

Data analysis

Andras Foldvari - andras.foldvari@edu.bme.hu

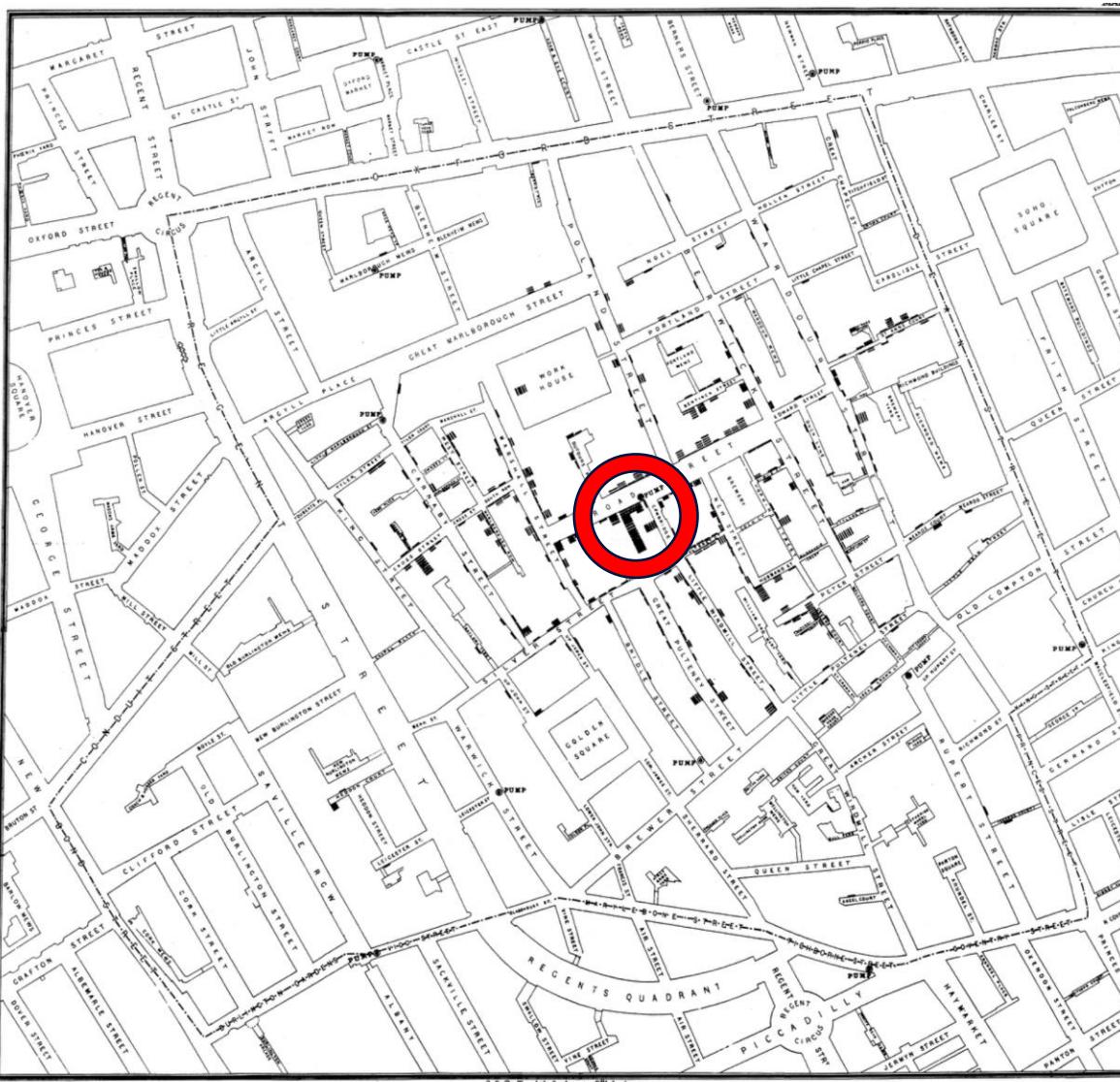
Sources: System Modeling, Empirical Systems Engineering



Budapest University of Technology and Economics
Department of Artificial Intelligence and Systems Engineering
ftsrg Research Group



Exploring Underlying Relations



Infographics

*Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.
Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite.*

Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui ont été en Russie, le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Léger, de Fezensac, de Chambray et le journal intérieur de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout, qui avaient été détachés sur Minsk en Mobilow et qui rejoignirent vers Orsha en Witebsk, avaient toujours marché avec l'armée.

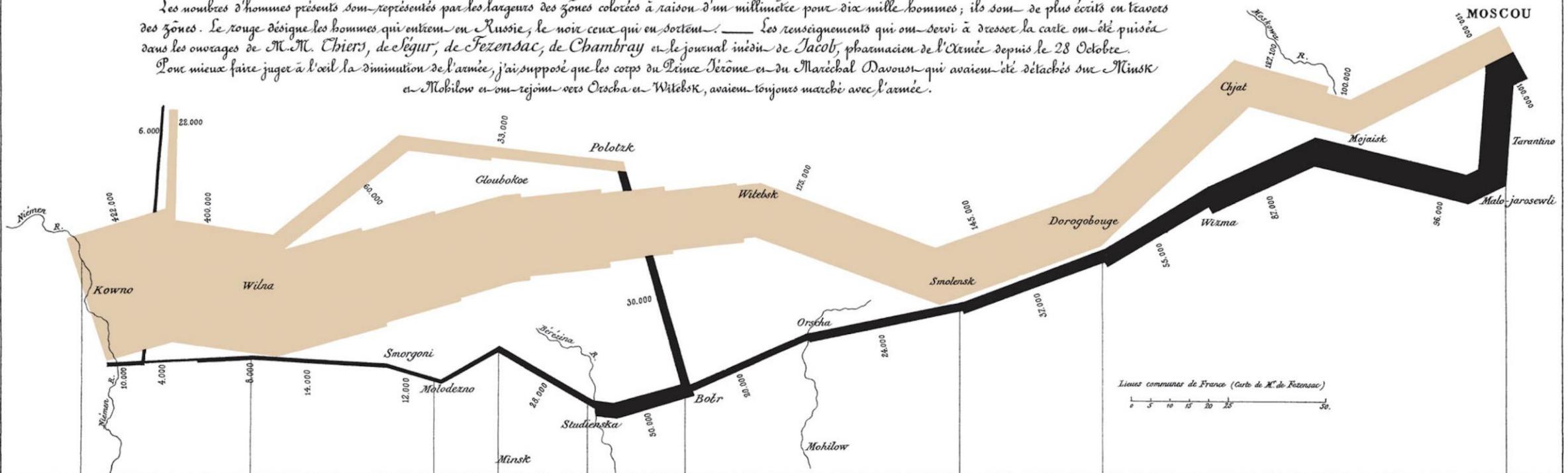
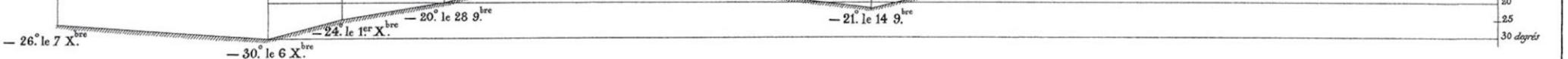


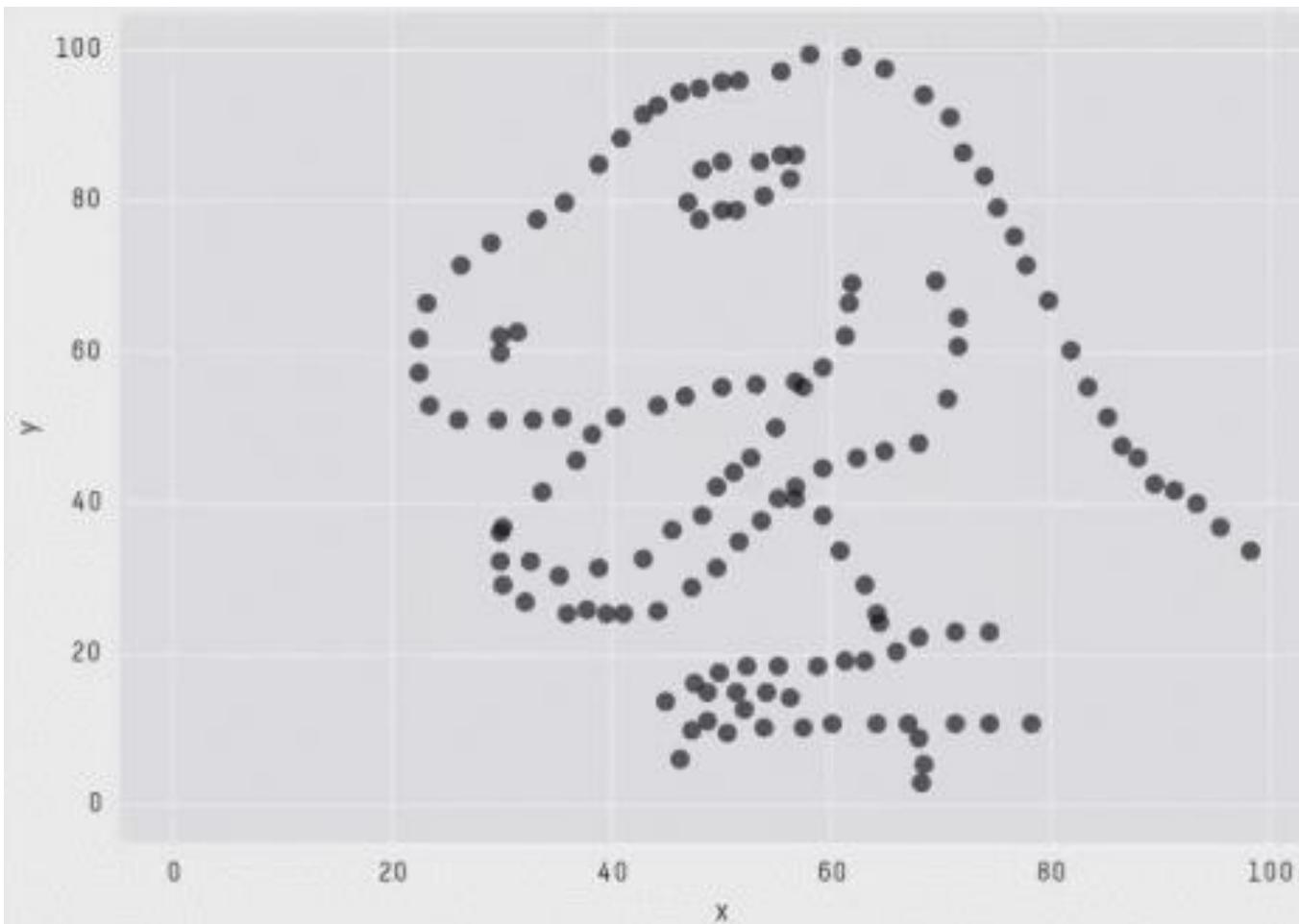
TABLEAU CRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les cosaques passent au galop
le Niemen gelé.





The datasaurus



X Mean: 54.2659224
Y Mean: 47.8313999
X SD : 16.7649829
Y SD : 26.9342120
Corr. : -0.0642526

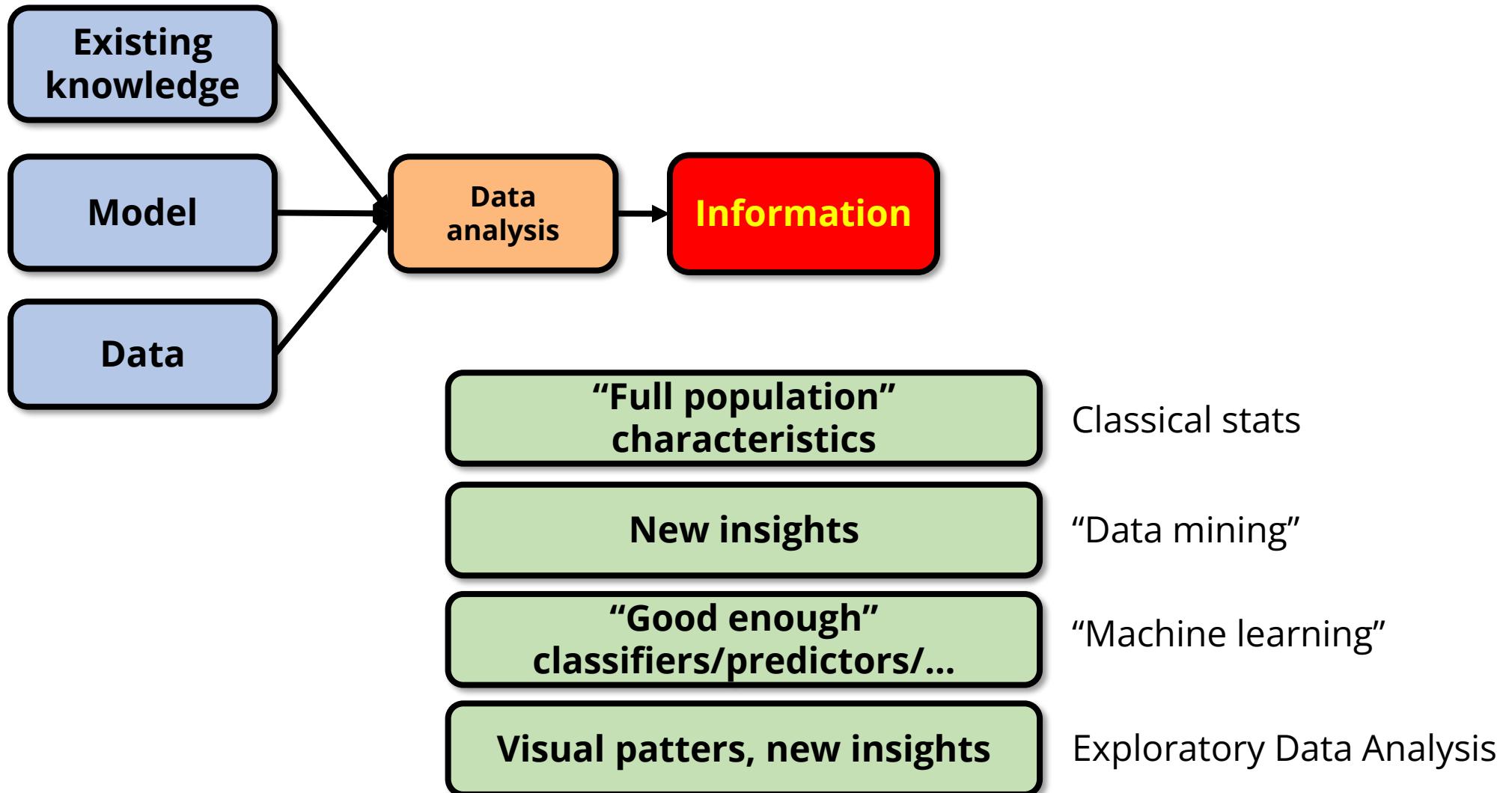
Alberto Cairo (2016): The datasaurus

CRISP-DM folyamatmodell



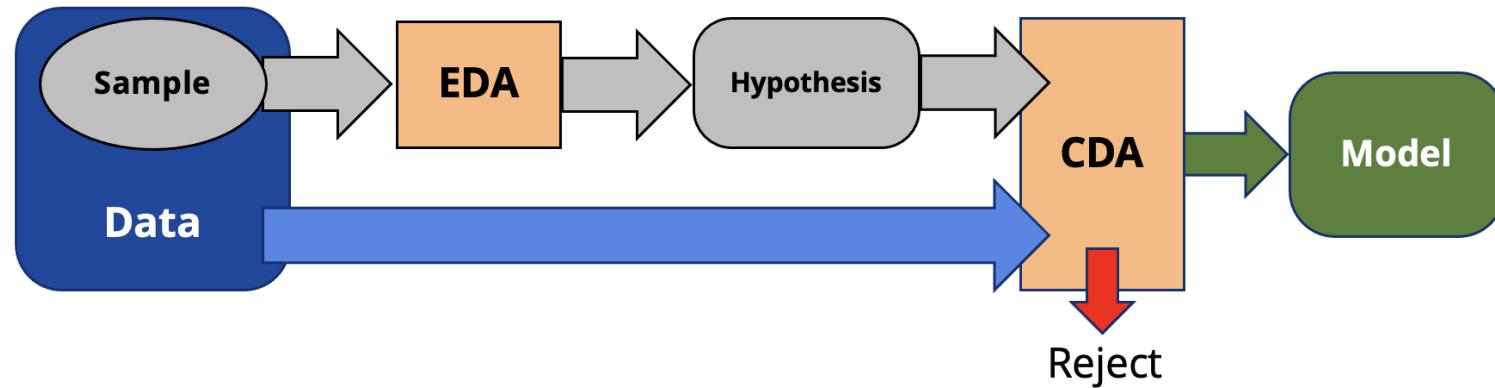
- **Cross-industry Standard Process for Data Mining**
- **IBM 2015**

Informal model of data analysis



Exploratory Data Analysis (EDA)

- Approach to **analyzing** data sets
 - Summary of **main characteristics**
 - Easy-to-understand,
 - Often: **visual graphs**
 - Without using **statistical models**
 - Or having formulated a hypothesis
- Prepares **Confirmatory Data Analysis**
 - Hypothesis testing
 - Statistical model building and checking
- (Dynamic) visualization delivers understanding: **outliers**, **trends** and **patterns**
- **Not** statistics-free - but “**inference**” sensitivity to errors should be managed
 - **Robust statistics** - low sensitivity to outliers
 - **Nonparametric statistics** – no a priori assumptions



Data models

Tabular Data

- Data is stored in rows and columns, in a tabular format.
- Columns represent attributes, and rows represent individual records.
- E.g., Pandas, MySQL, PostgreSQL

Graph data model

- Data is represented as nodes and edges.
- Useful for visualizing relationships and connections.
- E.g., Neo4j, GraphDB

Document data model

- Data is stored in structured document formats.
- Offers high flexibility, especially for handling unstructured or semi-structured data.
- JSON, XML
- E.g., MongoDB, CouchDB

Column data model

- Only relevant columns are read during queries.
- Commonly used for Big Data.
- E.g., Apache Cassandra, HBase

Further data models

- Key-Value
 - Redis
- Object-Oriented
- Semantic Models
 - RDF, OWL
- File-Based
 - AWS S3
- Time Series
 - InfluxDB

Tabular Data Model

- **Table row**= model element
- **Table column**= attribute

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

- Data analysis tools (e.g., R, Python): **dataframe**
 - Each row represents a measurement
 - Each column has a specific type

Data Cleaning

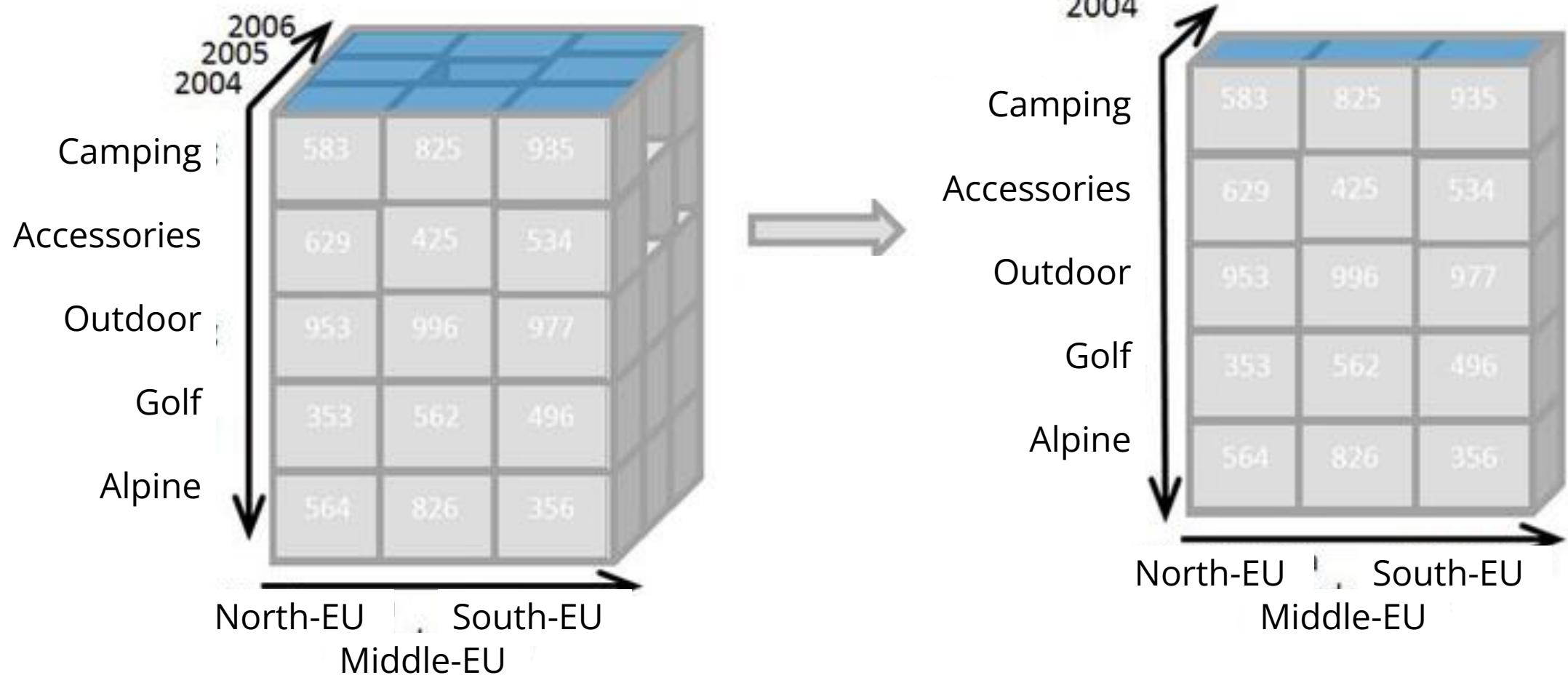
- The goal of data cleaning is to remove or correct errors, gaps, and inconsistencies that may distort analysis results.
- Handling missing values
 - They may indicate errors!
- Duplicates
- Erroneous data and outliers
 - Outliers can be valuable!
- Standardizing data formats
- Standardizing categories and labels
- Logical checks
- Tools: Python (pandas, NumPy), OpenRefine, Excel, Apache NiFi

Example data cleaning

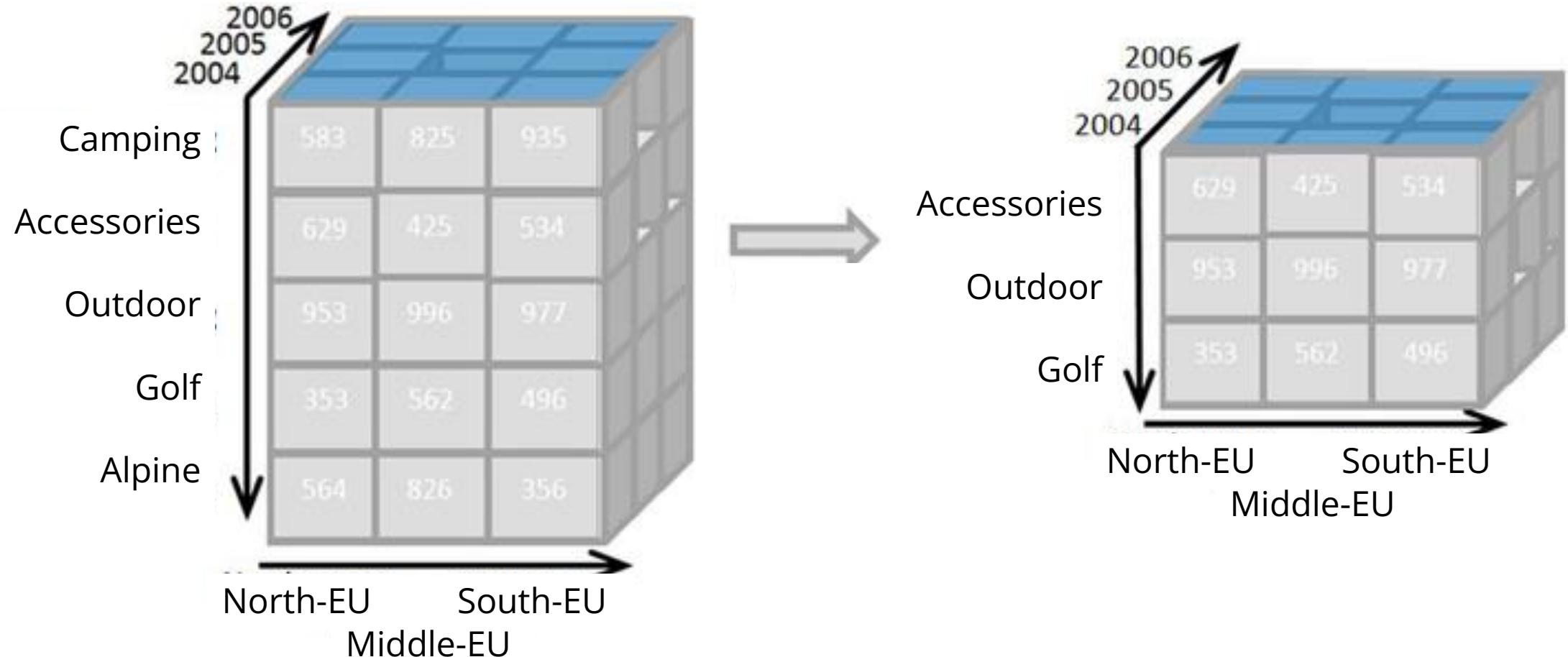
- Uniqueness
- Mixed columns (weight): kg and Unit are incompatible
- Multiple data in one cell (status)
 - Multiple statuses
 - State machine
- Uninterpreted textual description

Mission name	Organisation	Nation	Type (U/mass)	Launch date	Status	Mission description	Photo
TUBSAT-N	TU Berlin	Germany	8 kg	1998-07-07	Reentered, Was operational	Store and forward communication.	
TUBSAT-N1	Technische Universitat Berlin	Germany	2U	1998-07-07	Reentered, Was operational	Store and forward communication.	
Artemis JAK (MASat, Picosat 5)	Santa Clara University	US	0.2 kg	2000-01-27	No signal	Simple beacon transmitter.	

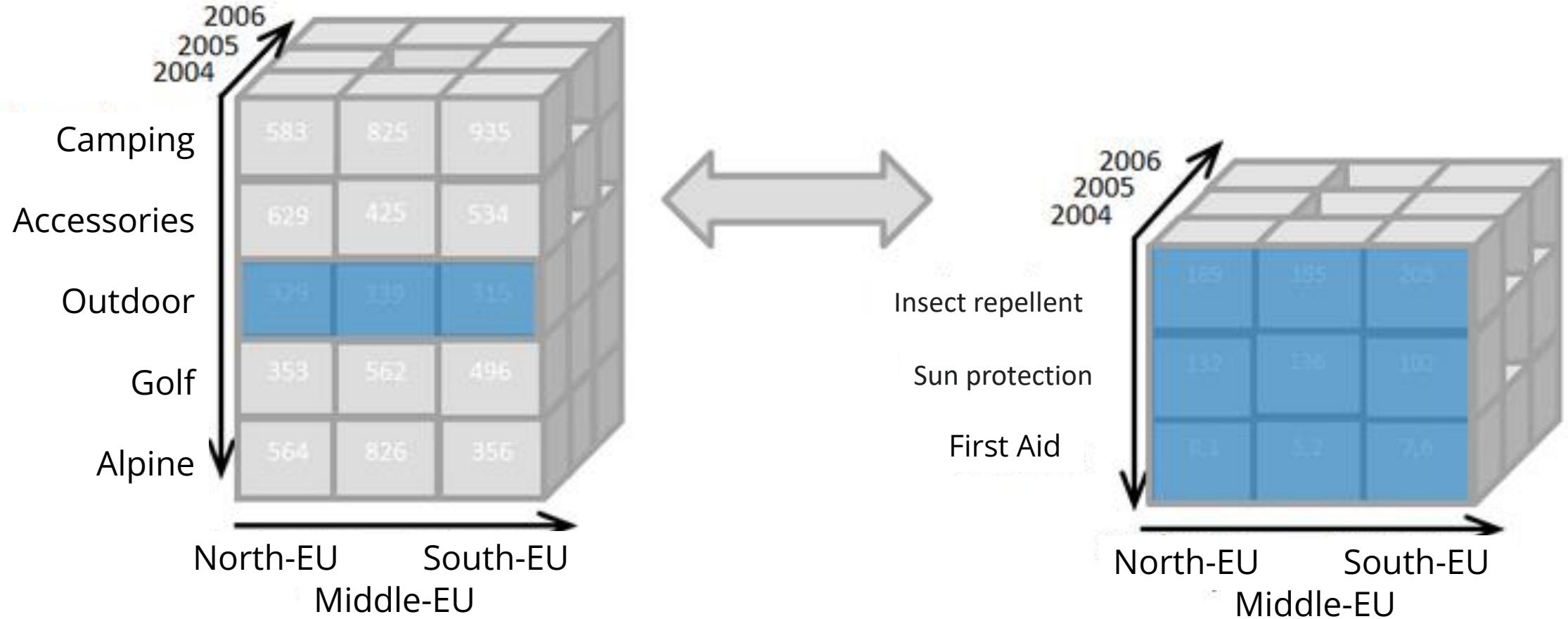
Slicing



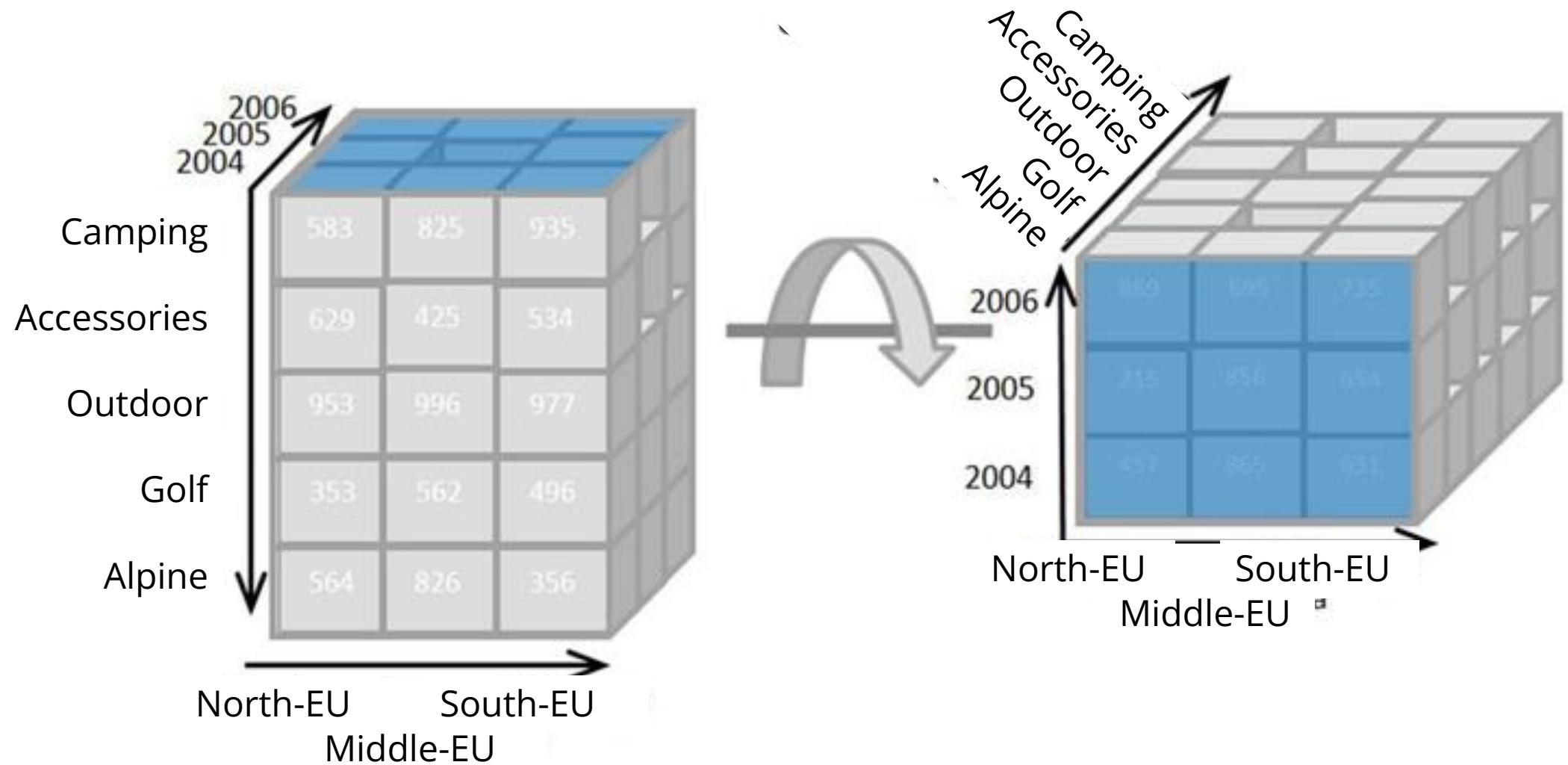
Dicing



Drill up & down



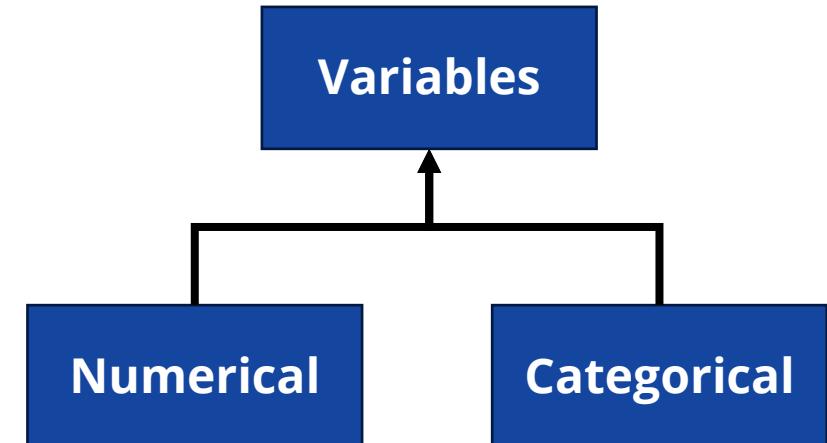
Pivoting



Numerical and Categorical Variables

- **Numerical**

- Basic arithmetic operations are interpretable.



- **Categorical**

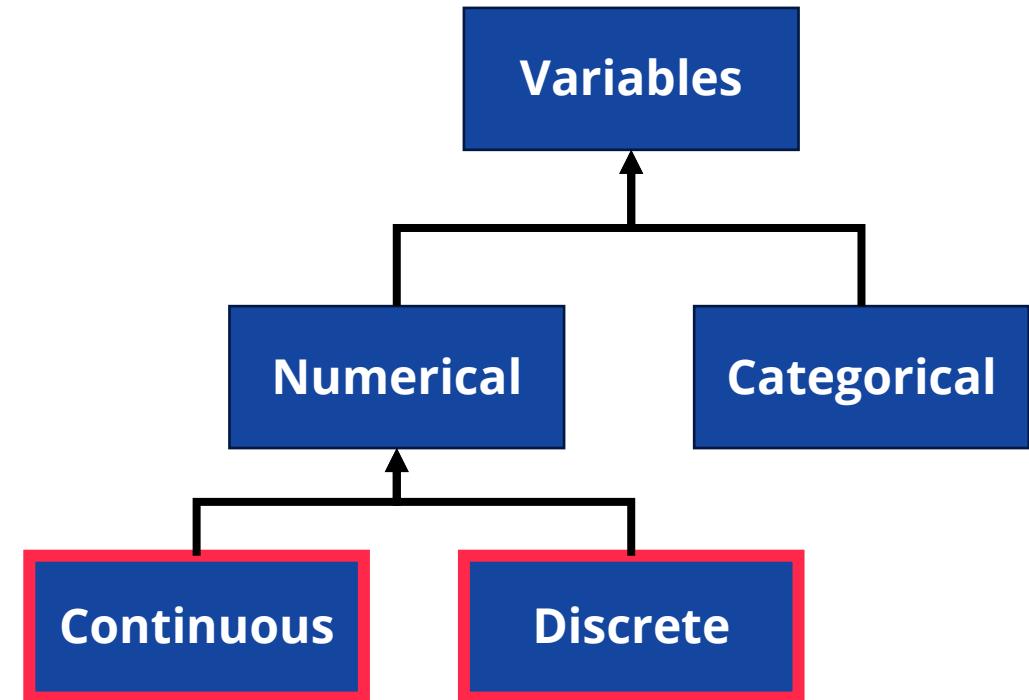
- Mathematical operations are not interpretable; at most, they can be ordered.



age, average temperature, phone number

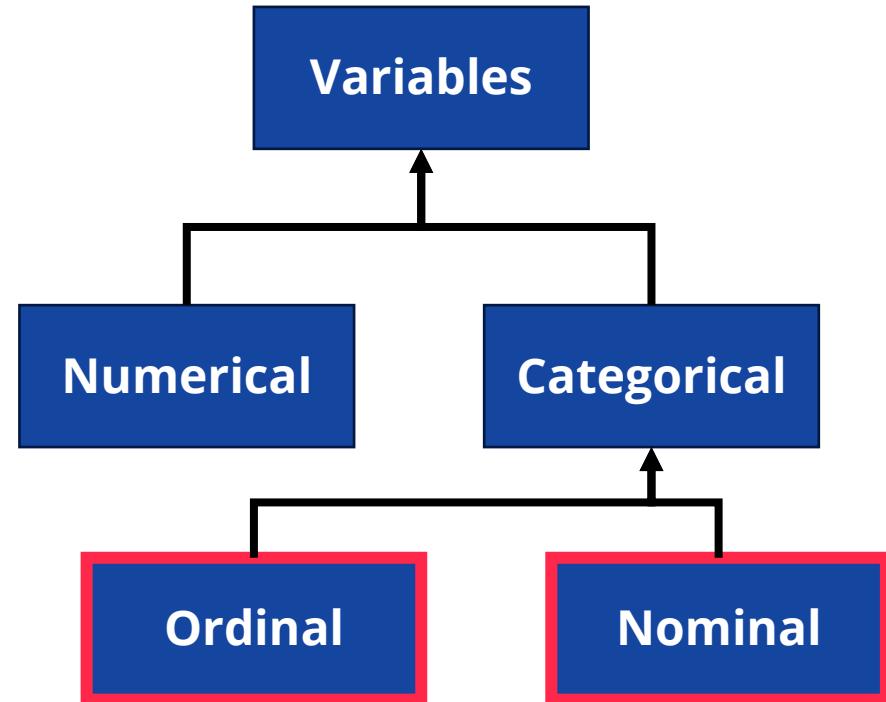
Numerical Variables

- Continuous
 - Measured – can take any value
 - within given range
 - with given precision
 - Pl. the average test score of the people in the room
- Discrete
 - Counted– can take finite number of values within a given range
 - Pl. the number of people attending the lecture



Categorical Variables

- Ordinal
 - Total ordering of values
 - E.g., stars of a hotel
- Nominal
 - E.g., colors, car brands





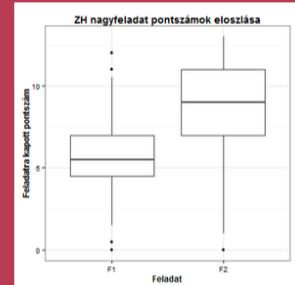
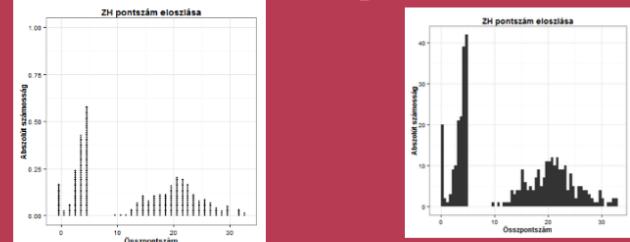
Data visualization

1 variable - distributions

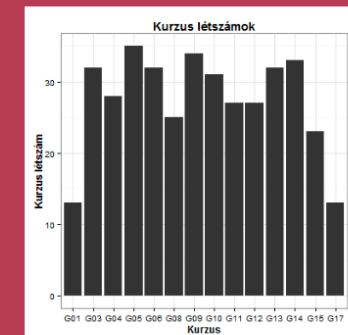
Variables

Numerical Categorical

ZH pontszám: [13, 15, 2, ...]



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



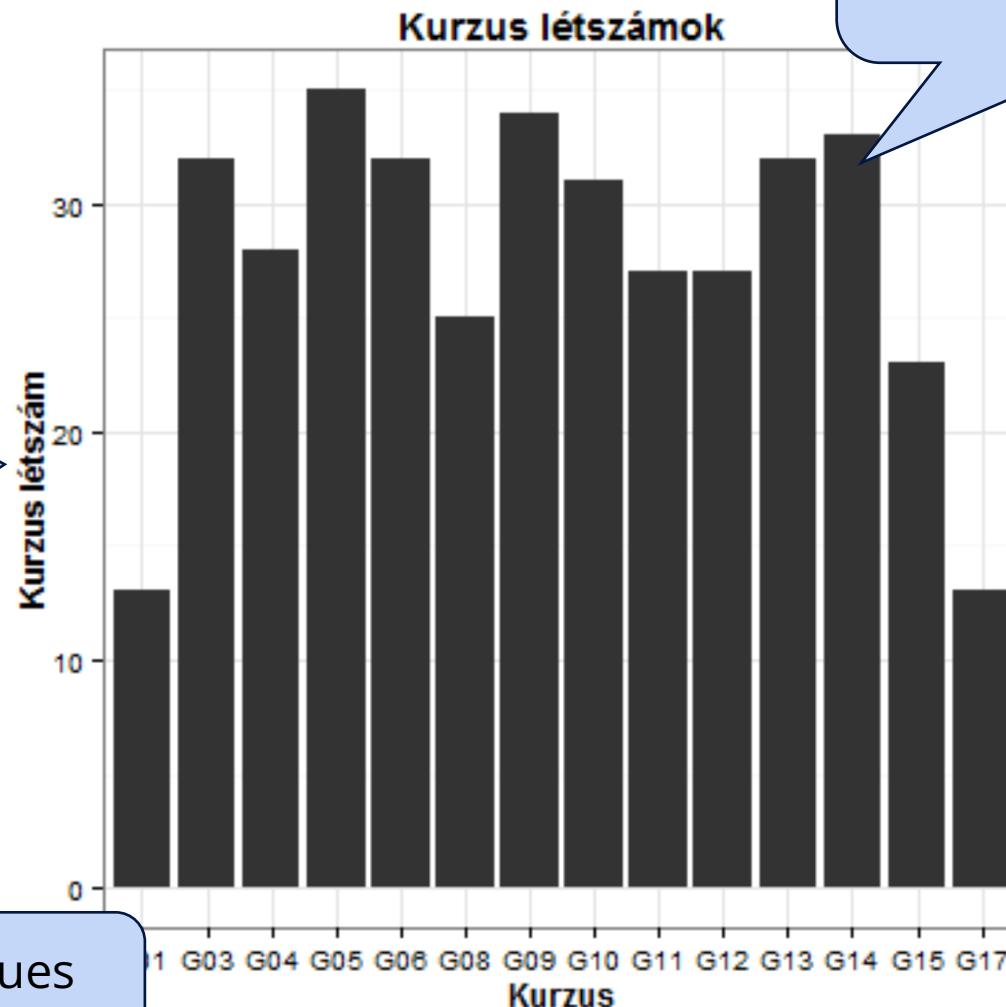
Bar chart

Input variable: course code

Question: how many students are enrolled in each course?

Absolute frequency

Design decision: breaking down the set of values
Example: Tuesday–Thursday–Friday



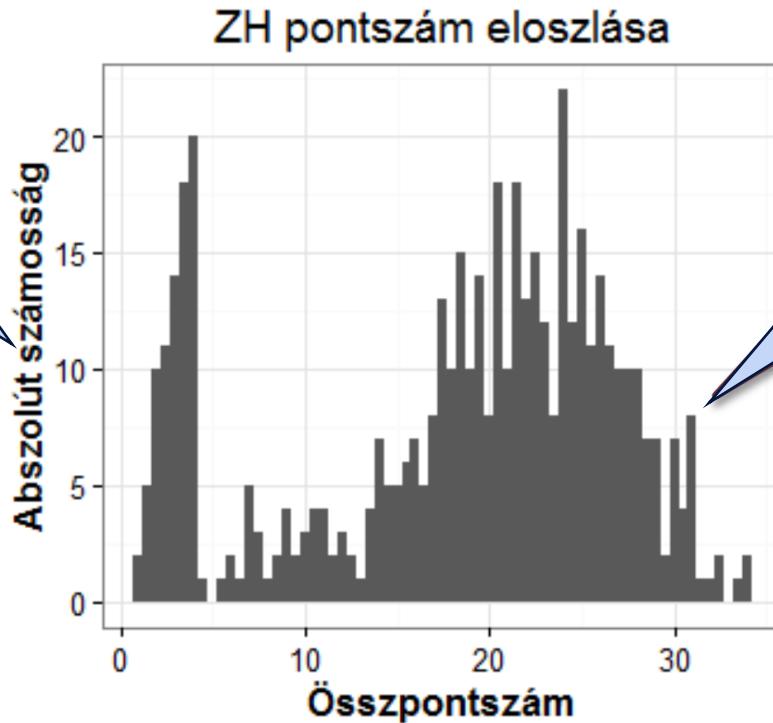
Bar height: frequency of a given value

Histogram

Input variable: Total score of midterm (ZH)

Question: How did the midterm scores turn out?

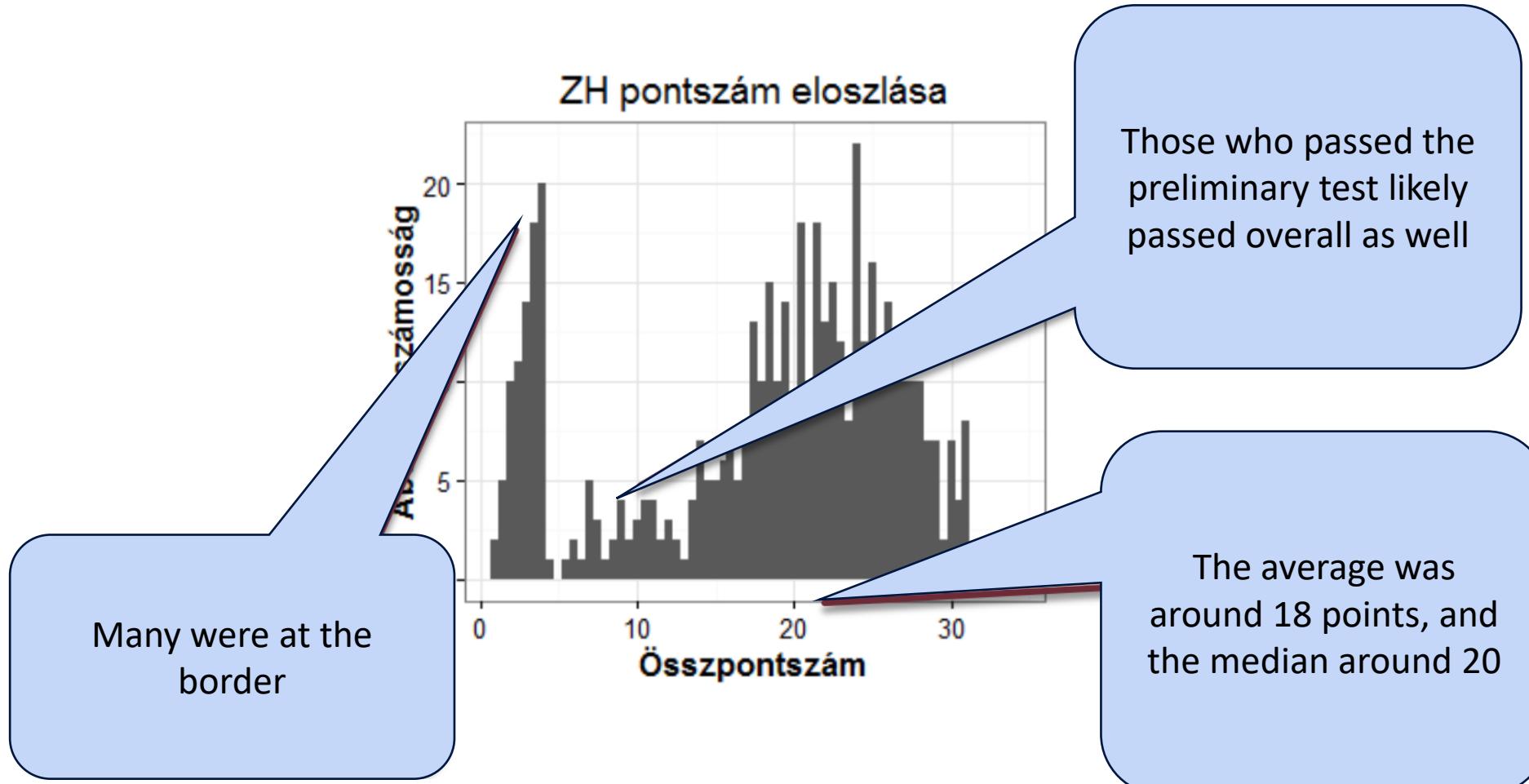
Absolute frequency



Bar height: count of a given interval

Design decision: how long should the interval be?
Example: is a resolution of 1 point enough, or should we go down to half-point intervals?

Histogram



Interesting titles



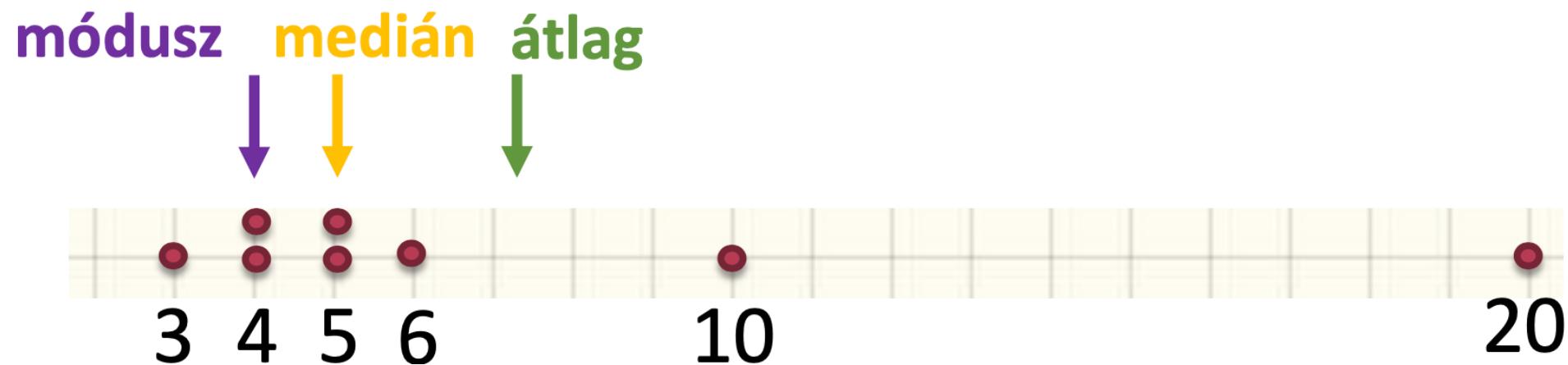
"Half of Hungarians earn less than the median wage."



Bonus question: What is wrong with the picture?

Mean, median, mode – the “middle”

- {3, 4, 4, 5, 5, 6, 10, 20}



Mode and Median

- If we arrange the values in ascending order, the middle value is the **median** of the dataset. If there is no middle value (in the case of an even number of values), the median is the average (arithmetic mean) of the two middle values.
- The **mode** is the value that occurs most frequently in the dataset. This is not necessarily unique; in such cases, we talk about multiple modes.

Percentile

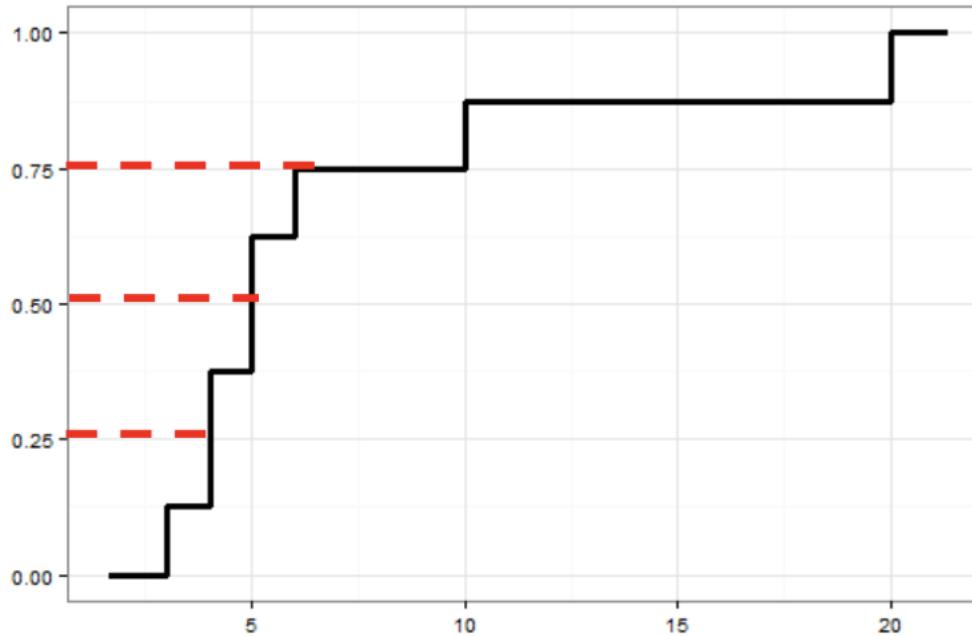
- At the n -th percentile, $n\%$ of the values are smaller.

- **Percentile**

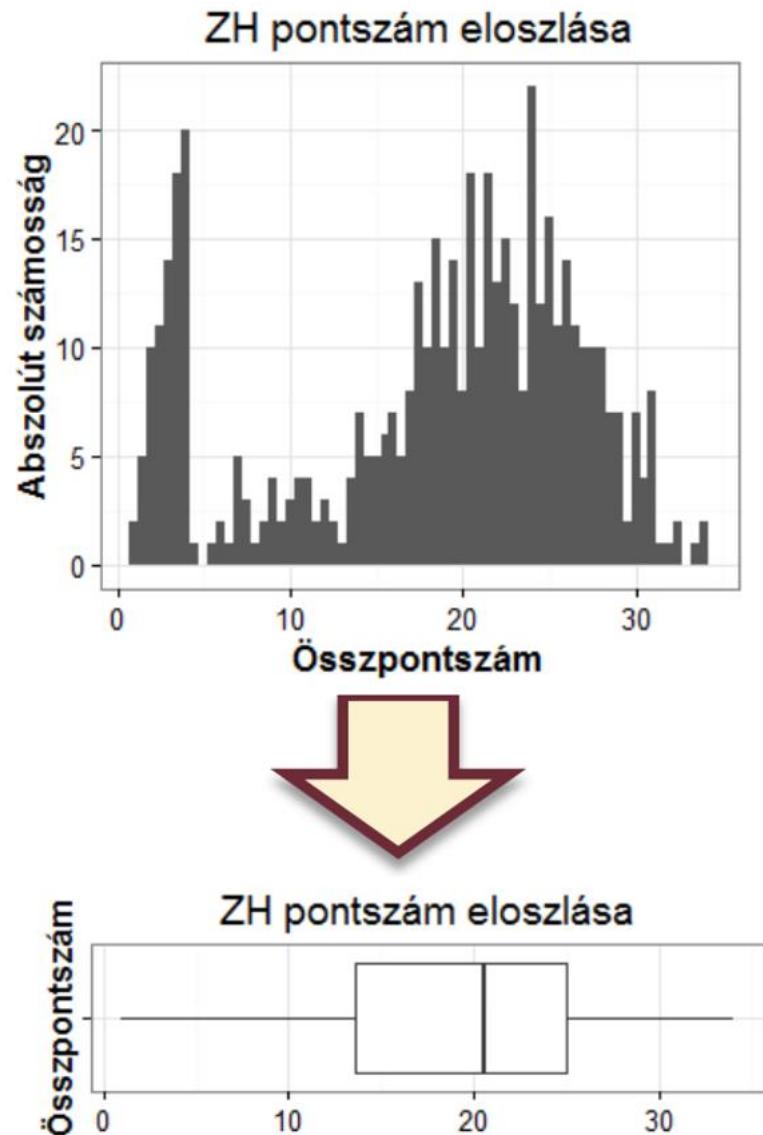
- $\{3, 4, 4, 5, 5, 5, 6, 10, 20\}$
 - 50th percentile: 5
 - 25th percentile: 4
 - 75th percentile: 6

- **Quartile**

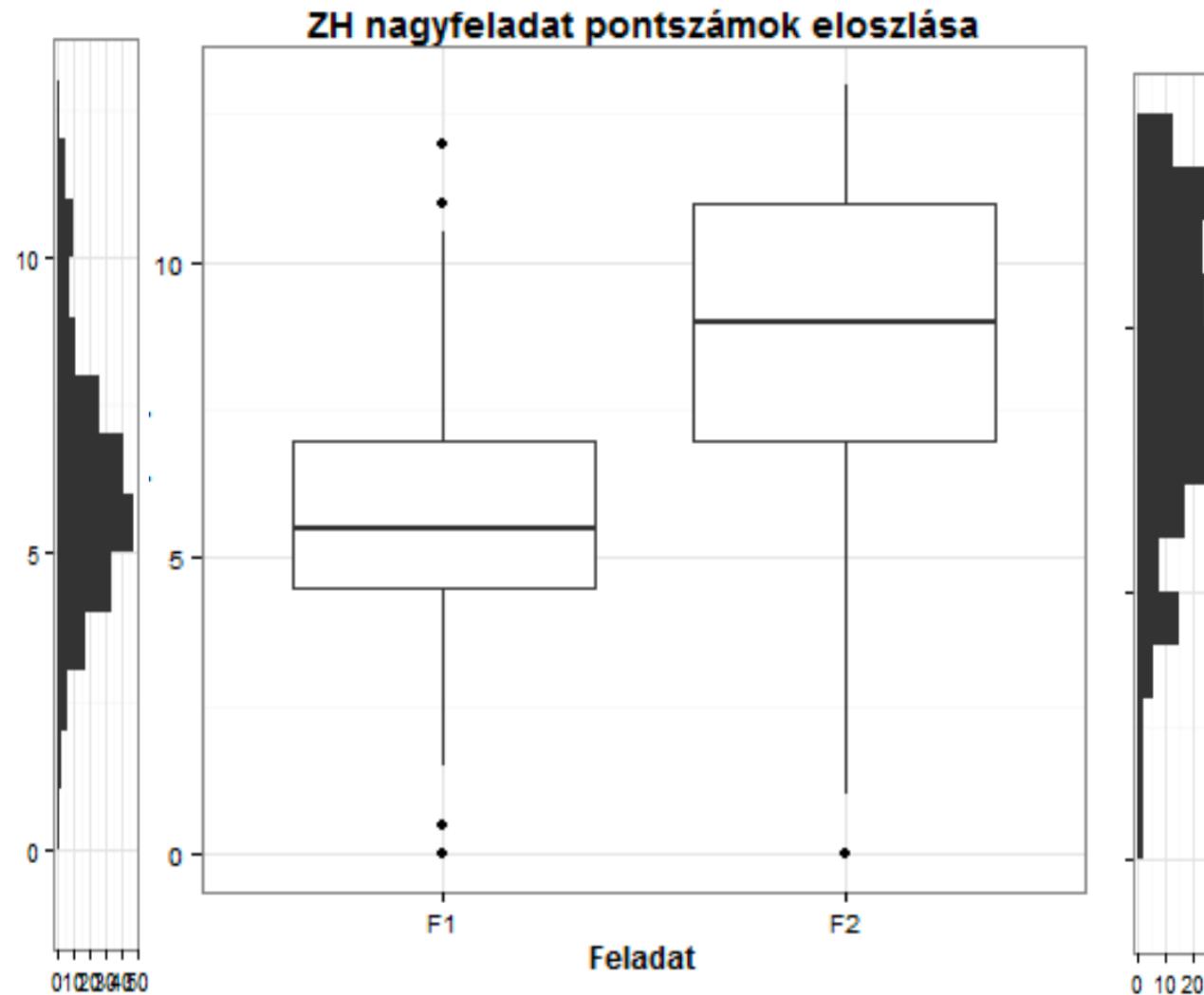
- Q1: 25th percentile
 - Q3: 75th percentile
 - Q2: median



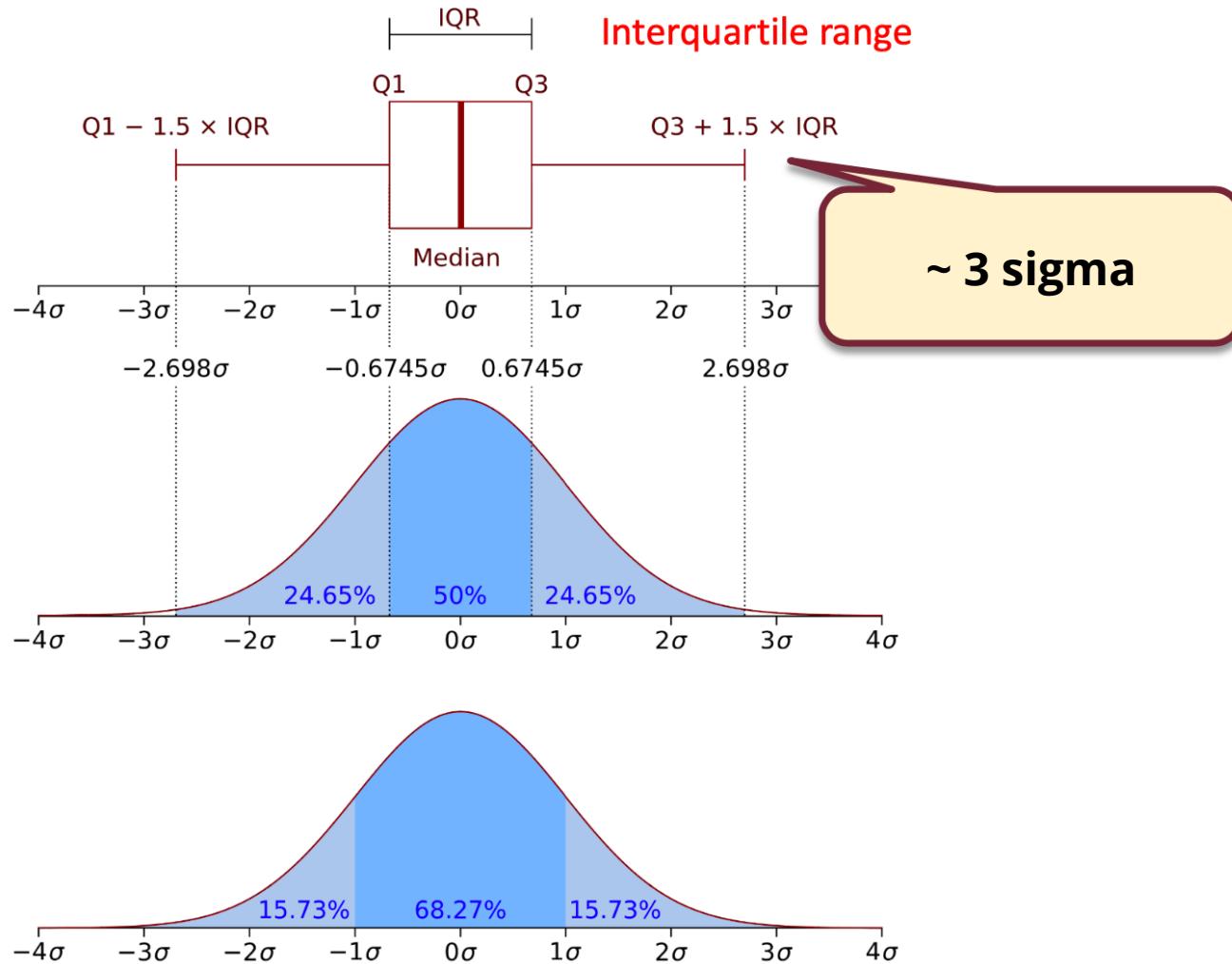
Boxplot



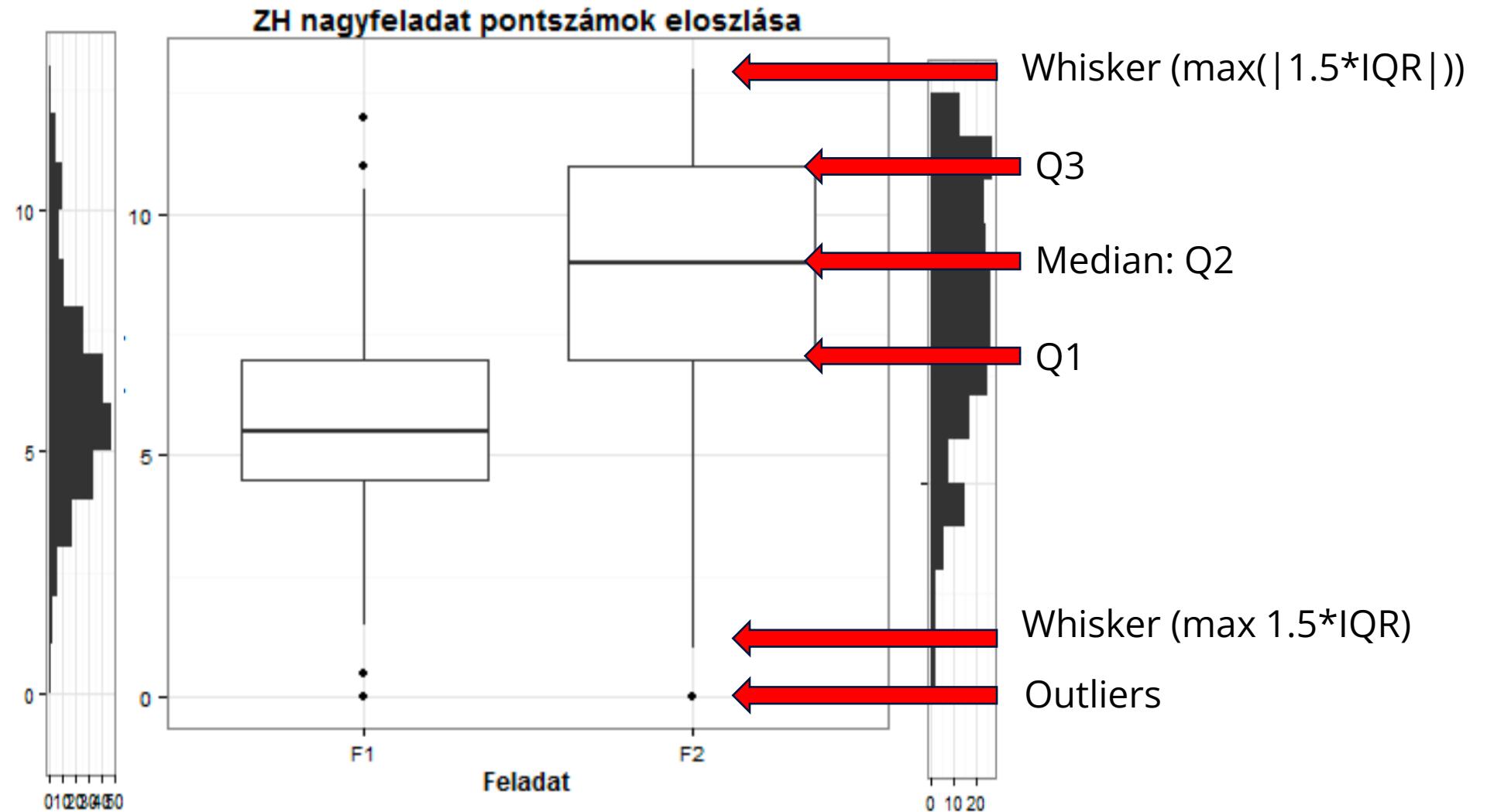
Boxplot



IQR

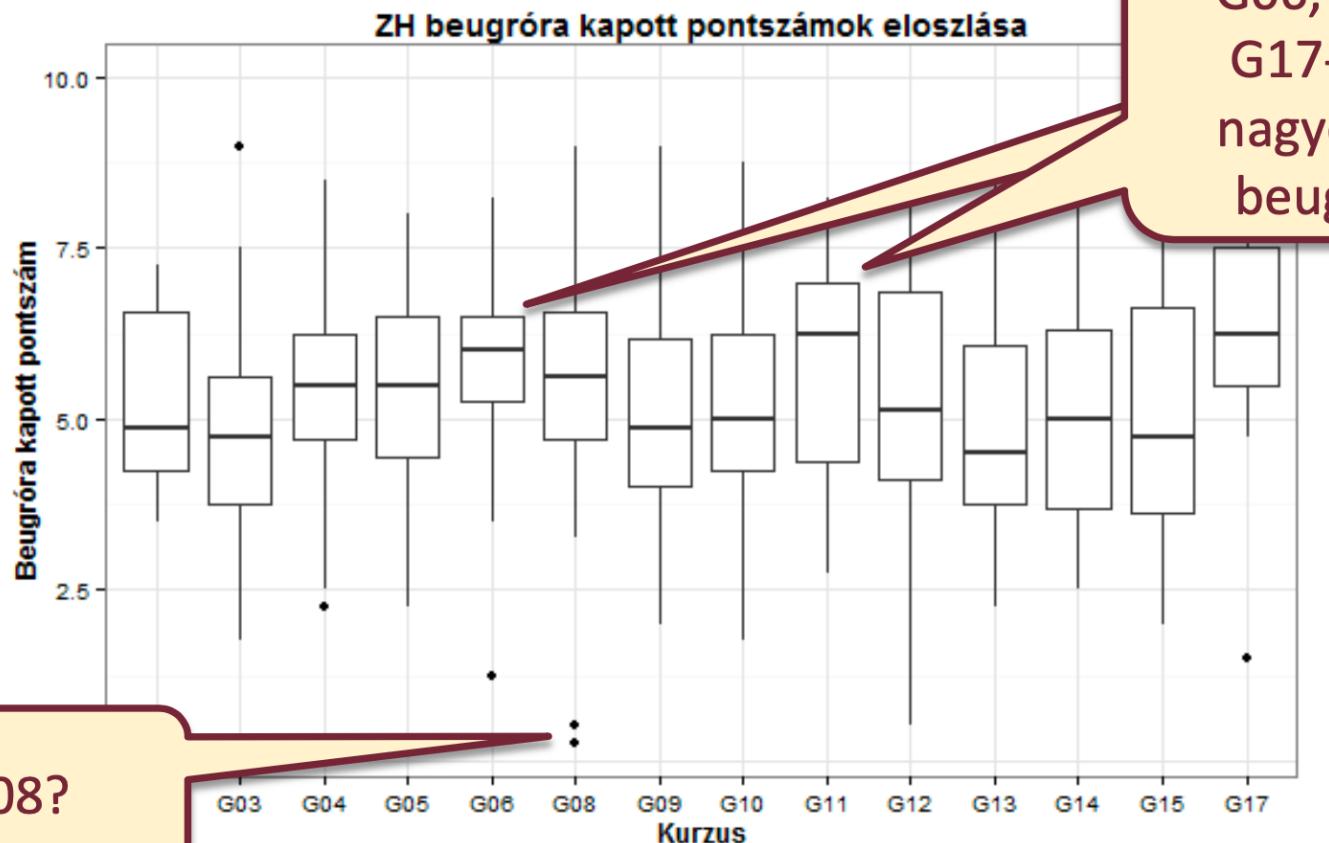


Boxplot



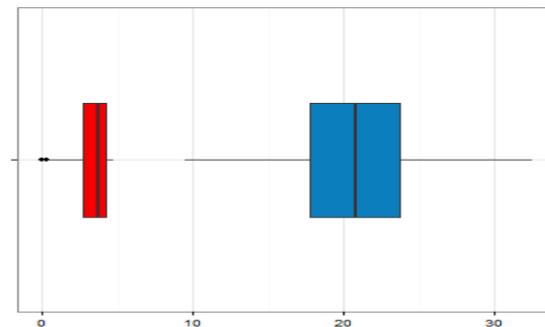
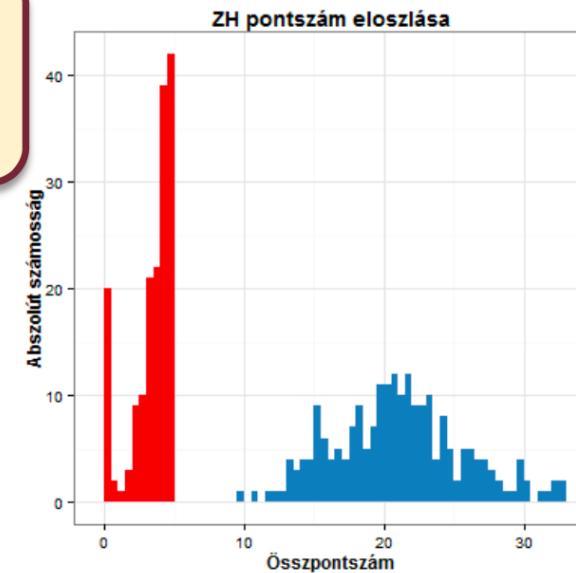
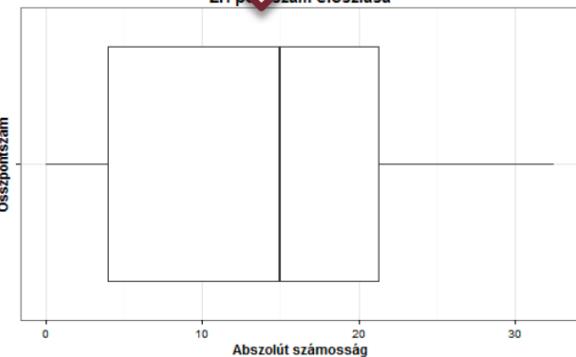
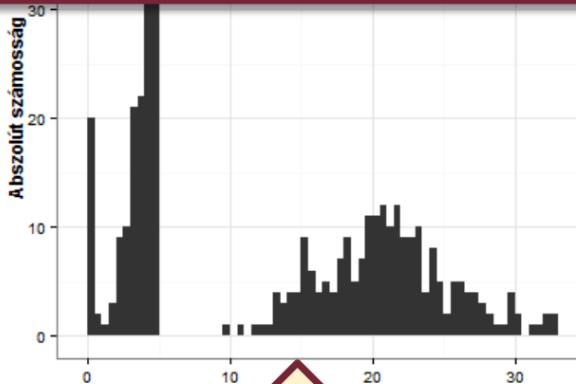
Boxplot example

- Melyik csoportban hogyan sikerültek a beugrók?



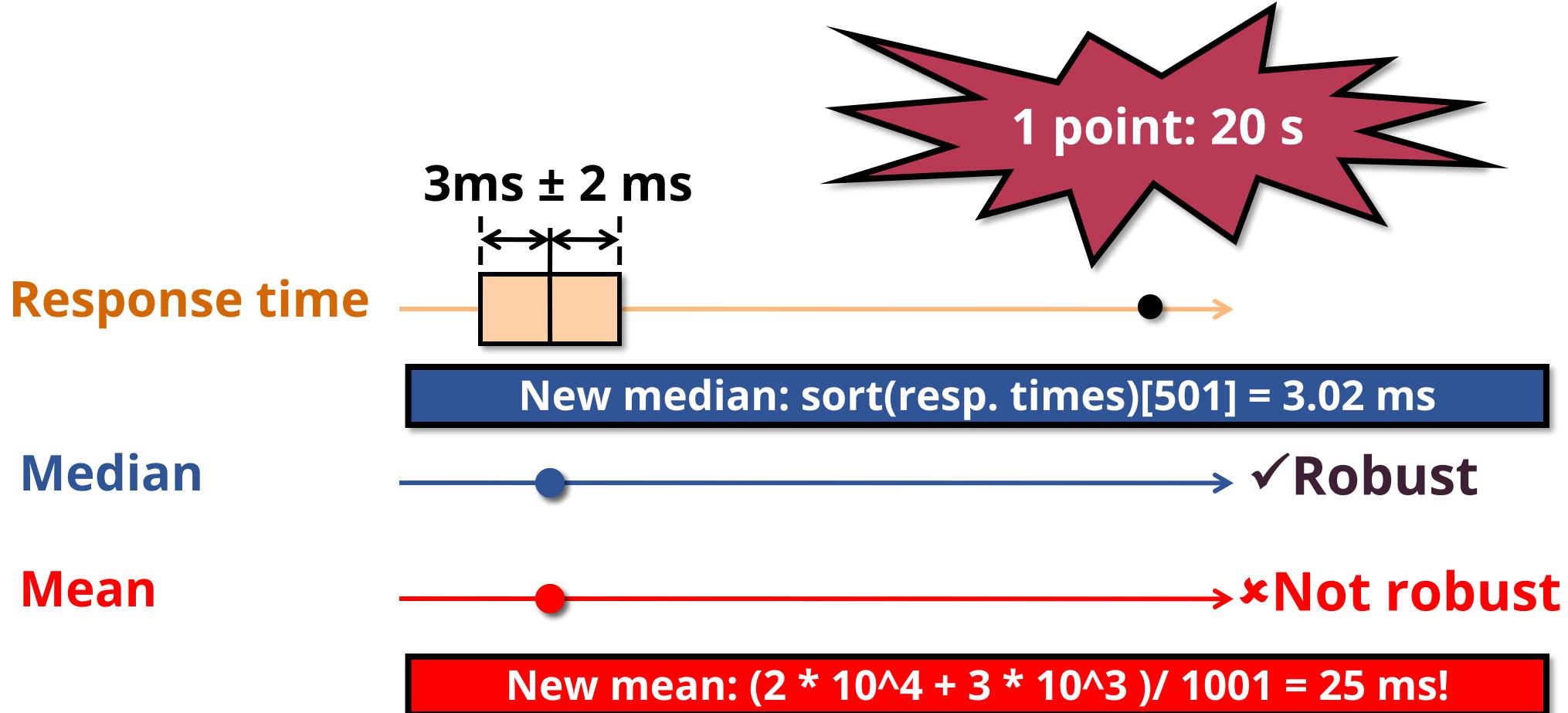
Boxplot abstraction

Abstraction: with a boxplot, we can lose important information!

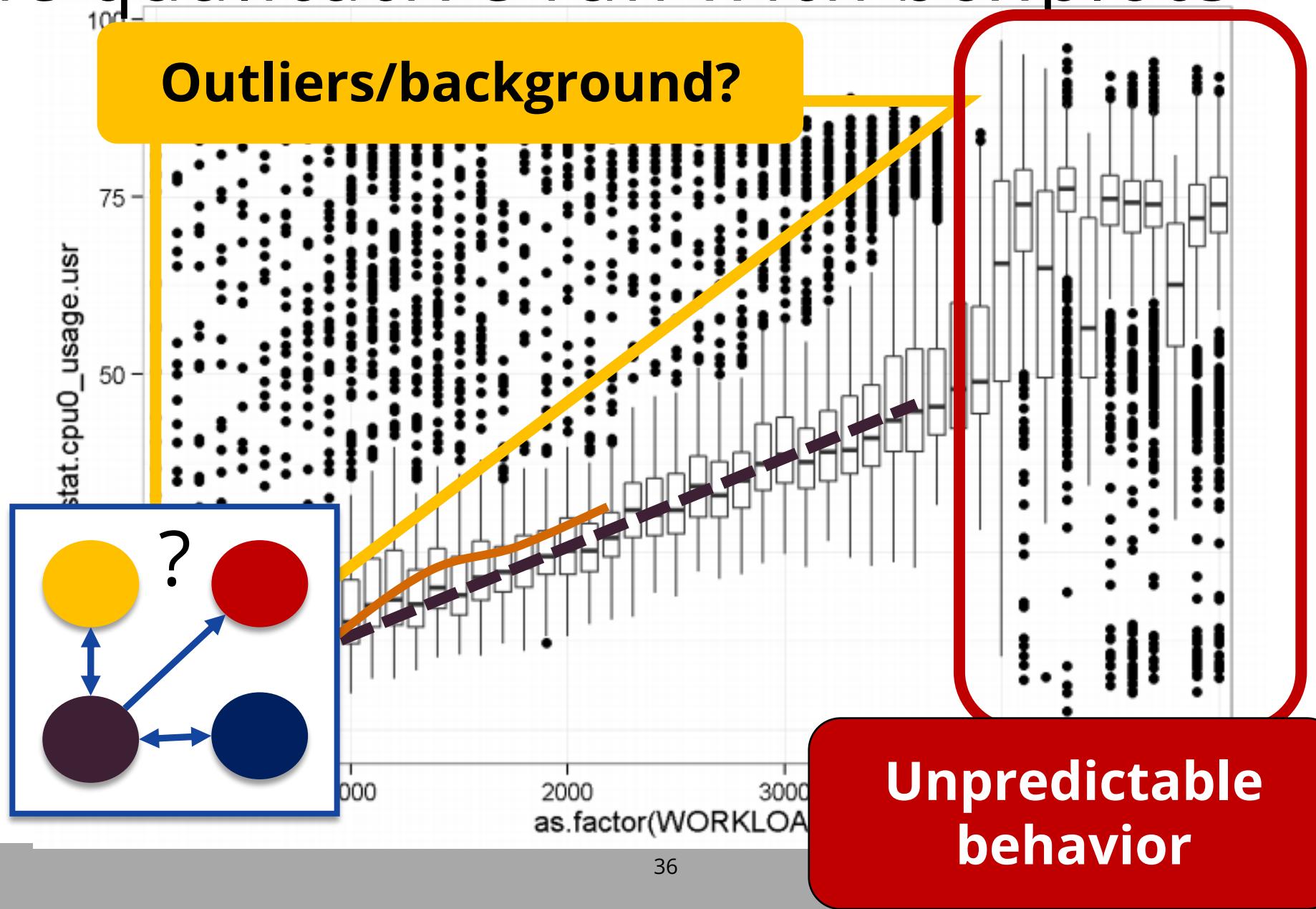


Robust statistics

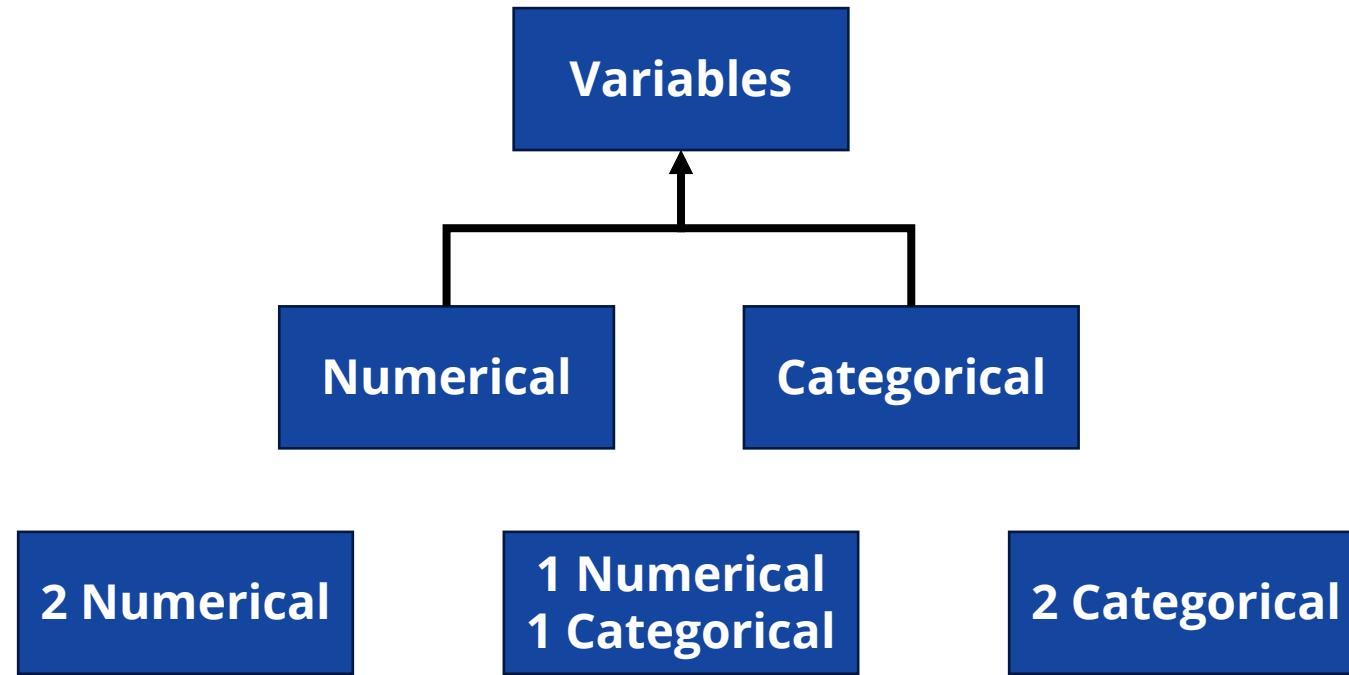
- Base population: 1000 samples from $U(1, 5)$



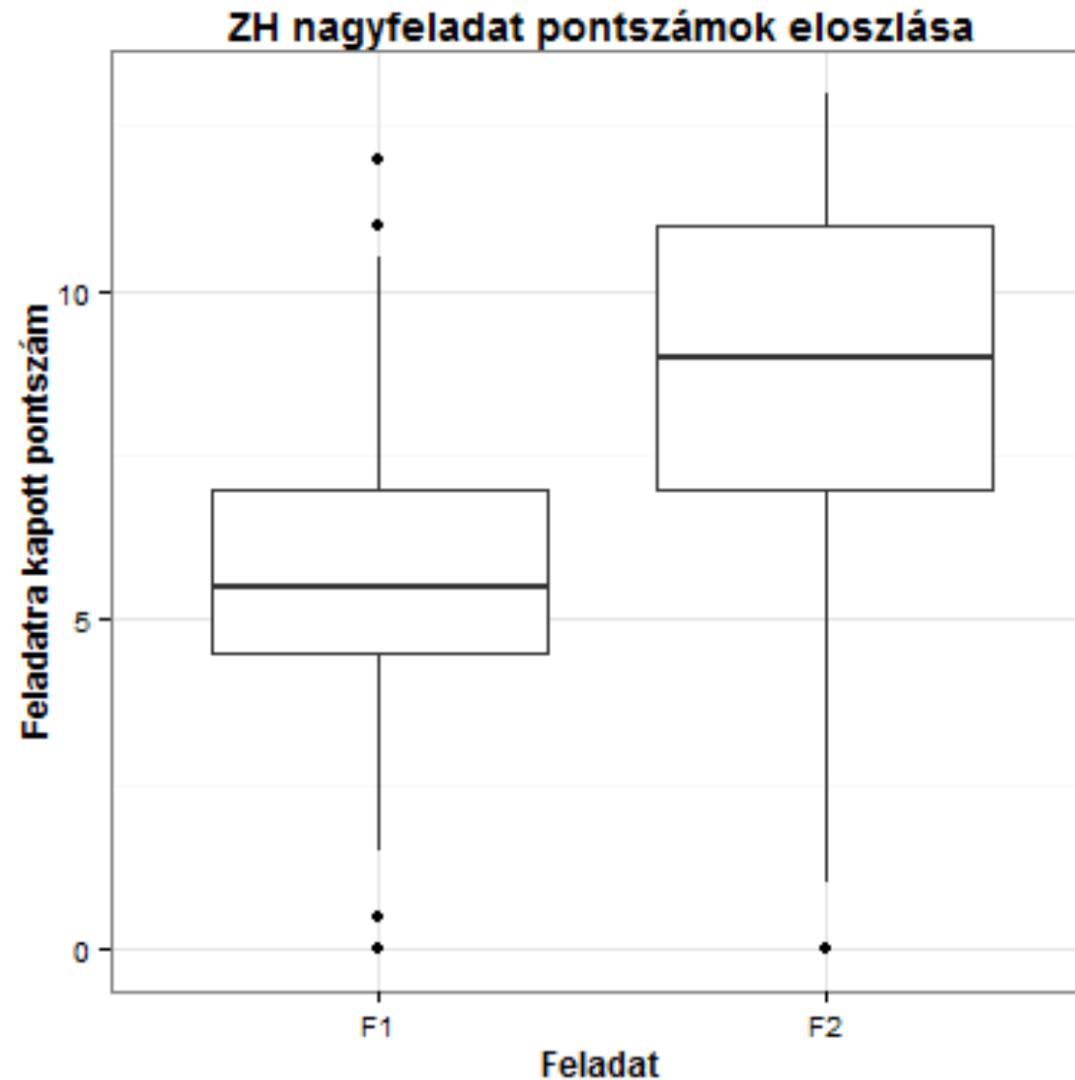
Some qualitative fun with boxplots



Visualization of 2 variables



Numerical by category

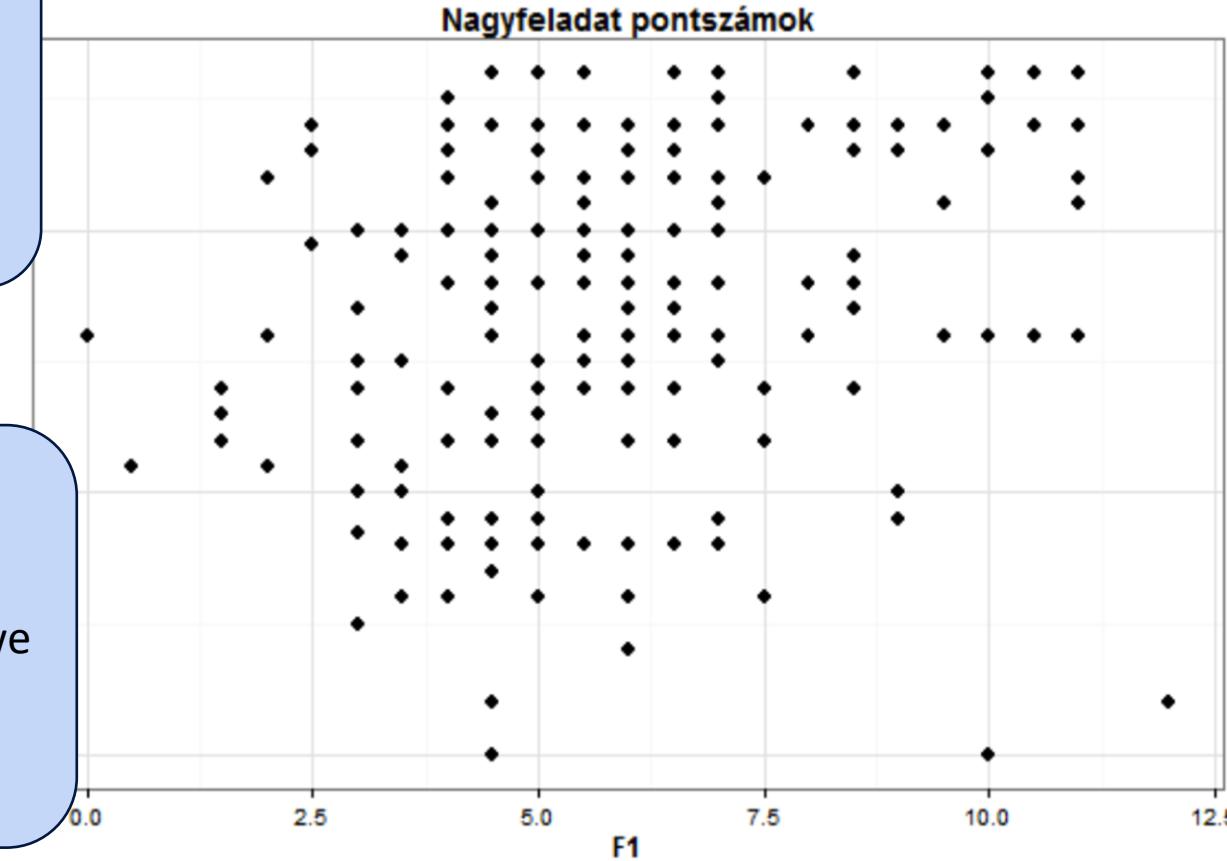


Scatterplot

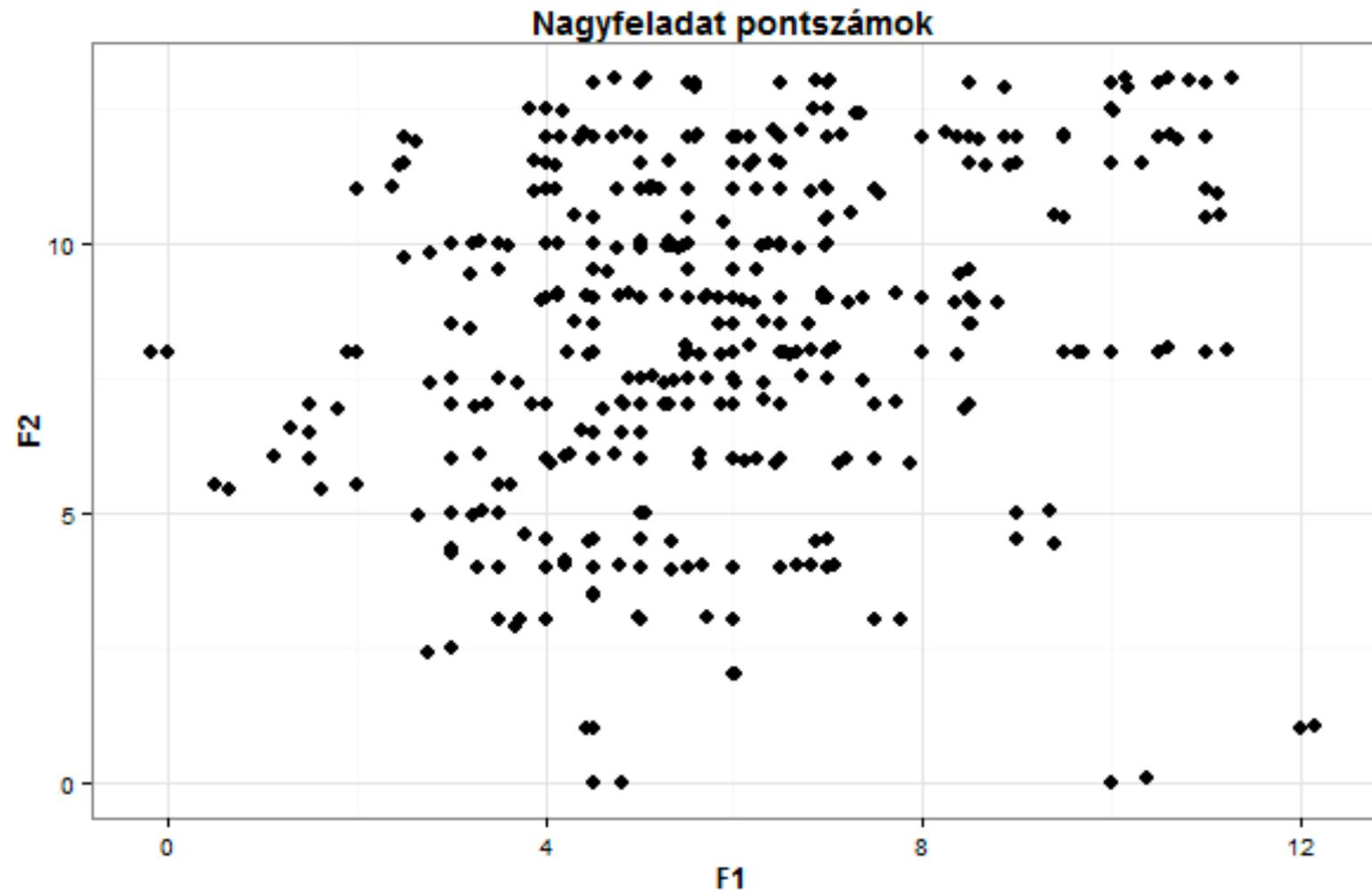
- **Input variable:** points received on tasks
- **Question:** how do they relate to each other?

We visualize pairs of points that occur together

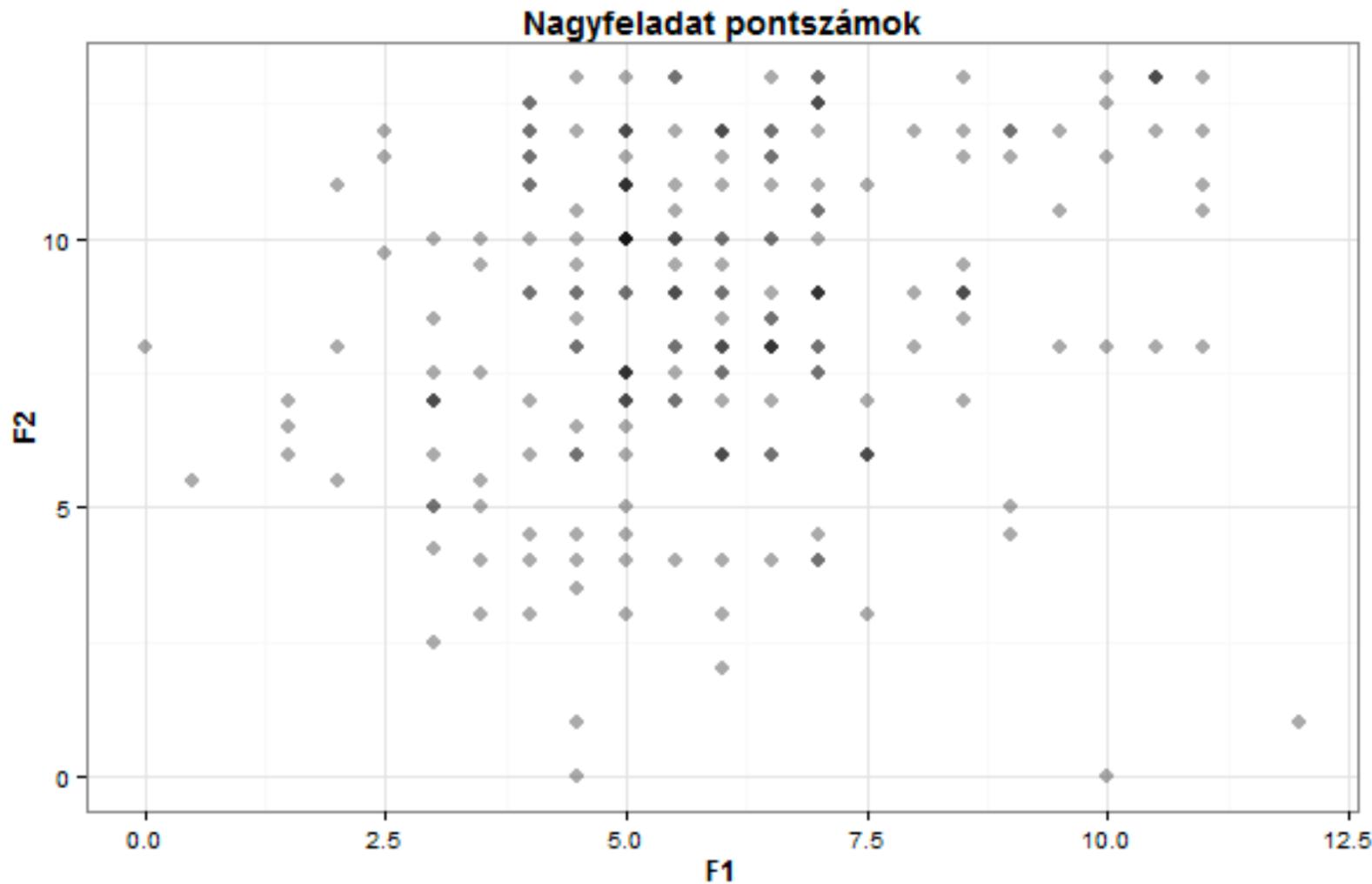
If the value of one variable is missing, we cannot plot it



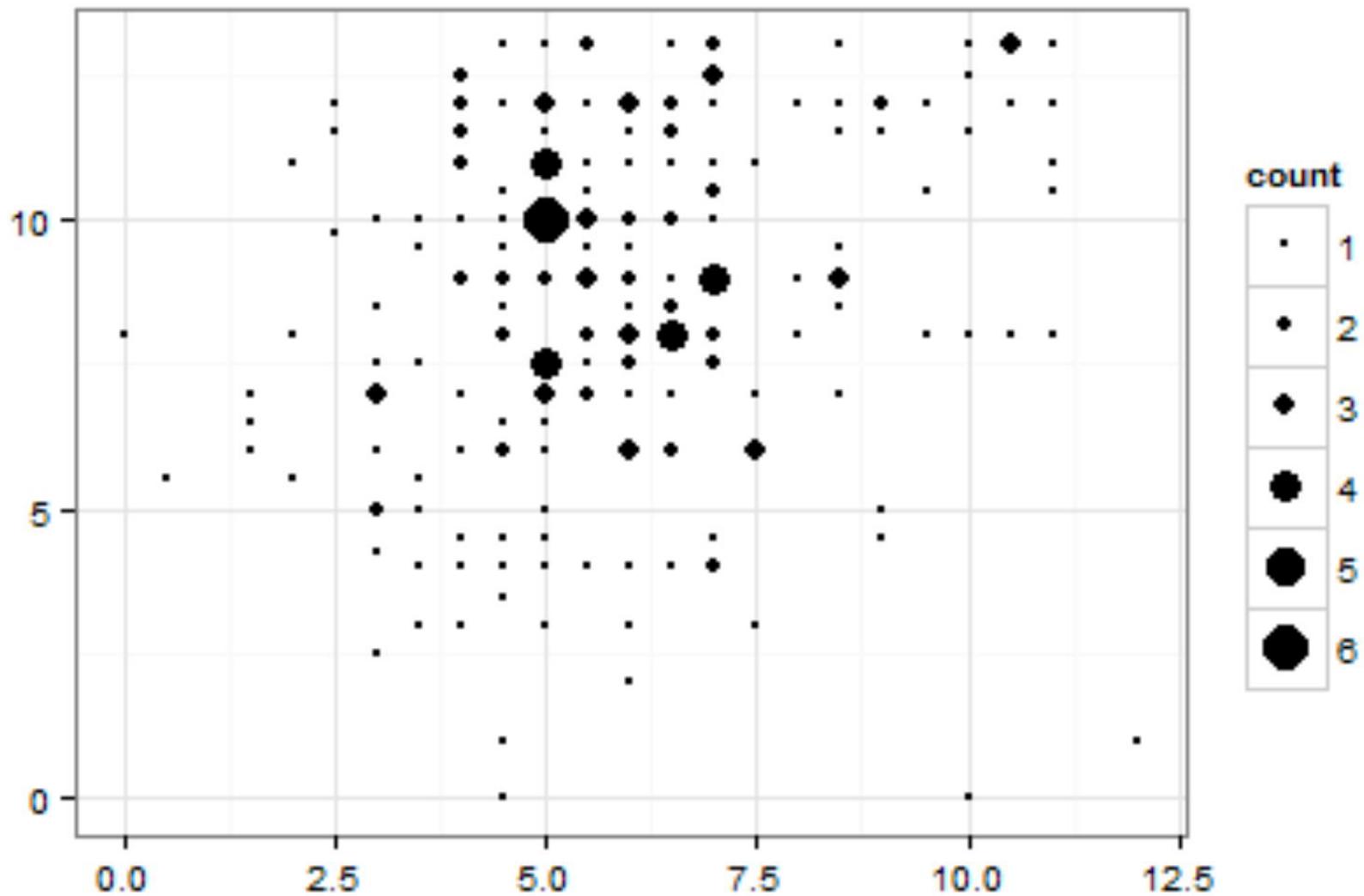
Jitter



Transparency

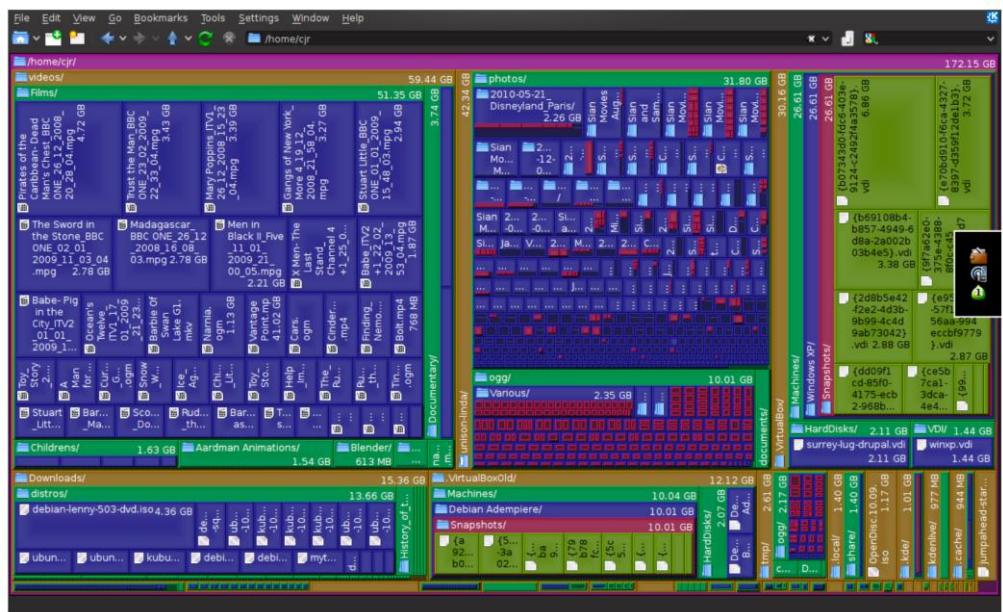
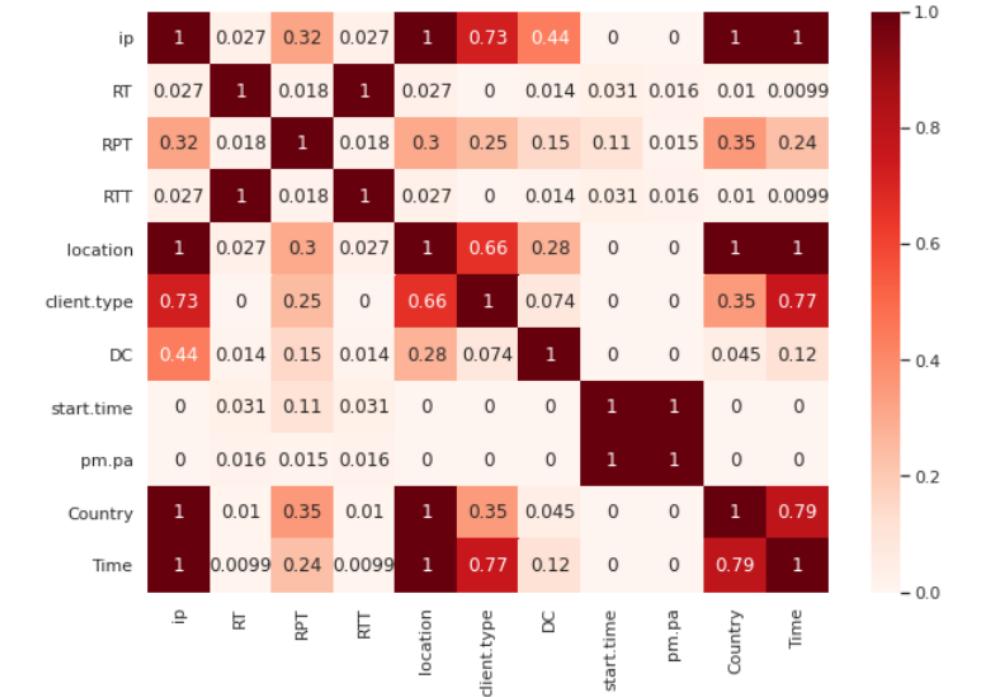


Size



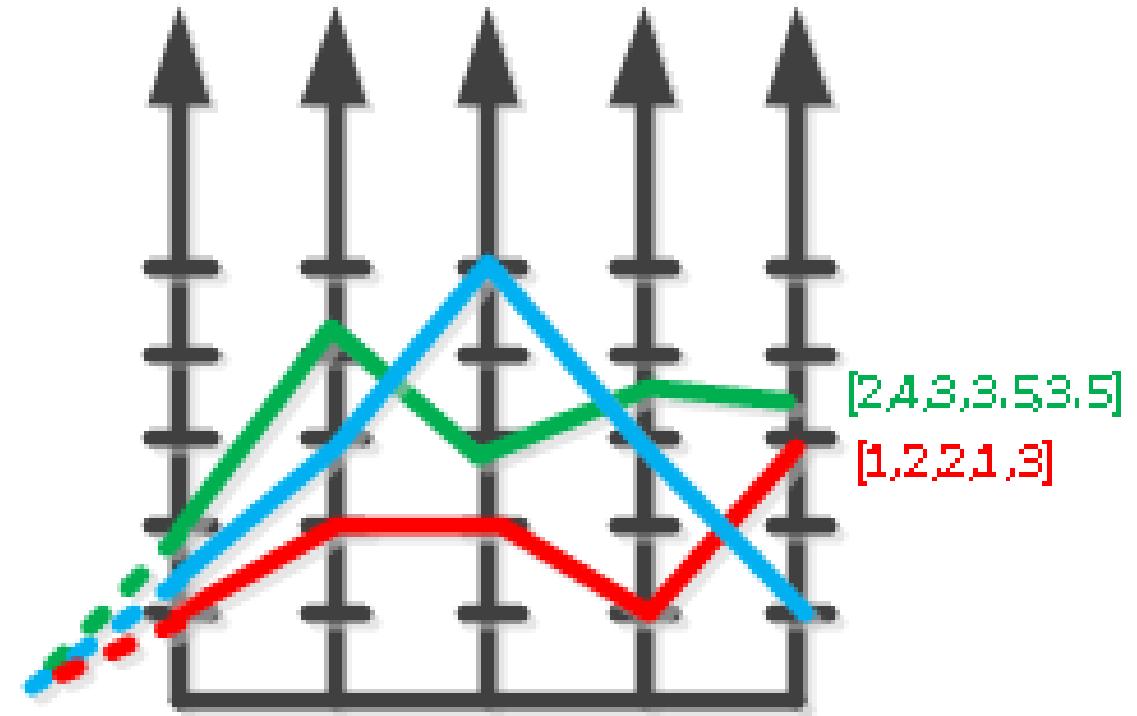
≥ 3 variable

- Modifying the attributes of graphical objects
 - Color
 - Size
 - Texture
 - Position – this may seem trivial, but it is important in a treemap
 - Pl. heatmap, treemap



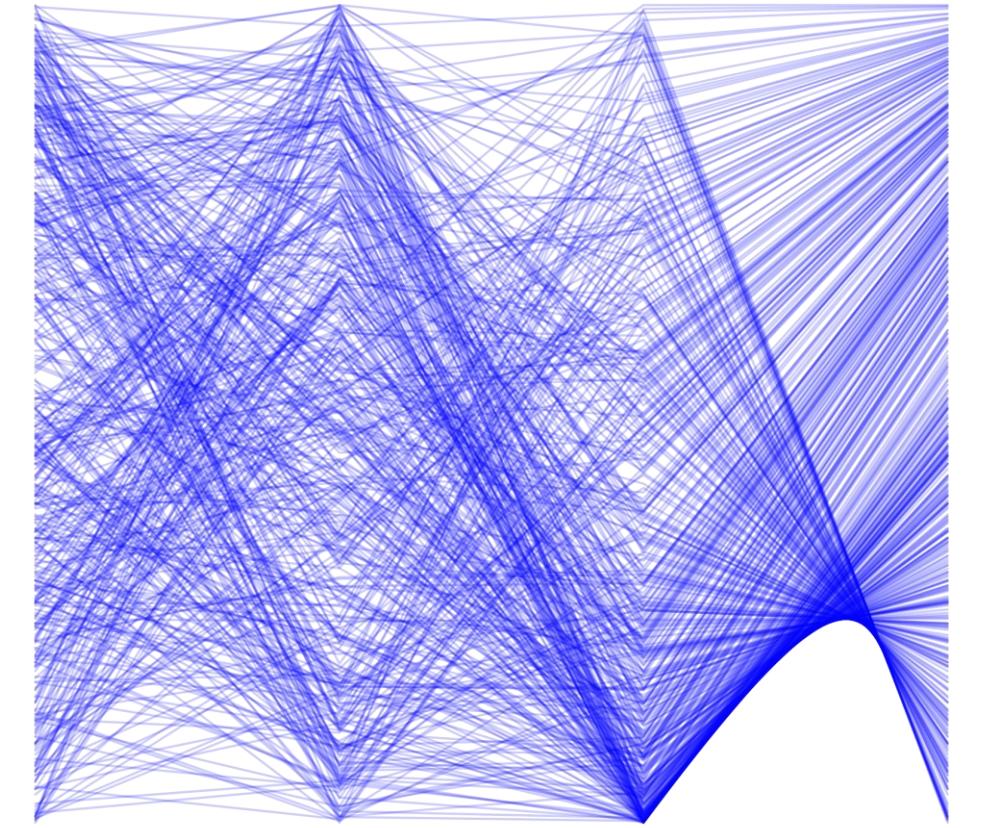
Parallel coordinates

- The trained eye can...
 - Remove redundancy
 - Recognize linear correlation
 - Coplanarity and other structures
 - Clustering
 - ...
- Variable order is crucial

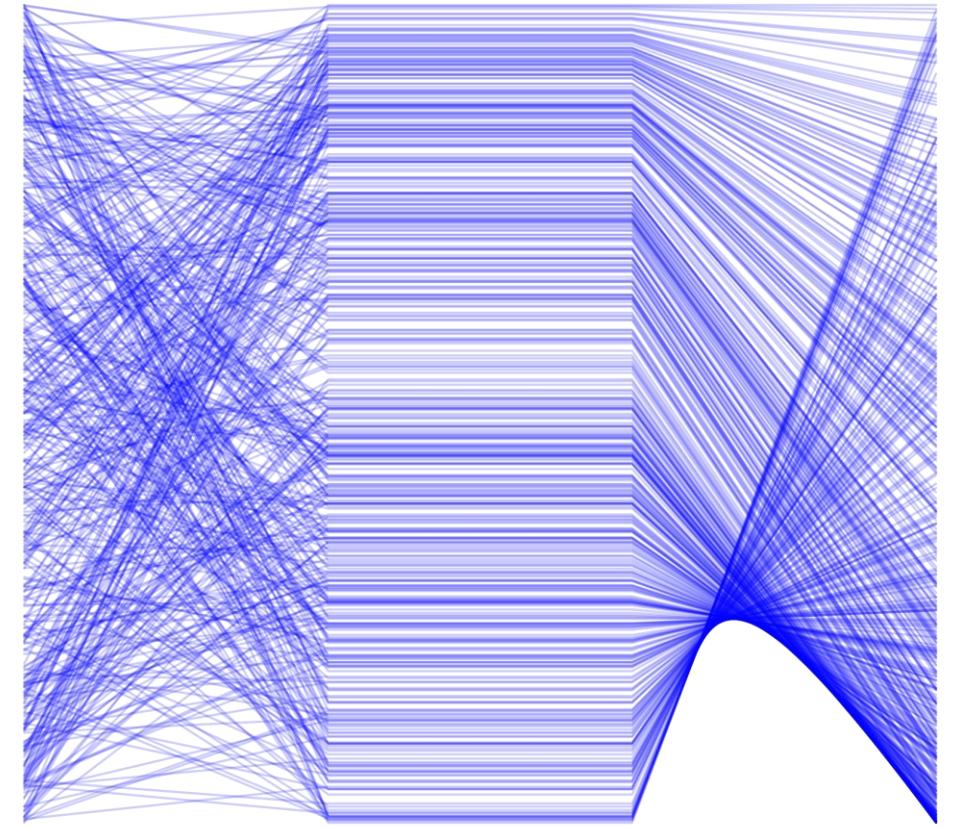


Inselberg, A: Parallel Coordinates, Visual Multidimensional Geometry and Its Applications, Springer 2009

The impact of axis ordering

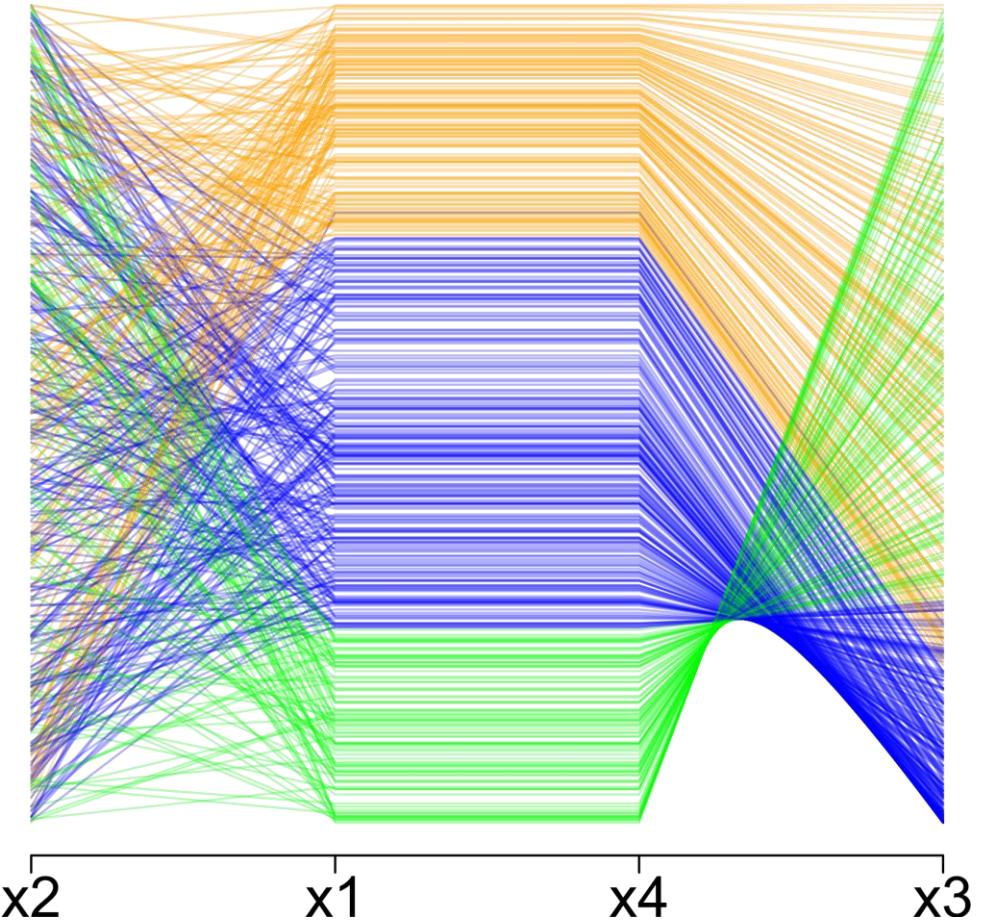
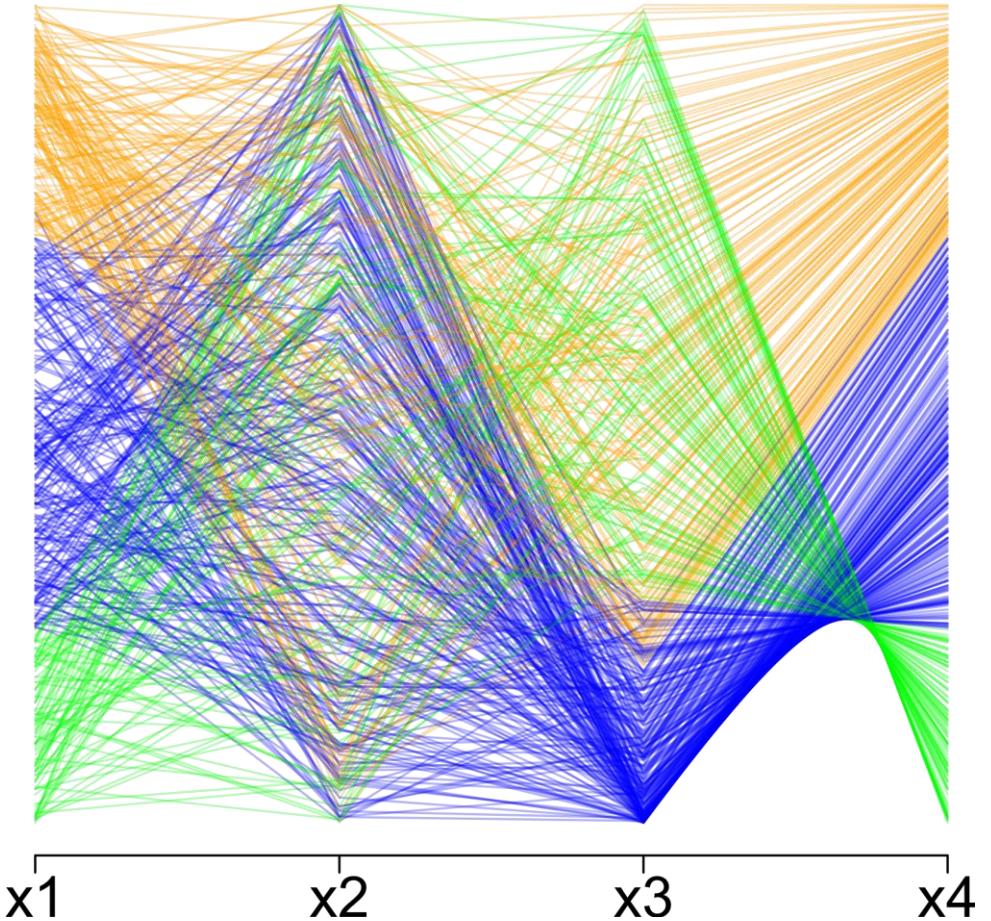


x1 x2 x3 x4

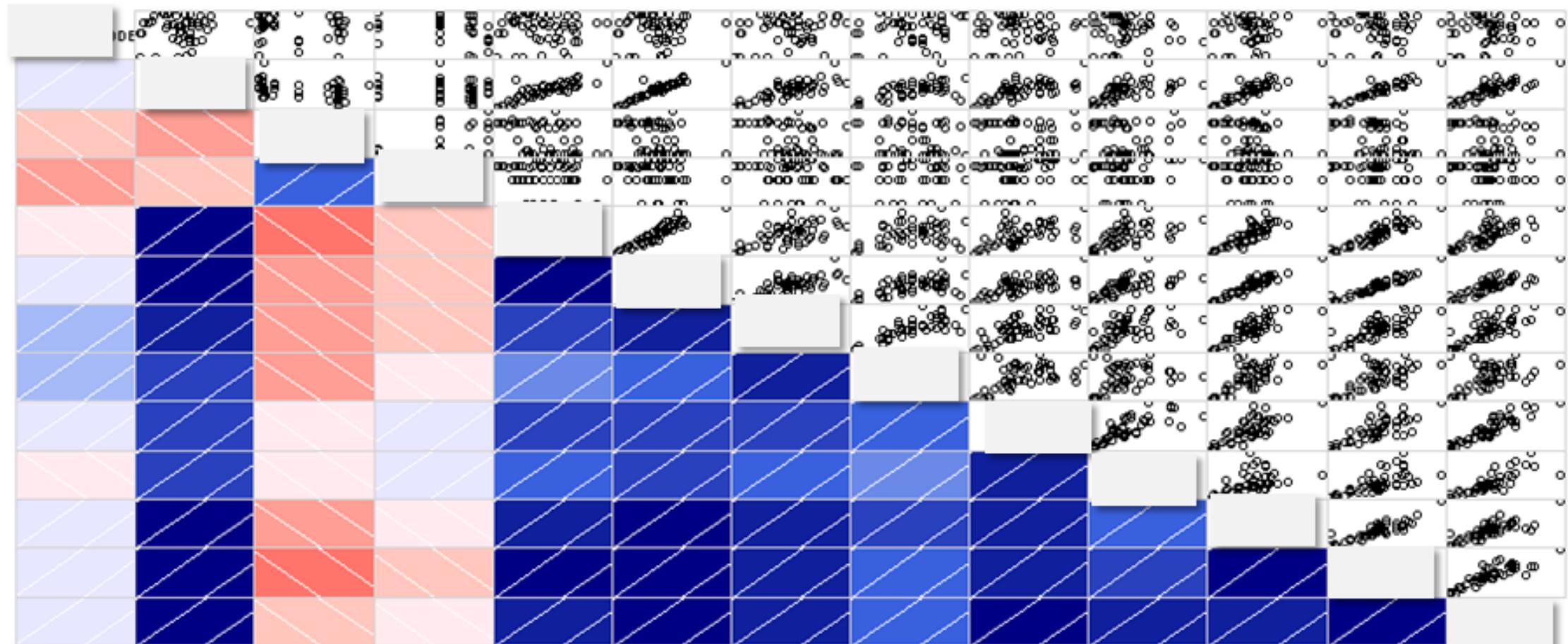


x2 x1 x4 x3

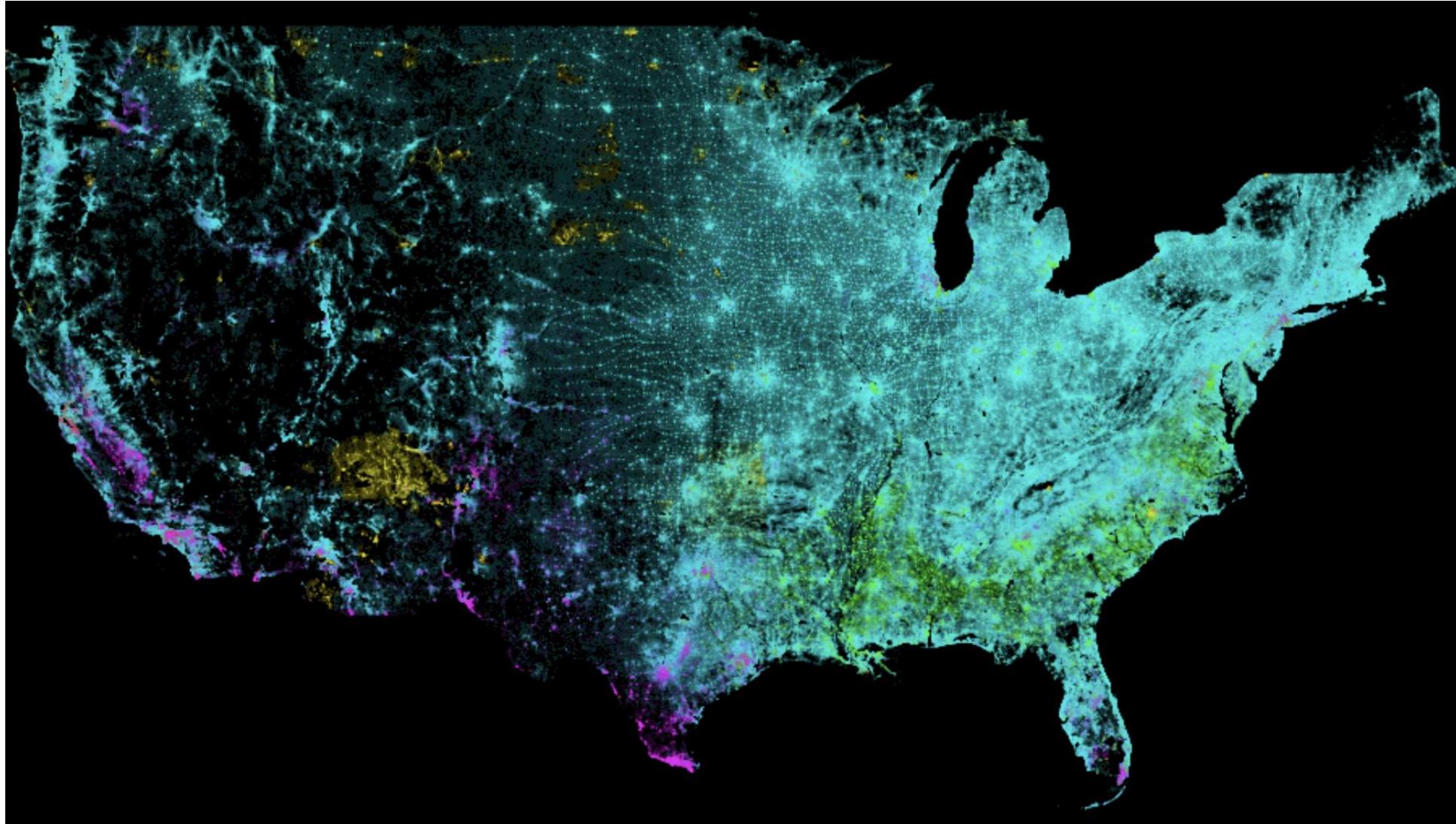
Model extraction



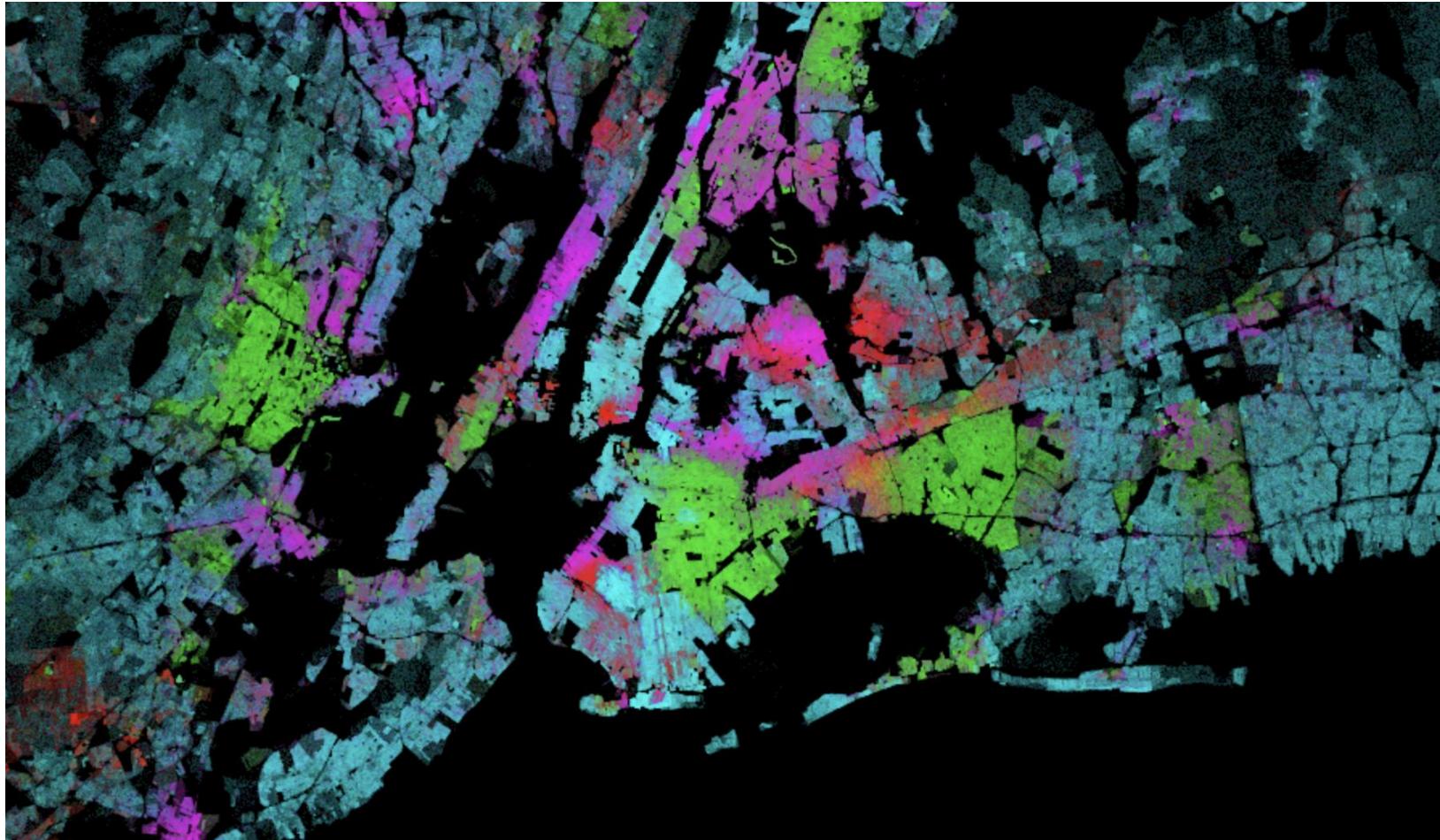
Correlation matrices/correlograms



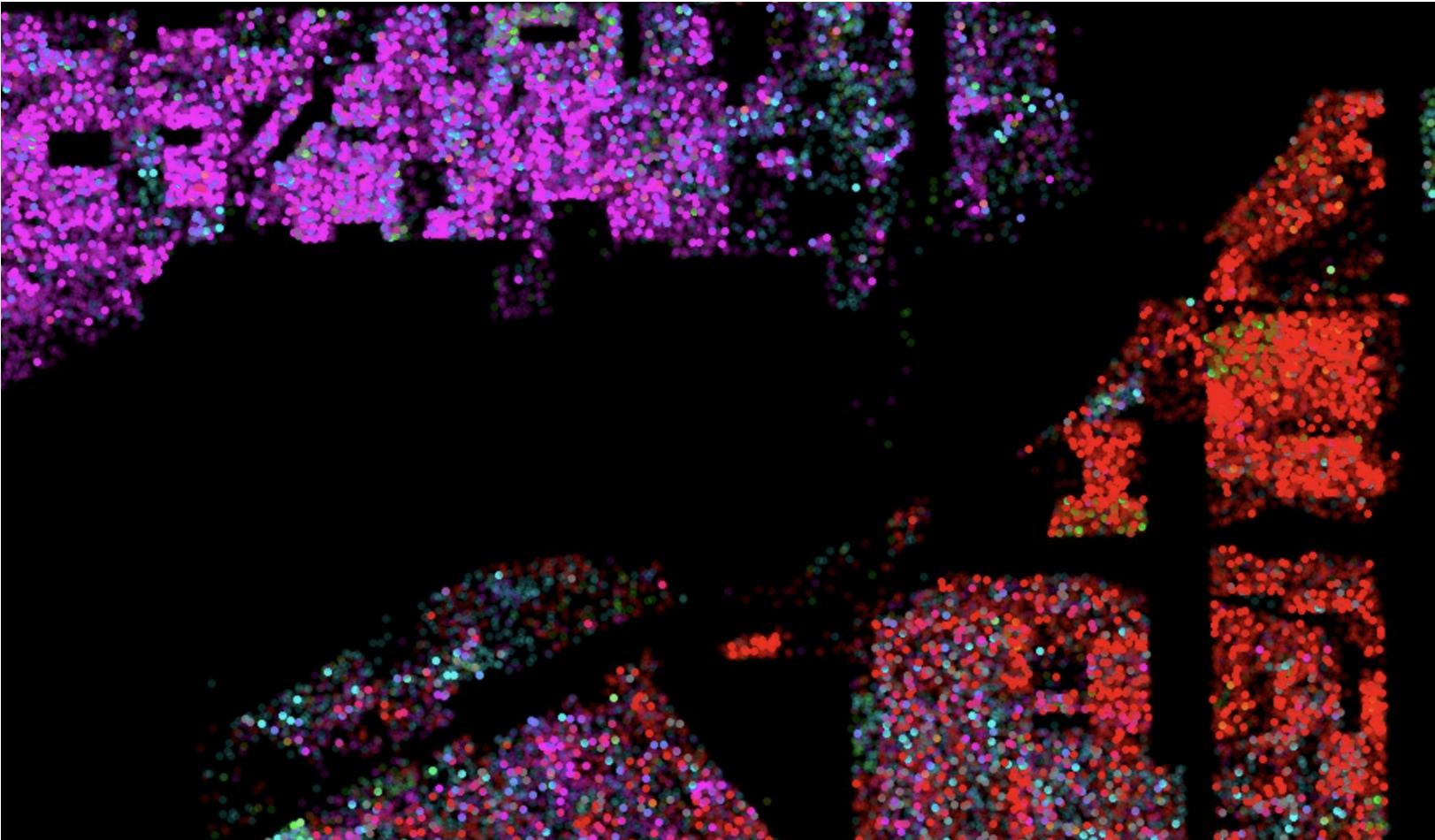
Big data



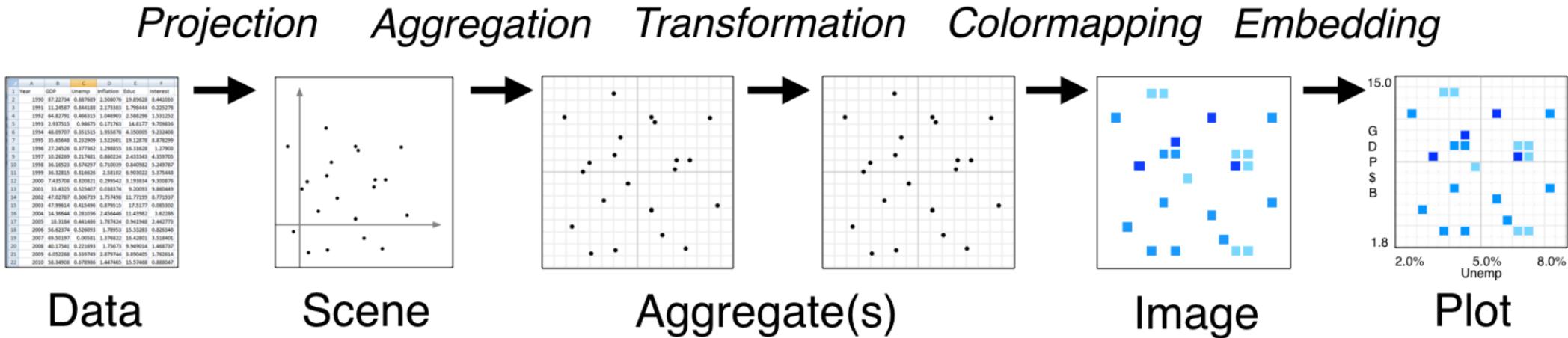
Big data



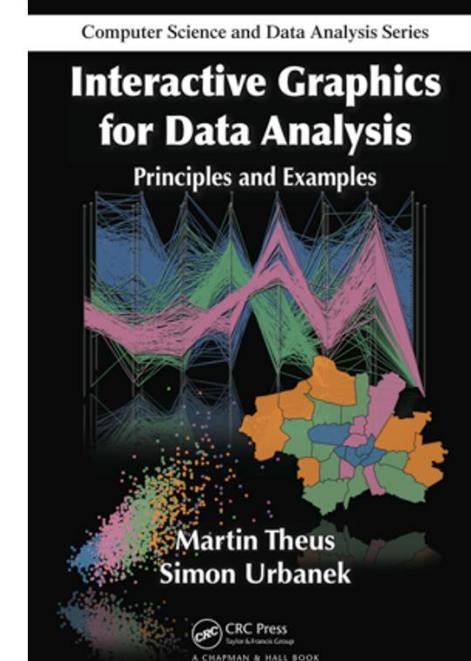
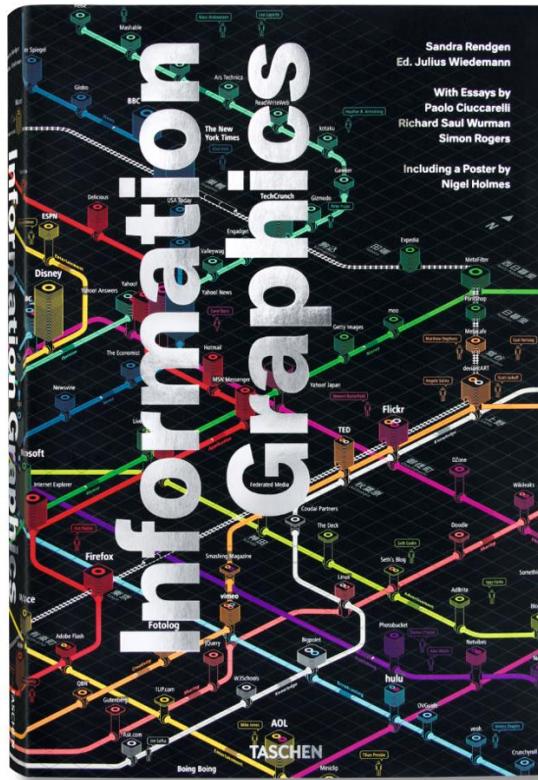
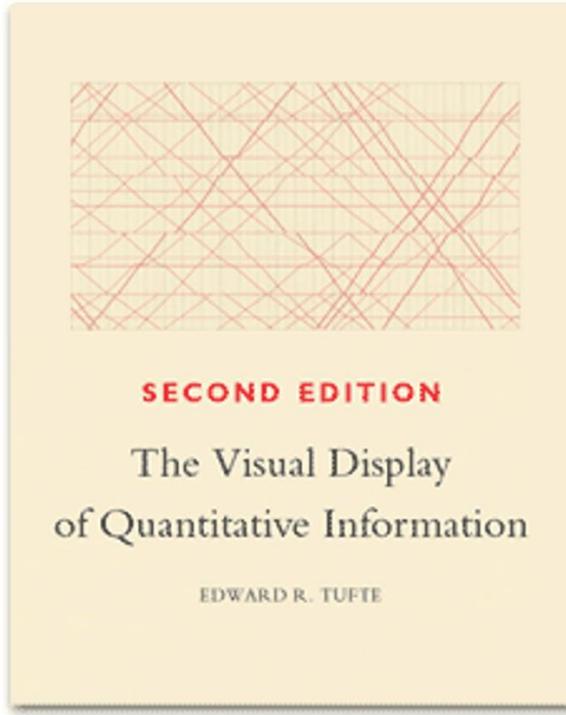
Big data



Big data visualization



Highly recommended learning sources



https://guides.library.duke.edu/datavis/vis_types

About Data Visualization

Visualization Types

Designing a Visualization

Helpful Tools and Tutorials

Learning About Visualization

Data analysis topics at ftsrg

- Model and knowledge-driven visual exploratory data analysis
- Extracting qualitative models from observations
- Data-driven model extraction (Process mining)
- Secure data sharing in data spaces
- Logic reasoning-based decision support
- Tool support for automated and visual data analysis
- LLM-based interactive data exploration with knowledge graph support