

**Attractiveness measures
for decision rules induced from data
-critical survey**

Izabela Brzezińska

Krytyczny przegląd miar atrakcyjności
reguł decyzyjnych wyindukowanych z danych

Plan

- n Introduction
- n Semantics of attractiveness measures
- n Knowledge representation semantics
 - n Brief attractiveness measures survey
 - n Concept of mining the Pareto-optimal support/certainty border
 - n Computational experiments
- n Conclusions
- n Further research plans
- n References

Introduction

- n $S = \langle U, A \rangle$ – *data table*, where U and A are finite, non-empty sets
 U – universe; A – set of attributes
- n $S = \langle U, C, D \rangle$ – *decision table*, where C – set of *condition attributes*,
 D – set of *decision attributes*, $C \cap D = \emptyset$
- n *Decision rule* or *association rule* induced from S
is a *consequence relation*: $F \textcircled{R} Y$ read as *if F , then Y*
where F and Y are condition and decision formulas expressed as attribute-value pairs

Introduction

n $\|\Phi\|$ is the set of all objects from U , having property F

n $\|\Psi\|$ is the set of all objects from U , having property Y

n *Support* of decision rule $F \rightarrow Y$:

$$\text{sup}(f \rightarrow y) = \text{card}(\|f \wedge y\|)$$

n *Certainty factor* for decision rule $F \rightarrow Y$ (Łukasiewicz, 1913): (called also *confidence*)

$$\text{cert}(f \rightarrow y) = \frac{\text{card}(\|f \wedge y\|)}{\text{card}(\|f\|)}$$

Semantics of attractiveness measures

- n In all practical applications, like *medical practice, market basket, customer satisfaction or risk analysis*, it is crucial to know how good the rules are for:
 - n knowledge representation
 - n prediction
 - n efficient intervention
- n “How good” is a question about *attractiveness measures* of rules
- n Review of literature shows that *there is no single measure which would be the best* for applications in all possible perspectives
(e.g. Bayardo and Agrawal 1999, Greco, Pawlak & Slowinski 2004, Yao & Zhong 1999, Hilderman and Hamilton 2001)
- n **Claim 1:** the adequacy of interestingness measures is dependent on their *semantics*

Knowledge representation semantics

- n Among commonly used attractiveness metrics are:
 - n support
 - n certainty (a.k.a. confidence)
 - n conviction
 - n lift
 - n laplace
 - n piatetsky-shapiro
 - n gini
 - n chi-squared
 - n gray-orlowska
 - n kamber-shinghal

- n Several algorithms are known to efficiently find the best rules according to one of these metrics (e.g. *Webb'95, Fukada et al.'96, Rastoni and Shim '98*).

List of all mentioned attractiveness metrics

$$\text{support}(f \rightarrow y) = \text{card}(\|f \wedge y\|)$$

$$\text{certainty}(f \rightarrow y) = \text{cert}(f \rightarrow y) = \frac{\text{card}(\|f \wedge y\|)}{\text{card}(\|f\|)}$$

$$\text{conviction}(f \rightarrow y) = \frac{\text{card}(U) - \text{card}(\|y\|)}{\text{card}(U) * (1 - \text{cert}(f \rightarrow y))}$$

$$\text{lift}(f \rightarrow y) = \frac{\text{card}(U) * \text{cert}(f \rightarrow y)}{\text{card}(\|y\|)} = \frac{\text{card}(U) * \text{card}(\|f \wedge y\|)}{\text{card}(\|y\|) * \text{card}(\|f\|)}$$

$$\text{laplace}(f \rightarrow y) = \frac{\text{card}(\|f \wedge y\|) + 1}{\text{card}(\|f\|) + k}, k = \text{number of classes}$$

$$\text{piatetsky - shapiro} = \text{card}(\|f \wedge y\|) - \frac{\text{card}(\|f\|) * \text{card}(\|y\|)}{\text{card}(U)}$$

List of all mentioned attractiveness metrics

$$\begin{aligned} \text{gini} = & 1 - \left(\frac{\text{card}(\|\psi\|)^2}{\text{card}(U)^2} + \frac{(\text{card}(U) - \text{card}(\|\psi\|))^2}{\text{card}(U)^2} \right) - \\ & - \frac{\text{card}(\|\phi\|)}{\text{card}(U)} * \left(1 - \left(\frac{\text{card}(\|\phi \wedge \psi\|)^2}{\text{card}(\|\phi\|)^2} + \frac{(\text{card}(\|\phi\|) - \text{card}(\|\phi \wedge \psi\|))^2}{\text{card}(\|\phi\|)^2} \right) \right) - \\ & - \frac{\text{card}(\|\sim \phi\|)}{\text{card}(U)} * \left(1 - \left(\frac{\text{card}(\|\sim \phi \wedge \psi\|)^2}{\text{card}(\|\sim \phi\|)^2} + \frac{(\text{card}(\|\sim \phi\|) - \text{card}(\|\sim \phi \wedge \psi\|))^2}{\text{card}(\|\sim \phi\|)^2} \right) \right) \end{aligned}$$

List of all mentioned attractiveness metrics

$$\begin{aligned}
 \text{chi-squared} = & \frac{\text{card}(\|\phi\|) * \left(\frac{\text{card}(\|\phi \wedge \psi\|)}{\text{card}(\|\phi\|)} - \frac{\text{card}(\|\psi\|)}{\text{card}(U)} \right)^2 - \text{card}(\|\sim \phi\|) * \left(\frac{\text{card}(\|\sim \phi \wedge \psi\|)}{\text{card}(\|\sim \phi\|)} - \frac{\text{card}(\|\psi\|)}{\text{card}(U)} \right)^2}{\frac{\text{card}(\|\psi\|)}{\text{card}(U)}} + \\
 & + \frac{\text{card}(\|\phi\|) * \left(\frac{\text{card}(\|\phi \wedge \sim \psi\|)}{\text{card}(\|\phi\|)} - \frac{\text{card}(\|\sim \psi\|)}{\text{card}(U)} \right)^2 - \text{card}(\|\sim \phi\|) * \left(\frac{\text{card}(\|\sim \phi \wedge \sim \psi\|)}{\text{card}(\|\sim \phi\|)} - \frac{\text{card}(\|\sim \psi\|)}{\text{card}(U)} \right)^2}{\frac{\text{card}(\|\sim \psi\|)}{\text{card}(U)}} +
 \end{aligned}$$

List of all mentioned attractiveness metrics

$$\text{gray - orlowska} = \left(\left(\frac{\text{card}(\|f \wedge y\|) * \text{card}(U)}{\text{card}(\|f\|) * \text{card}(\|y\|)} \right)^l - 1 \right) * \left(\frac{\text{card}(\|f\|) * \text{card}(\|y\|)}{\text{card}(U)^2} \right)^m$$

where l, k are parameters to weight the relative importance of the discrimination and support components.

$$\text{kamber-shinghal} = \frac{\text{card}(f \wedge y)}{\text{card}(f \wedge \sim y)} * \left(1 - \frac{\text{card}(\sim f \wedge y)}{\text{card}(\sim f \wedge \sim y)} \right)$$

The first component measures „how sufficient is Φ for Ψ ”.

The second component measures „how necessary is Φ for Ψ ”.

Conviction metric*

$$\text{conviction}(f \rightarrow y) = \frac{\text{card}(U) - \text{card}(\|y\|)}{\text{card}(U) * (1 - \text{cert}(f \rightarrow y))}$$

$$\text{conviction}(f \rightarrow y) = \frac{\text{card}(\|f\|) * \text{card}(\|\sim y\|)}{\text{card}(U) * \text{card}(\|f \wedge \sim y\|)}$$

$$\text{conviction}(f \rightarrow y) = \frac{\text{card}(\|\sim y\|)}{\text{card}(U) * \text{cert}(f \rightarrow \sim y)}$$

n The metric takes into account occurrence of objects with $\sim\Psi$ in the decision table.

n **Division by zero** occurs when there are no (Φ and $\sim\Psi$) objects.

It is a **strong disadvantage** of the metric!

*Bayardo, R.J.; Agrawal, R.; and Gunopulos, D. 1999. Constraint-Based Rule Mining in Large, Dense Databases. In *Proc. of the 15th Int'l Conf. on Data Engineering*, 188-197.

Piatetsky-Shapiro's metric*

$$piatetsky - shapiro = card(\|f \wedge y\|) - \frac{card(\|f\|) * card(\|y\|)}{card(U)}$$

- n This rule interest function is used to quantify the correlation between condition and decision attributes.
- n When $ps=0$, then F and Y are statistically independent and the rule is not interesting.
- n When $ps>0$ ($ps<0$), then F is positively (negatively) correlated to Y .
- n The metric does not take into account occurrence of objects with $\sim F$ nor $\sim Y$ in the decision table.
- n It can be transformed to be identical to *gain metric***.

*Piatetsky-Shapiro, G. 1991. Discovery , Analysis, and Presentation of Strong Rules. Chapter 13 of *Knowledge Discovery in Databases*, AAAI/MIT press, 1991.

**Fukada, T. et al. 1996. Data Mining using Two-Dimensional Optimized Association Rules: Scheme, Algorithms, and Visualization. In *Proc. of the 1996 ACM-SIGMOD Int'l Conf. on the Management of Data*, 13-23.

Lift metric*

$$\text{lift}(f \rightarrow y) = \frac{\text{card}(U) * \text{cert}(f \rightarrow y)}{\text{card}(\|y\|)} = \frac{\text{card}(U) * \text{card}(\|f \wedge y\|)}{\text{card}(\|y\|) * \text{card}(\|f\|)}$$

nThe metric is not straightforwardly influenced by number of objects with $\sim \Psi$ or $\sim \Phi$ in the decision table.

nMeasures the „independency“ of Φ and Ψ .

nLooks very similar to Horwitch'82 confirmation measure r.

$$r(y | f) = \log\left(\frac{\text{Pr}(y \wedge f)}{\text{Pr}(y) * \text{Pr}(f)}\right), \text{ where } \text{Pr}(X) = \frac{\text{card}(\|X\|)}{\text{card}(U)}$$

* International Business Machines, 1996. *IBM Intelligent Miner User's Guide*, Ver 1, Rel1.

Bayardo-Agrawal concept

- n Bayardo and Agrawal '99 introduced a concept involving a **partial order** on rules defined in terms of **support** and **certainty**.
- n They demonstrated that the set of rules that are optimal according to this partial order **includes all rules** that are **best** according to any of these metrics.

Optimized rule mining – problem statement

- n The input to the problem of mining optimized rules: $\langle K, U, \ell, L, N \rangle$
 - n K is a finite set of conditions;
 - n U is a data set;
 - n ℓ is a total order on rules;
 - n L is a condition specifying the rule consequent;
 - n N is a set of constraints on rules (e.g. minimum support, certainty).
- n **Optimized rule mining problem statement:**
Find a rule r_1 such that:
 1. r_1 satisfies the input constraints, and
 2. there exists no r_2 such that r_2 satisfies the input constraints and $r_1 < r_2$.

Mining optimized rules under partial order

- n With **partial order**, because some rules may be incomparable, there can be several **equivalence classes** containing optimal rules.
- n The previous problem statement requires an algorithm to identify **only a single rule** from one of these equivalence classes.
- n To mine at least one representative from each equivalence classes that contains an optimal rule, we need to rephrase the mining problem.
- n **Partial-order optimized rule mining problem statement:**
Find a set R of rules such that:
 1. *every rule r_i in R is optimal as defined by the optimized rule mining problem*
 2. *for every equivalence class of rules if the equivalence class contains an optimal rule, then exactly one member of this equivalence class is in R .*

Support-certainty optimality

n Consider the following partial order \leq_{sc} on rules. Given rules r_1 and r_2 , $r_1 \leq_{sc} r_2$ if and only if:

- $\text{sup}(r_1) \neq \text{sup}(r_2) \cup \text{cert}(r_1) < \text{cert}(r_2)$, or
- $\text{sup}(r_1) < \text{sup}(r_2) \cup \text{cert}(r_1) \neq \text{cert}(r_2)$.

Additionally, $r_1 =_{sc} r_2$ iff $\text{sup}(r_1) = \text{sup}(r_2)$ and $\text{cert}(r_1) = \text{cert}(r_2)$.

An optimal set of rules (optimized rule mining problem solutions) according to this partial order \leq_{sc} is regarded as *sc-upper border*. Intuitively, such a set of rules defines a *support-certainty border* above which no rule that satisfies the input constraints can fall.

Support-~ certainty optimality

n Consider the following partial order \leq_{s-c} on rules. Given rules r_1 and r_2 , $r_1 \leq_{s-c} r_2$ if and only if:

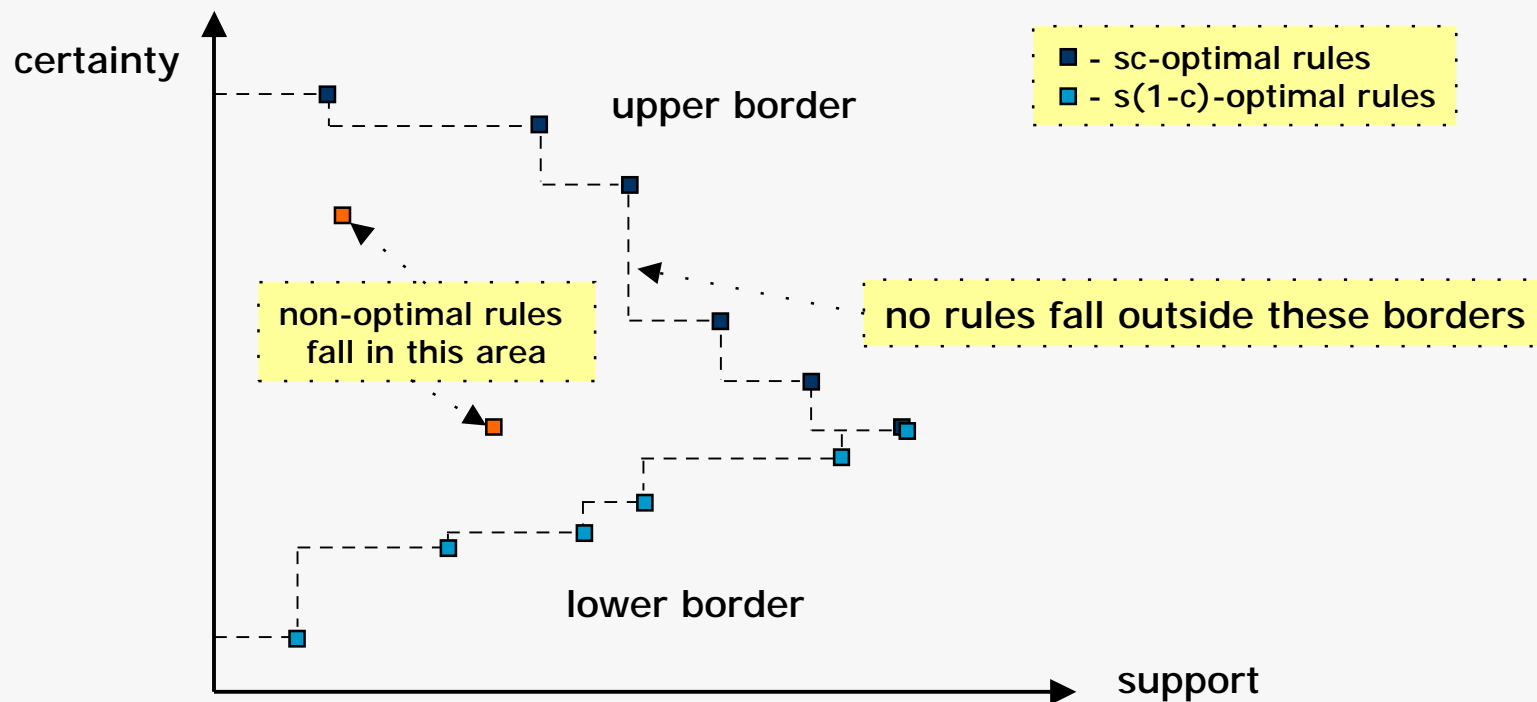
- $\text{sup}(r_1) \geq \text{sup}(r_2) \wedge \text{cert}(r_1) > \text{cert}(r_2)$, or
- $\text{sup}(r_1) < \text{sup}(r_2) \wedge \text{cert}(r_1) \geq \text{cert}(r_2)$.

Additionally, $r_1 =_{s-c} r_2$ iff $\text{sup}(r_1) = \text{sup}(r_2)$ and $\text{cert}(r_1) = \text{cert}(r_2)$.

An optimal set of rules according to this partial order \leq_{s-c} forms a **lower border**.

Support vs certainty pareto-optimal borders

- n Pareto optimal borders:
 - n support vs certainty => upper border
 - n support vs \sim certainty => lower border



Border mining

- n Mining the **upper support/certainty border** identifies optimal rules according to such interestingness metrics:
 - n support
 - n certainty (a.k.a. confidence)
 - n lift
 - n conviction
 - n laplace
 - n piatetsky-shapiro's rule-interest function (p-s)
- n If we also mine the **lower border**, such metrics will also be included:
 - n gini
 - n chi-squared

Theoretical implications – lemma 1

Definition 1:

We say that an intended to rank rules in order of interestingness

total order \leq_t is implied by partial order \leq_{sc} iff

$r_1 <_{sc} r_2 \text{ } \dot{\text{P}} \text{ } r_1 \not\leq_t r_2$, and $r_1 =_{sc} r_2 \text{ } \dot{\text{P}} \text{ } r_1 =_t r_2$.

Lemma 1:

Given the problem instance $I = \langle K, U, \ell_t, L, N \rangle$

such that \leq_t is implied by \leq_{sc} ,

an I -optimal rule is contained within any I_{sc} -optimal set

where $I_{sc} = \langle K, U, \ell_{sc}, L, N \rangle$.

Theoretical implications – proof 1

Proof 1:

Consider any rule r_1 that is not I_{sc} -optimal. Because r_1 is non-optimal, there must exist some rule r_2 that is optimal such that $r_1 <_{sc} r_2$. But then we also have that $r_1 \not\leq_t r_2$ since \leq_t is implied by \leq_{sc} .

This implies that any non- I_{sc} -optimal rule is either non- I -optimal, or it is equivalent to some I -optimal rule which resides in an I_{sc} -optimal equivalence class.

At least one I_{sc} -optimal equivalence class must therefore contain an I -optimal rule. Further, because $=_t$ is implied by $=_{sc}$, every rule in this equivalence class must be I -optimal. By definition, an I_{sc} -optimal set will contain one of these rules, and the claim follows. •

Theoretical implications – lemma 2

To identify the interestingness metrics that are implied by \preceq_{sc} we use:

Lemma 2:

The following conditions are sufficient for establishing that a total order \preceq_t defined over a rule value function $f(r)$ is implied by partial order \preceq_{sc} :

1. $f(r)$ is monotone in support over rules with the same certainty, and
2. $f(r)$ is monotone in certainty over rules with the same support.

Proof 2:

Suppose $r_1 \prec_{sc} r_2$, then consider a rule r where $\text{sup}(r_1) = \text{sup}(r)$ and $\text{cert}(r_2) = \text{cert}(r)$. By definition $r_1 \preceq_{sc} r$ and $r \preceq_{sc} r_2$.

If the total order has the monotonicity properties 1,2,

then $r_1 \preceq_t r$ and $r \preceq_t r_2$. Since total orders are transitive, we then have that $r_1 \preceq_t r_2$, which establishes the claim. •

Example - Laplace function

Consider a Laplace function, which is commonly used to rank rules for classification purposes:

Definition 2:

$$\text{laplace} (f \rightarrow y) = \frac{\text{card} (\|f \wedge y\|) + 1}{\text{card} (\|f\|) + k}$$

where k is a constant integer >1 , usually set to the number of classes when building a classification model.

Because:

$$\text{card} (\|f \wedge y\|) = \text{sup}(\|f \wedge y\|), \text{cert} = \frac{\text{card} (\|f \wedge y\|)}{\text{card} (\|f\|)}$$

we have:

$$\text{laplace} (f \rightarrow y) = \frac{\text{sup}(f \rightarrow y) + 1}{\frac{\text{sup}(f \rightarrow y)}{\text{cert}(f \rightarrow y)} + k}$$

Example - Laplace function

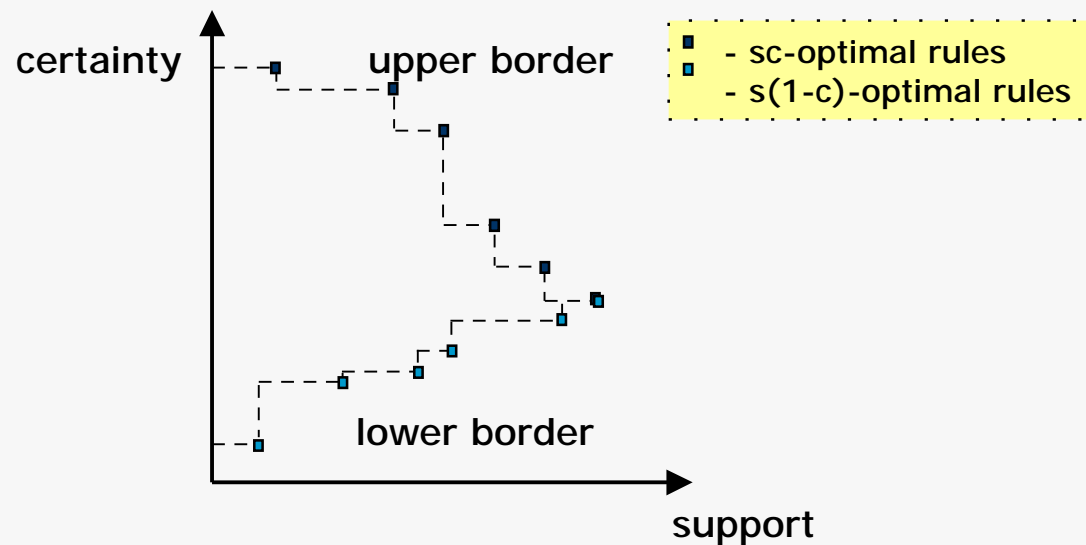
$$\text{laplace } (f \rightarrow y) = \frac{\text{sup}(f \rightarrow y) + 1}{\frac{\text{sup}(f \rightarrow y)}{\text{cert}(f \rightarrow y)} + k}$$

- n This expression is monotone in rule support since $k > 1$ and $\text{cert} \geq 0$.
- n It is also monotone in certainty among rules with equivalent support:
 - n note that if support is held constant, in order to raise the function value, we need to decrease the value of the denominator,
 - n the decrease of the denominator can only be achieved by increasing certainty.

Other attractiveness metrics included in the upper border

- n Bayardo and Agrawal '99 have also showed that the total orders listed below are also implied by partial order \mathcal{L}_{sc} :

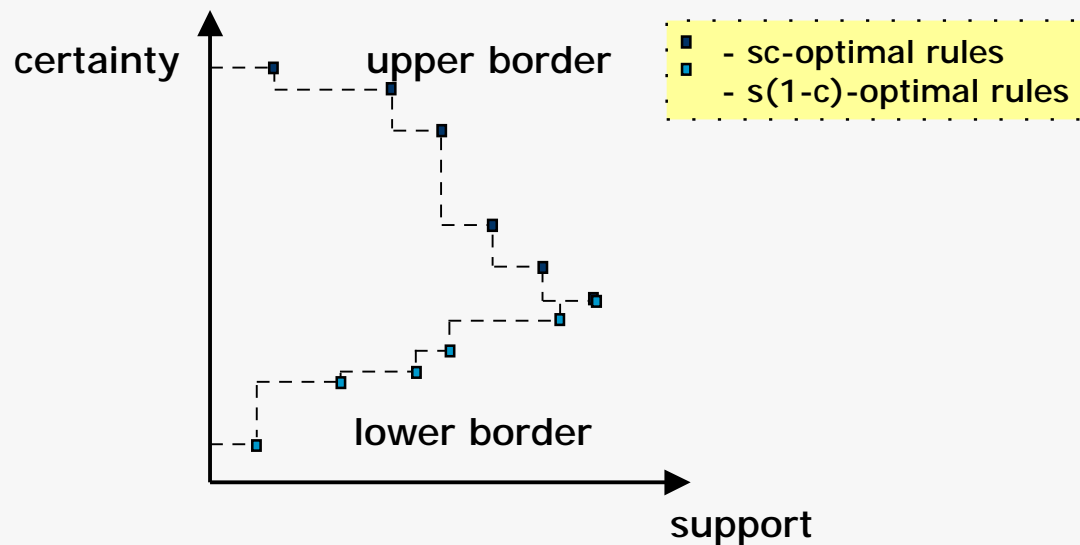
- n support
- n certainty
- n conviction
- n lift
- n laplace
- n piatetsky-shapiro



- n Thus, mining the upper support/certainty border identifies optimal rules according to these metrics.

Mining the lower border

- n The following metrics are **not** implied by \mathcal{L}_{sc} .
 - n gini
 - n chi-squared



- n However, Bayardo and Agrawal have shown that the optimal rules with respect to these metrics must reside on **either the upper or lower support/certainty border**.

Computational experiment

- n Decision rules were generated from **lower approximations** of preference-ordered decision classes defined according to **Variable-consistency Dominance-based Rough Set Approach (VC-DRSA)** (Greco, Matarazzo, Slowinski, Stefanowski 2001)

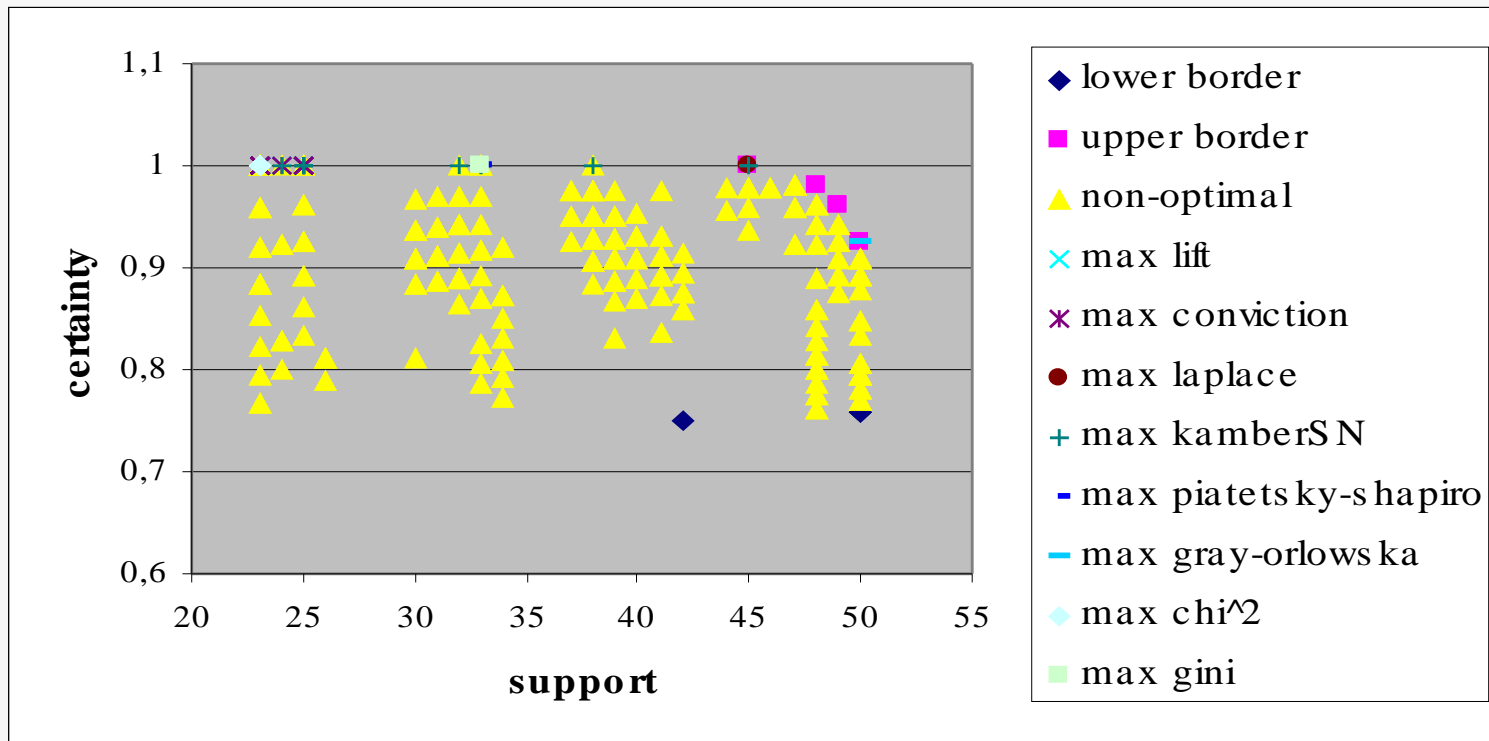
File	objects	atr+crit	classes	rules (alg)	consistency	length	coverage
<i>Buses</i>	76	0+8	3	266 (all)	≥ 0.75	≤ 3	≥ 0.9
<i>Nativity</i>	342	0+33	2	64 (mc)	≥ 0.75	no limit	no limit
<i>Urology</i>	500	18+9	3	186 (mc)	≥ 0.96	no limit	no limit

Rule induction algorithms: „all” = all rules algorithm (DOMAPRIORI)

„mc” = minimal-cover algorithm (DOMLEM)

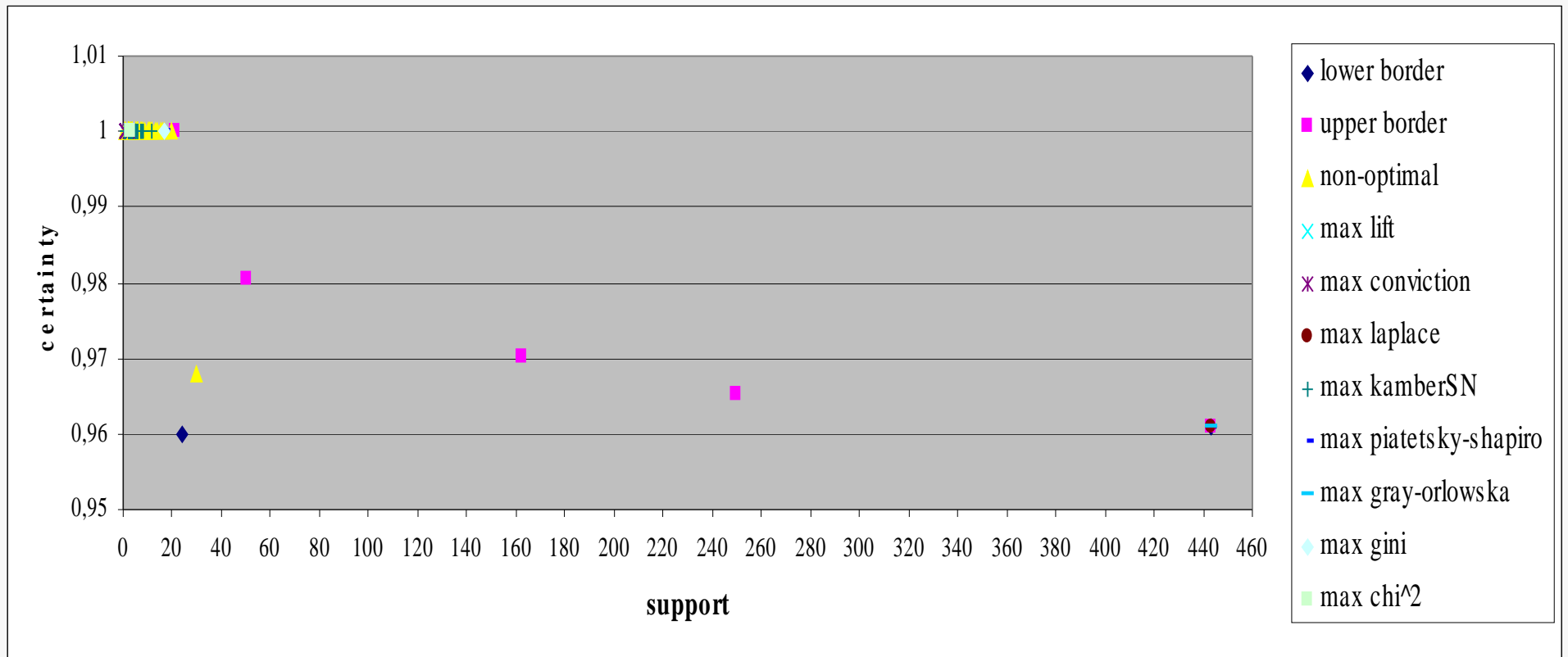
Computational experiment - Buses

n *Buses*



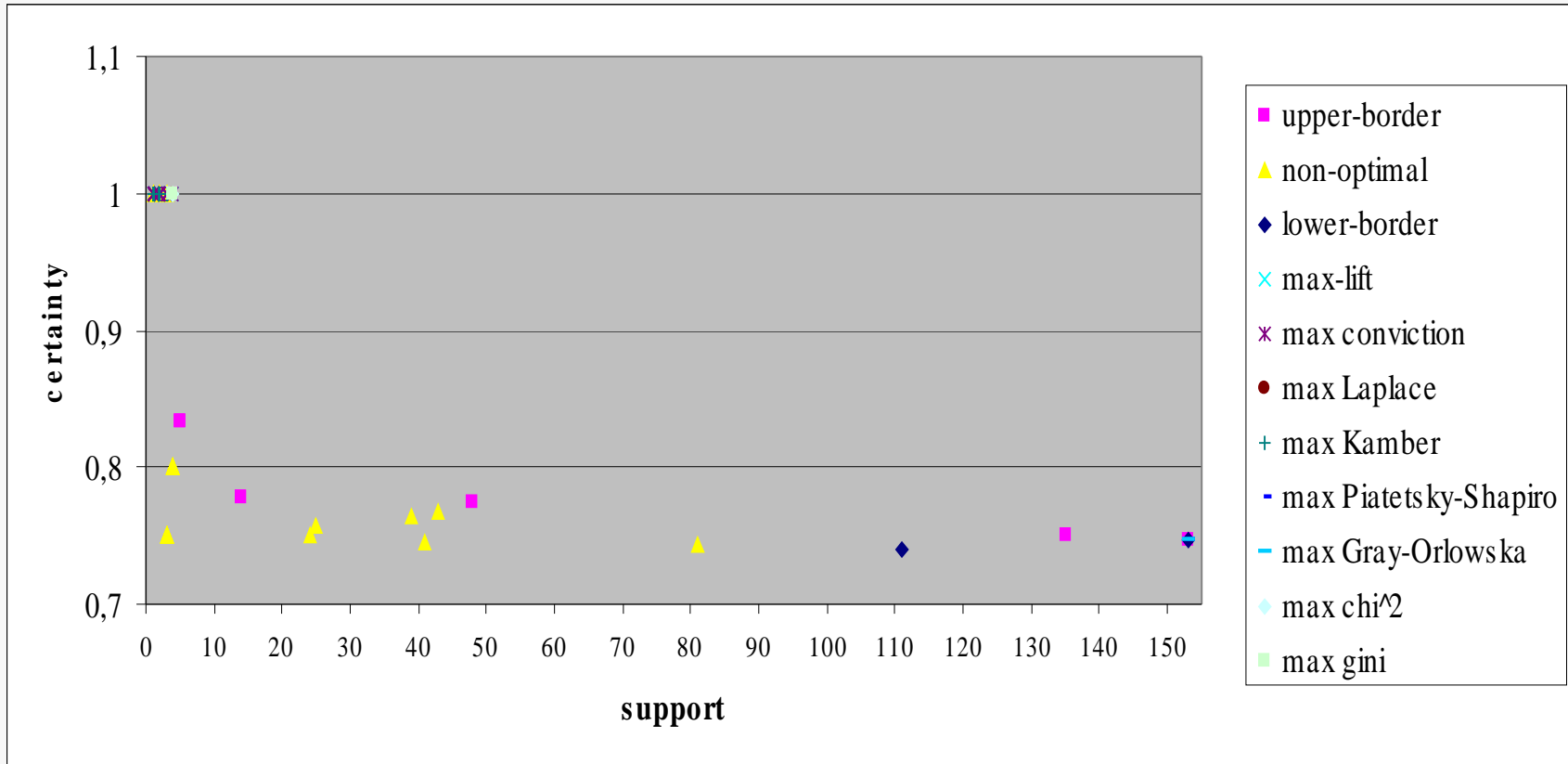
Computational experiment - Urology

n *Urology*



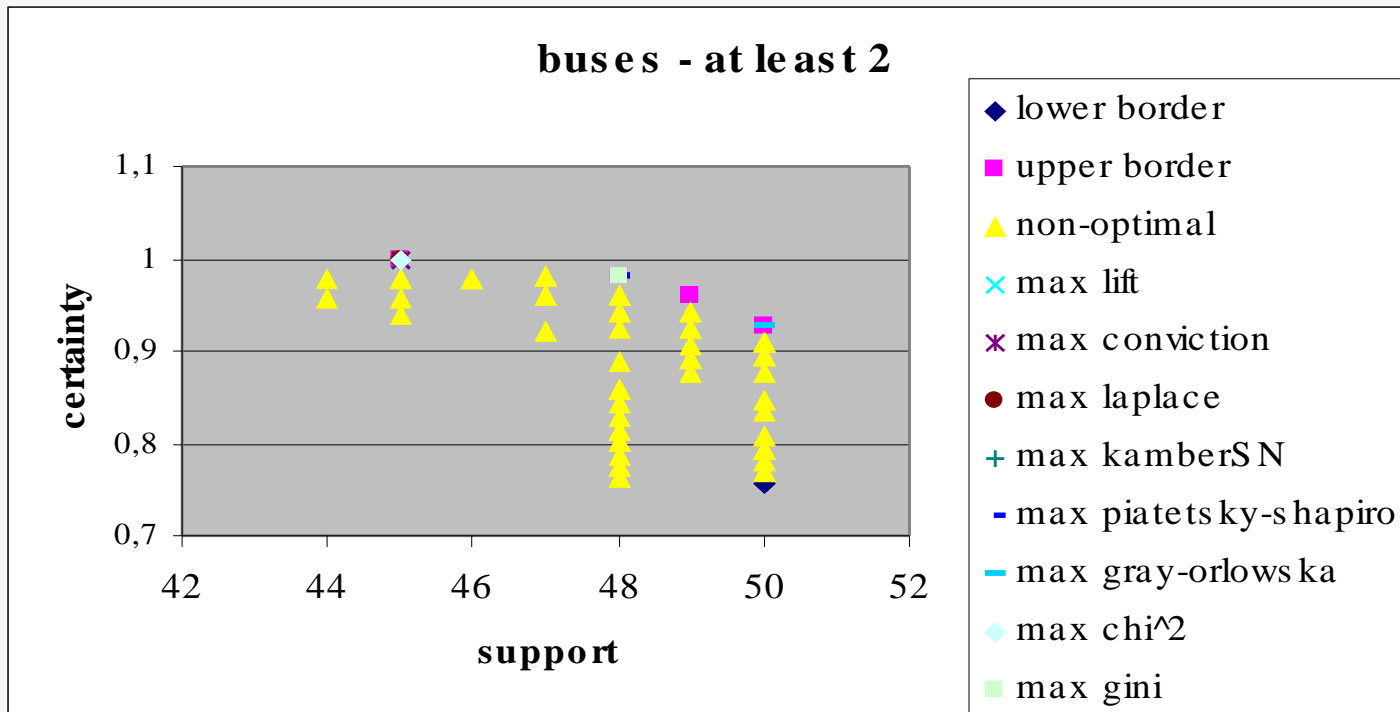
Computational experiment - Nativity

n *Nativity*



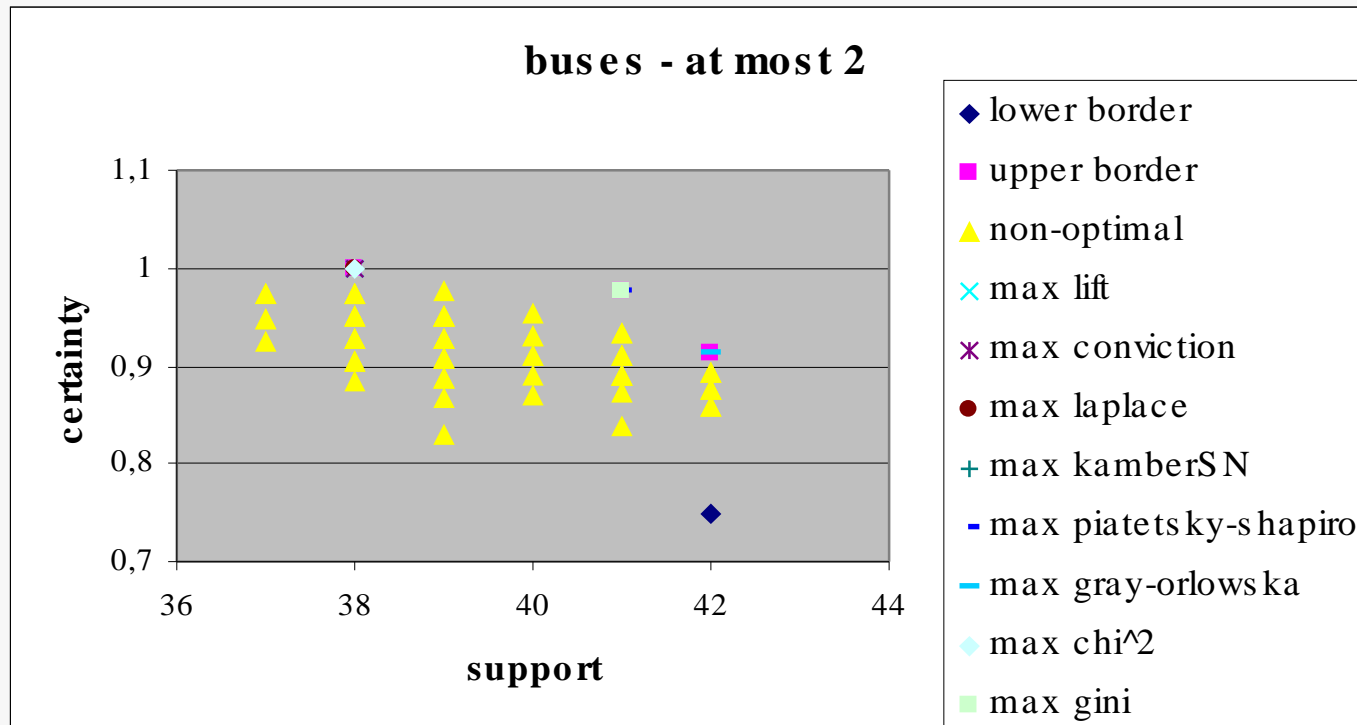
Computational experiment – Buses “at least 2”

n *Buses – for union of classes 1 and 2*



Computational experiment – Buses “at most 2”

n *Buses – for union of classes 2 and 3*



Conclusions – knowledge representation

- n A survey on attractiveness measures in knowledge representation aspect has been done.
- n By Bayardo and Agrawal a new optimized rule mining problem has been defined. It allows a partial order in place of the typical total order on rules.
- n Solving this optimized rule mining problem with respect to a particular partial order \leq_{sc} is guaranteed to identify a most-interesting rule according to several attractiveness metrics including: (support, certainty, laplace, conviction, piatetsky-shapiro, lift, gini, chi-squared).
- n The computational experiment has expressed that indeed Pareto optimal support/certainty border contains rules optimal with respect to any of those metrics.
- n Moreover, the computational experiment placed in upper support/certainty border also rules optimal according to kamber-shinghal and gray-orlowska metrics. However, an analytical proof is required.

Further research

- n Computational experiments have placed rules optimal according to gray-orlowska metric and kamber-shinghal metric in the upper support/certainty border. Can it be analitically verified whether these total orders are implied by partial order \leq_{sc} ?
- n Is it possible to imply discussed total orders (like: support, certainty, laplace, etc.) by partial order other than support/certainty?
- n Are some metrics (eg. lift) confirmation metrics? Analitical proof of posessing hypothesis symmetry and monotonicity properties.
- n From decision to association rules...

References

- n Bayardo, R.J., Agrawal, R.: Mining the most interesting rules. [In]: Proc. of the Fifth ACM SIGKDD Intl Conf. on Knowledge Discovery and Data Mining, (1999) 145-154
- n Greco, S., Matarazzo, B., Pappalardo, N., Słowiński, R.: Measuring expected effects of interventions based on decision rules. *Journal of Experimental and Applied Artificial Intelligence*, 17 (1-2) (2005) 103-118
- n Greco, S., Pawlak, Z., Słowiński, R.: Can Bayesian confirmation measures be useful for rough set decision rules? *Engineering Applications of Artificial Intelligence*, 17 (4) (2004) 345-361
- n Hilderman, R.J., Hamilton, H.J.: *Knowledge Discovery and Measures of Interest*. Kluwer Academic Publishers, Boston (2001)
- n Yao, Y.Y., Zhong, N.: An analysis of quantitative measures associated with rules. [In]: *Proceedings of the 3rd Pacific-Asia Conference on Knowledge Discovery and Data Mining*, LNAI 1574, Springer-Verlag, Berlin (1999) 479-488