The Use of Rough Set Methods in Knowledge Discovery in Databases

Marcin Szczuka The University of Warsaw Poland

Why KDD?

"We are drowning in the sea of data, but what we really want is **knowledge**"

PROBLEM: How to retrieve useful information (knowledge) from massive data sources?

SOLUTION: Data Warehouses + Data Mining

- Acquisition of data in real time
- Discovery of interesting knowledge (patterns, dependences, rules, regularities, ...) from large data sets.



Data mining

- Data mining = the iterative and interactive process of discovering non-trivial, implicit, previously unknown and potentially useful(interesting) information or patterns from data in large databases.
- Wikipedia: Data mining (the analysis step of the Knowledge Discovery in Databases process, or KDD), a relatively young and interdisciplinary field of <u>computer science</u>, is the process of extracting patterns from large <u>data sets</u> by combining methods from <u>statistics</u> and <u>artificial intelligence</u> with <u>database</u> <u>management</u>.

Data Mining



Challenges and Competitions

- Since 1997 ACM KDD Cup attracted numerous researchers in the area.
- In 2009 Netflix Prize (\$1 million) has opened the gate for really serious business-oriented competitions.
- Several platforms for on-line KDD/DM competitions are now existing: <u>kaggle.com</u>, <u>tunedit.org</u>, ...
- Crowdsourcing of highly advanced DM tasks.
- There are several opportunities to win thousands of dollars (or Euros), and that does not include ...

Heritage Health Prize

- The goal of the prize is to develop a predictive algorithm that can identify patients who will be admitted to the hospital within the next year, using historical claims data.
- \$3 million (US) Grand Prize.
- Six Milestone Prizes totalling \$230,000, which are awarded in varying amounts at three designated intervals during the Competition.
- Ends in April 2013

http://www.heritagehealthprize.com/

Rough Set Methodology

What Rough Sets are all about and how to make them work for you?

Rough Set Methodology

- Rough set methodology is based on indiscernibility between objects.
- Rough set methods utilize the comparison between elements, e.g., discernibility, indiscernibility, similarity, ...
- In KDD, rough set methodology can be applied to:
 - implement efficient methods for mining interesting templates from data: data reduction, minimal decision rules, decomposition, hierarchical learning ...
 - cooperate and improve existing methods in data mining like decision trees, association rules, clustering, kNN, neural networks, Bayesian networks...
 - design practical KDD projects for complex application domains like: search engines, medical data analysis, bioinformatics, ...

A simple example

| Id. | Age | LEM | Walk | | | | |
|-----------------------|-------|-------|------|--|--|--|--|
| u 1 | 16-30 | 50+ | Yes | | | | |
| <i>u</i> ₂ | 16-30 | 0 | No | | | | |
| <i>u</i> ₃ | 31-45 | 1-25 | No | | | | |
| u ₄ | 31-45 | 1-25 | Yes | | | | |
| <i>u</i> ₅ | 46-60 | 26-49 | No | | | | |
| u ₆ | 16-30 | 26-49 | Yes | | | | |
| u ₇ | 46-60 | 26-49 | No | | | | |

Information system: Rows = objectsColumns = attributes (features) If decision (attribute) is present the information system becomes **decision table**. We usually denote it by: S=(U,A), where *U* – the set (universe) of objects A – the set of attributes

(In)Discernibility

In our example the information is not precise (complete) enough. We cannot discern precisely walking from nonwalking patients. So, the concepts of walking and nonwalking patient are partly *indiscernible*.

| Age/LEM | 0 | 1-25 | 26-49 | 50+ |
|---------|-----------------------|---|---|------------|
| 16-30 | u ₂ | | u ₆ | u 1 |
| 31-45 | | u ₃ u ₄ | | |
| 46-60 | | | <i>u</i> ₅ <i>u</i> ₇ | |

Approximations

Lower approximation corresponds to certainty about the object belonging to a concept (set). They **definitely** belong to the set.



Approximations

Upper approximation corresponds to possibility of the object belonging to a concept (set). It is possible (likely, feasible) that they belong to the set, i.e., they **roughly** are in the set.

| Age/LEM | 0 | 1-25 | 26-49 | 50+ |
|---------|----------------|---|---|------------|
| 16-30 | u ₂ | | u ₆ | u 1 |
| 31-45 | | u ₃ u ₄ | | |
| 46-60 | | | <i>u</i> ₅ <i>u</i> ₇ | |

Approximations

Upper approximation corresponds to possibility of the object belonging to a concept (set). It is possible (likely, feasible) that they belong to the set, i.e., they **rougly** are in the set.

| Age/LEM | 0 | 1-25 | 26-49 | 50+ |
|---------|-----------------------|---|---|-----|
| 16-30 | U ₂ | | u ₆ | u, |
| 31-45 | | u ₃ u ₄ | | |
| 46-60 | | | <i>u</i> ₅ <i>u</i> ₇ | |

Boundary

The boundary region represents the uncertain portion of our data set. We do not have enough information to definitely establish the status of objects in this area.



The donut



Rough Set Techniques in DM

- Over the years work of several "rough setters" resulted in creation of algorithmic methods for analysis of various kinds of data.
- These methods utilize fundamental rough set notions such as discernibility, approximation, information function and reduct.

These methods are dealing with:

- Reduction of data size and complexity.
- Discovery of frequent patterns and decision rule discovery.
- Continuous attributes' discretization.
- Data decomposition.
- Others including new feature construction, instance-based learning, ...

Reduction

- Do we need all attributes?
- Do we need to store the entire data?
- Is it possible to avoid a costly test?
- **Reducts** are subsets of attributes that preserve the same amount of information. They are, however, (NP-)hard to find.
- Efficient and robust heuristics exist for reduct construction task.
- Searching for reducts may be done efficiently with the use of , e.g., evolutionary computation.
- Overfitting can be avoided by considering several reducts, pruning rules and lessening discernibility constraints.

Data reduction with RS

What is a reduct?

Reducts are minimal subsets of attributes which contain a necessary portion of information from the set of all attributes.

Given an information system S = (U,A) and a monotone evaluation function:

$$\mu_S: \mathscr{P}(A) \to \mathfrak{R}_+$$

The set $B \subset A$ is called μ -*reduct* of A, iff:

 $\mathbf{I}. \qquad \mu(B) = \mu(A),$

2. for any proper subset $B_o \subset B$ we have $\mu(B_o) < \mu(B)$. The set $B \subset A$ is called *approximate reduct*, iff:

$$\mu(B) \leq \mu(A) - \varepsilon,$$

2. for any proper subset ...

Some types of reducts

- Information reduct:
- μ_I(B) = number of pairs of objects discerned by B
 Decision-oriented reduct:
 - $\mu_D(B)$ = number of pairs of **conflicting objects** discerned by *B*
- Object-oriented reduct:
 - $\mu_x(B)$ = number of objects discerned from *x* by *B*
- Frequent reducts;
- α-reducts;

Pattern and rule discovery

By examining the structure of indiscernibility classes and reducts one can summarize information carried by objects using patterns of the form:

$$(a_{i1} = v_1) \wedge \ldots \wedge (a_{ik} = v_k)$$

These patterns may be further converted into associations. In classification problems, we can produce decision rules of the form:

$$(a_{i_1} = v_1) \wedge \ldots \wedge (a_{i_k} = v_k) \Longrightarrow d = v_d$$

Patterns and rules can be filtered, pruned, generalized, and composed. That permits management of discovered knowledge.

Discretization

Attributes that have many different values, e.g.., real-valued, may pose a technical problem for some algorithmic DM methods. **Discretization** (quantization) of attribute values can be done using rough set framework.

- We consider all pairs of objects. Then we consider all possible cuts on attributes in discourse. We choose the cut that induces best split w.r.t. the number of objects from different decision classes that are discerned by this split.
- This is called **Maximal Discernibility (MD)** heuristic. Various modifications, especially concerning the choice of best cut, exist.









Data decomposition with RS

Large data sets may not be possible to process as-is. The ability to decompose data set into smaller chunks is a requirement. These fragments, after decomposition represented as leafs in decomposition tree, are supposed to be more uniform and easier to cope with decision-wise.



Rough Sets vs. Others

How RS methods fit the bigger picture and how they compare with some popular DM techniques?

Rough Sets

in Decision Tree Construction

Main problem: Search for minimal decision tree compatible with a given decision table. This problem is NP-hard

Heuristics:

- Decision tree are constructed from a given set of candidate partitions;
- Best-first searching strategy, a quality measure must be defined, e.g..,
 - Entropy gain;
 - Gini's index;
- Rough set based measure: *discernibility measure = number of resolved conflicts*
- In our research: decision trees constructed by discernibility measure have many interesting properties.

Association Rule Generation

Association rule generation methods consist of two steps:

- 1. Generate as many large templates as possible: $T = D_1 \wedge D_2 \wedge ... \wedge D_k$ i.e, support(T) > s and $support(T \wedge D) < s$ for any descriptor D)
- 2. For any template *T*, search for a partition $T = P \land Q$ such that:
 - support(P) < support(T) /c
 - *P* is the smallest template satisfying the previous condition

Surprisingly: the second step can be solved by rough set methods (using α -reducts).

Rough Sets vs. k-NN

- How to define the measure function or neighbourhood?
- Similarity relation can be learned from data using rough set methods.
- A simple idea:

the more reducts an attribute appears in the more important this attribute is.

• Filtering methods in rule based classifiers can simulate both decision tree and k-NN classifiers.

Multivariate Analysis

- Multiple Regression: analyzes the relationship between several attributes and the decision;
- Principal Components Analysis and Factor Analysis A linear dimensionality reduction technique, which:
 - identifies orthogonal directions of maximum variance in the original data,
 - projects the data into a lower-dimensionality space formed of a sub-set of the highest-variance components.
- Discriminant Analysis: searching for the best set of attributes that discriminates objects from two or more decision classes.
- Cluster Analysis: grouping similar objects.

Multivariate Analysis & RS

- PCA and clustering methods can be applied as pre-processing step to rough set methods.
 e.g., to extract new features from the original set of attributes.
- Experimental results: they can improve the quality of rough set classifiers.

Rough Sets

in Network Construction

- A Bayesian network is an acyclic directed graph that models probabilistic dependencies between the domain variables:
- **Q:** How to construct Bayesian networks from data? In many cases, the problem is NP-hard.
- **A:** Searching for structure + probability distribution;
- **RS:** Structure can be reconstructed by calculating frequent reducts!
- Rough sets vs. neural networks:
- Rough sets vs. Petri net
- Rough sets vs. Belief (cause-and-effect) networks.

Software tools

How to do your own experiments with Rough Set methods?

Tools overview

- The past:
 - Rosetta,
 - RSES,
 - ROSE, 4eMka
- The present:
 - Rseslib et consortes
 - jRS, jMAF, jAMM
 - Debellor
 - TunedIT
- Related tools and perspectives for the future (one of possible scenarios)

Rosetta

http://www.lcb.uu.se/tools/rosetta/

Supported OS:



Partly Supported OS:



| Rosetta - cleveland.ros | | _ | | | | | _ | | | | | | | | | ا _ |
|--|-----------------------|--|---|-----------------|-------------------------------|------------------|---------------------------------|-----------|--|-------------------------------|------------|---------------|-----------------|-----------------|-----------|---------|
| e <u>E</u> dit ⊻lew <u>Wi</u> ndow <u>H</u> elp | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | |
| Project | | 🗖 tra | ining set, disci | etized | | | | | | | | | | | | |
| P Structures | | | ср | trestbps | chol | fbs | restecg | thalach | n exa | ang | oldpeak | slope | ca | thal | num | dis |
| E D cleveland | | 1 | Asymptomati | [*, 127) | [233, 295] |) False | LV hypertrop | [151,*) | 1 N | NO | [*, 0.6) | Upsloping | [*, 1) | Normal | 0 | |
| D training set | | 2 | Typical angin | [*, 127) | [198, 233] |) False | LV hypertrop | [112, 151 | I) Yi | 'es | 0.6, 1.7) | Upsloping | [1,*) | Normal | 0 | |
| Itraining set, discretized | | 3 | Atypical angi | [*, 127) | [198, 233] |) False | Normal | [112, 151 | I) N | No | [*, 0.6) | Flat | [1,*) | Reversable d | 0 | |
| B K reducts genetic app | proximate | 4 | Asymptomati | [*, 127) | [233, 295] |) False | LV hypertrop | [112, 151 | I) N | 10 I | 0.6, 1.7) | Upsloping | [*,1) | Reversable d | 1 | Y |
| R approximate rule | 15 | 5 | Typical angin | [*, 127) | [233, 295] |) False | Normal | [112, 151 | I) N | NO I | 0.6, 1.7) | Flat | [*, 1) | Reversable d | 1 | , |
| R johnson rules, ex | xact | 7 | i ypical angin | [141, *] | [198, 233] | Faise | LV nypertrop | [151,*] | | 100 | 0.6, 1.7) | Flat | [^,1) [4, A) | Reversable d | 0 | |
| E c:temp/cuts txt | View | <u> </u> | Asymptomati | [141, 1] | [233, 295] | False | Normai | [*, 112] | | es | 0.6, 1.7) | Fiat | [1,") [1, t) | Reversable d | 1 | |
| E D testing set | Remove | 9 | Asymptomati | [* 127] | [130, 235] |) False | LV hypertrop | [112, 13] | | 10 | [* 0.6) | Linsioning | [1 Å] | Normal | 2 | |
| e D testing set, discretized | Duplicate | 10 | Asymptomati | [* 127) | [295 1] | False | Normal | [112 151 | n y | 'es | [1,7,4] | Flat | [1.1] | Reversable d | 2 | |
| Confusion matrix | | 11 | Non-anginal p | [127, 141] | [198, 233] |) False | LV hypertrop | [151,*) | N | NO | [1.7.*) | Flat | [1,1) | Normal | 0 | |
| | Save | 12 | Asymptomati | [127, 141] | [198, 233] |) True | LV hypertrop | [151,*) | I Yi | 'es | [1.7,*) | Downsloping | [*,1) | Reversable d | 1 | |
| C: /temp/roc.txt | Save as | 13 | Asymptomati | [*, 127) | [233, 295] |) False | Normal | [112, 151 | Y) | 'es | [1.7,*) | Flat | [1,*) | Reversable d | 2 | |
| Algorithms | Load | 14 | Asymptomati | [*, 127) | [233, 295] |) False | LV hypertrop | [*, 112) | i Yi | 'es | [1.7,*) | Downsloping | [1,*) | Normal | 3 | |
| | Export + | 15 | Asymptomati | [*, 127) | [233, 295] |) False | Normal | [*, 112) | i Yi | 'es | [1.7,*) | Flat | [1,*) | Reversable d | 1 | |
| onrovimate rulee | Ciller, N | Denie 4 | | 1 | | | | | 1 | | | 1 _ 1 | | 1 | | |
| | riller | Dasie I | itterning | | | | | | | | - | | | le. | | - |
| | Execute > | Quainy | mittering | | | LHS Support | RHS Support | RHS | Accuracy | LHS | Coverage | RHS COV | erage | RHS Stabilit LI | HS Len(| - |
| oldpeak((*, 0.6)) AND slope(Upslop | | Quality | tiltering loop | (Yes) | | 43 | 36,7 | 0.837209 | 9, 0.162791 | 0.282 | 895 | 0.433735, 0.1 | 01449 1 | 1.0, 1.0 3 | [| - |
| restecg(Normal) AND thalach([151 | Statistics | > diseas | e(No) OK diseas | e(Yes) | | 42 | 36,6 | 0.85/143 | 5, 0.142857 | 0.276 | 316 | 0.433735, 0.0 | 86957 1 | 1.0,1.0 3 | H | - |
| restecg(Normal) AND trialacri([151 theleok(151_A)) AND eldpook(151_0) | Information | => dises | (se(res) OR disc | Base(NU) | | 39 | 9,00 | 0.102564 | 0.03/430 | 0.250 | 797 | 0.057971,0.4 | 121007 1 | 1.0,1.0 3 | | |
| evang(No) AND oldpeak([* 0.6)) A | Report | -> diseas | e(No) OR disea | e(Ves) | | 41 | 35.6 | 0.853600 | | | | | L | | | |
| oldpeak([1,7,*)) => disease(Yes) OR disease(No) | | | | | | 39 | 29.10 | 0.743 | c:/temp | | | | | | | _ |
| 52 cp(Asymptomatic) AND exang(Yes) AND slope(Flat) => disease(Yes) OR disease(No) | | | | 31 | 29, 2 | 0.9354 | | | | c:/te | mp/class. | log | | 1 | | |
| 53 cp(Asymptomatic) AND restecg(LV hypertrophy) => disease(No) OR disease(Yes) | | | | | 38 | 9,29 | 0.2368 2 | 4 | | P | anking = | (0.68975 | 53) Yes (1) | 220 rule(| s) | |
| 54 oldpeak((*, 0.6)) AND ca((*, 1)) AND thal(Normal) => disease(No) OR disease(Yes) | | | | 38 | 34, 4 | 0.894; 2 | 5 | | | | (0.31024 | 47) No (0) 2 | 01 rule(s |) | | |
| 55 restecg(Normal) AND ca([*, 1)) AND thal(Normal) => disease(No) OR disease(Yes) | | | | 38 | 34, 4 | 0.8947 2 | 6 0Ъ | ject 4: | ERROR A | ctual = | Yes (1) | | | | | |
| sex(Male) AND slope(Flat) AND ca(| ([1,*)) => disease(No |) OR dise | ease(Yes) | | | 30 | 2,28 | 0.0666 2 | 7 | | P | redicted = | No (0) | | | |
| sex(Male) AND cp(Asymptomatic) | AND exang(Yes) => | disease() | Yes) OR disease | (No) | | 31 | 28,3 | 0.9032 2 | 18 | | P | lanking = | (0.56015 | 59) No (0) 2 | 22 rule(s |) |
| | sable detect) => dise | asel Yes, |) OR disease(No, |) | | 29 | 28,1 | 0.965: 2 | 9 | | | | (0.43984 | 41) Yes (1) | 226 rule(| s) |
| | | | | | | | | | <u>и</u> ив | ject s: | ок а | ctual = | NO (U) | | | |
| onfusion matrix | | X | iohnson rules | . exact | | | -1 | o xi lă | 10 | | , r 1 | contring = | NO (0) | 12) No (0) 2 | 02 milo/c | |
| Predicted | | 712 | | | Rul | e | | | 13 | | | | (0.48878 | 37) Yes (1) | 182 rule(| , ≲) |
| hie | Vee | ┛║┾ | cm(Asymp | dometic) AND | avang(Vac) | AND co(1 t)) - | -> disease(Vec) | | 4 оъ | nect 6: | ok À | ctual = | No (0) | | | |
| No Yes I cp(Asymptomatic) AND exang Yes No 74 7 0.01259 2 thelech(112, 151)) AND exang Yes | | |) exang(165) | AND that (Reve | rsable defect) = | <u>,</u> 13 | 15 | | р | redicted = | No (0) | | | | | |
| al Yes 17 | 53 0.757143 | 3 | age(157, 62)) AND sex(Male) AND cp(Asymptomatic) => disease(Yes) | | | | | 3 | 16 | | 3 | tanking = | (0.70745 | 51) No (0) 2 | 14 rule(s |) |
| 0.813187 0.8 | 383333 0.84106 | 4 | 4 age((*, 45)) AND cp(Atypical angina) => disease(No) | | | | | | 7 | | | | (0.29254 | 49) Yes (1) | 192 rule(| s) |
| Class Yes | | 5 cp(Atypical angina) AND trestbps((127, 141)) => disease(No) | | | | | | 3 | в оъ | ject 7: | ok A | ctual = | Yes (1) | | | |
| Area 0.919489 | | 6 age([*, 45)) AND trestbps([127, 141)) AND exang(No) => disease(No) | | | | | 0) 3 | 19 | | P | redicted = | Yes (1) | | | | |
| C Std. error 0.024094 | | 7 sex(Female) AND cp(Non-anginal pain) => disease(No) | | | | N0) | 4 | 10 | | P | anking = | (0.75373 | 31) Yes (1) | 285 rule(| s) | |
| Thr. (0, 1) 0.428 | | 8 | age([57,6 | i2)) AND exan | g(Yes) AND | ca([1,*)) => dis | ease(Yes) | 4 | 1 | | | | (0.24626 | 59) No (0) 2 | 40 rule(s |) |
| Thr. acc. 0.428 | | 9 | age([49, 5 | 7)) AND oldpe | eak((*, 0.6)) / | \ND ca([*, 1)) ≕ | disease(No) | _ | i2 0b | ject 8: | ok Å | ctual = | No (0) | | | |
| | | - 10 | age([57, 62)) AND cp(Asymptomatic) AND oldpeak([1.7, *)) => disease | | | | | ise 4 | 13 | | | redicted = | No (U) | | | |
| ress | × | 11 | age([*, 45)) AND cp(Non-anginal pain) => disease(No) | | | | _ | 14 | | * | anging = | (0.82107 | /2) NO (U) 2 | 229 mule(s | , ~> | |
| | Clear | 12 | 2 ca([1, *)) AND thal(Fixed defect) => disease(Yes) | | | | | _ 16 | 16 0h | inct 9. | olr à | ctual = | Vec (1) | 10) IES (1) | 235 Ture(| |
| | | | 3 sex(Female) AND cp(Asymptomatic) AND thal(Reversable defect) => d | | | | | > d | 17 | geco s. | 0A A | redicted = | Yes (1) | | | |
| Applying JohnsonHeducer to training set, discretized | | 14 | age([*, 45)) AND exang(Yes) => disease(Yes) | | | | - | 18 | Ranking = (0.804212) Yes (1) 258 rule(s) | | | | | | | |
| Done applying JohnsonReducer to training set, discretized | | | sex(male) | (7)) AND treet | 5,)) AND 8. http://127_44 | ti) AND thak®a | verseble defect | | 19 | (0.195788) No (0) 199 rule(s) | | | | | | , · |
| lication took 00:00:01 | | 17 | ser(Ferna | ie) AND cn(At | vnical angins | a) => disease(N | non cable defect | - 5 | 0 0ъ | ject 10 | : ok | Actual : | No (0) | | | |
| Jying BatchLlassifier to testing set, disci Carea – 0.919489 | retized | 15 | Cp(Non-ar | nginal pain) AN | D chol(1233 | 2951) AND elor | ~/ ne(Linsloping) => | - di 5 | 1 | | | Predicted : | No (0) | | | |
| C std. error = 0.024094 | | 19 | ager(157_F | (2)) AND trests | bps/[127.14 | 1)) AND slope(F | lat) => disease(| Ye 🔳 5 | 12 | | | Ranking | (0.7300 | 053) No (0) | 243 rule(| s) |
| | | 1 | | | | | | | | | | | | | | |

RSES <u>http://logic.mimuw.edu.pl/~rses/</u>



Supported OS:







ROSE 2 & 4eMka2

http://idss.cs.put.poznan.pl/site/software.html



Rseslib <u>http://rsproject.mimuw.edu.pl</u>



- Open Source (GNU GPL) library of RS and related classifications algorithms and data structures, written in Java.
- Accepts CSV+, ARFF (Weka), and RSES formats.
- Several small GUI-based tools to make usage easier : Visual Rseslib, Qmak, Trickster.

jRS, jMAF, jAMM <u>http://idss.cs.put.poznan.pl/site/software.html</u>

- jRS is a Java library implementing methods of analysis provided by the Dominance-based Rough Set Approach and Variable Consistency Dominance-based Rough Set Approach.
- jMAF is a Rough Set Data Analysis Framework written in Java language. It is based on jRS library.
- JAMM is a GUI decision support tool designed to guide the user through analysis and solving of multiple-criteria classification problems using jRS.



Debellor http://www.debellor.org



- Debellor is an open source framework for scalable data mining and machine learning.
- The unique feature of Debellor is *data streaming*, which enables efficient processing of massive data and is essential for scalability of algorithms.
- You may implement your own algorithms in Debellor's architecture and combine them with pre-existing ones (over 100), in particular:
 - Rseslib algorithms;
 - Most of WEKA algorithms.

TunedIT platform http://www.tunedit.org/



- <u>TunedIT Research</u> delivers web-based tools to help data mining scientists conduct repeatable experiments and easily evaluate data-driven algorithms.
- <u>TunedIT Challenges</u> platform for hosting data competitions - for educational, scientific and business purposes.
- Repository of data sets, methods and algorithms, including:
 - UCI, StatLog, KDD Cup, PROMISE, and other data sets;
 - Algorithms libraries: Debellor ,Weka, and Rseslib;
 - TunedTester framework and various examples.
- Knowledge base gathering various experimental results.

RoughICE

<u>http://www.mimuw.edu.pl/~bazan/roughice/</u>

- The Rough Set Interactive Classification Engine
- Software platform supporting the approximation of spatio-temporal complex concepts in the given concept ontology.
- Support for the dialogue with the user (domain expert).
- Design classifiers and knowledge flow for complex concepts.



Infobright http://www.infobright.org/

- Column-oriented database (warehouse) engine.
- Suitable for large data thanks to compression methods.

INFOBR GHT

- Supporting analytic processing of large amounts of data.
- Build around rough set concepts of approximation and supporting approximate SQL query answering.
- Based on MySQL architecture.
- Both Open Source (ICE) and commercial (IEE) versions available.

Thank you! Questions and comments are welcome.

Bibliography and links

Where to look for more information?

- W. Frawley, G. Piatetsky-Shapiro, C. Matheus: Knowledge Discovery in Databases: An Overview. AI Magazine, Fall 1992, pgs 213-228.
- D. Hand, H. Mannila, P. Smyth: *Principles of Data Mining*. MIT Press, Cambridge, MA, 2001
- Z. Pawlak: *Rough Sets: theoretical aspects of reasoning about data*. Kluwer AP, Dordrecht, 1991
- L. Polkowski, A. Skowron (eds.): Rough Sets in Knowledge Discovery, part 1 and 2 Physica-Verlag, Heidelberg, 1998
- S.K. Pal, L. Polkowski, A. Skowron (eds.): *Rough-Neural Computing*. Springer-Verlag, Heidelberg, 2004

- A. Wojna: Analogy-based Reasoning in Classifier Construction. Ph. D. Thesis, Department of Mathematics, Informatics and Mechanics, The University of Warsaw, Warsaw, 2004
- J. Bazan, Nguyen H. S., Nguyen S. H., P. Synak, J. Wróblewski: *Rough set algorithms in classification problem*. In: L. Polkowski, S. Tsumoto, T.Y. Lin (eds.), *Rough Set Methods and Applications*, Physica-Verlag, Heidelberg, 2000, pp. 49–88
- Nguyen H.S., Nguyen S.H.: Rough Sets and Association Rule Generation. Fundamenta Informaticae, Vol. 40/4, 1999, pp. 310-318

Links

- International Rough Set Society http://www.roughsets.org
- Rough Set Database System (RSDS) <u>http://rsds.univ.rzeszow.pl/</u>
- KDNuggets <u>http://www.kdnuggets.com/</u>

Software

- <u>Rosetta http://www.lcb.uu.se/tools/rosetta/</u>
- <u>RSES http://logic.mimuw.edu.pl/~rses/</u>
- <u>ROSE, 4eMka, jRS, JAMM, etc. -</u> <u>http://idss.cs.put.poznan.pl/site/software.html</u>
- <u>Rseslib http://rsproject.mimuw.edu.pl</u>
- <u>Debellor http://www.debellor.org</u>
- <u>TunedIT-http://www.tunedit.org/</u>
- Infobright (ICE) http://www.infobright.org