# Learning Rules With Numerical and Categorical Attributes from Linked Data Sources

Andre de Oliveira Melo

Saarland University

*andresony@gmail.com*

March 18, 2013

# Overview

# Semantic Web

### Semantic Web

*"provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries"*
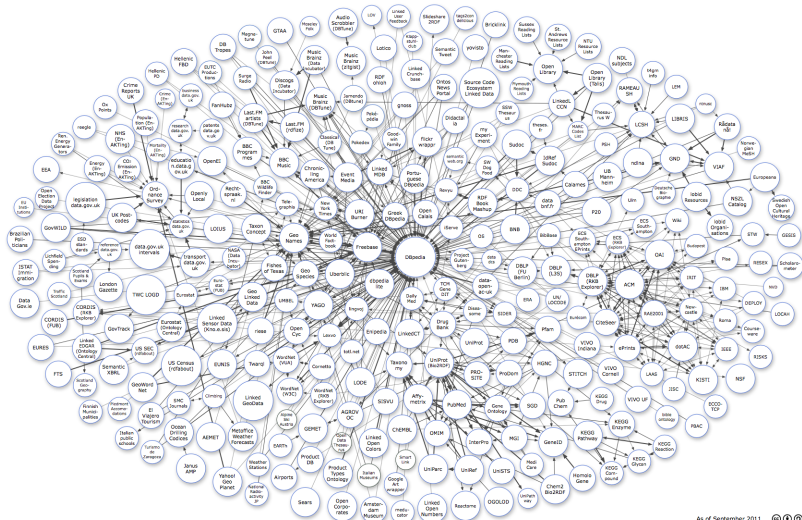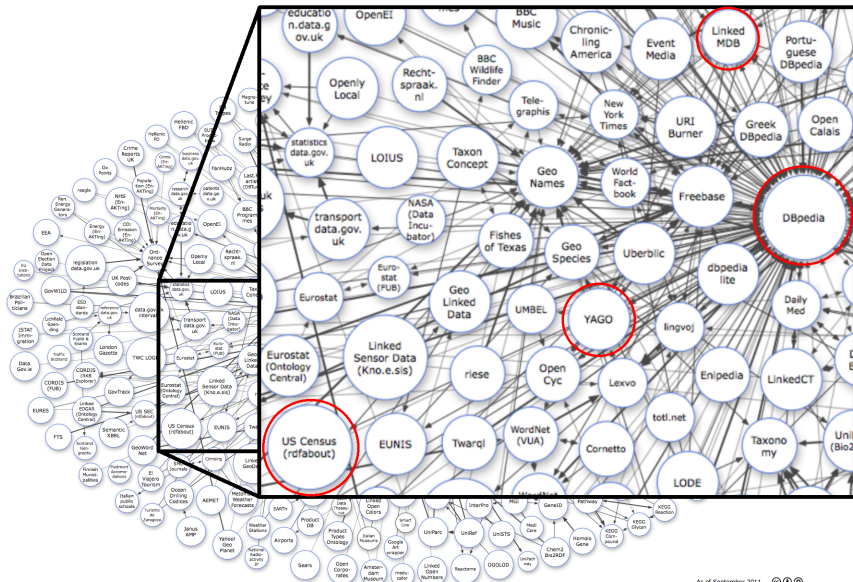
# Semantic Web

### Semantic Web

*"provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries"*

### Linked Data

*"collection of interrelated datasets on the Web"*

*"recommended best practises for exposing, sharing, and connecting pieces of data, information and knowledge on the Semantic Web"*

As of September 2011

As of September 2011

## Motivation

Learn Datalog rules from data:

$$\underbrace{livesIn(X, Y)}_{head} \text{ :- } \underbrace{isMarriedTo(X, Z), livesIn(Z, Y)}_{body}$$

Support and confidence thresholds

- Support: $supp(head \text{ :- } body) = supp(head \wedge body)$
- Confidence: $conf(head \text{ :- } body) = \dfrac{supp(head \wedge body)}{supp(body)}$

## Rules with constants

Refining rules with constants is relevant

$$speaks(X, Z) :- livesIn(X, W)$$

Searching constants for $Z$ and $W$ we can learn:

$$speaks(X, englsih) :- livesIn(X, australia)$$
$$speaks(X, spanish) :- livesIn(X, argentina)$$
$$speaks(X, portuguese) :- livesIn(X, brasil)$$

What about numerical constants?

$$speaks(X, english) :- hasIncome(X, \$3.71Billion)$$
$$speaks(X, portuguese) :- livesIn(X, W), hasPopulation(W, 193946886)$$

# Refining rules with numerical intervals

*maritalStatus(X, single) :- age(X,Y)* [conf=0.40]
*maritalStatus(X,married) :- age(X,Y)* [conf=0.46]
*maritalStatus(X,widowed) :- age(X,Y)* [conf=0.06]

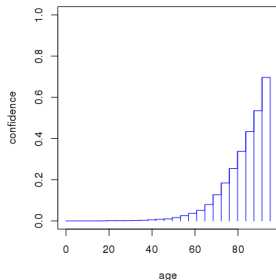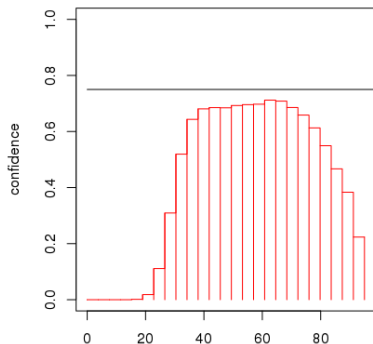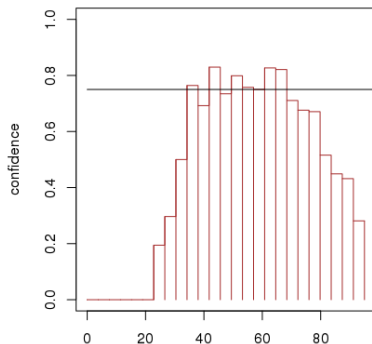Single

Married

Widowed

# Combine with categorical constants

For *maritalStatus*($X$, *married*), *minConf* = 0.75 is not satisfied.
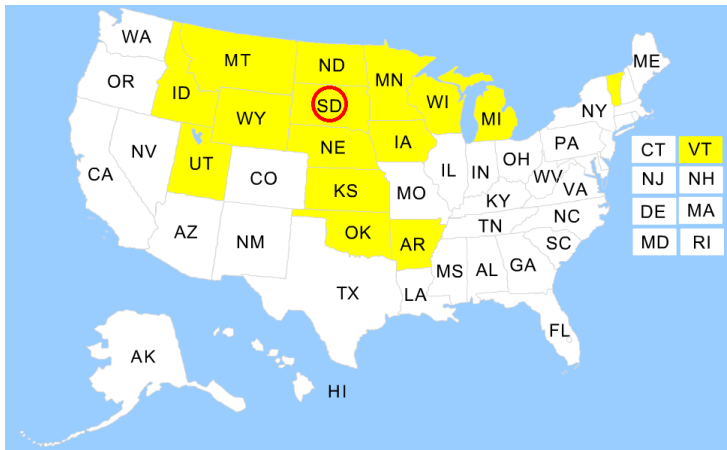Refine by State?

USA

South Dakota



$\Rightarrow$

# Combine with categorical constants

We can find intervals for $Y$ that satisfy $minConf = 0.75$ for
$maritalStatus(X,widowed) :- livesIn(X,sd), age(X,Y)$ and...

## Base-rule and Refined-rule

- **Base-rule**: Numerical argument with no constant
  $r_1 : marritalStatus(X, single) :- livesIn(X, sd), age(X, Y)$
  [conf=0.49,supp=2368]
- **Refined-rule**: Base-rule with restricted numerical variable
  $r_2 : marritalStatus(X, single) :- livesIn(X, sd), age(X, Y), Y \in [33, 67]$
  [conf=0.77,supp=1092]

We are interested in refinements that bring a significant confidence gain:
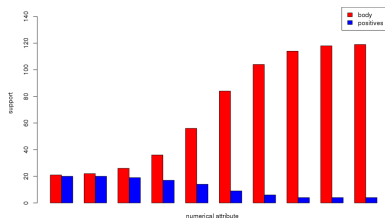
$$gain_{r_{ref}, r_{base}} = \frac{conf(r_{ref})}{conf(r_{base})} \qquad (1)$$

For our example: $gain_{r_2, r_1} = \frac{0.77}{0.49} = 1.57$

What base-rules have refined-rules with significant confidence gain?

- Satisfy support threshold
- Do not necessarily satisfy confidence threshold
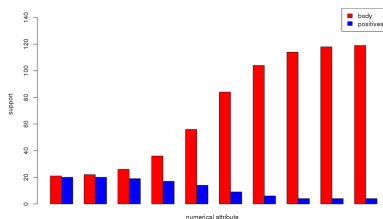- Divergent body and positives (body∧head) probability distributions

Frequency histograms

What base-rules have refined-rules with significant confidence gain?

- ▶ Satisfy support threshold
- ▶ Do not necessarily satisfy confidence threshold
- ▶ Divergent body and positives (body∧head) probability distributions

Frequency histograms

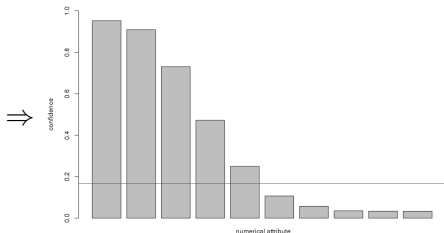Confidence distribution



$\Rightarrow$

What base-rules have refined-rules with significant confidence gain?

- ▶ Satisfy support threshold
- ▶ Do not necessarily satisfy confidence threshold
- ▶ Divergent body and positives (body∧head) probability distributions

Body and Positives distributions

Confidence distribution



$\Rightarrow$

## Motivation

Problem?

▶ Search space grows exponentially with the number of predicates and constants

▶ Querying support and confidence distributions is very expensive

Idea:

▶ Analyze combinations of numerical and categorical properties

▶ Measure their level of interestingness

▶ Extend top-down ILP to detect and suggest interesting combinations

# Logic Programming Concepts

- ▶ Literal: predicate symbol with bracketed n-tuple, e.g:
  $L = livesIn(X, Y)$
- ▶ Clause: a disjunction of literals (negated or not), e.g:
  $c = (L_1 \vee L_2 \vee \ldots \vee \neg L_{m-1} \vee \neg L_m)$
- ▶ Safe Datalog Rule: every variable in the head appear in a non-negated literal in the body, negated literal variables in the body should appear in some positive literal in the body, e.g.:
  $speaks(X, Y) :- wasBornIn(X, Z), hasOfficialLanguage(Z, Y)$
- ▶ Hypothesis: a set of clauses $\mathcal{H}$
  - ▶ Completeness: $\mathcal{H}$ covers all positive examples
  - ▶ Consistency: $\mathcal{H}$ covers no negative examples

# Inductive Logic Programming (ILP)

Inductive Logic Programming: Finds a hypothesis $\mathcal{H}$ that covers all positive, and no negative examples

$positiveExamples + negativeExamples + backgroundKnowledge \rightarrow hypothesis$

| Training Examples | Background Knowledge |
|---|---|
| daughter(mary,ann) + | parent(ann,mary) |
| daughter(eve,tom) + | parent(ann, tom) |
| daughter(tom,ann) - | parent(tom,eve) |
| daughter(eve,ann) - | parent(tom,ian) |
| | female(ann) |
| | female(mary) |
| | female(eve) |

$\mathcal{H} = daughter(X, Y) :\text{-} female(X), parent(Y, X)$

# Inductive Logic Programming (ILP)

Approaches

- ▶ Bottom-up: Start with least general $\mathcal{H}$ then perform generalizations
- ▶ **Top-down**: Start with most general $\mathcal{H}$ then perform specializations
    - ▶ Specialization loop: adds literals to a clause and ensures consistency
    - ▶ Covering loop: adds clauses to the hypothesis and ensures completeness
    - ▶ Apriori-style pruning

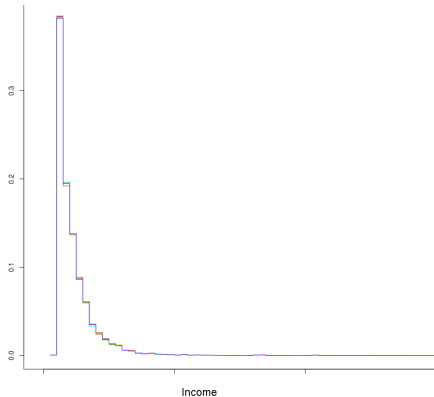What about large, noisy, and incomplete LOD datasets such as YAGO and DBpedia?:

- ▶ Sample data to reduce size
- ▶ Restrict the number of literals in a clause
- ▶ Tolerate a certain level of inconsistency and incompleteness

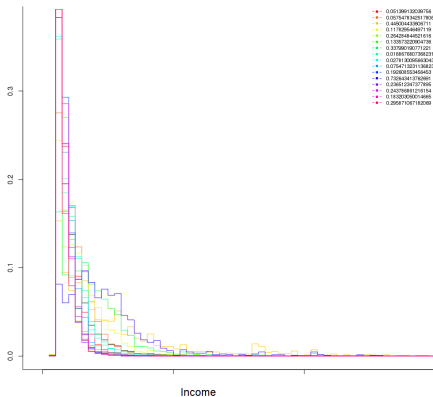    Expected Accuracy: $A(c) = P(e \in \mathcal{E}^+|c) = \dfrac{n^+(c)}{n^+(c) + n^-(c)}$

# Correlation between Literals

Let's say we want to refine a clause with $hasIncome(X, Y)$ with an interval for $Y$. Refine by $quarterOfBirth$ or $hasEducation$?
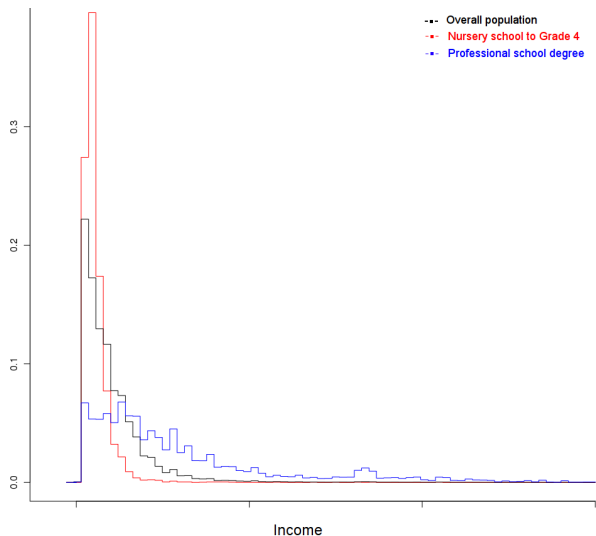
$hasIncome(X, Y), quarterOfBirth(X, Z)$

$hasIncome(X, Y), hasEducation(X, Z)$

# Correlation between Literals

# Interestingness Measure

How to measure the interestingness of adding a literal $l$ to a clause $c$?

- ▶ Extract the frequency histograms of $\{c\}$ and $\{c \wedge l\}$ over a numerical attribute $Y$
- ▶ Normalize the histograms to obtain their probability distributions, and measure their divergence (e.g., with Kullback-Leibler)

But, divergence alone isn't a good idea because:

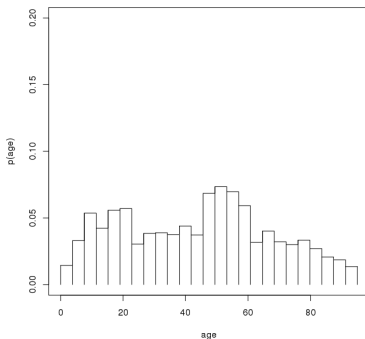- ▶ Lower support histograms are more likely to have a divergent distribution (*sampling error*)
- ▶ Rules with high support are still interesting

Then combine both measures: divergence*support

# Interestingness Measure
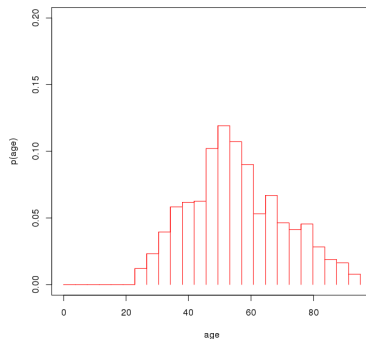
$$\underbrace{maritalStatus(X, married)}_{l=head} \; :\text{-} \; \underbrace{livesIn(X, sd), age(X, Y)}_{c=body}$$



age(X,Y),livesIn(X,sd)

age(X,Y),livesIn(X,sd),marriedStatus(X,married)

# Correlation Lattice

- Build a lattice similar to an *itemset lattice*
- Numerical property $r(X, Y)$ as root
- The "items" are literals that can be joined with the root's non-numerical variable $X$
- Root's numerical attribute $Y$ is discretized in $k$ buckets $\{b_1, \ldots, b_k\}$
- Each node $x$ has a frequency histogram $h(x) = < h_1(x), \ldots, h_k(x) >$ from its clause support distribution

    where $h_i(x) = supp(x|Y \in b_i)$    and    $|h(x)|_1 = supp(x)$

# Correlation Lattice

$hasIncome(X, Y)$

# Correlation Lattice

$hasIncome(X, Y)$

$\cdots$

$hasIncome(X, Y),$
$quarterOfBirth(X, q1)$

$hasIncome(X, Y),$
$quarterOfBirth(X, q4)$

# Correlation Lattice



$hasIncome(X, Y)$

. . .          . . .

$hasIncome(X, Y),$
$quarterOfBirth(X, q1)$

$hasIncome(X, Y),$
$quarterOfBirth(X, q4)$

$hasIncome(X, Y),$
$hasEducation(X, nursery)$

$hasIncome(X, Y),$
$hasEducation(X, phd)$

# Correlation Lattice



$hasIncome(X, Y)$

$\ldots$     $\ldots$     $\ldots$

$hasIncome(X, Y),$
$quarterOfBirth(X, q1)$

$hasIncome(X, Y),$
$quarterOfBirth(X, q4)$

$hasIncome(X, Y),$
$hasEducation(X, nursery)$

$hasIncome(X, Y),$
$hasEducation(X, phd)$

# Correlation Lattice



$hasIncome(X, Y)$

$\cdots$

$\cdots$

$\cdots$

$hasIncome(X, Y),$
$quarterOfBirth(X, q1)$

$hasIncome(X, Y),$
$quarterOfBirth(X, q4)$

$hasIncome(X, Y),$
$hasEducation(X, nursery)$

$hasIncome(X, Y),$
$hasEducation(X, phd)$
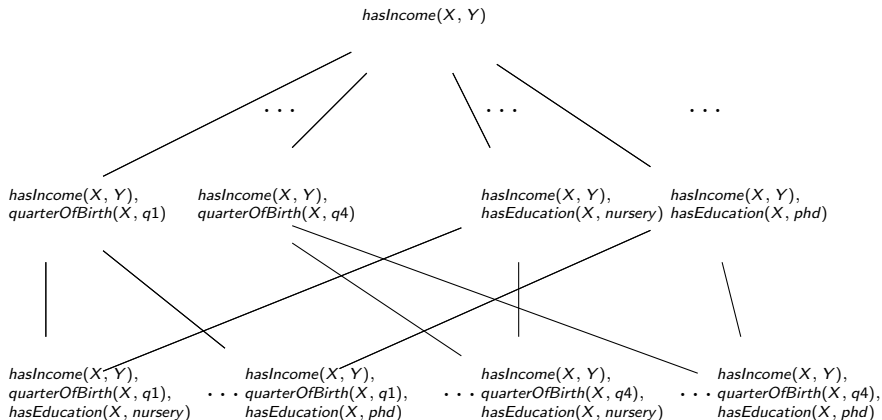
$hasIncome(X, Y),$
$quarterOfBirth(X, q1),$
$hasEducation(X, nursery)$

# Correlation Lattice

# Correlation Lattice

# Correlation Lattice

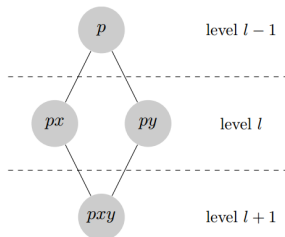- Number of nodes in a lattice with $\ell$ levels $n$ properties and $m$ constants per property:

$$\sum_{i=1}^{\ell} \binom{nm}{i} \tag{2}$$

- Too expensive, we need to reduce size
  - Prune by support (safe)
  - Restrict $\ell$ to the maximum clause size allowed in the core-ILP
  - Restrict the literals added to the lattice in order to reduce $n$ and $m$
  - Prune by interestingness or independence (heuristics)

# Independence checks

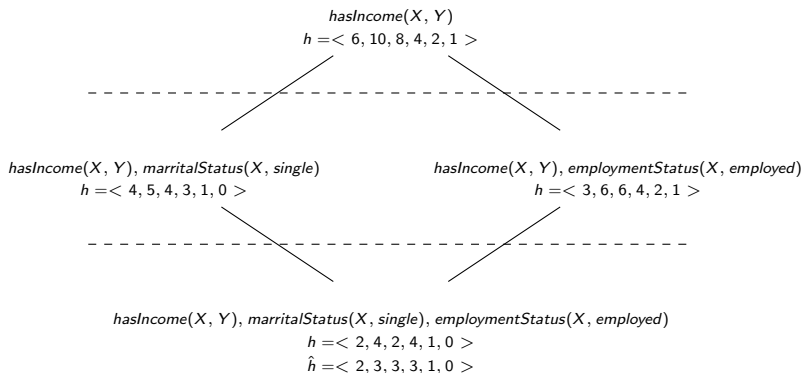▶ Checks if a pair of nodes joining nodes are independent given their
  common parent



(where $p$ is a clause, $x$ and $y$ are literals, s.t. $x \neq y$ and $x, y \notin p$)

▶ Estimate $\hat{h}(pxy)$ assuming independence of $x$ and $y$ given $p$
▶ Query actual $h(pxy)$ and perform a Pearson's chi-squared test
  $H_0 = x$ and $y$ are independent given $p$
  $H_1 = x$ and $y$ are dependent given $p$

# Independence checks

$hasIncome(X, Y)$
$h = < 6, 10, 8, 4, 2, 1 >$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

$hasIncome(X, Y), marritalStatus(X, single)$
$h = < 4, 5, 4, 3, 1, 0 >$

$hasIncome(X, Y), employmentStatus(X, employed)$
$h = < 3, 6, 6, 4, 2, 1 >$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

$hasIncome(X, Y), marritalStatus(X, single), employmentStatus(X, employed)$
$h = < 2, 4, 2, 4, 1, 0 >$
$\hat{h} = < 2, 3, 3, 3, 1, 0 >$

$$\chi^2 = \sum_{i=1}^{k} \frac{(h_i - \hat{h}_i)^2}{\hat{h}_i} = 1 \quad \Rightarrow \quad p\text{-value} = 0.96$$
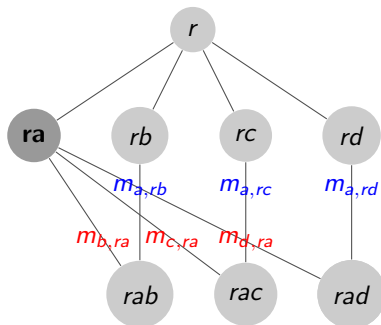
# Independence checks

- If there's not enough evidence of dependence, we assume independence, then:

  $x \text{ :- } p, y \equiv x \text{ :- } p$

  $y \text{ :- } p, x \equiv y \text{ :- } p$

- The lower the p-value (greater $\chi^2$), the greater the evidence that $x$ and $y$ are dependent given $p$, therefore the more interesting it is to join the nodes $py$ and $px$

- As heuristics, we can set a maximum *p-value* threshold to prune independent nodes

# Refinement Suggestions

In the ILP refinement loop, the clauses have a fixed head while the body is refined. Assuming we have $a$ as head literal, $r$ as root and $b$, $c$, $d$ as possible new literals:
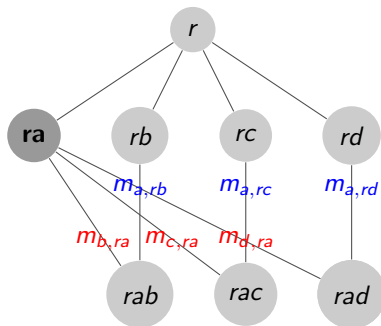


$$a \left| \begin{array}{l} b \ [m_{a,rb}] \\ c \ [m_{a,rc}] \\ d \ [m_{a,rd}] \end{array} \right.$$

What literal is more interesting to add to the clause $a:\text{-}r$?

# Refinement Suggestions

In the ILP refinement loop, the clauses have a fixed head while the body is refined. Assuming we have $a$ as head literal, $r$ as root and $b$, $c$, $d$ as possible new literals:



$$a \left|\; \begin{array}{l} b\; [m_{a,rb}] \\ c\; [m_{a,rc}] \\ d\; [m_{a,rd}] \end{array} \right.$$

What literal is more interesting to add to the clause $a$:-$r$?

$argmax_{i \in \{\; b,c,d\}} m_{a,ri}$
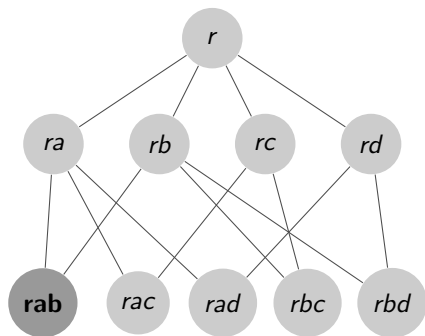
# Search in the Lattice

What has to be done?

- ▶ Search the node with body literals
- ▶ For each child of such node check head literal can be further added, if so collect the new literal and the interestingness value of adding the head
- ▶ Sort the possible new literals by interestingness

Alternative?

- ▶ Create mapping in every node with the possible head literals as key and sorted literals to be added to body as value
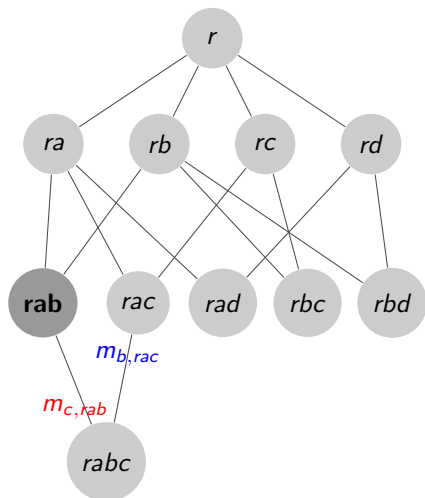- ▶ Only add entry if head and new literal not independent given body

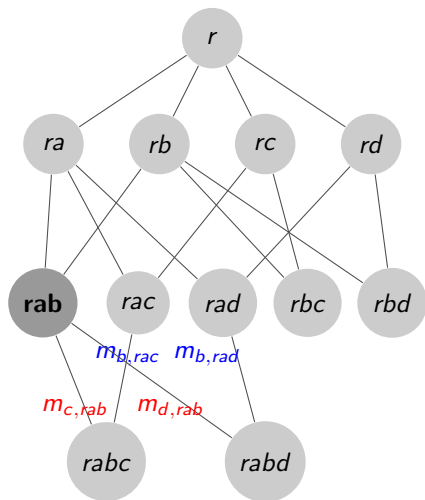# Refinement Suggestions



Suggestions Map

# Refinement Suggestions



Suggestions Map
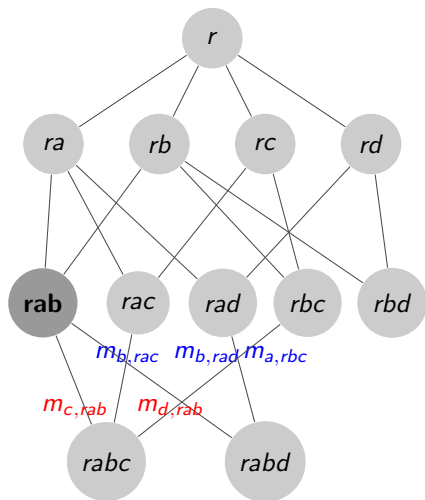
$$b \quad \mid \quad c \ [m_{b,rac}]$$

# Refinement Suggestions



Suggestions Map

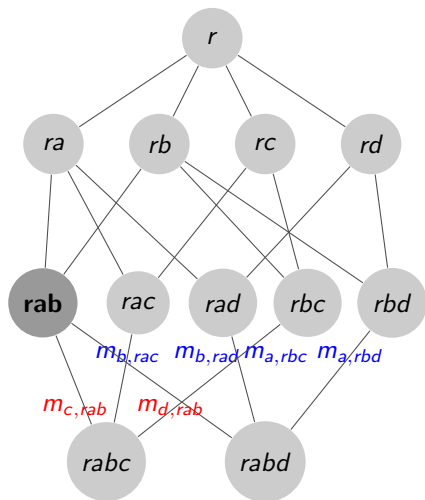$$b \quad \left| \begin{array}{l} c \ [m_{b,rac}] \\ d \ [m_{b,rad}] \end{array} \right.$$

# Refinement Suggestions



Suggestions Map

| $b$ | $c$ $[m_{b,rac}]$ |
| | $d$ $[m_{b,rad}]$ |
| --- | --- |
| $a$ | $c$ $[m_{a,rbc}]$ |

# Refinement Suggestions



Suggestions Map

| $b$ | $c$ $[m_{b,rac}]$ |
| --- | --- |
| | $d$ $[m_{b,rad}]$ |
| $a$ | $c$ $[m_{a,rbc}]$ |
| | $d$ $[m_{a,rbd}]$ |

# Incorporating the Lattice in the Core-ILP

In the refinement step, we detect clauses with body containing a lattice root

- ▶ If clause satisfies support threshold and does not satisfy confidence threshold
- ▶ Then search in lattice for body literals and head
- ▶ Check the interestingness of adding the head to the body and analyze whether to search for numerical intervals
- ▶ Query the lattice for suggestions of interesting literals to be added to the clause

# Experiments

Overall Settings:

- ▶ We compare 4 interestingness measures:
    1. ● *supp*: Support Only
    2. ■ *klsupp*: KL-divergence*Support
    3. ● *kldiv*: KL-divergence Only
    4. ⋆ *jssupp*: JS-divergence*Support
- ▶ Thresholds:
    - ▶ *minConf* = 0.75
    - ▶ *minSupp* = 25
    - ▶ *minGain* = 1.25

$1^{st}$ Experiment: evaluation of the Correlation Lattice

- ▶ All data joined by person only (anonymized)
- ▶ All properties categorical (categories as literals)
- ▶ Create a lattice for *hasIncome* property

$2^{nd}$ Experiment: evaluation of the ILP extension

- ▶ All data joined by person only (anonymized)
- ▶ All properties categorical (categories as literals)

# 1$^{st}$ Experiment

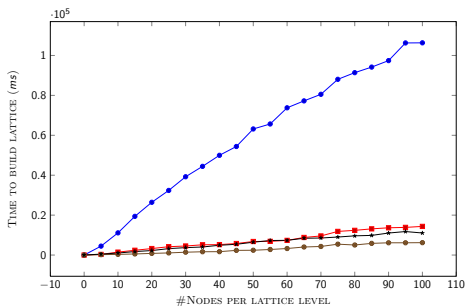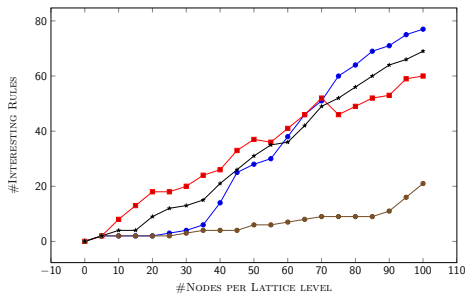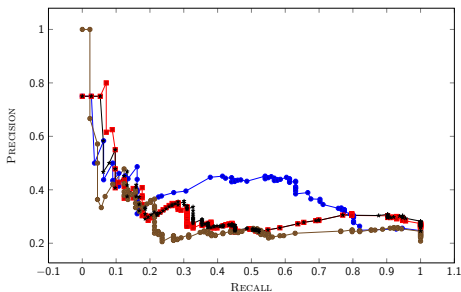Figure: Build time per lattice size



Figure: Interesting rules per lattice size



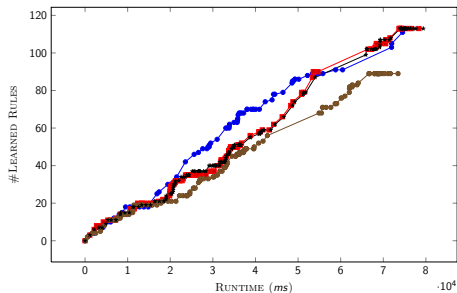Legend: [●*supp* ■ *klsupp* ● *kldiv* ⋆ *jssupp* ]

# $2^{nd}$ Experiment

Figure:  Precision-Recall graph from interestingness predictions (rules with *runtime* attribute)
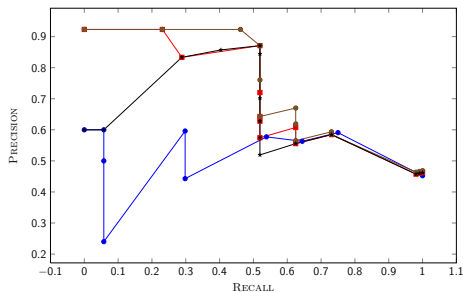
Figure:  Interesting rules per runtime (rules with attribute *runtime*)



Legend: [●*supp* ■ *klsupp* ● *kldiv* ⋆ *jssupp* ]
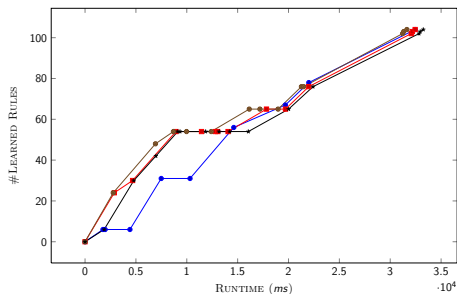
# $1^{st}$ Experiment

Figure: Precision-Recall graph from interestingness predictions (rules with *budget* attribute)

Figure: Interesting rules per runtime (rules with *budget* attribute)



Legend: [●*supp* ■ *klsupp* ● *kldiv* ⋆ *jssupp* ]

# Thank you