

Cluster analysis of time series data

Tomáš Bartoň

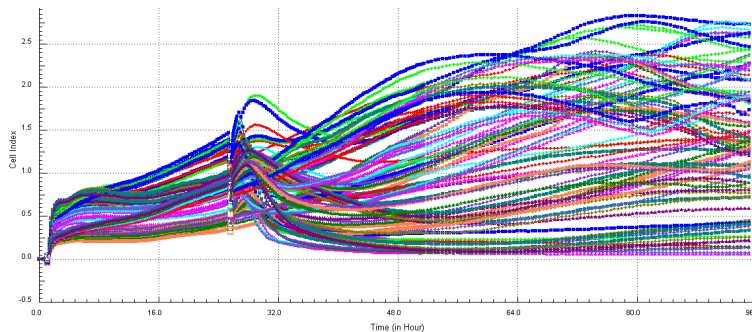
Supervisor: **Ing. Pavel Kordík, Ph.D.**

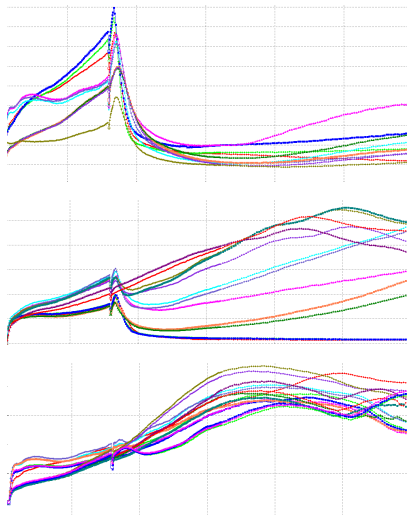
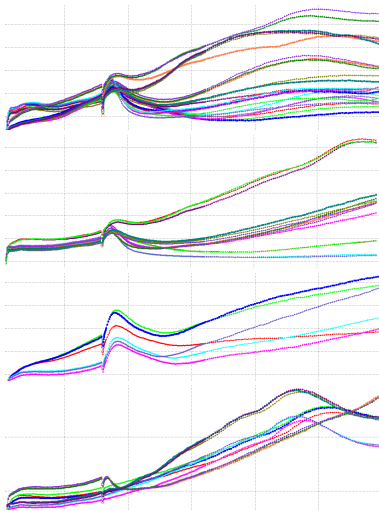
Department of Theoretical Computer Science
Faculty of Information Technology
Czech Technical University in Prague

January 5, 2012

The Problem

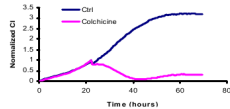
- find suitable method for identification of patterns
- assign samples into (unknown) groups



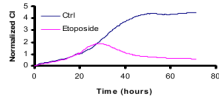


Goals

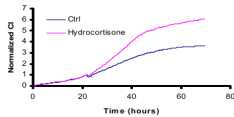
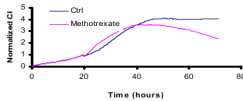
- capture global trends
- absolute values (sometimes) doesn't matter
- signals are not periodical
- discover unknown patterns



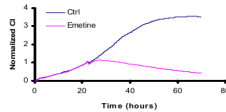
Anti-mitotic



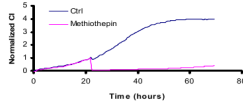
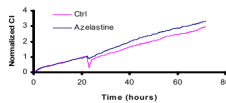
DNA Damaging



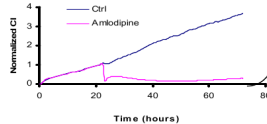
Nuclear Receptor modulator



Protein Synthesis Inhibitors



Calcium Modulators



Phases of clustering process

- 1 data cleaning
- 2 data integration
- 3 data selection
- 4 data transformation
- 5 clustering
- 6 pattern evaluation
- 7 knowledge representation

Clustering

No “correct” clustering exists

Definition

“Those methods concerned in some way with the identification of homogeneous groups of objects”

[Arabie et al., 1996]

Definition

“A cluster is a set of entities that are alike, and entities from different clusters are not alike” [Everitt, 1993]

- clustering can be used for understanding data

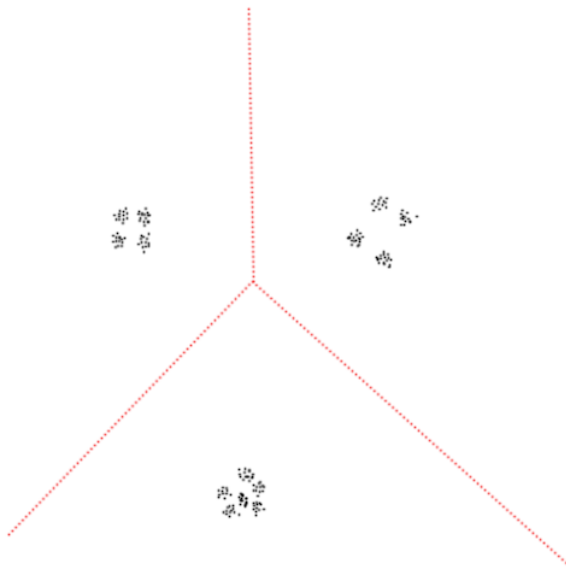


- to perform clustering you need to understand data

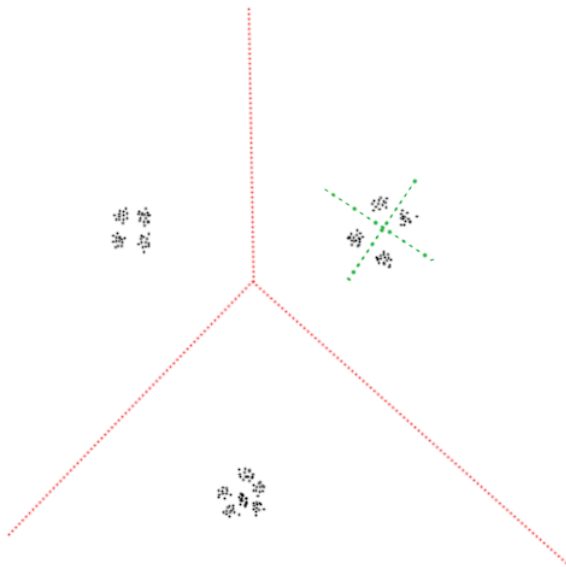


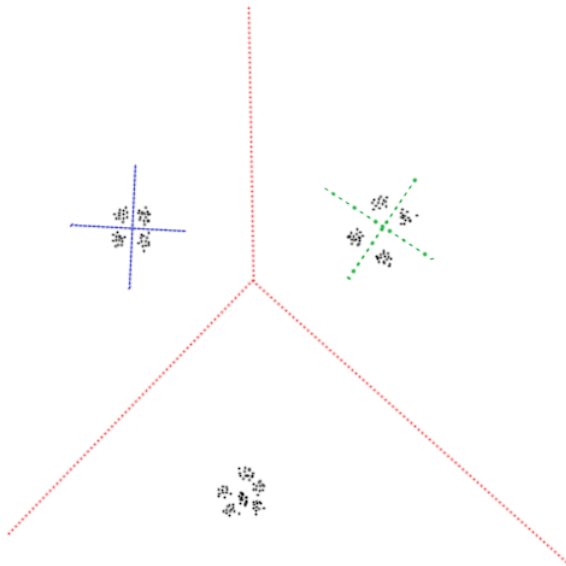
Determine number of clusters



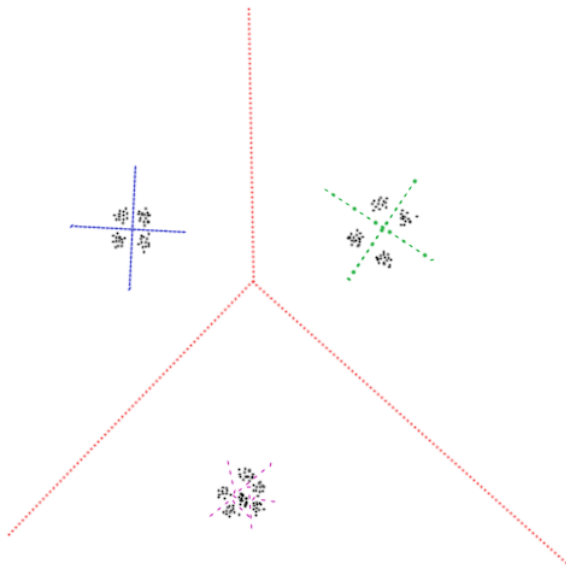
$k = 3$ 

$$k = 6$$

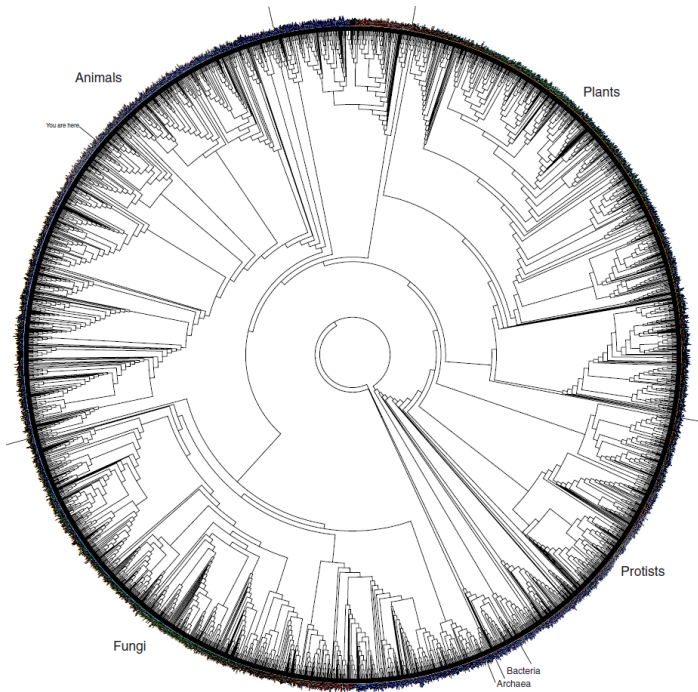


$k = 9$ 

$$k = 14$$



Determining the number of clusters in a data set is challenging [Mufti et al., 2005]





from Chinese encyclopedia Heavenly Emporium of Benevolent Knowledge. Animals are divided into [Borges, 1952]:

- those that belong to the emperor
- embalmed ones
- those that are trained
- suckling pigs
- mermaids
- fabulous ones
- those that are included in this classification
- innumerable ones
- etcetera

Clustering is ill-defined [Caruana et al., 2006]

All we care about is the “usefulness” of the clustering for achieving our final goal [Guyon et al., 2009]

Time series

Problem

- sensitive to small changes
- sum of distance does not capture shape of curve
- computationally expensive
- redundant information

Autoregressive model

- predict an output of a system based on the previous outputs

$$X_t = c + \sum_{i=1}^p \varphi_i \cdot X_{t-i} + \epsilon_t$$

- φ_i – parameters of the AR model
- X_t – amplitude of the signal
- ϵ_t – white noises

Moving average

$$X_t = \mu + \epsilon_t + \sum_{i=1}^q \phi_i \cdot \epsilon_{t-i}$$

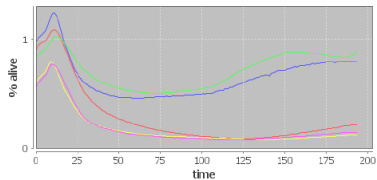
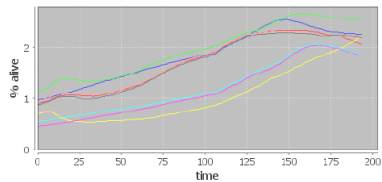
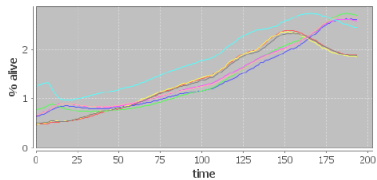
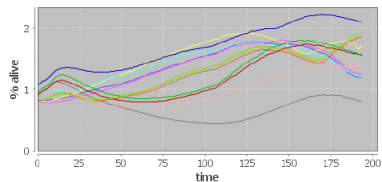
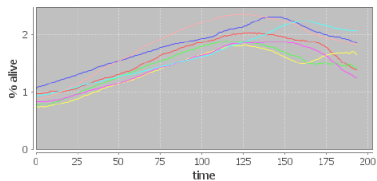
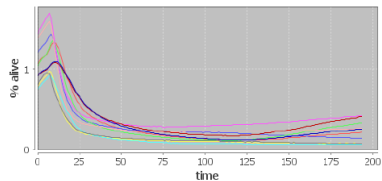
- ϕ_i – parameters of the AR model
- μ – expectations of X_t (often assumed to equal 0)
- ϵ_t – white noises

Autoregressive–moving-average model

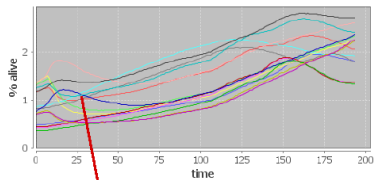
putting all together:

$$X_t = c + \epsilon_t + \sum_{i=1}^p \varphi_i \cdot X_{t-i} + \epsilon_t + \sum_{i=1}^q \phi_i \cdot \epsilon_{t-i}$$

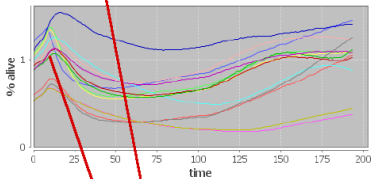
- ARMA(p, q) refers to the model with p autoregressive terms and q moving-average terms
- in Matlab function `armax[Time-domain, data object]`

Cluster 1**Cluster 4****Cluster 2****Cluster 5****Cluster 3****Cluster 6**

Cluster 7

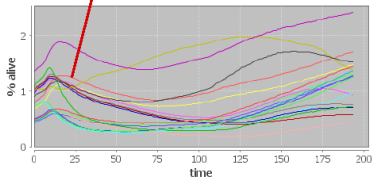


Cluster 8

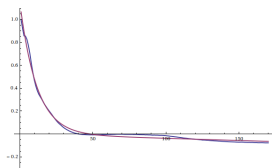
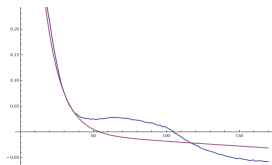
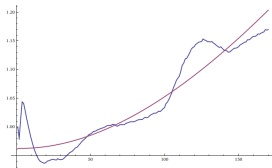


these should be in one cluster

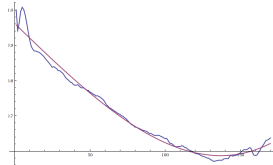
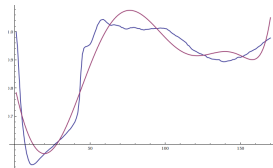
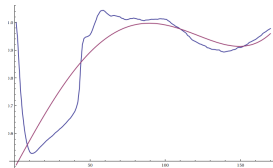
Cluster 9



Exponential

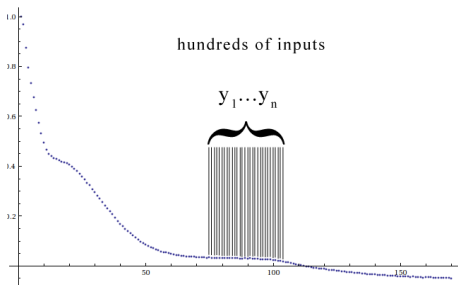


Polynomial



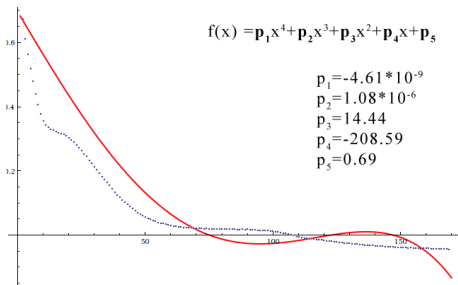
Representation of inputs

Measured values



- too many inputs
- does not represent patterns

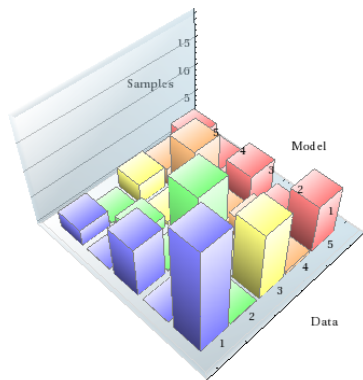
Approximated model



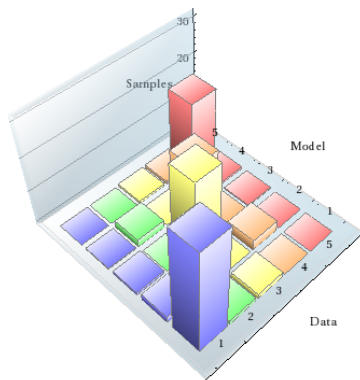
- only 5 parameters describing whole curve
- represent patterns

How many parameters do we need?

2 parameters



5 parameters



Which parameters to choose?

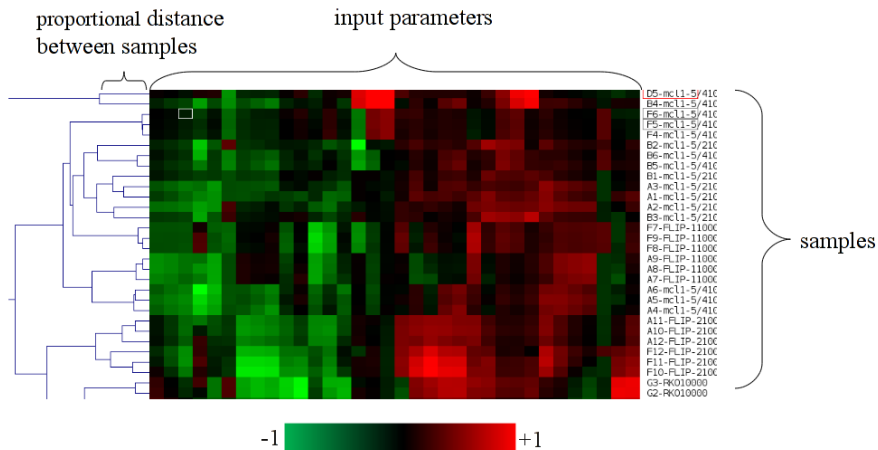
on previous slide input parameters were following:

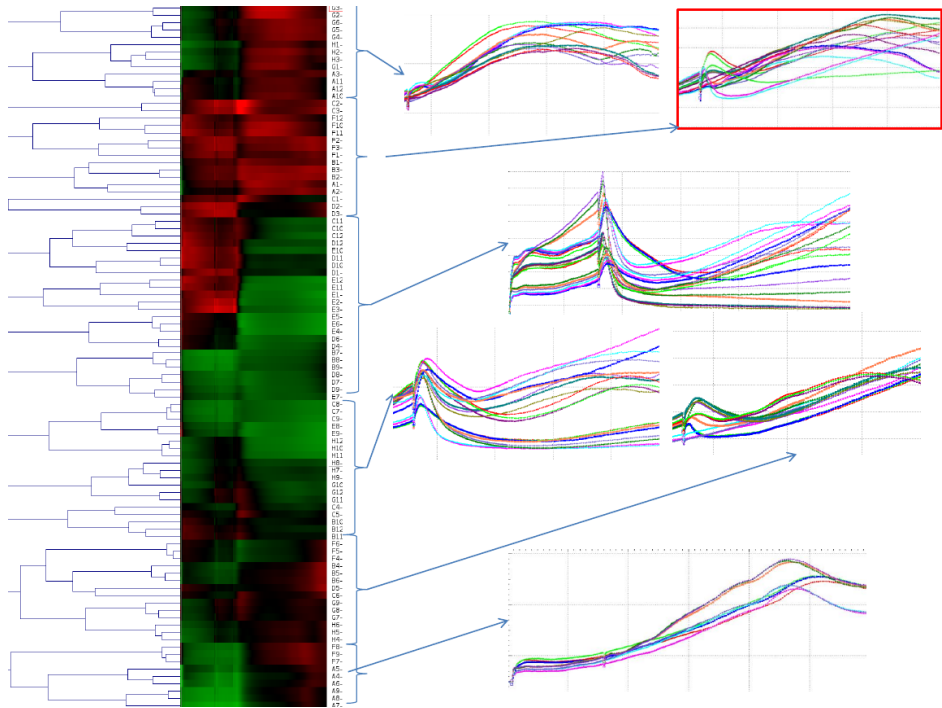
- mean
- minimum
- maximum
- linear coefficient
- quadratic coefficient

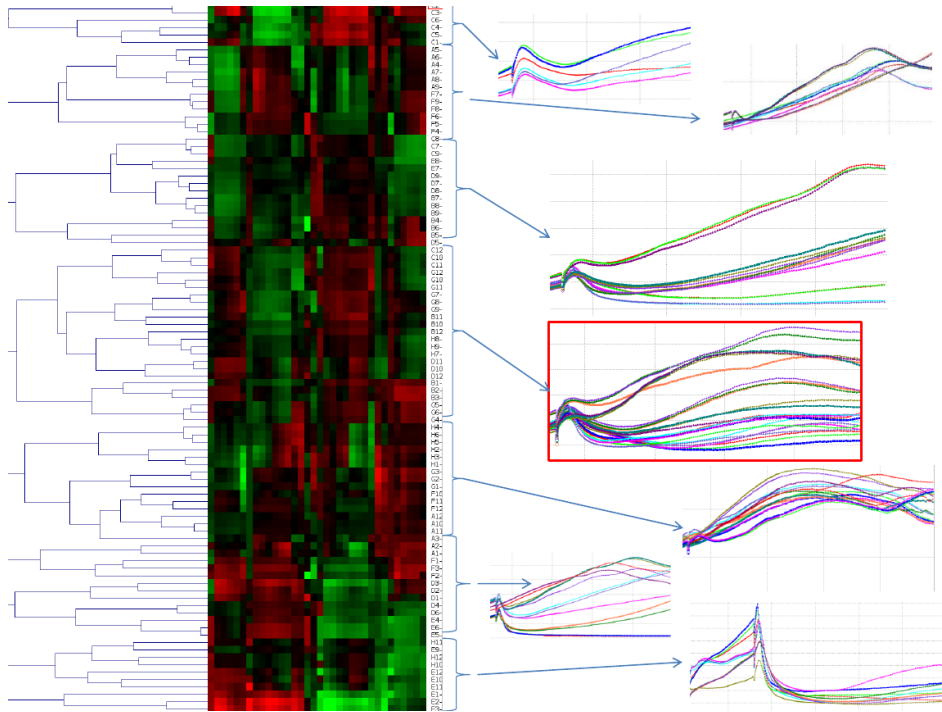
for EEG clustering is [Siuly et al., 2011] using:

- minimum
- maximum
- mean
- median
- modus
- first quartile
- third quartile
- inter-quartile range
- standard deviation

Dendrogram







C-Index

The C-index was reviewed in Hubert and Levin [1976]

$$p_{c-index} = \frac{d_w - \min(d_w)}{\max(d_w) - \min(d_w)}$$

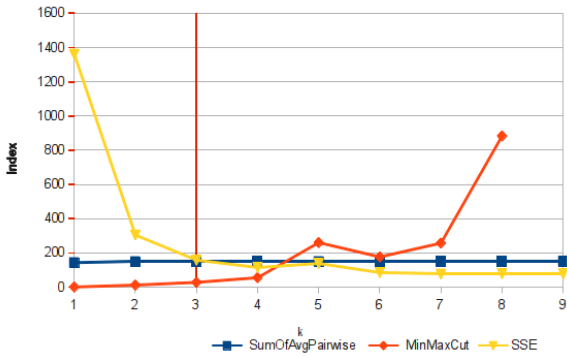
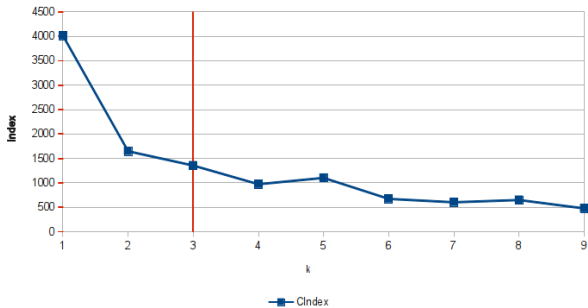
where d_w is the sum of the within cluster distances.

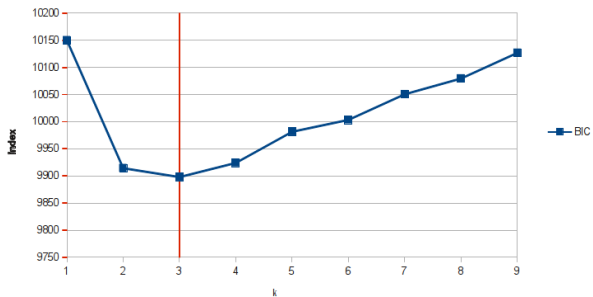
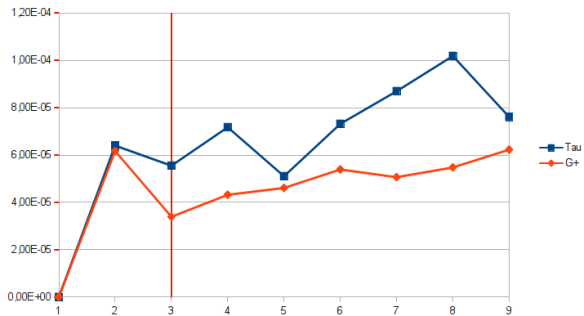
Gamma

$$p_{gamma} = \frac{s(+) + s(-)}{s(+) - s(-)}$$

where $s(+)$ represents the number of consistent comparisons involving between and within cluster distances, and $s(-)$ represents the number of inconsistent outcomes Milligan and Cooper [1985]

An Iris dataset





The strive for objectivity, repeatability, testability etc. is perfectly right attitude as long as their proper place in the “hierarchy of aims” is maintained, but becomes very harmful if these tools dominate over the purpose of scientific research. [Holynski, 2005, p. 487]

Questions?

Thanks for your attention!

`tomas.barton@fit.cvut.cz`

References I

P. Arabie, L. J. Hubert, and G. D. Soete. Clustering and classification. World Scientific, 1996.

J.L. Borges. El idioma analítico de John Wilkins. Obras completas, 1952.

Rich Caruana, Mohamed Elhawary, Nam Nguyen, and Casey Smith. Meta clustering. In Proceedings of the Sixth International Conference on Data Mining, ICDM '06, pages 107–118, Washington, DC, USA, 2006. IEEE Computer Society. ISBN 0-7695-2701-9. doi: 10.1109/ICDM.2006.103. URL <http://dx.doi.org/10.1109/ICDM.2006.103>.

B. S. Everitt. Cluster Analysis. Edward Arnold, 1993.

References II

- I. Guyon, U. Von Luxburg, and R.C. Williamson. Clustering: Science or art. In NIPS 2009 Workshop on Clustering Theory, 2009.
- Roman B. Holynski. Philosophy of science from a taxonomist's perspective. Genus, 16(4):469–502, 2005.
- L.J. Hubert and J.R. Levin. A general statistical framework for assessing categorical clustering in free recall. Psychological Bulletin, 83(6):1072, 1976.
- Glenn W. Milligan and Martha C. Cooper. An examination of procedures for determining the number of clusters in a dataset. Psychometrika, 50(2):159–179, June 1985.

References III

- G. Bel Mufti, P. Bertrand, and L. El Moubarki. Determining the number of groups from measures of cluster stability. In Proceedings of International Symposium on Applied Stochastic Models and Data Analysis, pages 404–412, 2005.
- Siuly, Yan Li, and Peng (Paul) Wen. Clustering technique-based least square support vector machine for eeg signal classification. Computer Methods and Programs in Biomedicine, 104(3): 358–372, 2011. URL <http://dblp.uni-trier.de/db/journals/cmpb/cmpb104.html#SiulyLW11>; <http://dx.doi.org/10.1016/j.cmpb.2010.11.014>; <http://www.bibsonomy.org/bibtex/2976ac83c8e51dd3ff108ee52da22902d/dblp>.