



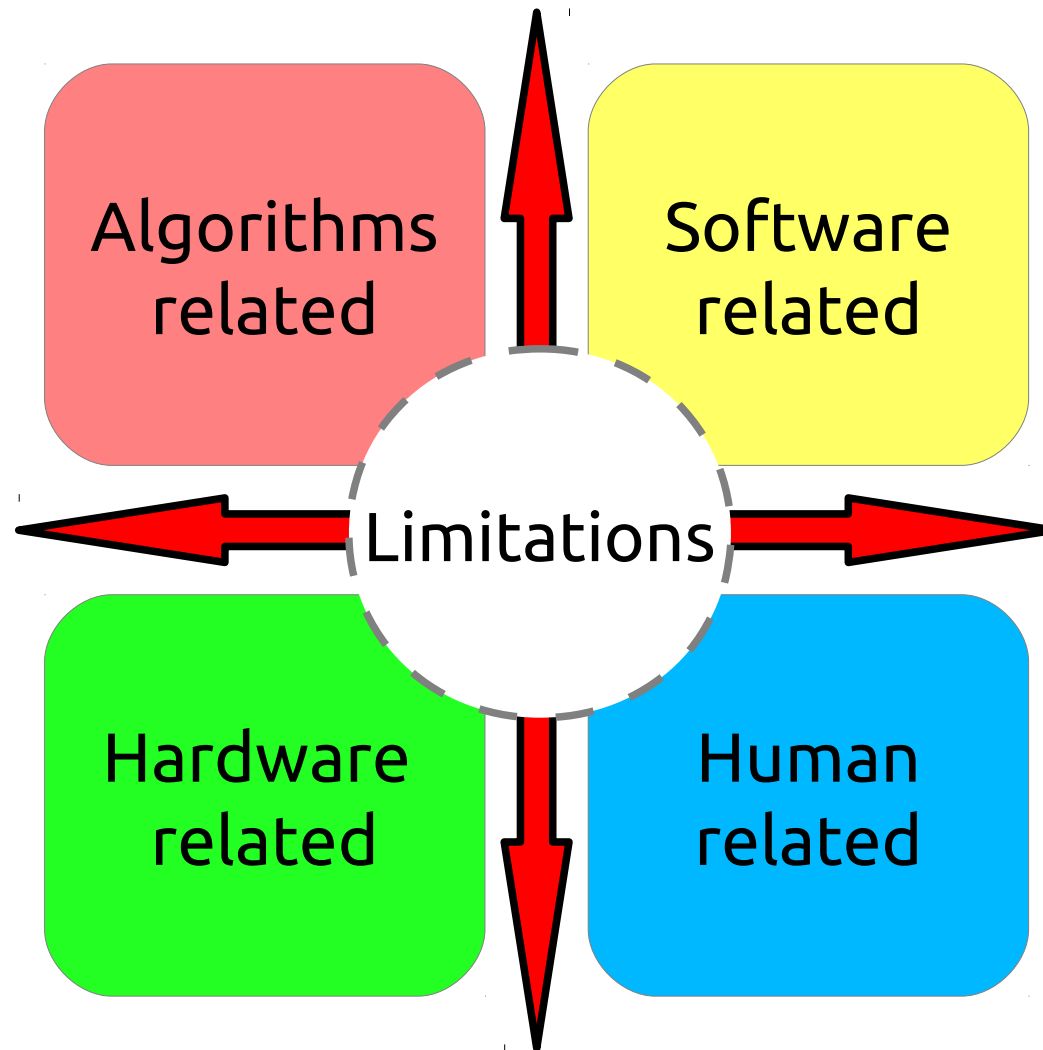
**WIESŁAW PIETRUSZKIEWICZ**

**WEST POMERANIAN UNIVERSITY OF TECHNOLOGY IN SZCZECIN**

# Agenda

- Data processing limits
- Ways to ease processing limits
- Evaluation of applied DM
- Examples
- DM applied in software
- Conclusions
- Future research

# Data processing limits



# Data processing limits

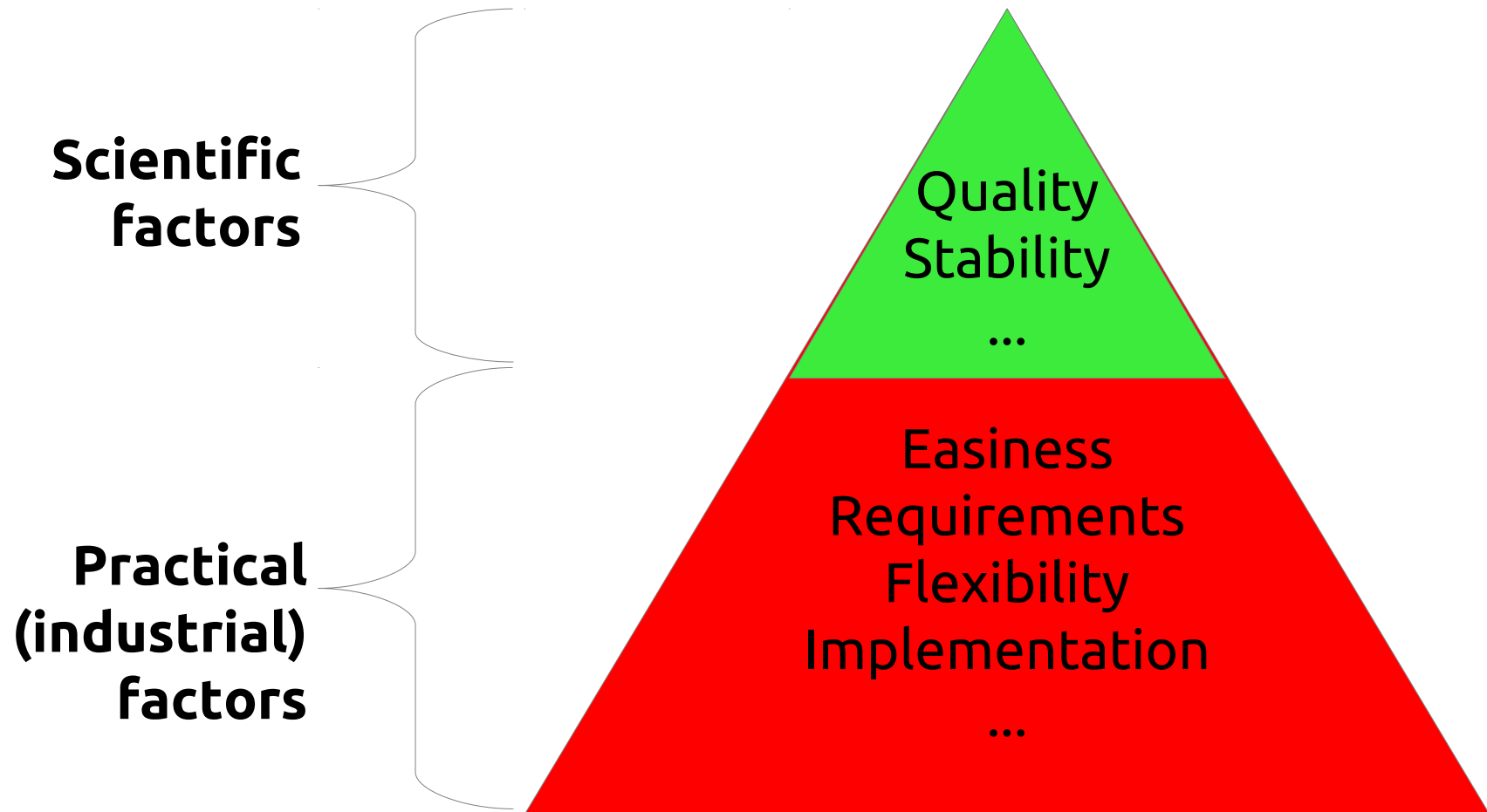
- Data processing limits for applied DM influence:
  - Capabilities
  - Flexibility
  - Usability
  - Economics (resources usage)
- Limits usually reveal themselves for applied DM, more often than for research-oriented DM

# Easing processing limits

We can ease or omit limits by:

- Creating better algorithms (easier, less resources-consuming ...)
- Using better resources (faster, distributed, parallel ...)
- Deploying advanced software technologies (e.g. allowing to scale application)
- Using specialised supporting software

# Evaluation of DM outcomes



## Evaluation of Data Mining

# Scientific evaluation

- Usually concerned about the quality of outcomes
- Neglects other important factors (crucial for applied DM)
- A real-life requires us to make compromises e.g., details vs. speed

# Practical (industrial) evaluation

## Feature selection

- Stability
- Unbiasedness
- Scalability
- Flexibility

## Classification

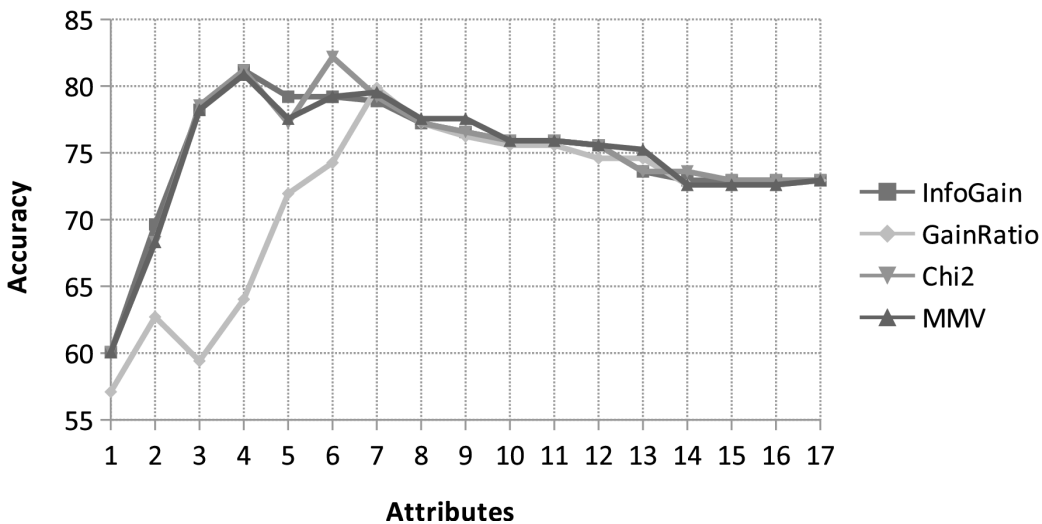
- Speed of learning
- Number of adjustments
- Speed of simulation
- Demand on resources



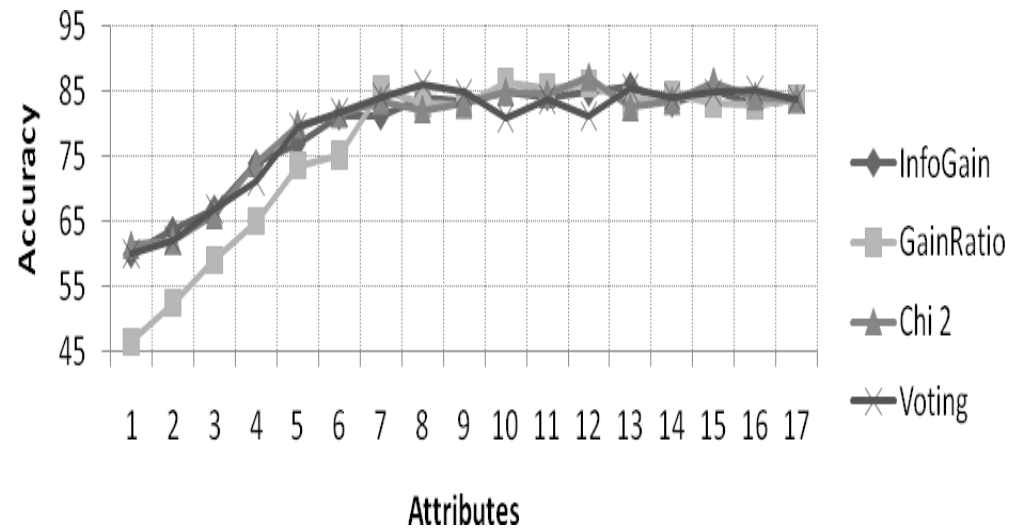
# Example I – Attributes selection

- Check unbiasedness for feature selection compare filters with MMV (examined unbiased multi-measure algorithm – Pietruszkiewicz 2010)
- We'd rather expect good results on average (not best only by a chance)
- Examined using two dataset:
  - Household **bankruptcy** dataset (Rozenberg&Pietruszkiewicz 2008)
  - **Images** UCI dataset

# Example I – Attributes selection

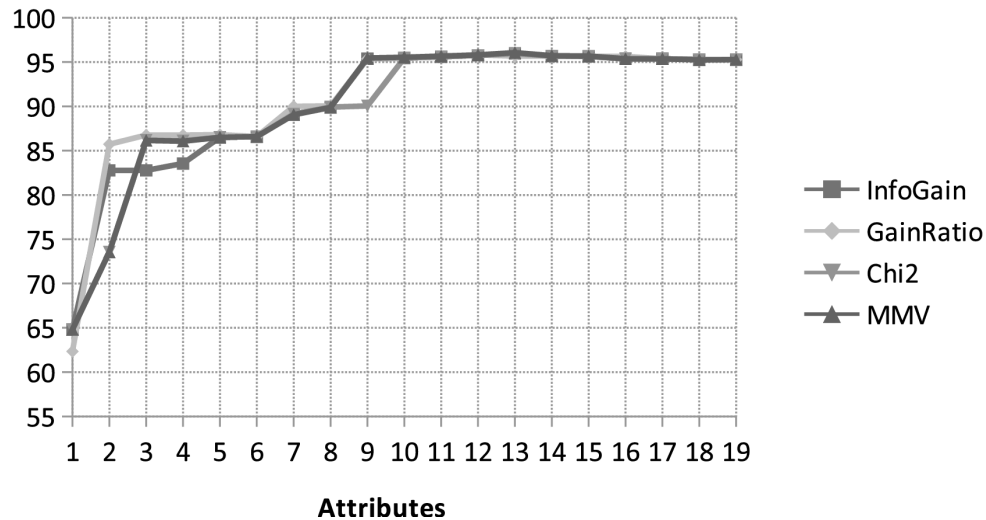


Attributes  
Decision trees  
for **Bankruptcy** dataset

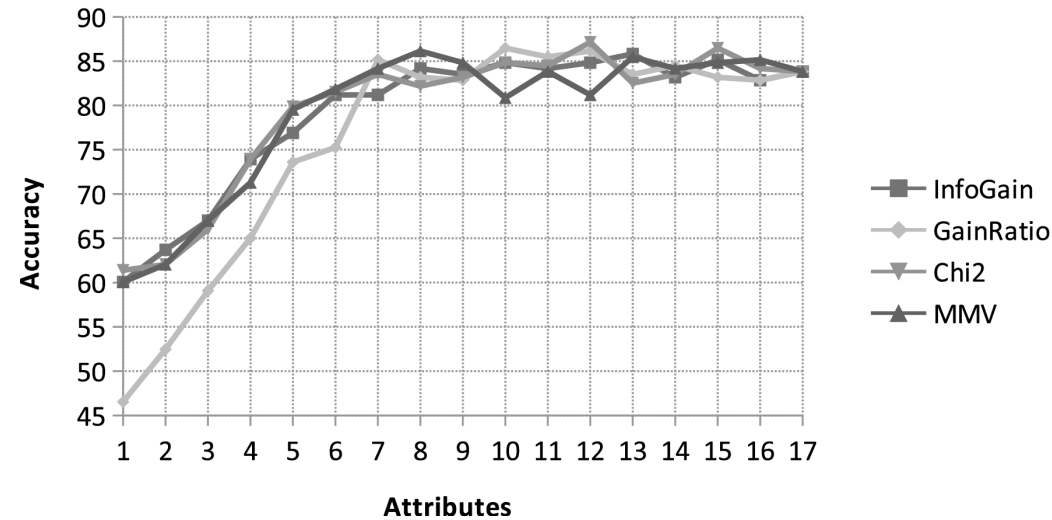


Attributes  
Neural networks  
for **Bankruptcy** dataset

# Example I – Attributes selection



Decision trees  
for **Images** dataset



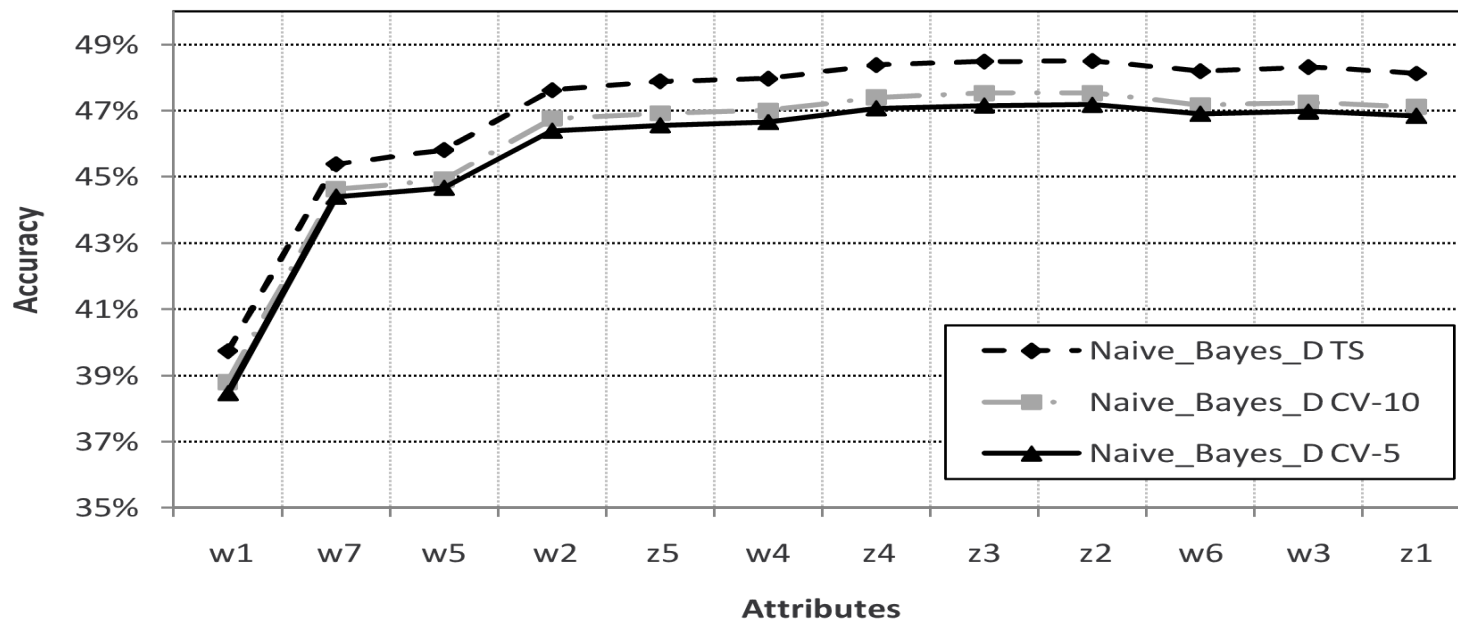
Neural networks  
for **Images** dataset

# Example I – Outcomes

- Stable results
- Unbiased selection (not in favour of any method)
- Easy to use
- Could be developed further by voting or adding other filters

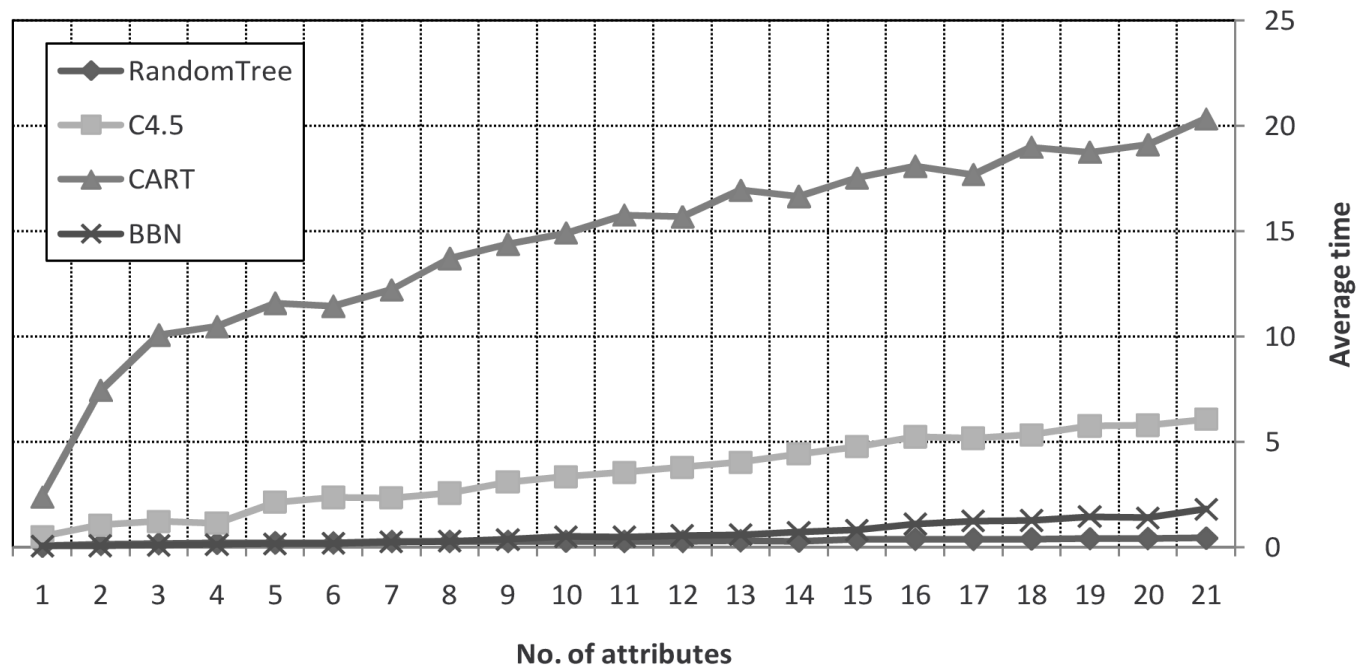
# Example II - Classification

Done for one (the most difficult to classify) of four datasets for OS projects estimation (Dzega 2008).



# Example II - Classification

Time is another important factor – large differences between methods.



## Example II – Outcomes

- More is not always better (quality), but is usually (internal feature selection done by decision trees) more costly&complex
- Methods differ very much in the time of learning and simulation
- Time as a resource, more time per experiment = less experiments (variants tested)

# DM in software systems

- Features important for DM deployed in software systems:
  - Understandability of results
  - Easy implementation
  - Demand for resources
- In development phase they influence time&costs
- For production-ready systems they influence their quality



# DM in software systems

Form of knowledge	Main advantages	Main disadvantages
Decision rules	<ul style="list-style-type: none"><li>• Easy software implementation</li><li>• Obvious even for non-skilled users</li></ul>	<ul style="list-style-type: none"><li>• Delivers not-optimal structure in terms of size and speed</li><li>• Possible incoherency of rules</li></ul>
Decision tree	<ul style="list-style-type: none"><li>• Fast functioning (learning and simulation)</li><li>• Intuitive visual representation</li></ul>	<ul style="list-style-type: none"><li>• Can be used only in classification tasks (or simplified quasi-regression)</li><li>• Non-resistant to noisy observations at run-time</li></ul>
Causal (Bayesian) networks	<ul style="list-style-type: none"><li>• Deal with partial information at run time</li><li>• Work with partial observations</li></ul>	<ul style="list-style-type: none"><li>• Results depend on particular algorithms deployed</li><li>• Suitable only to classification problems</li></ul>
Mathematical models	<ul style="list-style-type: none"><li>• Flexible form suitable to different tasks</li><li>• Usually fast at run-time</li></ul>	<ul style="list-style-type: none"><li>• Even slightly complex results might be unclear for people</li><li>• Very broad space of search for optimal model class</li></ul>

# DM in software systems - enhancements

- Simplification (e.g. selection of less complicated algorithms, reduction of outcomes complexity)
- GPU-supported processing (high power, low costs, **green/eco-friendly** solution)
- Distributed processing (clusters, clouds)

# Conclusions

- Evaluation of DM outcomes should involve much more than scientific features
- Introduction of practical features of quality
- Knowledge representation forms influence further intelligent software development
- Introduction of software-relevant DM features

# Future research

- Set of benchmarks to measure introduced quality factors (UCI benchmarks are oriented on quality of prognosis)
- Practical guideline for usage of algorithms (many features of algorithms quality are dataset-independent)

Thank you!

Questions and comments are welcomed ...

If you would like to contact:  
[wieslaw@pietruskiewicz.com](mailto:wieslaw@pietruskiewicz.com)