

Mining Implications from Lattices of Closed Trees

José L. Balcázar, Albert Bifet and Antoni Lozano

Universitat Politècnica de Catalunya

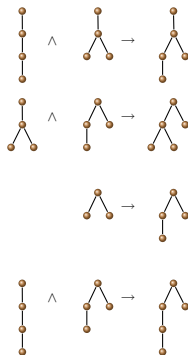
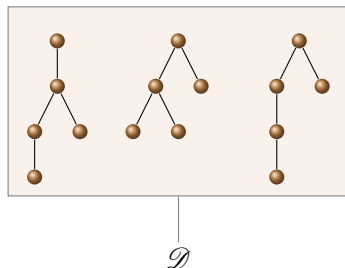
Extraction et Gestion des Connaissances EGC'2008
2008 Sophia Antipolis, France



Introduction

Problem

Given a dataset \mathcal{D} of rooted, unlabelled and unordered trees, find a “basis”: a set of rules that are sufficient to infer all the rules that hold in the dataset \mathcal{D} .

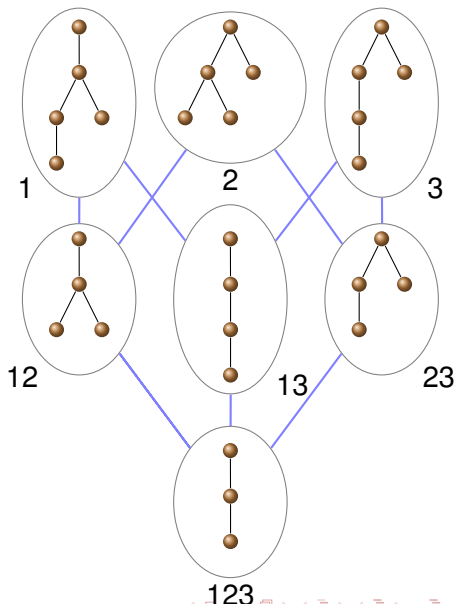


Introduction

Set of Rules:

$$A \rightarrow \Gamma_{\mathcal{D}}(A).$$

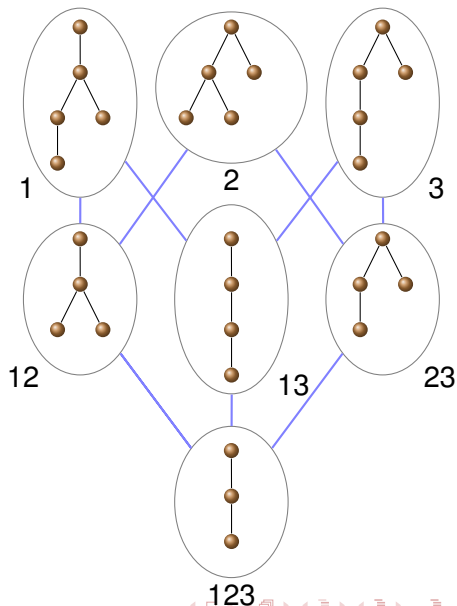
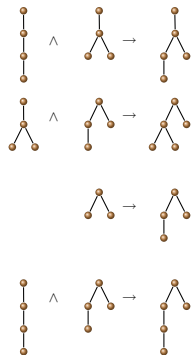
- antecedents are obtained through a computation akin to a hypergraph transversal
- consequents follow from an application of the closure operators



Introduction

Set of Rules:

$$A \rightarrow \Gamma_{\mathcal{D}}(A).$$



Trees

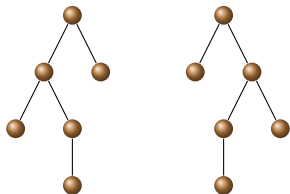
Our trees are:

- Rooted
- Unlabeled
- Unordered

Our subtrees are:

- Induced
- Top-down

Two different ordered trees
but the same unordered tree



Deterministic association rules

- Logical implications are the traditional mean of representing knowledge in formal AI systems. In the field of data mining they are known as **association rules**.

M	a	b	c	d
m_1	1	1	0	1
m_2	0	1	1	1
m_3	0	1	0	1

$$a \rightarrow b, d$$

$$d \rightarrow b$$

$$a, b \rightarrow d$$

- Deterministic association rules are implications with 100% confidence.
- An advantage of deterministic association rules is that they can be studied in purely logical terms with **propositional Horn logic**.

Propositional Horn Logic

M	a	b	c	d
m_1	1	1	0	1
m_2	0	1	1	1
m_3	0	1	0	1

$$a \rightarrow b, d \quad (\bar{a} \vee b) \wedge (\bar{a} \vee d)$$

$$d \rightarrow b \quad \bar{d} \vee b$$

$$a, b \rightarrow d \quad \bar{a} \vee \bar{b} \vee d$$

- Assume a finite number of variables.
 - $V = \{a, b, c, d\}$
- A clause is **Horn** iff it contains at most one positive literal.
 - $\bar{a} \vee \bar{b} \vee d$ $a, b \rightarrow d$
- A **model** is a complete truth assignment from variables to $\{0, 1\}$.
 - $m(a) = 0, m(b) = 1, m(c) = 1, \dots$
- Given a set of models M, the Horn theory of M corresponds to the conjunction of all Horn clauses satisfied by all models from M.

Theorem

Given a set of models M , there is exactly one minimal Horn theory containing it. Semantically, it contains all the models that are intersections of models of M . This is sometimes called the **empirical Horn approximation**.

We propose

- Closure operator
- translation of tree set of rules to a specific propositional theory

Closure Operator

- \mathcal{D} : the finite input dataset of trees
- \mathcal{T} : the (infinite) set of all trees

Definition

We define the following the Galois connection pair:

- For finite $A \subseteq \mathcal{D}$
 - $\sigma(A)$ is the set of subtrees of the A trees in \mathcal{T}

$$\sigma(A) = \{t \in \mathcal{T} \mid \forall t' \in A (t \preceq t')\}$$

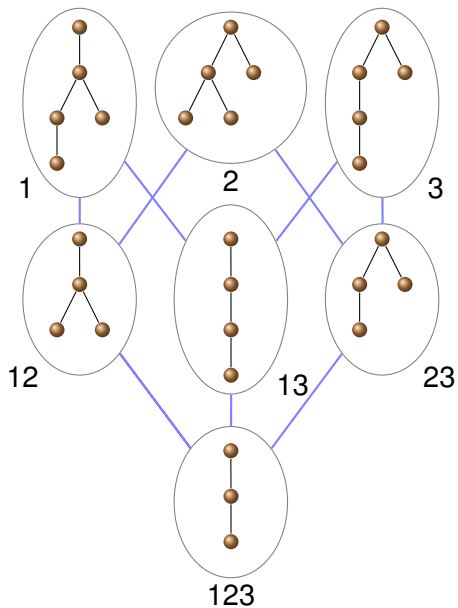
- For finite $B \subset \mathcal{T}$
 - $\tau_{\mathcal{D}}(B)$ is the set of supertrees of the B trees in \mathcal{D}

$$\tau_{\mathcal{D}}(B) = \{t' \in \mathcal{D} \mid \forall t \in B (t \preceq t')\}$$

Closure Operator

The composition $\Gamma_{\mathcal{D}} = \sigma \circ \tau_{\mathcal{D}}$ is a closure operator.

Galois Lattice of closed set of trees



Intuition

- One propositional variable v_t is assigned to each possible subtree t .
- A set of trees A corresponds in a natural way to a model m_A .
- Let m_A be a model: we impose on m_A the constraints that if $m_A(v_t) = 1$ for a variable v_t , then $m_A(v_{t'}) = 1$ for all those variables $v_{t'}$ such that $v_{t'}$ represents a subtree of the tree represented by v_t .

$$\mathcal{R}_0 = \{v_{t'} \rightarrow v_t \mid t' \preceq t, t \in \mathcal{U}, t' \in \mathcal{U}\}$$

Theorem

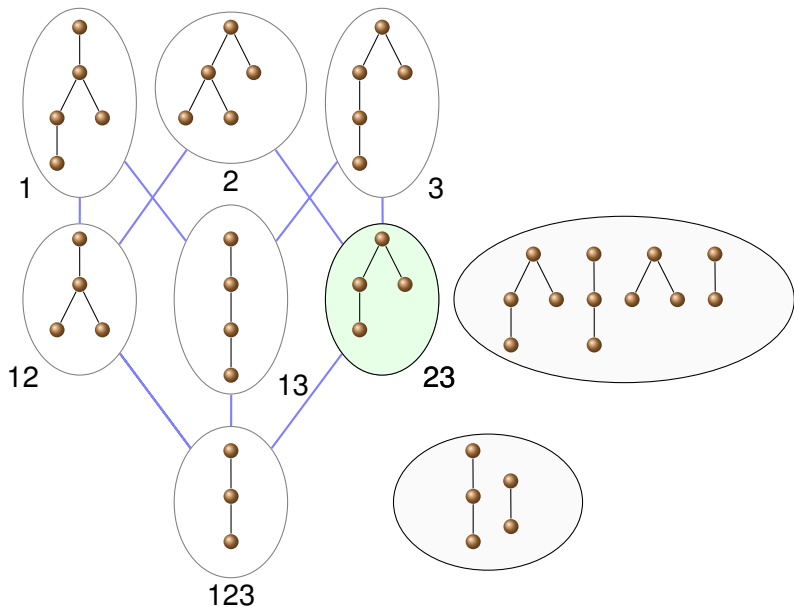
The following propositional formulas are logically equivalent:

- the conjunction of all the Horn formulas that are satisfied by all the models m_t for $t \in \mathcal{D}$
- the conjunction of \mathcal{R}_0 and all the propositional translations of the formulas in $\mathcal{R}'_{\mathcal{D}}$

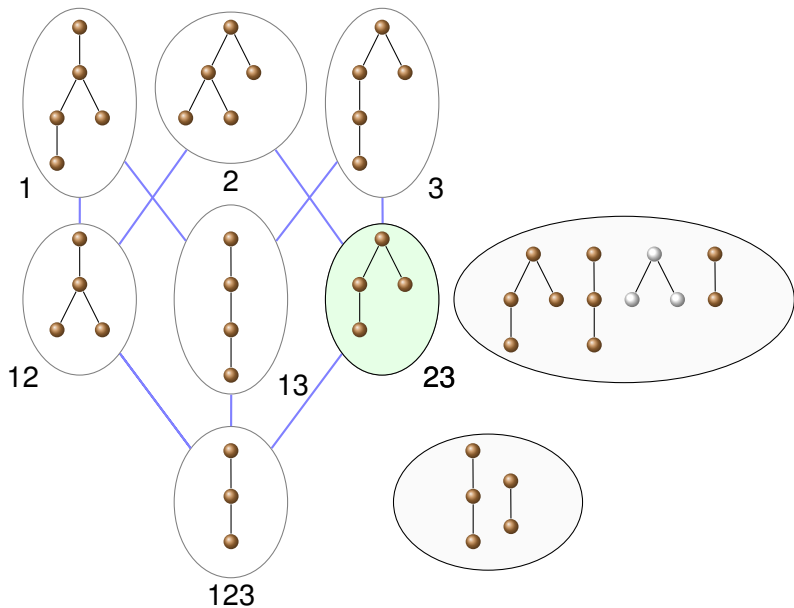
$$\mathcal{R}'_{\mathcal{D}} = \bigcup_{\mathcal{C}} \{A \rightarrow t \mid \Gamma_{\mathcal{D}}(A) = \mathcal{C}, t \in \mathcal{C}\}$$

- the conjunction of \mathcal{R}_0 and all the propositional translations of the formulas in a subset of $\mathcal{R}'_{\mathcal{D}}$ obtained transversing the hypergraph of differences between the nodes of the lattice.

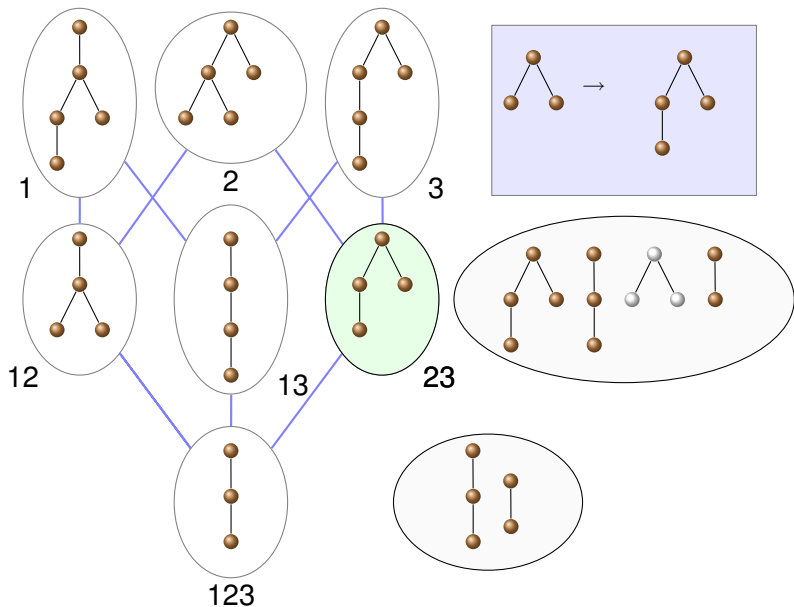
Association Rule Computation Example



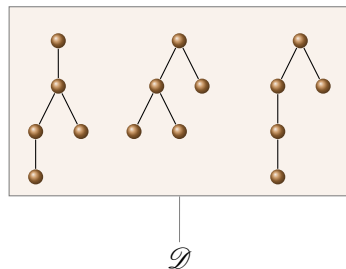
Association Rule Computation Example



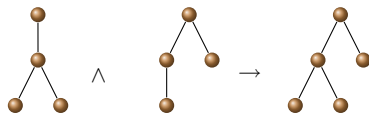
Association Rule Computation Example



Implicit rules



Implicit Rule

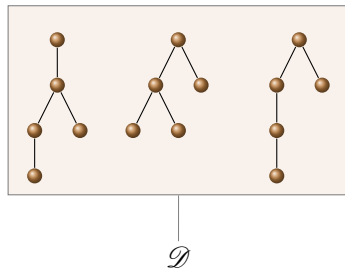


Given three trees t_1 , t_2 , t_3 , we say that $t_1 \wedge t_2 \rightarrow t_3$, is an *implicit Horn rule* (abbreviated, an *implicit rule*) if for every tree t it holds

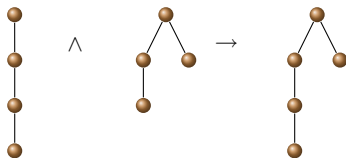
$$t_1 \preceq t \wedge t_2 \preceq t \leftrightarrow t_3 \preceq t.$$

t_1 and t_2 have implicit rules if $t_1 \wedge t_2 \rightarrow t$ is an implicit rule for some t .

Implicit rules



Implicit Rule

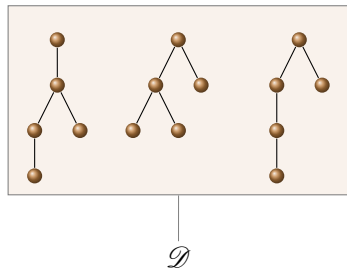


Given three trees t_1 , t_2 , t_3 , we say that $t_1 \wedge t_2 \rightarrow t_3$, is an *implicit Horn rule* (abbreviately, an *implicit rule*) if for every tree t it holds

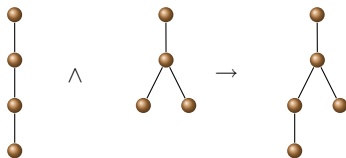
$$t_1 \preceq t \wedge t_2 \preceq t \leftrightarrow t_3 \preceq t.$$

t_1 and t_2 have implicit rules if $t_1 \wedge t_2 \rightarrow t$ is an implicit rule for some t .

Implicit rules



NOT Implicit Rule

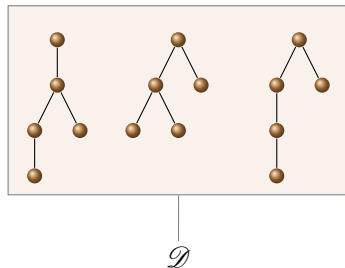


Given three trees t_1 , t_2 , t_3 , we say that $t_1 \wedge t_2 \rightarrow t_3$, is an *implicit Horn rule* (abbreviated, an *implicit rule*) if for every tree t it holds

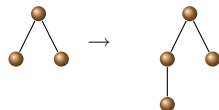
$$t_1 \preceq t \wedge t_2 \preceq t \leftrightarrow t_3 \preceq t.$$

t_1 and t_2 have implicit rules if $t_1 \wedge t_2 \rightarrow t$ is an implicit rule for some t .

Implicit rules



NOT Implicit Rule



Given three trees t_1 , t_2 , t_3 , we say that $t_1 \wedge t_2 \rightarrow t_3$, is an *implicit Horn rule* (abbreviately, an *implicit rule*) if for every tree t it holds

$$t_1 \preceq t \wedge t_2 \preceq t \leftrightarrow t_3 \preceq t.$$

t_1 and t_2 have implicit rules if $t_1 \wedge t_2 \rightarrow t$ is an implicit rule for some t .

Implicit Rules

Theorem

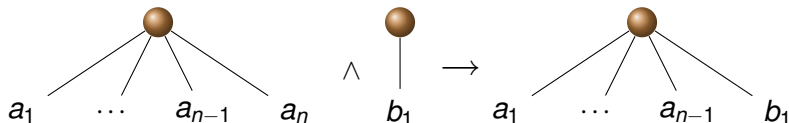
All trees a, b such that $a \preceq b$ have implicit rules.

Theorem

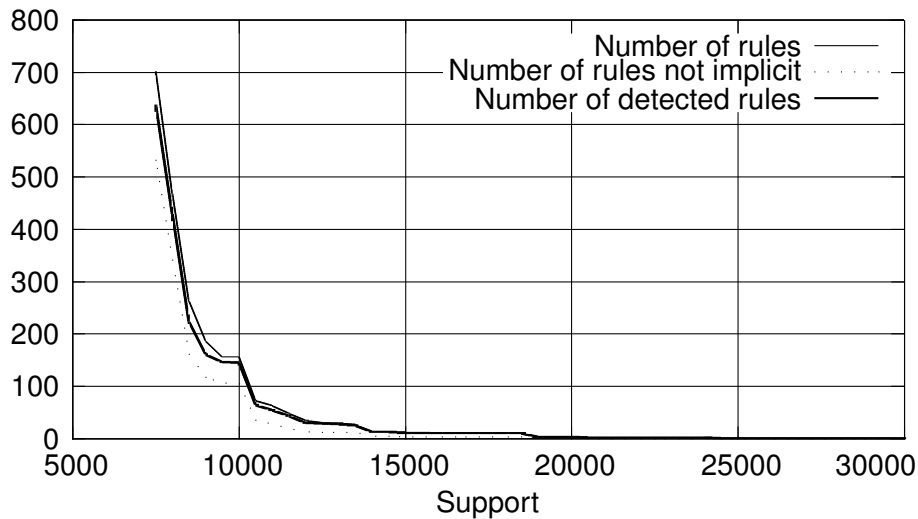
Suppose that b has only one component. Then they have implicit rules if and only if a has a maximum component which is a subtree of the component of b .

- for all $i < n$

$$a_i \preceq a_n \preceq b_1$$



Experimental Validation: CSLOGS



Conclusions

- A way of extracting high-confidence association rules from datasets consisting of unlabeled trees
 - antecedents are obtained through a computation akin to a hypergraph transversal
 - consequents follow from an application of the closure operators
- Detection of some cases of **implicit rules**: rules that always hold, independently of the dataset