



<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

## Multivariate data visualization

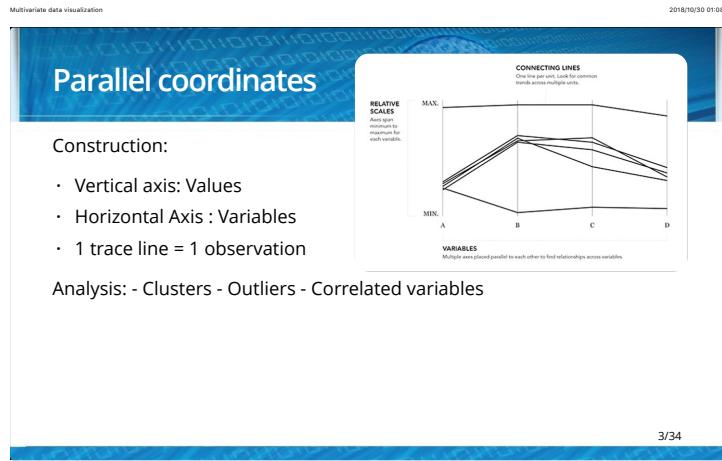
Continous variables involved in the following:

- Parallel coordinate plots
- Heatmaps
- Star charts

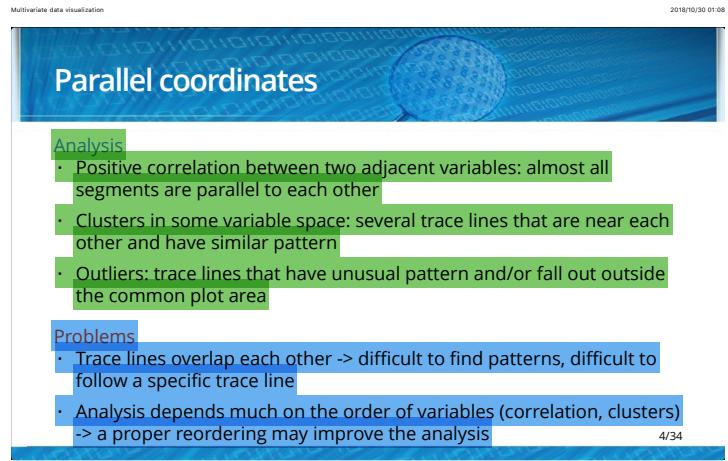
2/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 2 页 (共 34 页)



<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

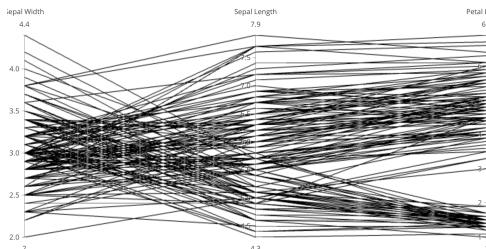


<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 4 页 (共 34 页)

## Parallel coordinates

**Example:** Iris dataset - How many clusters do you see?



5/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 5 页 (共 34 页)

## Ordering problem

- Problem of ordering (variables, observations) is one of the key problems in multidimensional visualization
  - Sometimes has a huge impact on perception (heatmaps)
  - A lot of approaches exist

**Problem formulation:** Given data set  $\chi = (x_{ij} | i = 1, \dots, n, j = 1, \dots, p)$

- Select order  $\Psi = i_1, \dots, i_p$  that optimizes visual perception (analysis) -> this defines reordering of data columns  $\Psi : \chi \rightarrow \chi'$

**Note:**  $p!$  possible orderings exist...

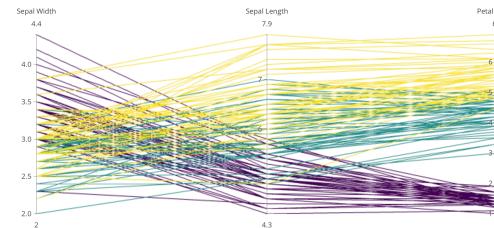
7/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 7 页 (共 34 页)

## Parallel coordinates

- Sometimes clusters overlap with categories given by some variable
  - Non-mixing groups is not the same as clustering!



6/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 6 页 (共 34 页)

## Ordering problem

### Solution

- early approaches (for ex. Ankerst et al. 1998):
  - Choose a distance (proximity) matrix  $D = \{d_{ij} = d(x^i, x^j)\}$  between variables (columns)
    - Euclidian distance on scaled columns
    - 1- correlation
  - This defines graph with vertices  $1, \dots, p$  and edge weights  $d_{ij}$  -> Hamiltonian path (Traveling Salesman Problem)

$$\min_{\Psi} \sum_{j=1}^{p-1} d'_{j,i+1}$$

TSP is NP-complete  $\rightarrow$  Approximate solutions are used.

8/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 8 页 (共 34 页)



## Ordering problem

**Solution:** modern approaches

- Based on:
  - Decreasing visual clutter
  - Clustering data points/dimensions
  - Outlier detection
  - Dimensionality reduction (for ex. MDS)
  - ...
- **Note:** most of these can be applied both for ordering observations and ordering variables
  - Just transpose the data matrix...

9/34



## Gradient measures

**Aim:** distances should increase from diagonal

$$d_{ik} \leq d_{ij} \text{ for } 1 \leq i < k < j \leq n$$

$$d_{kj} \leq d_{ij} \text{ for } 1 \leq i < k < j \leq n$$

**Objective function:**

$$L(D) = \sum_{i < k < j} f(d_{ik}, d_{ij}) + f(d_{kj}, d_{ij})$$

where

$$f(z, y) = \text{sign}(z - y) \text{ or } f(z, y) = z - y$$

11/34



## Ordering problem

**Objective functions:**

- Gradient measures (anti-Robinson)
- Hamiltonian path length
- Least squares
- ...

They based on  $\min_{\Psi} L(\Psi(D))$

**Optimization algorithms:**

- Partial enumeration
- Traveling salesman solvers

*Hierarchical clustering*

10/34



## Other objectives

**Hamiltonian path length:**

$$L(D) = \sum_{i=1}^{n-1} d_{i,i+1}$$

**Least squares criterion (PCA)**

- Solution is similar to first PCA component

$$L(D) = \sum_i \sum_j (d_{ij} - |i - j|)^2$$

12/34

# Optimization algorithms

## Partial enumeration methods

- Ex: Branch and bounds and dynamic programming
- Constructing solutions by parts

## TSP solver

- Suitable for hamiltonian path objective
- Find shortest path by dynamic programming or heuristics

13/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

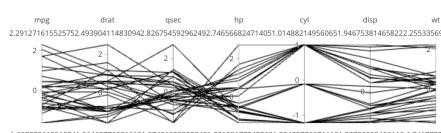
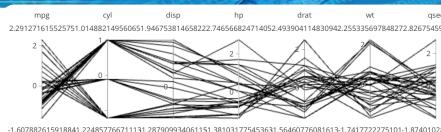
第 13 页 (共 34 页)

14/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 14 页 (共 34 页)

# Effect of ordering



15/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 15 页 (共 34 页)

16/34

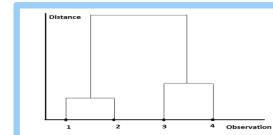
<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 16 页 (共 34 页)

# Optimization algorithms

## Hierarchical clustering

- Observations are joined into clusters
- Clusters are joined in larger clusters
- Until only one cluster left
- Leaves and branches are permuted to minimize given objective



14/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 14 页 (共 34 页)

# Heatmaps

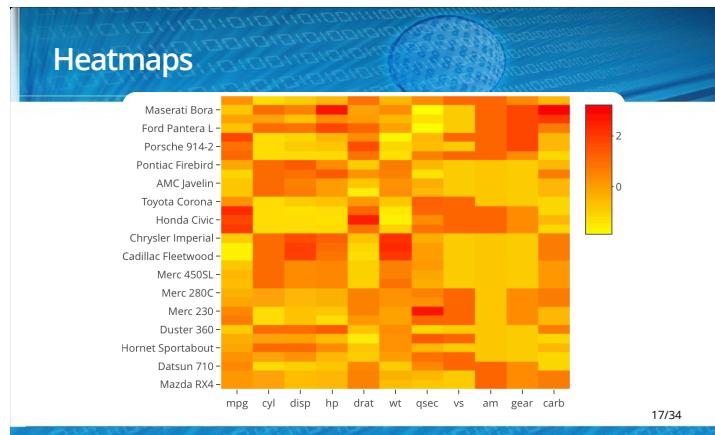
A heat map visualizes a matrix [ n x m ]

- Normally rows=observations, columns= parameters
- Heatmap has the corresponding size
- Each cell of the matrix corresponds a cell in the heatmap
- High values correspond intense colors in this map (or visa versa for other color schemes!)
- Names of variables and observations are shown

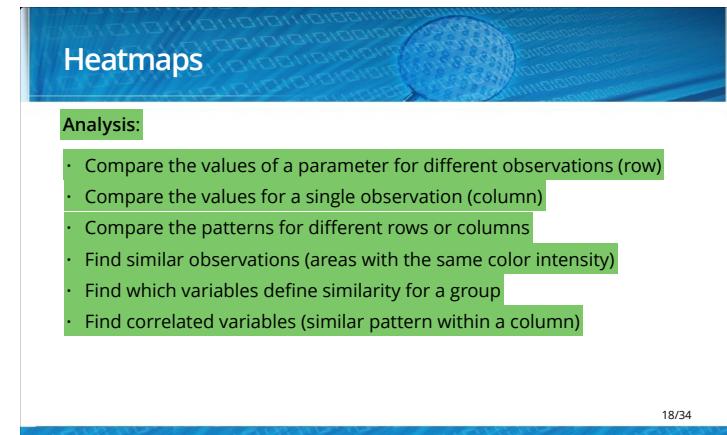
16/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 16 页 (共 34 页)



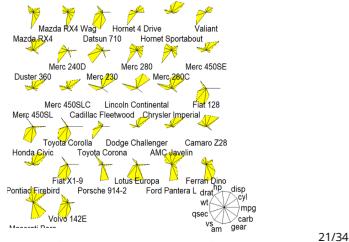
- Exercise (last picture):
  - How many clusters do you see?
  - Which variables define clusters?
  - Which variables are correlated?



- Gradient measure objective used
- See new analysis possibilities

## Radar charts

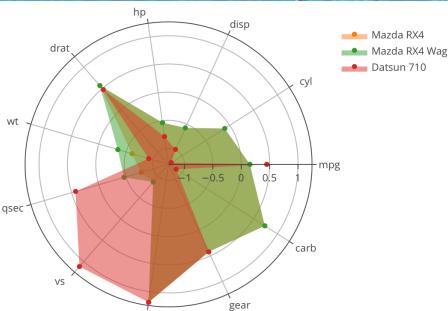
- Use polar coordinate system
- Map column value as a coordinate in certain direction



<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 21 页 (共 34 页)

## Radar charts



<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 23 页 (共 34 页)

## Radar charts

If juxtaposed, analyse:

- Clusters
- Outliers
- Outlying directions

If superimposed,

- Comparing variable length
- Seeing similar and outlying observations

22/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 22 页 (共 34 页)

## Radar charts

Problems:

- Difficult to judge orientations
- Number of dimensions are observations is very limited
  - Number of observations is extremely limited if superimposed
- More close radar charts easier to compare
- Perception is much affected by observation ordering

Ordering:

- Same as before plus
- Dimensions can be sorted to promote more symmetric charts

24/34

<https://www.ida.liu.se/~732A98/info4/Lecture4.html#1>

第 24 页 (共 34 页)



- Now with reordering by Gradient Measures



25/34



Idea:

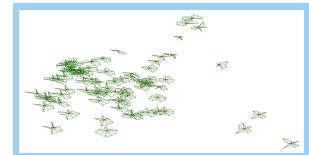
1. Make same kind of plot for subsets of data
2. Plot together
3. See patterns/differences

Analogy: cutting a sausage

27/34

## Radar charts

Other positioning possible - PCA/MDS

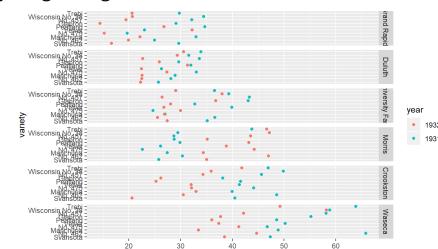


26/34

## Trellis plots / facets

- Example: Barley data

- Anything strange?



28/34



- Faceting = one more aesthetics
- What can be analysed?
  - Patterns within/between plots
  - Conditional dependence  $Y \sim X|Z$
  - Variable interaction, additivity
- > Useful tool for modeling!
- Compare : 3D- scatter plots

29/34

30/34



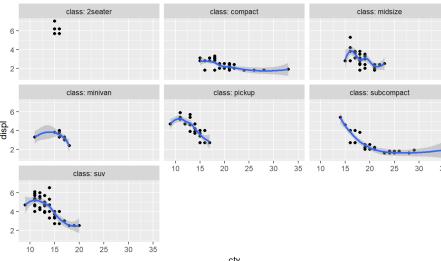
- Design issues:
    - How to order rows/columns in trellis?
      - A: X=one var, Y=another var (`facet_grid`)
      - B: independently of aes (`facet_wrap`)
    - How to handle categorical vars?
      - One value/panel
      - Group
      - Ordering? (R: decide factor levels)
    - How to handle real-valued vars?
      - Split equal size/length
- Shingles*

31/34

32/34



- Another car data: is there additivity?

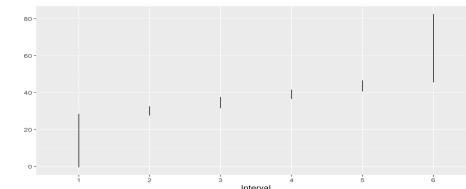


30/34

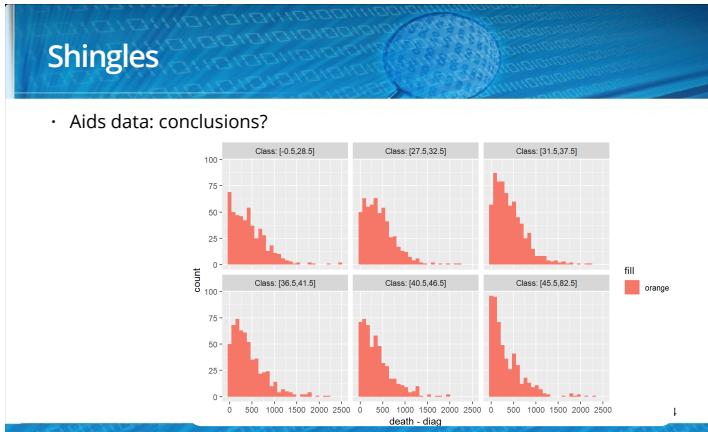


- Creates overlap
- To avoid boundary effects

**Example:** AIDS data (Age, Time of Death, Time of Diag)



32/34



## Read at home

- Chapter 5
- Paper "Hahsler, M., Hornik, K., & Buchta, C. (2008). Getting things in order: an introduction to the R package seriation. *Journal of Statistical Software*, 25(3), 1-34".
- (Browse through) paper "Ankerst, M., Berchtold, S., & Keim, D. A. (1998, October). Similarity clustering of dimensions for an enhanced visualization of multidimensional data. In *Information Visualization, 1998. Proceedings. IEEE Symposium on* (pp. 52-60). IEEE."
- Becker, R. A., Cleveland, W. S., & Shyu, M. J. (1996). The visual design and control of trellis display. *Journal of computational and Graphical Statistics*, 5(2), 123-155.