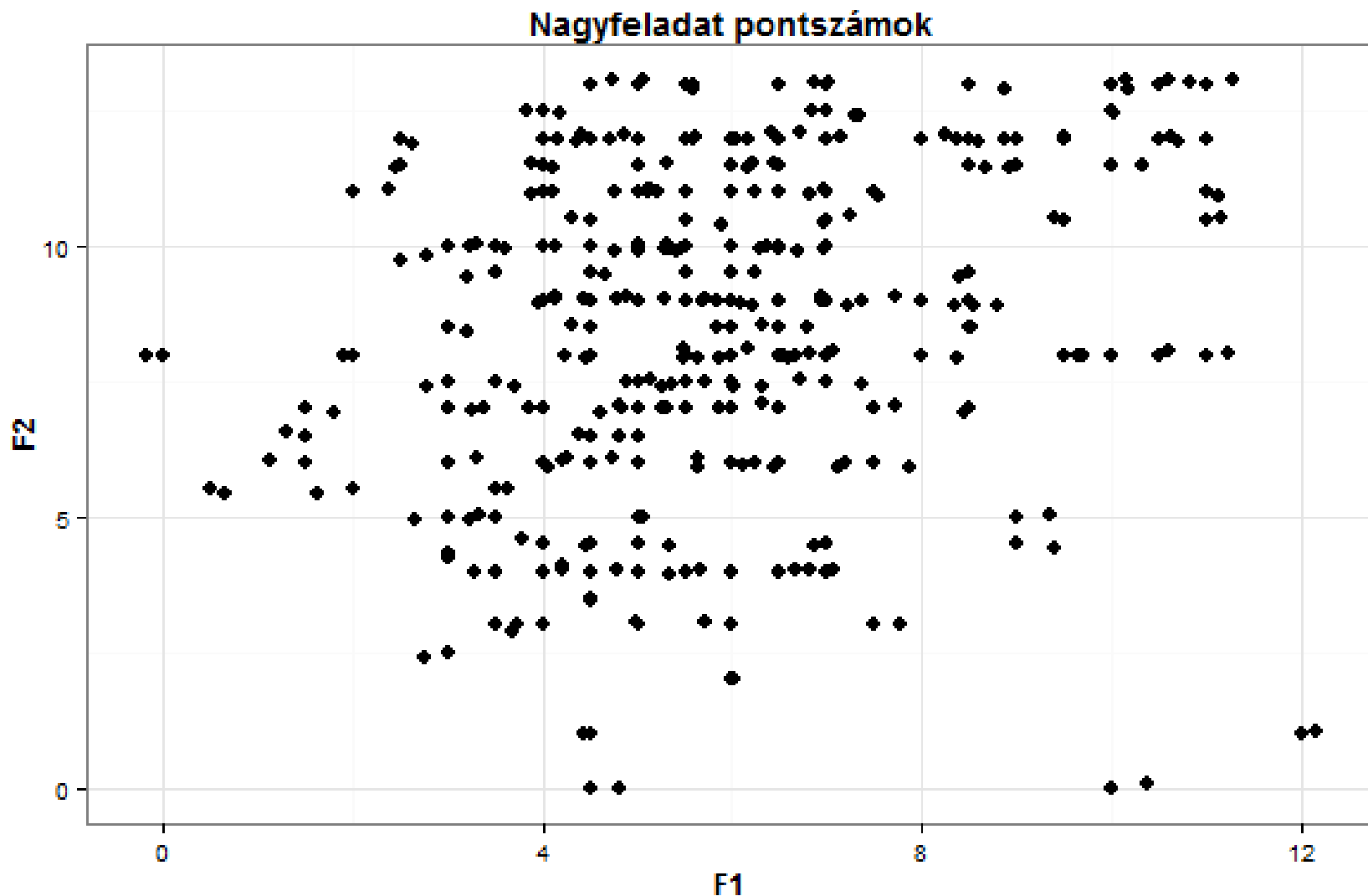


Hogyan kezeljük a takarásokat?

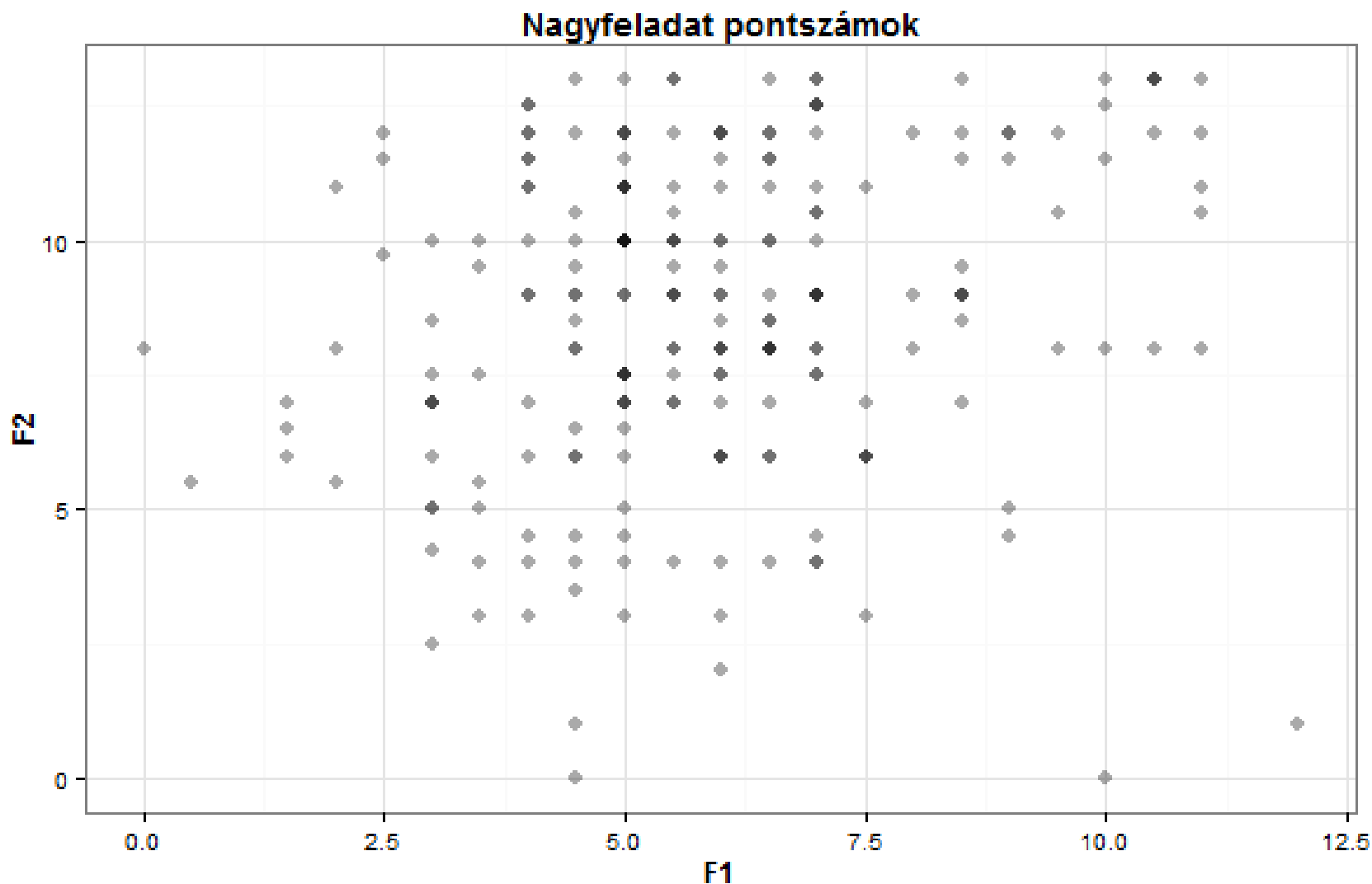
Overplotting



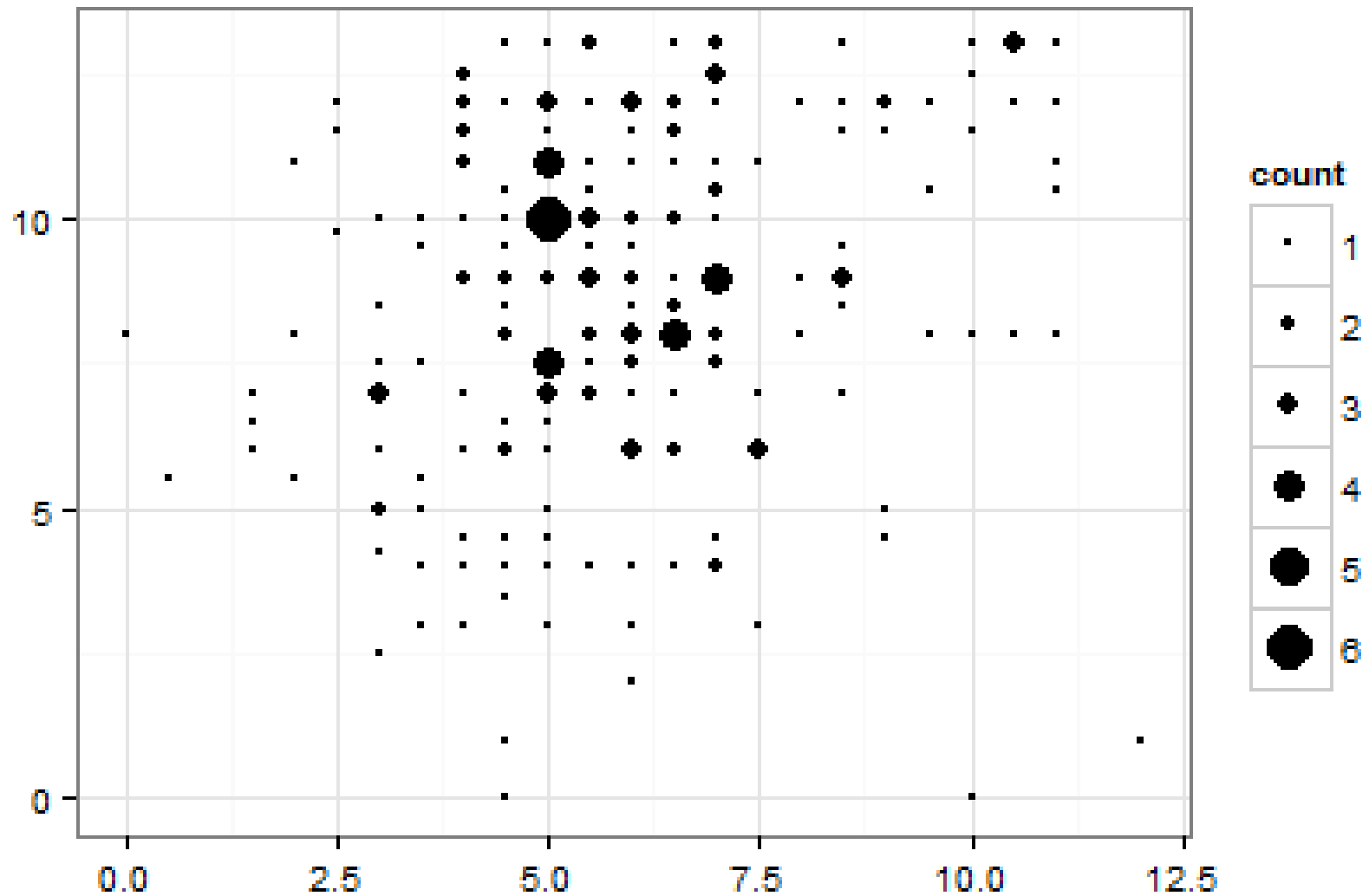
Overplotting megoldások 1: jitter



Overplotting megoldások 2: átlátszóság



Overplotting megoldások 3: méret

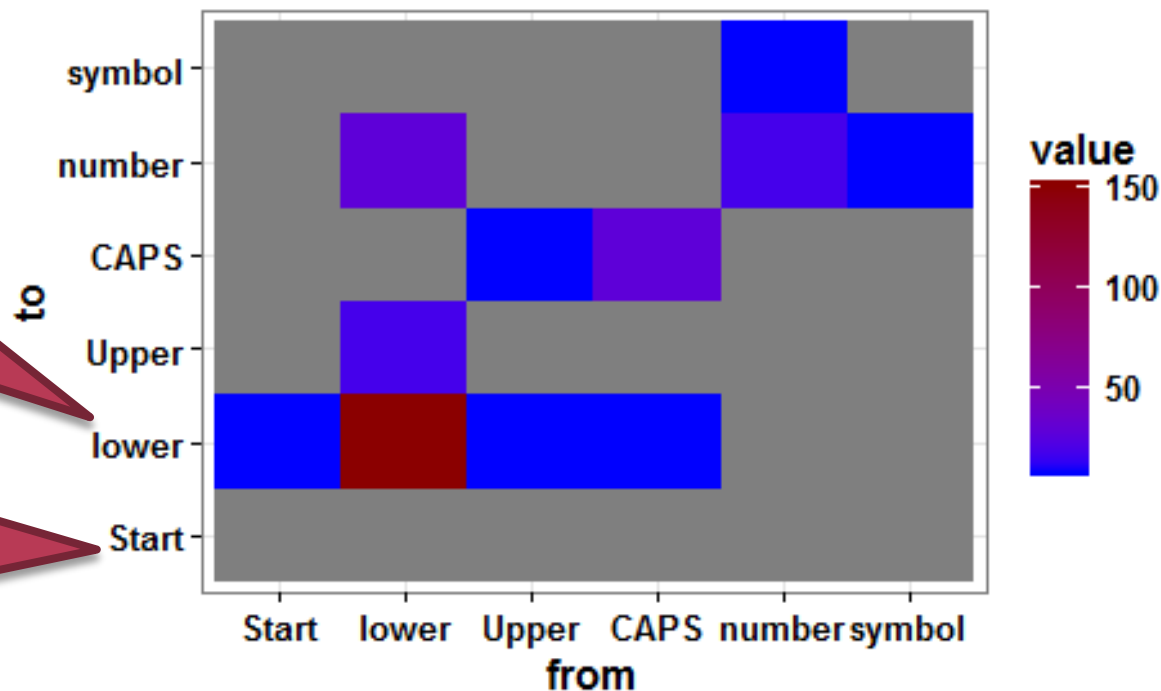
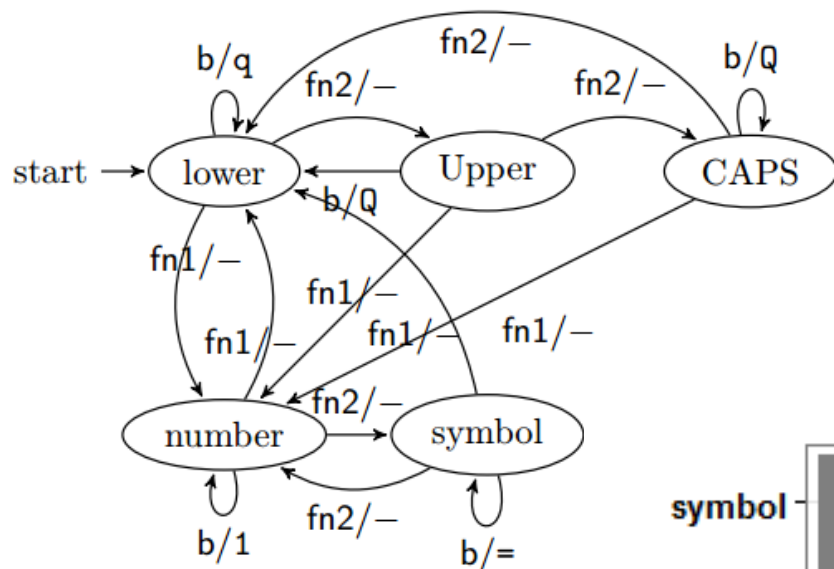


SOK VÁLTOZÓ

≥ 3 változó

- A grafikai objektumok attribútumait változtatom
 - Szín
 - Méret
 - Textúra
 - Hely – ez triviálisnak tűnik, de a treemapnél van jelentősége
- Pl. heatmap, treemap

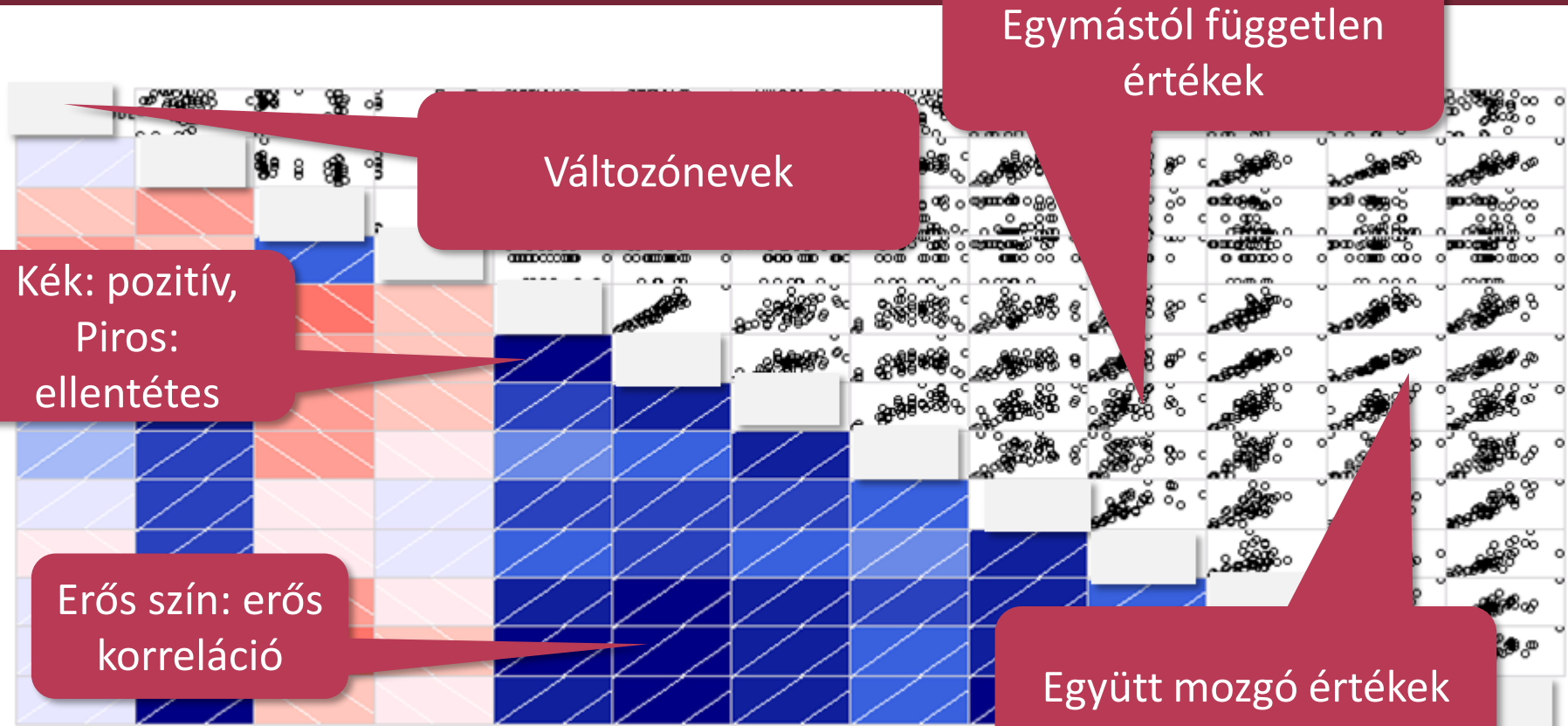
Heatmap: lefutási statisztikák



Inkább csak sima szöveget írunk

A Startban mindig csak kezdünk, oda nem jutunk vissza

Kitekintés: több érték páronkénti korrelációja



R statisztikai szoftver „corrgram” csomagjával előállítva.

Korreláció (ld. Valószínűségi számítás):

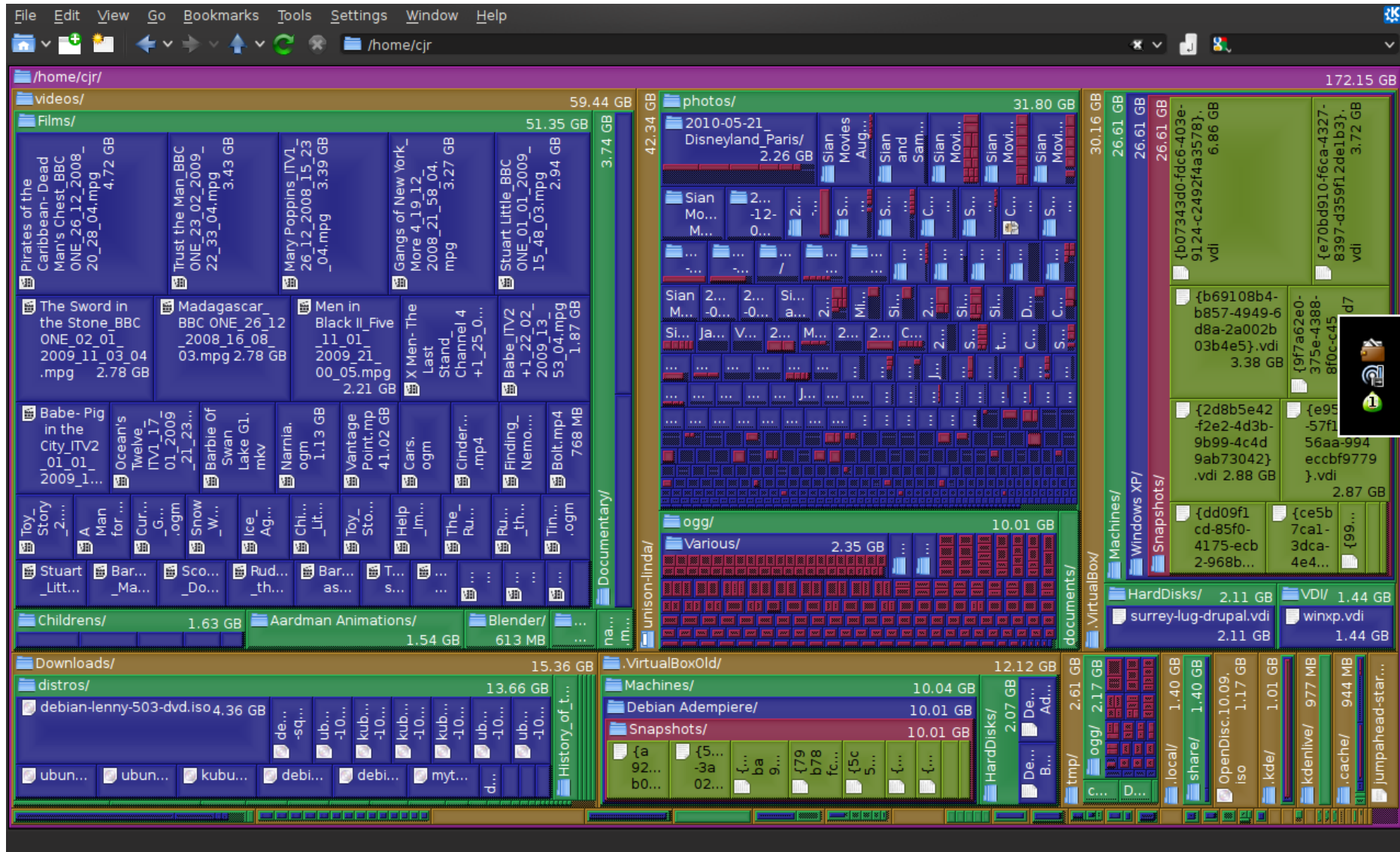
két érték közti lineáris kapcsolat erőssége és iránya

Átló felett: **scatterplot mátrix**

Cél: együtt mozgó értékek kiszűrése, **kiugró értékek (outlierek)** azonosítása.

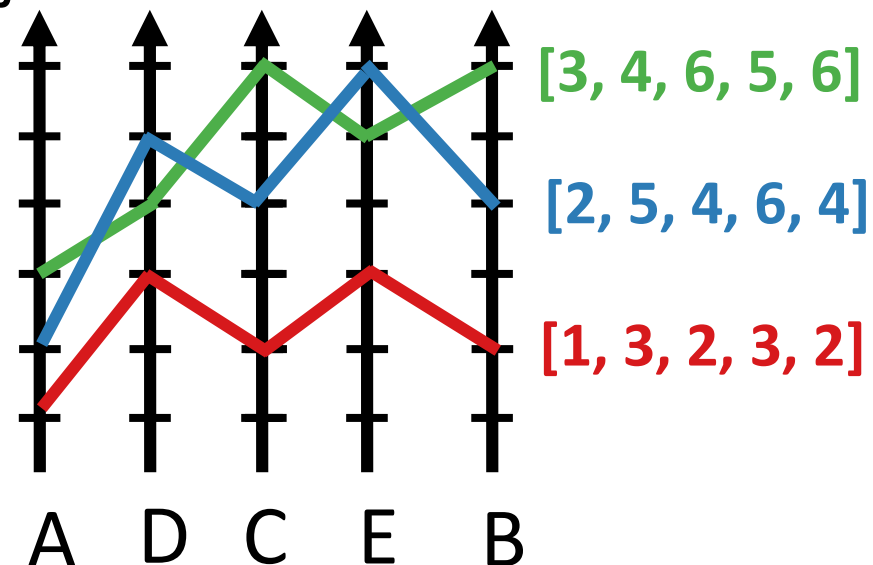
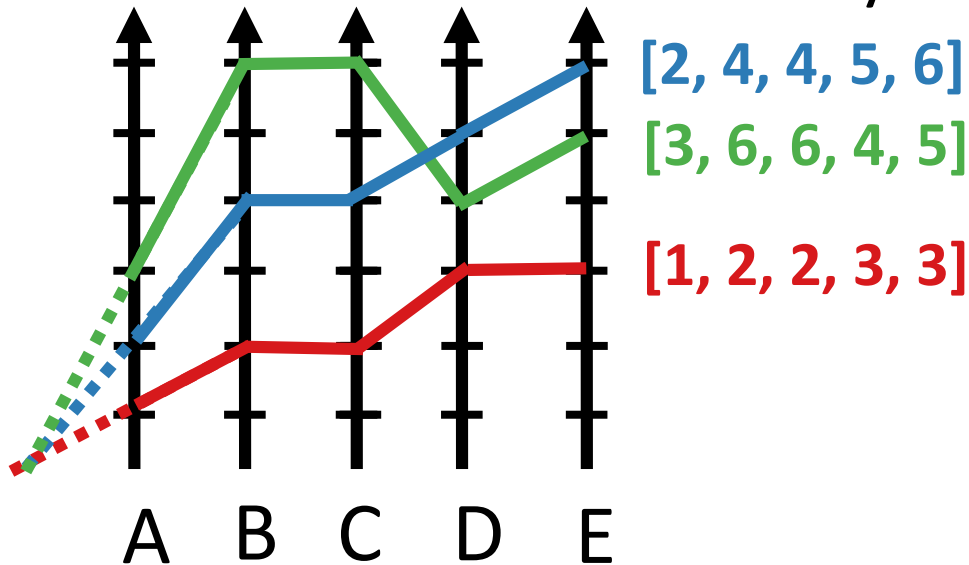
→ Mik a terhelés/előrejelzés szempontjából lényeges változók?

Treemap: állományrendszer



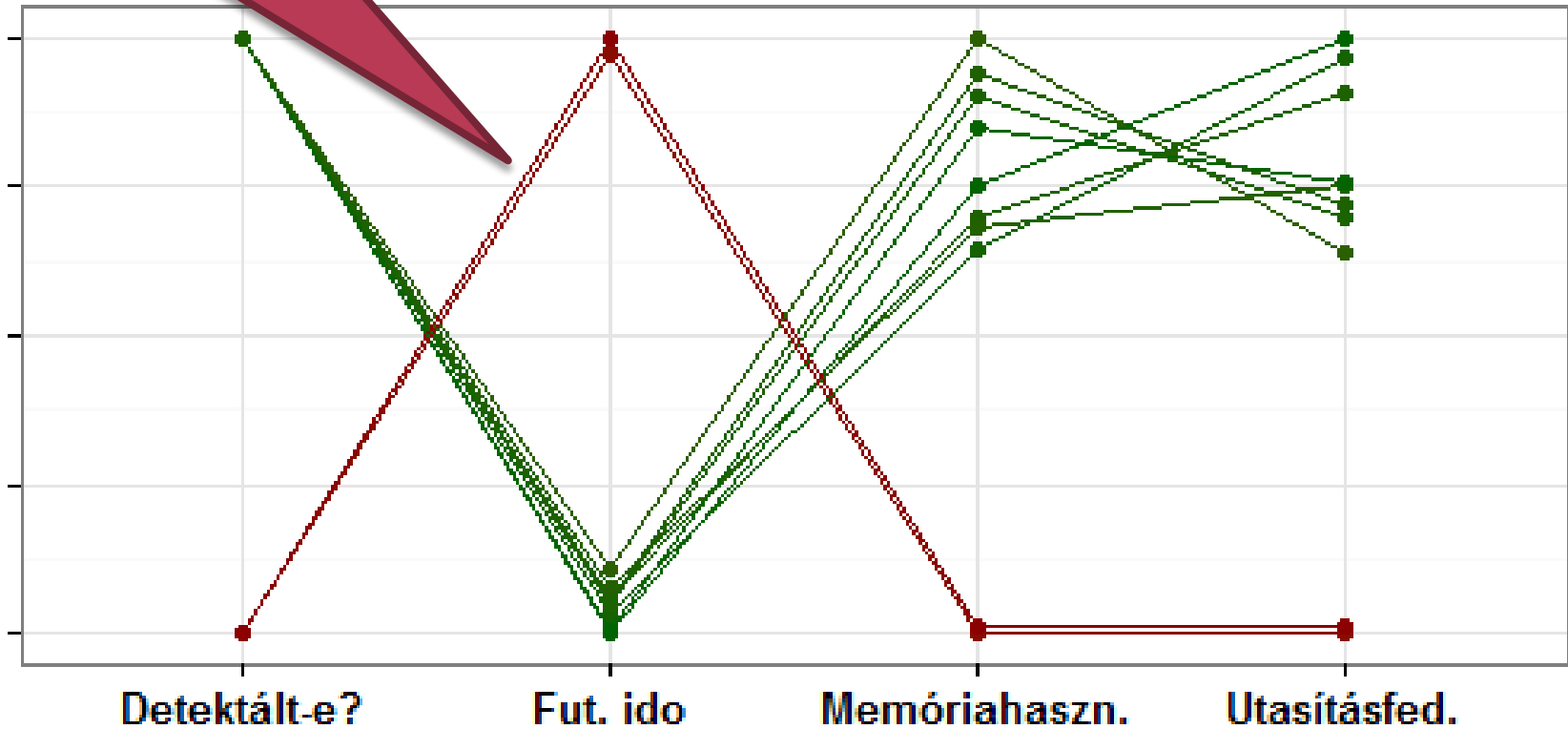
Párhuzamos koordináták

- Tengelyek: dimenziók/koordináták
 - tetszőleges számú
 - tetszőleges skála
- Egy vonal egy mérés (darabszám?)
- Kompakt és skálázható
- Koordináta sorrend befolyásolja a kiértékelést



Párhuzamos koordináták: tesztesetek elemzése

1 teszteset 1 törött vonal

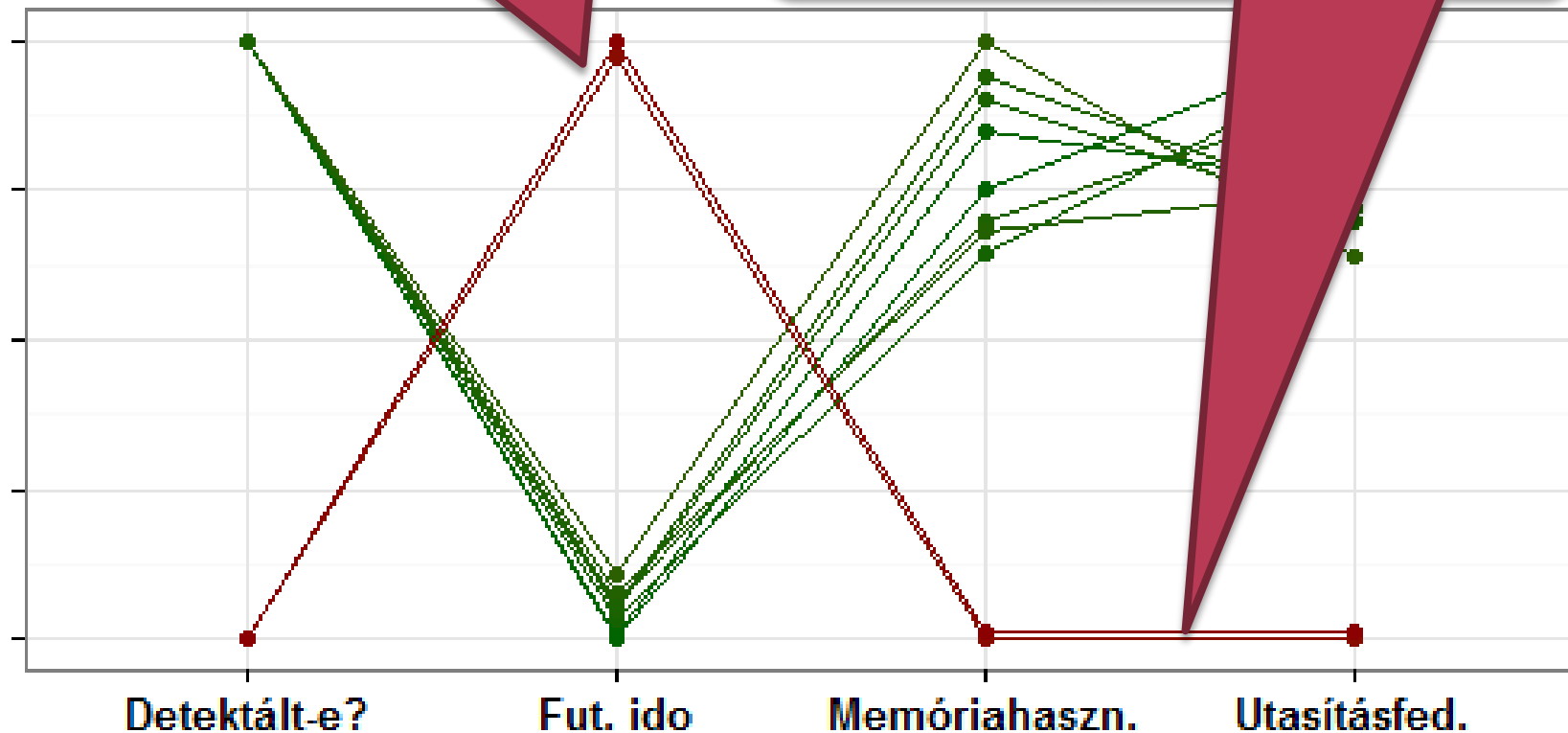


A változók az x tengelyen jelennek meg

Párhuzamos koordináták: tesztesetek elemzése

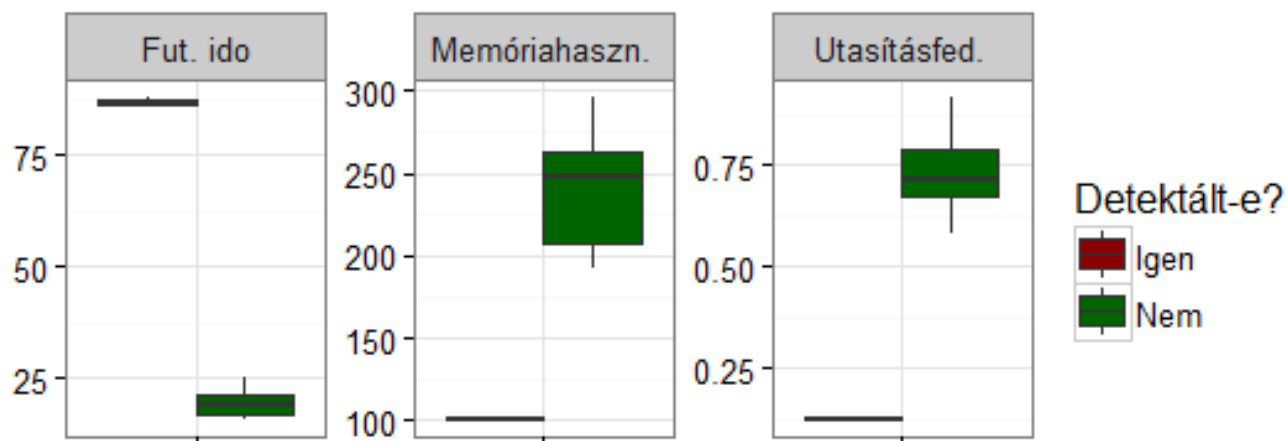
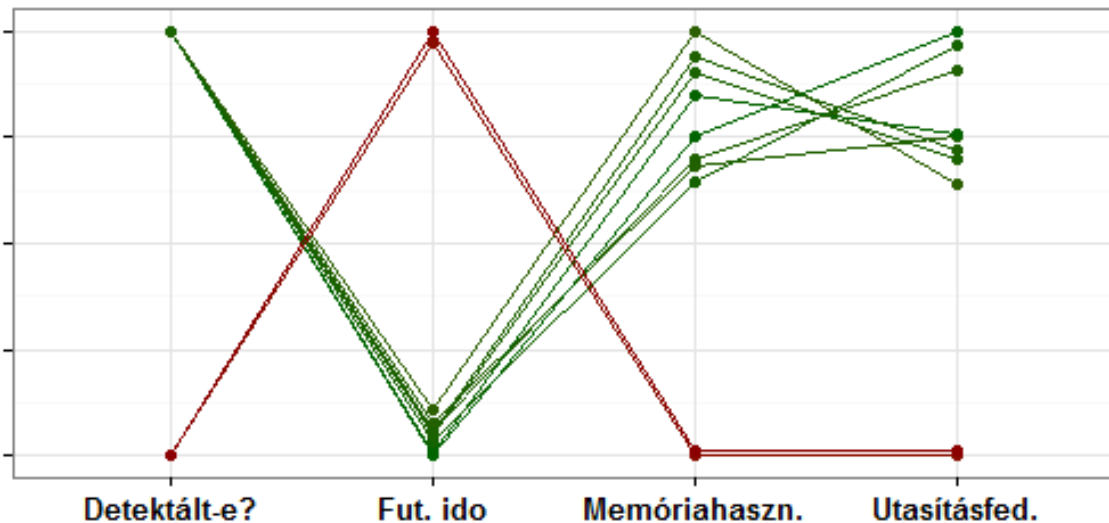
Timeout?

A hibát detektálók az érdemi számításig valószínűleg el sem jutnak

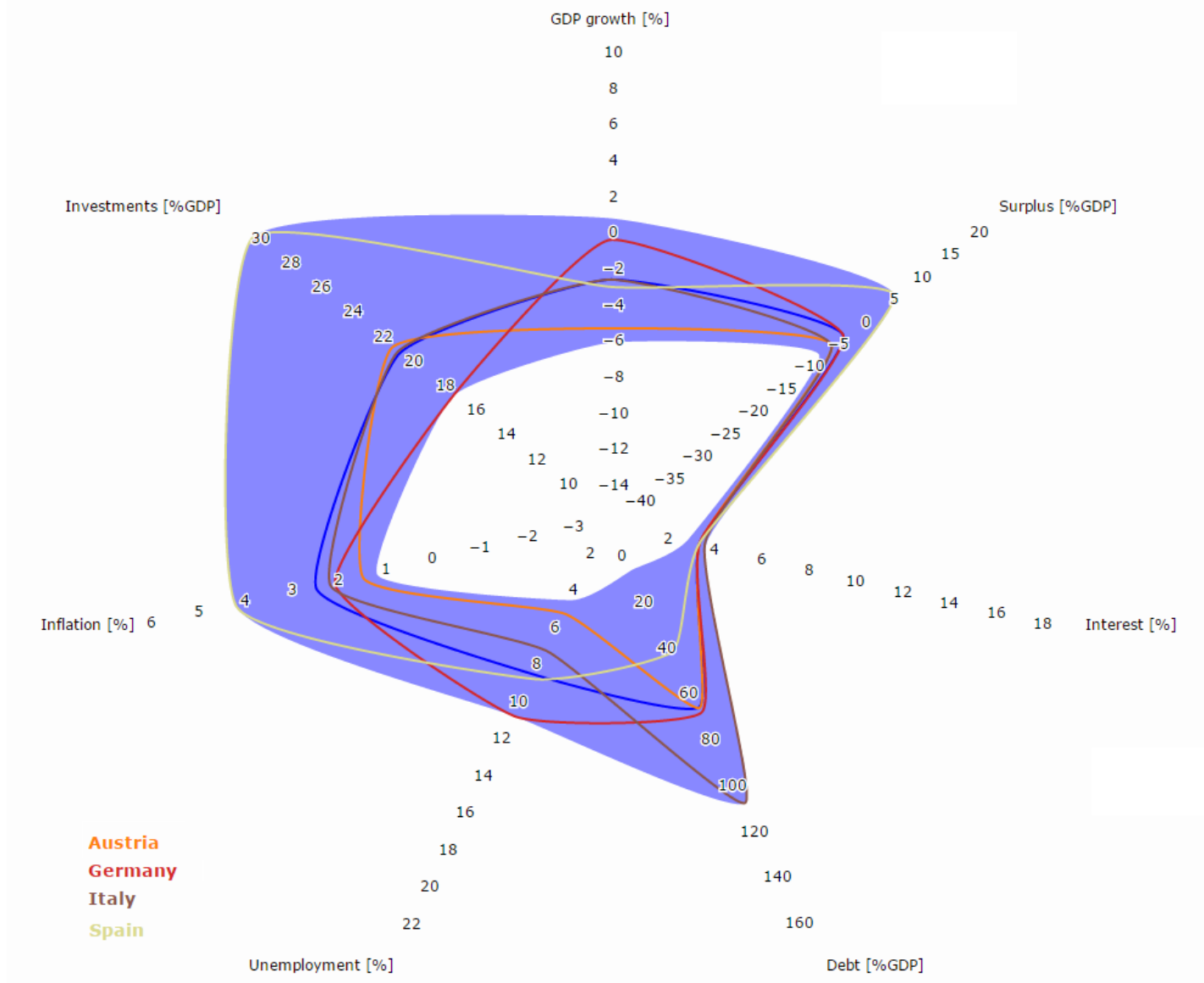


A futási idő és a memóriahasználat valószínűleg pozitív kapcsolatban állnak (sikeres teszteknel)

Párhuzamos koordináták: viz. alternatívák



Radar chart: egy párhuzamos koord. kiterjesztés



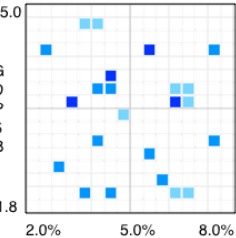
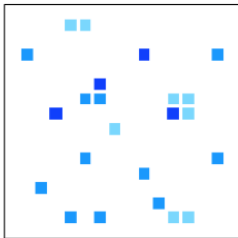
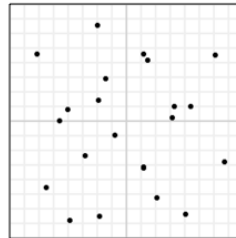
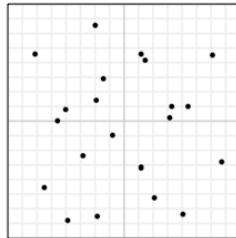
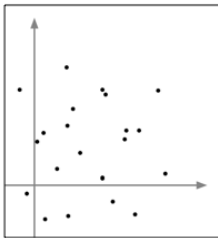
AGGREGÁLÁS + VIZUÁLIS MEGJELENÉS

Datashader

- Gyakran nagy az adat (1000000+)
- DE nem az egyedi adatvektor, hanem az értékek eloszlása érdekes

Projection *Aggregation* *Transformation* *Colormapping* *Embedding*

Year	GOP	Unemp	Inflation	Educ	Interest
1990	87.22774	0.88769	2.50876	19.89628	8.442063
1991	11.24687	0.84458	2.17383	17.96444	8.22579
1992	64.82791	0.44615	1.04903	2.588296	1.531252
1993	2.827515	0.56875	0.171763	14.8277	9.798936
1994	48.09787	0.51513	1.55378	8.50005	9.324268
1995	15.85648	0.25209	1.52801	18.12878	8.878299
1996	27.24226	0.37782	1.28835	16.16028	1.27903
1997	50.26089	0.21761	0.86024	2.433843	8.193702
1998	16.16523	0.87427	0.750219	0.84092	5.249787
1999	16.12853	0.83626	2.88202	6.902022	1.575446
2000	7.415708	0.82021	0.299542	1.193834	9.300876
2001	13.42125	0.525407	0.088174	9.20093	9.880449
2002	41.62787	0.306719	1.27498	11.77219	8.772817
2003	47.99834	0.415496	0.879515	17.5177	0.085302
2004	14.39844	0.28558	2.40446	1.47892	1.62286
2005	18.1334	0.46458	1.782424	0.942948	2.62773
2006	16.62174	0.520093	1.78933	15.33283	0.828348
2007	19.50517	0.02081	1.978622	16.42061	1.518401
2008	40.17541	0.221893	1.75473	9.949314	1.468717
2009	9.052268	0.189748	2.879744	1.89005	1.762614
2010	58.34968	0.47896	1.48360	15.57068	0.888847



Data

Scene

Aggregate(s)

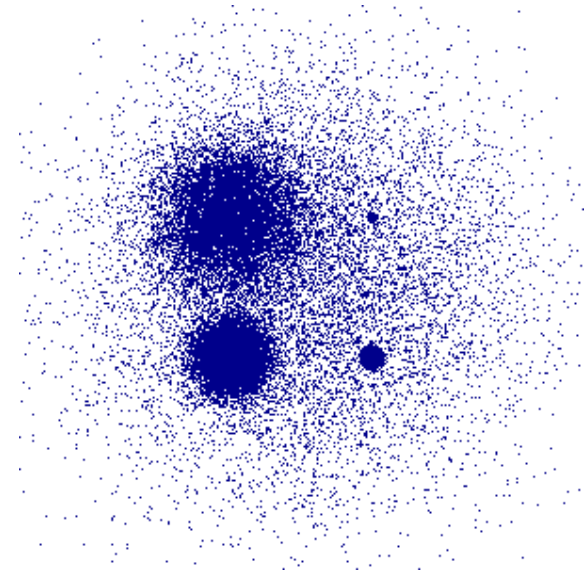
Image

Plot

http://datashader.org/getting_started/2_Pipeline.html

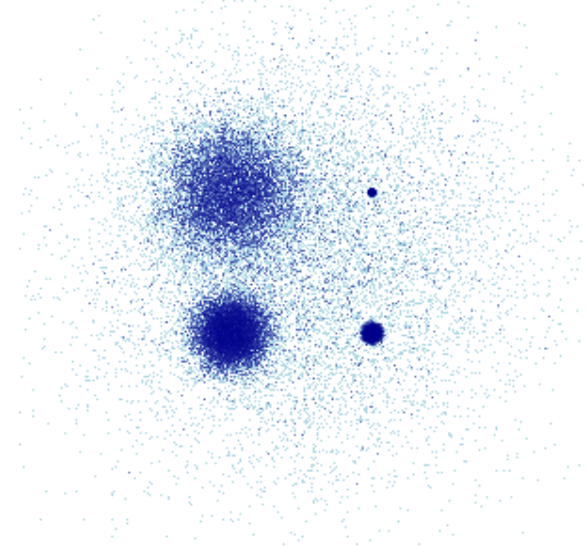
- Adat
- for **x**, **y**, **val**, **cat** in
 [(**2**, **2**, **10**, "d1"),
 (**2**, **-2**, **20**, "d2"),
 (**-2**, **-2**, **30**, "d3"),
 (**-2**, **2**, **40**, "d4"),
 (**0**, **0**, **50**, "d5")]

Előfordulás
 Cellaintenzitás :
 Van-e adatpont a
 cellában?



Gyakoriság
 Cellaintenzitás :
 Cellabeliek száma

 Összes minta

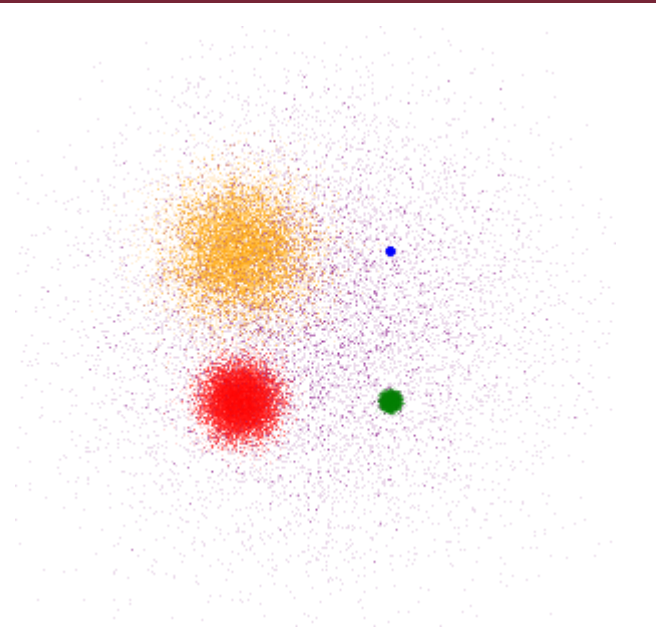


Szűrés, színezés

- Adat
- for **x**, **y**, **val**, **cat** in
[(**2**, **2**, **10**, "d1"),
(**2**, **-2**, **20**, "d2"),
(**-2**, **-2**, **30**, "d3"),
(**-2**, **2**, **40**, "d4"),
(**0**, **0**, **50**, "d5")]

Színezés
kategóriánként

```
color_key = dict(  
    d1='blue',  
    d2='green',  
    d3='red',  
    d4='orange',  
    d5='purple')
```



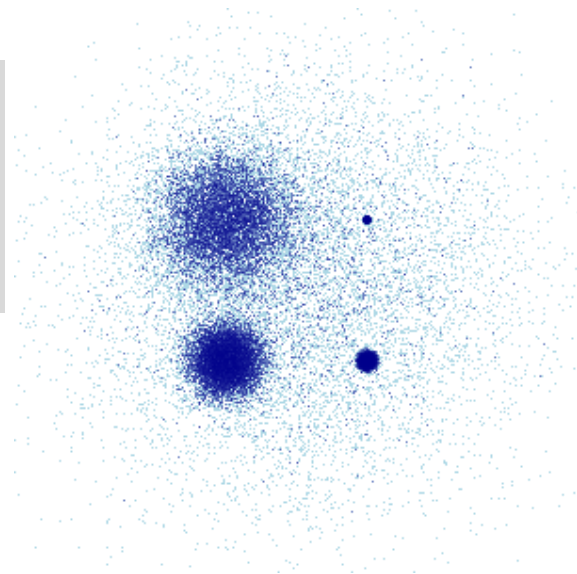
Szűrt előfordulás
Cellaintenzitás
Cellabeliek száma<--
Cellabeliek száma | d3)



Súlyozás

- Adat
- for x , y , val , cat in
[(2, 2, 10, "d1"),
(2, -2, 20, "d2"),
(-2, -2, 30, "d3"),
(-2, 2, 40, "d4"),
(0, 0, 50, "d5")]

Gyakoriság
Cellaintenzitás :
Gyakoriság
négyzete



Cellaintenzitás
Cellabeliek
kivéve kiesők

