



# Fast Theta-Subsumption with Constraint Satisfaction Algorithms

JÉRÔME MALOBERTI

MICHÈLE SEBAG

*Laboratoire de Recherche en Informatique, CNRS UMR 8623, Bât 490, Université Paris-Sud, F-91405 Orsay*

jerome.maloberti@lri.fr

michele.sebag@lri.fr

**Editor:** Peter Flach

**Abstract.** Relational learning and Inductive Logic Programming (ILP) commonly use as covering test the  $\theta$ -subsumption test defined by Plotkin. Based on a reformulation of  $\theta$ -subsumption as a binary constraint satisfaction problem, this paper describes a novel  $\theta$ -subsumption algorithm named *Django*,<sup>1</sup> which combines well-known CSP procedures and  $\theta$ -subsumption-specific data structures. *Django* is validated using the stochastic complexity framework developed in CSPs, and imported in ILP by Giordana et Saïtta. Principled and extensive experiments within this framework show that *Django* improves on earlier  $\theta$ -subsumption algorithms by several orders of magnitude, and that different procedures are better at different regions of the stochastic complexity landscape. These experiments allow for building a control layer over *Django*, termed *Meta-Django*, which determines the best procedures to use depending on the order parameters of the  $\theta$ -subsumption problem instance. The performance gains and good scalability of *Django* and *Meta-Django* are finally demonstrated on a real-world ILP task (emulating the search for frequent clauses in the mutagenesis domain) though the smaller size of the problems results in smaller gain factors (ranging from 2.5 to 30).

**Keywords:** relational learning, constraint satisfaction, phase transition, meta-learning,  $k$ -locality

## 1. Introduction

Supervised learning relies heavily on the generality relation, stating whether a hypothesis covers a (training) example. As the assessment of a hypothesis depends on its coverage, the implementation of the generality relation, or covering test, must be efficient.

The complexity of the covering test is among the greatest challenges facing supervised learning from relational examples, or *Inductive Logic Programming* (ILP) (Muggleton & De Raedt, 1994; Nienhuys-Cheng & de Wolf, 1997). ILP systems commonly use as generality relation, a decidable restriction of logical implication named  $\theta$ -subsumption (Plotkin, 1970), with exponential worst-case complexity in the size of the candidate hypothesis. This drawback has motivated the use of learning biases, more precisely syntactic biases and/or search biases favoring short hypotheses (Nédellec et al., 1996; De Raedt & Bruynooghe, 1993; Quinlan, 1990; Muggleton, 1995; Blockeel & De Raedt, 1998; Srinivasan, 2001). In parallel, efficient implementations of  $\theta$ -subsumption have been proposed (Kietz & Lübbe, 1994; Scheffer, Herbrich, & Wyszotzki, 1997); and ILP systems based on an approximation of  $\theta$ -subsumption have been developed, see for instance (Sebag & Rouveirol, 2000; Kramer, Lavrac, & Flach, 2001).

This paper presents a new correct and complete algorithm implementing  $\theta$ -subsumption, termed *Django*. The originality of the approach is that it formalizes  $\theta$ -subsumption as a Constraint Satisfaction problem (Tsang, 1993), and thoroughly exploits efficient CSP algorithms in order to resolve  $\theta$ -subsumption. *Django* is evaluated along an original framework developed by Giordana and Saitta (2000), which transposes in the context of ILP the principled framework developed in the Constraint Satisfaction literature since the early 90's (Cheeseman, Kanefsky, & Taylor, 1991; Hogg, Huberman, & Williams, 1996).

In this framework, the focus is shifted from a worst-case complexity analysis to a statistical complexity analysis. This new complexity paradigm is based on the definition of order parameters. Typically, the order parameters associated to  $\theta$ -subsumption are the number of variables and predicates in the hypothesis, the number of literals and constants in the example. Computational complexity is thus viewed as a random variable, ranging over the empirical complexity of all problem instances with same order parameters.

The  $\theta$ -subsumption landscape associated to these order parameters has been investigated experimentally and intensively (Giordana & Saitta, 2000). In accordance with earlier findings in CSPs (Hogg, Huberman, & Williams, 1996), three complexity regions appear. In a wide region, corresponding to under-constrained, satisfiable CS problems (the hypothesis is “short”, or general, relatively to the example)  $\theta$ -subsumption almost always succeeds. In another wide region, corresponding to over-constrained, unsatisfiable CS problems (the hypothesis is “long”, or specific, relatively to the example),  $\theta$ -subsumption almost always fails. These two wide regions are respectively referred to as YES and NO regions; they are separated by a narrow region, termed phase transition (PT), where the probability for  $\theta$ -subsumption to succeed abruptly drops from almost 1 to almost 0.

The main goal of the paper is to provide a  $\theta$ -subsumption algorithm with good performance in the PT region, without sacrificing the performance in the YES and NO regions. This goal is motivated by two remarks. Firstly, after the CSP literature, the PT region concentrates the hardest problem instances on average, i.e. the most challenging problems; therefore the scalability of a novel  $\theta$ -subsumption algorithm must be investigated in this particular region. Secondly, experiments on artificial relational learning problems suggest that this PT region is particularly relevant to ILP (Botta et al., 2003), for it includes most of the sought hypotheses. Indeed, relevant hypotheses are unlikely to belong to the YES or NO regions since hypotheses in these regions tend to cover all or none of the examples.

Experimental results show that the *Django* algorithm improves on previous  $\theta$ -subsumption algorithms (Kietz & Lübbe, 1994; Scheffer, Herbrich, & Wysotzki, 1997) by one or several orders of magnitude in all three regions of the  $\theta$ -subsumption landscape. These results are confirmed by experimentations in the real-world mutagenesis domain (King, Srinivasan, & Sternberg, 1995).

Lastly, as could have been expected, there is nothing like a universally efficient CSP procedure: some procedures are more efficient than others in particular regions of the  $\theta$ -subsumption landscape. The respective efficiency of the procedures is demonstrated by their *competence map*, summarizing principled experimentations inspired by Giordana and Saitta (2000). Based on these competence maps, a control layer termed *Meta-Django* is defined; it achieves the built-in selection of the best *Django* procedures depending on the problem instance at hand (Section 7). This way, *Meta-Django* benefits from the most

efficient procedures implemented in *Django*, and does not suffer from their computational overhead in the regions where these procedures are not adequate. *Meta-Django* is validated on the well-known ILP benchmark, the mutagenesis domain (King, Srinivasan, & Sternberg, 1995).

The paper is organized as follows. The next section briefly introduces  $\theta$ -subsumption and reviews existing  $\theta$ -subsumption algorithms (Kietz & Lübbe, 1994; Scheffer, Herbrich, & Wysotzki, 1997). Section 3 presents the Constraint Satisfaction framework and the main algorithms used to solve CSPs. Section 4 presents the *Django* algorithm, first described in (Maloberti & Sebag, 2001). *Django* first transforms the  $\theta$ -subsumption problem referred to as *primal problem*, into a binary CSP referred to as *dual problem*; the dual problem is then handled using diverse procedures developed for binary CSPs and  $\theta$ -subsumption-specific data structures. The experimental setting inspired by Botta et al. (2003), is described and discussed in Section 5, together with the experimental goals. Empirical results are reported in Section 6 and discussed with respect to the state of the art. Section 7 is devoted to *Meta-Django*. Some related works are discussed in Section 8 and the paper ends with some perspectives for further research.

## 2. Theta-subsumption, definition and algorithms

After defining  $\theta$ -subsumption, this section reviews the main algorithms proposed in the ILP literature for optimizing the  $\theta$ -subsumption test.

### 2.1. Definition and standard subsumption test

In the rest of the paper,  $\mathcal{C}$  and  $Ex$  denote Horn clauses in Datalog language (including no function symbols other than constants).

*Definition 1* (Theta-subsumption). Clause  $\mathcal{C}$   $\theta$ -subsumes  $Ex$  iff there exists a substitution  $\theta = \{X_i/V_{i_j}\}$ , where  $X_i$  ranges over the variables in  $\mathcal{C}$  and  $V_{i_j}$  ranges over the variables and constants in  $Ex$ , such that  $\mathcal{C}\theta \subseteq Ex$ .

*Example 1.*

$$\begin{aligned} \mathcal{C}: tc(X_0) \leftarrow atm(X_0, X_1), atm(X_0, X_2), atm(X_0, X_3), bond(X_1, X_2), \\ bond(X_1, X_3) \\ Ex: tc(m) \leftarrow atm(m, m_1), atm(m, m_2), bond(m_1, m_2) \end{aligned}$$

$\mathcal{C}$   $\theta$ -subsumes  $Ex$  according to substitution  $\theta = \{X_0/m, X_1/m_1, X_2/m_2, X_3/m_2\}$ .

The natural search space for  $\theta$ -subsumption therefore is the set of substitutions, mapping the set of variables in  $\mathcal{C}$  noted  $\text{arg}(\mathcal{C})$ , onto the set of variables and constants in  $Ex$ , noted  $\text{arg}(Ex)$ .

A more computationally efficient approach is to consider the set of mappings from the set of literals in  $\mathcal{C}$  onto the set of literals in  $Ex$  (built on the same predicate symbol), referred to as literal mappings. In the above toy example, each one of the three literals  $atm(X_0, X_1)$ ,  $atm(X_0, X_2)$ ,  $atm(X_0, X_3)$  is mapped onto either  $atm(m, m_1)$  or  $atm(m, m_2)$ , and both literals  $bond(X_1, X_2)$  and  $bond(X_1, X_3)$  are mapped onto  $bond(m_1, m_2)$ .

Through a literal mapping, each variable in  $\mathcal{C}$  is associated to a set of variables and constants in  $Ex$ . The literal mapping is said to be consistent if it associates each variable in  $\mathcal{C}$  to a single variable or constant in  $Ex$ .

*Definition 2* (Consistency of a literal mapping). A mapping  $m = \{p_i/p'_i\}$  from the literals in  $\mathcal{C}$  onto the literals in  $Ex$  is said to be *consistent* if there exists a substitution  $\theta$  such that, for each literal  $p_i$  in  $\mathcal{C}$ ,  $p'_i = p_i\theta$ .

The standard algorithm for  $\theta$ -subsumption is based on Prolog SLD resolution (Robinson, 1965). Within SLD resolution, all mappings from the literals in  $\mathcal{C}$  onto the literals in  $Ex$  built on the same predicate symbol are constructed and checked for consistency. The order of construction depends on the order of the literals in  $\mathcal{C}$  and  $Ex$ ; as known by all Prolog programmers, this order has a significant impact on the SLD efficiency.

## 2.2. Determinate matching

Kietz and Lubbe have proposed a first improvement over SLD resolution, known as determinate matching (Kietz & Lübbe, 1994). In determinate matching, the literals in  $\mathcal{C} = p_1 \dots p_K$  are reordered in such a way that, if possible, for each  $p_i$  in  $\mathcal{C}$  there exists a single candidate literal  $p'_i$  in  $Ex$  which is consistent with the previous assignments, i.e. such that  $\{p_1/p'_1, \dots, p_i/p'_i\}$  is consistent;  $p_i$  is then termed determinate literal. After all determinate literals in  $\mathcal{C}$  have been mapped onto literals in  $Ex$ , and if necessary, the search resumes using SLD resolution.

Clause  $\mathcal{C}$  is said to be determinate *with respect to*  $Ex$  if it only contains determinate literals wrt  $Ex$ .

*Example 2.*

$$\mathcal{C}: tc(X_0) \leftarrow p(X_0, X_1), p(X_1, X_2), p(X_2, X_3)$$

$$Ex: tc(a_0) \leftarrow p(a_0, a_1), p(a_1, a_2), p(a_2, a_3)$$

It is clear that  $\mathcal{C}$   $\theta$ -subsumes  $Ex$  according to substitution  $\theta = \{X_0/a_0, X_1/a_1, X_2/a_2, X_3/a_3\}$ : determinate matching succeeds for clauses  $\mathcal{C}$  and  $Ex$ .

The determinacy property is desirable as it guarantees the polynomial complexity of the  $\theta$ -subsumption test. The limitation is that this property depends on the example considered.

*Example 3.*

$$\mathcal{C}: tc(X_0) \leftarrow p(X_0, X_1), p(X_1, X_2), p(X_2, X_3)$$

$$Ex: tc(a_0) \leftarrow p(a_0, a_1), p(a_1, a_2), p(a_2, a_3), p(a_0, a_4)$$

Adding literal  $p(a_0, a_4)$  to the body of  $Ex$  causes the determinacy test to fail as  $p(X_0, X_1)$  now admits  $p(a_0, a_1)$  and  $p(a_0, a_4)$  as candidate literals.

### 2.3. Graph context

In order to circumvent this problem, Scheffer, Herbrich, and Wysotzki (1997) propose a method similar to the *vertex invariants* (Read & Corneil, 1977) used in graph isomorphism to prune candidate literals and thus, extend the scope of determinate matching. The idea is the following: when graph isomorphisms are required to preserve some attributes for the vertices or nodes in the graph, such as the degree and label of a node, or the labels of the adjacent vertices, these attributes can be used to prune and partition the candidate set into equivalence classes.

Scheffer, Herbrich, and Wysotzki (1997) build two graphs respectively associated to  $\mathcal{C}$  and  $Ex$ , where the vertices correspond to literals, and the edges connect all pairs of literals that share a variable in  $\mathcal{C}$  (resp. a variable or a constant in  $Ex$ ).

To each literal  $p$  in  $\mathcal{C}$  (resp. in  $Ex$ ) is associated its neighborhood or graph context, recursively defined as follows. The 1-neighbors of  $p$  are all literals in  $\mathcal{C}$  (resp. in  $Ex$ ) adjacent to  $p$ . The  $i$ -th neighbors are recursively constructed, as 1-neighbors of  $(i - 1)$ -neighbors.

One easily shows that a necessary condition for a consistent literal mapping  $m$  is that it maps every literal  $p$  in  $\mathcal{C}$  onto literal  $p'$  in  $Ex$ , such that all predicate symbols occurring in  $p'$   $i$ -neighbors also are  $i$ -neighbors of  $p$ .

*Example 4.* Let  $\mathcal{C}$  and  $Ex$  be defined as in Example 3, and let us consider their associated *graph context of depth 1*, noted 1-GC.

The 1-GC of  $\mathcal{C}$  gives for each literal and each variable  $X$  in this literal, the set of pairs (predicate  $p$ , argument  $k$ ) such that  $X$  also appears in a literal built on predicate  $p$  as  $k$ -th argument:

$$\begin{aligned} tc(X_0): & \{(1/p.1)\} \\ p(X_0, X_1): & \{(1/tc.1), \langle 2/p.1 \rangle\} \\ p(X_1, X_2): & \{(1/p.2), \langle 2/p.1 \rangle\} \\ p(X_2, X_3): & \{(1/p.2)\} \end{aligned}$$

The 1-GC of literal  $tc(X_0)$ , i.e.  $1/p.1$ , indicates that literal  $tc(X_0)$  is linked by its first variable to the first variable of a literal built on predicate symbol  $p$ .

Likewise, the 1-GC of  $Ex$  (Example 3) is given as:

$$\begin{aligned} tc(a_0): & \{(1/p.1)\} \\ p(a_0, a_1): & \{(1/tc.1), \langle 1/p.1 \rangle, \langle 2/p.1 \rangle\} \\ p(a_0, a_4): & \{(1/tc.1), \langle 1/p.1 \rangle\} \\ p(a_1, a_2): & \{(1/p.2), \langle 2/p.1 \rangle\} \\ p(a_2, a_3): & \{(1/p.2)\} \end{aligned}$$

This example shows that the use of graph contexts extends the scope of determinate matching. Literal  $p(X_0, X_1)$  now admits a single candidate,  $p(a_0, a_1)$ ; candidate  $p(a_0, a_4)$  is pruned as its 1-GC does not contain that of  $p(X_0, X_1)$ .

As graph contexts help to reduce the set of candidate literals, their use extends the set of  $\theta$ -subsumption problems that can be resolved in polynomial time. Nevertheless, the determinacy of the  $\theta$ -subsumption test still depends on the clause  $Ex$  at hand.

*Example 5.*

$$\begin{aligned} C: tc(X_0) &\leftarrow p(X_0, X_1), p(X_1, X_2), p(X_2, X_3) \\ Ex: tc(a_0) &\leftarrow p(a_0, a_1), p(a_0, a_2), p(a_1, a_2), p(a_2, a_3) \end{aligned}$$

In this example, literal  $p(a_0, a_2)$  cannot be pruned on the basis of its 1-GC. Graph contexts of depth 2 would be required to ensure determinate matching.

In many cases, the matching test can be made deterministic by developing the graph context up to some depth  $d$ ; but the construction of the graph context is exponential in  $d$ .

#### 2.4. Other approaches

Scheffer, Herbrich, and Wysotzki (1997) further define a *substitution graph*-based  $\theta$ -subsumption test. The nodes of the substitution graph are all pairs of literals  $(p/p')$  where  $p$  and  $p'$  respectively are literals in  $C$  and  $Ex$  built on the same predicate symbol. Two nodes  $(p/p')$  and  $(q/q')$  are linked by an edge iff mapping  $(p/p', q/q')$  is consistent (Scheffer, Herbrich, & Wysotzki, 1997).

In Scheffer, Herbrich, and Wysotzki (1997), the SLD search is replaced by a maximal clique search in the substitution graph. The worst-case complexity remains exponential, but the advantage is that the consistency check is performed only once.

Another approach presented by Kietz and Lubbe (1994) proceeds by decomposing the substitution graph into mutually independent components (*k-locality*). Whenever such a decomposition applies, it significantly reduces the complexity of the  $\theta$ -subsumption test.

Indeed, all above heuristics can be viewed as particular cases of Constraint Satisfaction algorithms. Since we aim at a principled exploitation of CSP algorithms for  $\theta$ -subsumption, the next section is devoted to a brief presentation of CS.

### 3. Constraint satisfaction problem

This section introduces Constraint Satisfaction problems (CSPs) together with the main resolution algorithms; the reader is referred to (Tsang, 1993) for a comprehensive presentation.

### 3.1. Overview of CSPs

*Definition 3* (Constraint Satisfaction Problem). A CSP involves:

- A set of  $n$  variables  $\mathcal{X} = \{X_1, \dots, X_n\}$ ; to each variable  $X_i$  is associated its variable domain  $\text{dom}(X_i)$ , including all possible values for  $X_i$ ;
- A set of  $m$  constraints, specifying the simultaneously admissible values of the variables.

Each constraint  $r$  is defined on some variables in  $\mathcal{X}$ :

- it can be thought of as a  $k$ -ary predicate  $r(X_{i_1}, \dots, X_{i_k})$ ;
- the constraint scope noted  $\text{arg}(r)$  is the set of variables  $\{X_{i_1}, \dots, X_{i_k}\}$ ;
- the set of admissible values for the variables in the constraint scope is termed the constraint relation, noted  $\text{rel}(r)$ . Constraint relation  $\text{rel}(r)$  can be thought of as a set of  $N$  literals  $\{r(a_{1,i_1}, \dots, a_{1,i_k}), \dots, r(a_{N,i_1}, \dots, a_{N,i_k})\}$ , with  $a_{l,i_j}$  a value in  $\text{dom}(X_{i_j})$  for  $l = 1..N$  and  $j = 1..k$ .

Although a constraint relation can be infinite in the general case (e.g. a numerical constraint on real-valued variables), only finite constraint relations are considered in the rest of the paper.

A CSP is *satisfiable* iff it admits a solution, assigning to each variable  $X_i$  a value  $a_i$  in  $\text{dom}(X_i)$  such that all constraints are satisfied. Such a solution can be viewed as a mapping  $\theta = \{X_i/a_i\}$  such that for each constraint  $r(X_{i_1}, \dots, X_{i_k})$ ,  $r\theta = r(a_{i_1}, \dots, a_{i_k})$  belongs to  $\text{rel}(r)$ . In other words, the CSP defined by constraints  $r_1, \dots, r_m$  is *satisfiable* iff conjunction  $r_1, \dots, r_m$   $\theta$ -subsumes the conjunction of all literals in  $\text{rel}(r_1), \dots, \text{rel}(r_m)$ .

*Example 6.* Let us consider the CSP defined as:

- $\mathcal{X} = \{X_0, X_1, X_2, X_3\}$ , where each variable has domain  $\{m, m_1, m_2\}$ ;
- six constraints are defined as follows:  $tc$  is defined on  $X_0$ , with  $\text{rel}(tc) = \{m\}$ ;  $atm_1$  (respectively  $atm_2$  and  $atm_3$ ) is defined on  $(X_0, X_1)$  (resp.  $(X_0, X_2)$  and  $(X_0, X_3)$ ), with domain  $\{\langle m, m_1 \rangle, \langle m, m_2 \rangle\}$ ;  $bond_1$  and  $bond_2$  are respectively defined on  $(X_1, X_2)$  and  $(X_1, X_3)$  with domain  $\{\langle m_1, m_2 \rangle\}$ .

The above CSP is equivalent to the  $\theta$ -subsumption test in Example 1, and its solution is given by  $\theta = \{X_0/m, X_1/m_1, X_2/m_2, X_3/m_2\}$ .

Two CSPs are *equivalent* iff they are defined on the same variables and admit the same solutions. As any CSP can be embedded into a binary CSP, i.e. involving only binary constraints, most CS algorithms consider binary CSPs. In the following, we assume with no loss of generality that at most one constraint is set on any variable pair.

**Complexity.** The CSP complexity is exponential in the number  $n$  of variables, and linear in the number  $m$  of constraints. More precisely, assuming that all variables have the same domain size  $L$ , the CSP resolution complexity is  $\mathcal{O}(L^n \times m)$ . A first way of decreasing the complexity thus is by decomposing the CSP into independent or scarcely related subproblems—hierarchizing the set of variables (Gyssens, Jeavons, & Cohen, 1994; Gottlob, Leone, & Scarcello, 2000) or the set of constraints (Dechter & Pearl, 1989)—in the same spirit as  $k$ -locality (Kietz & Lübke, 1994).

CS algorithms basically combine two kinds of procedures: reduction procedures and search procedures. Reduction procedures reduce the variable domains, thereby transforming a CSP into an equivalent one of lower complexity. Search procedures are concerned with the search strategy, and how to best explore the search space.

### 3.2. Reduction procedures

The simplest reduction algorithm, termed *node consistency*, is concerned with unary constraints. Let  $r(X)$  be a unary constraint defined on variable  $X$ , and let  $\text{rel}(r)$  denote its set of admissible values. With no loss of generality, the domain of  $X$  can be restricted to those values appearing in  $\text{rel}(r)$ . Conversely,  $r$  is trivially satisfied if  $\text{dom}(X)$  is included in  $\text{rel}(r)$ .

Node consistency is extended based on the *2-consistency value* definition.

**Definition 4 (2-Consistency of a value).** Let  $X$  be a constrained variable and let  $a$  denote a value in  $\text{dom}(X)$ . Value  $a$  is *2-consistent* (or consistent if there is no ambiguity from the context) if, for every variable  $Y$  such that there exists a constraint  $r(X, Y)$ , there exists some value  $b$  in  $\text{dom}(Y)$  such that  $r(a, b)$  belongs to  $\text{rel}(r)$ ;  $b$  is termed *support* of  $a$  wrt  $r$ .

*Example 7.*

$$\begin{aligned} p(X, Y) \text{ rel}(p) &= \{\langle a, b \rangle, \langle c, d \rangle\} \\ q(X, Z) \text{ rel}(q) &= \{\langle a, e \rangle, \langle f, g \rangle\} \end{aligned}$$

Value  $a$  is consistent for  $X$  (supported by  $b$  for  $Y$  and  $e$  for  $Z$ ); value  $c$  is not consistent.

Clearly, if  $a$  has no support value with respect to some constraint  $r$  involving  $X$ , then  $X$  cannot be mapped onto  $a$ . Therefore, all non 2-consistent values can soundly be removed from  $\text{dom}(X)$ .

The above pruning is achieved by *arc consistency*; this algorithm considers each value in the domain of every variable, and removes all values that are not 2-consistent.

Node and arc consistency can be generalized using  $k$ -consistency:

**Definition 5 ( $k$ -Consistency of value tuples).** A partial assignment  $\theta$  assigning a value to a subset  $\mathcal{X}_\theta$  of the constrained variables  $\mathcal{X}$ , is consistent if it violates no constraint, i.e. it satisfies all constraints defined on (a subset of)  $\mathcal{X}_\theta$ .

A CSP is  $k$ -consistent iff for any consistent assignment  $\theta$  over  $k - 1$  variables, for any variable  $X$  not in  $\mathcal{X}_\theta$  there exists a value  $a_X$  such that  $\theta' = \theta \cup \{X/a_X\}$  is consistent.



A CSP is strongly  $k$ -consistent if it is  $j$ -consistent for all  $j$  less than  $k$ .

If a consistent partial solution over some  $k - 1$  variables of a CSP cannot be extended to another variable, then this partial solution can soundly be removed (with all its specializations) from the assignment search space. This pruning algorithm is termed  $k$ -consistency.

Particular cases of  $k$ -consistency are given below:

|               |   |                  |
|---------------|---|------------------|
| 1-consistency | = | node consistency |
| 2-consistency | = | arc consistency  |
| 3-consistency | = | path consistency |

For complexity and implementation reasons, arc consistency is generally preferred to  $k$ -consistency with  $k > 2$ . On one hand, checking  $k$ -consistency is exponentially complex in  $k$ . On the other hand,  $k$ -consistency aims at pruning  $(k - 1)$ -tuples of values, which requires specific data structures for  $k > 2$ ; managing these structures might partially or totally offset the computational gains for some problem instances.

**Complexity.** The best known worst-case complexity of arc consistency algorithms is  $\mathcal{O}(mL^2)$ , with  $m$  being the number of constraints and  $L$  the variable domain size.

Note that in the particular case of tree-structured CSPs, arc consistency is a necessary and sufficient condition for the CSP to be satisfiable. More generally, a CSP with an acyclic hyper graph can be solved with polynomial complexity (Dechter & Pearl, 1989; Gottlob & Leitsch, 1985).

### 3.3. Search procedures

CSP algorithms construct an assignment solution  $\{X_i/a_i\}$  through a depth first exploration of the assignment space, referred to as substitution tree. Each node corresponds to a variable  $X_i$ ; the edge attached to this node is the candidate value  $a_i$  tentatively assigned to  $X_i$ . On each assignment, consistency is checked; on failure, another candidate value for the current node is considered; if no other value is available, the search backtracks.

Several approaches have been proposed in order to improve: (i) the backtracking procedure (*look-back* algorithms); (ii) the choice of the next variable and candidate value to consider (*look-ahead* algorithms).

Look-back algorithms aim at preventing the repeated exploration of a same substitution subtree on backtracking (*thrashing*). For instance, Conflict directed BackJumping (CBJ) (Prosser, 1993) registers all variable conflicts occurred during the exploration, which allows for backtracking directly to the appropriate tree level. On the other hand, it may happen that the overhead due to maintaining the conflict registers offsets the look-back advantages for some particular CSP instances.

Look-ahead algorithms aim at minimizing the number of assignments considered. The best known look-ahead procedure is constraint propagation, pruning all candidate values

which are inconsistent with the current assignment, also known as arc-consistency. This way, inconsistencies are detected earlier and less nodes are visited; on the other hand, the assignment operation is computationally heavier as it involves the constraint propagation check.

Forward checking (FC) employs a limited propagation, only checking the variables connected to the assigned variable (partial arc-consistency). Maintaining Arc Consistency (MAC) extends FC by checking arc-consistency on adjacent variables if a value has been removed from their domain. Again, the overhead due to constraint propagation might offset its advantages on medium-size weakly constrained CSPs. Currently, the CSP resolution algorithms considered as most efficient combine FC and CBJ procedures.

The variable order can be optimized too. The ordering is achieved either statically (once for all), or dynamically (the yet unassigned variables are reordered on each assignment). Dynamic variable ordering is generally more efficient than static variable ordering. Variable ordering is most often based on the First Fail Principle (Bessi re, 1996), favoring the variable with the smallest domain. This way, failures will occur sooner rather than later. Last, the candidate values can also be ordered; the value with less conflicts (greatest support) with the other variable domains is commonly preferred.

#### 4. CS resolution for theta-subsumption

In this section, we first discuss how to best transform a  $\theta$ -subsumption problem into a binary CSP, in order to delegate the  $\theta$ -subsumption resolution to efficient CSP procedures. Based on this transformation, the *Django* algorithm is presented; it combines standard CSP procedures with  $\theta$ -subsumption specific data structures and procedures.

##### 4.1. Standard binary representation

The equivalence of  $\theta$ -subsumption and CS problems follows from their definitions in Sections 2 and 3. However,  $\theta$ -subsumption generally involves  $n$ -ary predicates, while most CSP procedures consider binary constraints. A specific change of representation is thus required to allow  $\theta$ -subsumption to exploit standard CS algorithms.

The usual way of binarizing an  $n$ -ary CSP replaces each  $n$ -ary constraint  $r(X_{i_1}, \dots, X_{i_n})$  using an additional variable  $Y_r$  and  $n$  additional binary constraints  $r_1, \dots, r_n$  linking  $X_{i_1}, \dots, X_{i_n}$  to  $Y_r$ . Constraint  $r_j(X_{i_j}, Y_r)$  expresses that  $X_{i_j}$  is the  $j$ -th argument of the initial constraint  $r$ .

Let us illustrate the binarization procedure on an example.

*Example 8.*

$$\begin{aligned} \mathcal{C}: tc(X_0) &\leftarrow p(X_0, X_1), q(X_0, X_2, X_3) \\ \text{Ex: } tc(a_0) &\leftarrow p(a_0, a_1), p(a_1, a_2), q(a_0, a_2, a_3), q(a_0, a_1, a_3) \end{aligned}$$

The  $\theta$ -subsumption problem is first rewritten into a CSP problem:

- The set of variables  $\mathcal{X}$  is  $\{X_0, X_1, X_2, X_3\}$ ,
- The domain variables are defined and reduced using arc consistency (Section 3.2). The domain of  $X_0$  is reduced to  $a_0$  by consistency with the unary constraint  $tc$ .

$$\text{dom}(X_1) = \text{dom}(X_2) = \text{dom}(X_3) = \{a_0, a_1, a_2, a_3\};$$

- Two constraints  $r_1$  and  $r_2$  respectively account for literals  $p(X_0, X_1)$  and  $q(X_0, X_2, X_3)$ :

$$r_1(X_0, X_1) \text{ is associated to } \text{rel}(r_1) = \{\langle a_0, a_1 \rangle, \langle a_1, a_2 \rangle\}$$

$$r_2(X_0, X_2, X_3) \text{ is associated to } \text{rel}(r_2) = \{\langle a_0, a_2, a_3 \rangle, \langle a_0, a_1, a_3 \rangle\}.$$

This CSP is represented as a hyper graph (figure 1).

In order to binarize the above CSP, an additional variable is associated to each literal in  $\mathcal{C}$ . (Though literal  $p(X_0, X_1)$  is binary, we still introduce variable  $Y_p$  for the sake of pedagogy). The resulting binary CSP (figure 2) is given as:

- Two additional variables  $Y_p$  and  $Y_q$  are considered besides variables  $X_0$  to  $X_3$ ,
- The variable domains associated to  $Y_p$  and  $Y_q$  are respectively defined as

$$\text{dom}(Y_p) = \{p(a_0, a_1), p(a_1, a_2)\}, \text{ and } \text{dom}(Y_q) = \{q(a_0, a_2, a_3), q(a_0, a_1, a_3)\},$$

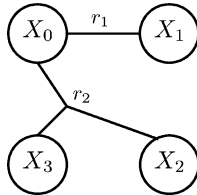


Figure 1. Hyper graph representation of Example 8.

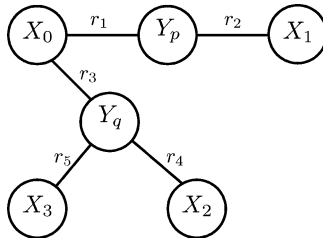


Figure 2. Binary graph representation of Example 8.

– The set of constraints now includes 5 binary constraints:

$$\begin{aligned}
 r_1(X_0, Y_p) \text{ is associated to } \text{rel}(r_1) &= \{\langle a_0, p(a_0, a_1) \rangle, \langle a_1, p(a_1, a_2) \rangle\}; \\
 r_2(X_1, Y_p) \text{ is associated to } \text{rel}(r_2) &= \{\langle a_1, p(a_0, a_1) \rangle, \langle a_2, p(a_1, a_2) \rangle\}; \\
 r_3(X_0, Y_q) \text{ is associated to } \text{rel}(r_3) &= \{\langle a_0, q(a_0, a_2, a_3) \rangle, \\
 &\quad \langle a_0, q(a_0, a_1, a_3) \rangle\}; \\
 r_4(X_2, Y_q) \text{ is associated to } \text{rel}(r_4) &= \{\langle a_2, q(a_0, a_2, a_3) \rangle, \\
 &\quad \langle a_1, q(a_0, a_1, a_3) \rangle\}; \\
 r_5(X_3, Y_q) \text{ is associated to } \text{rel}(r_5) &= \{\langle a_3, q(a_0, a_2, a_3) \rangle, \\
 &\quad \langle a_3, q(a_0, a_1, a_3) \rangle\};
 \end{aligned}$$

This standard binarization creates two types of constrained variables: the actual constrained variables ( $X_0$  to  $X_3$  in the above example) and the additional variables standing for the literals in  $\mathcal{C}$  ( $Y_p$  and  $Y_q$  in the above example). As the  $\theta$ -subsumption algorithm must handle both types of variables in different ways, standard binarization is not compatible with a straightforward use of CSP algorithms.

#### 4.2. Dual constrained representation

For this reason, we devised another transformation of a  $\theta$ -subsumption problem into a binary CSP, making a clear distinction between literals and variables in  $\mathcal{C}$ .

Concretely, each  $\theta$ -subsumption problem  $(\mathcal{C}, Ex)$  is transformed into a binary CSP termed *dual CSP*, where:

- (i) Each literal  $p$  in  $\mathcal{C}$  gives rise to a constrained variable  $Y_p$  termed *dual variable*, as opposed to the variables in  $\mathcal{C}$  referred to as *primal variables*.
- (ii) The domain of a dual variable  $Y_p$  is the set of all literals in  $Ex$  built on predicate symbol  $p$ .
- (iii) The constraints in the dual CSP, termed *dual constraints*, implements the consistency requirement (Definition 1): literals in  $\mathcal{C}$  must be mapped onto literals in  $Ex$  in such a way that it induces the mapping of (primal) variables in  $\mathcal{C}$  onto a single constant or variable in  $Ex$ . Formally,

*Definition 6* (Dual representation). The dual representation of a  $\theta$ -subsumption problem  $(\mathcal{C}, Ex)$  is defined as follows.

- To each literal in  $\mathcal{C}$  is associated a dual constrained variable. Specifically, to the  $i$ -th literal built on predicate symbol  $p$  in  $\mathcal{C}$  noted  $p_i$  is associated the dual constrained variable  $Y_{p,i}$ . Subscript  $i$  will be omitted for readability when there is a single literal built on the predicate symbol  $p$ .
- The domain of a dual constrained variable  $Y_{p,i}$  involves all literals in  $Ex$  built on the predicate symbol  $p$ . All dual variables related to the same predicate symbol  $p$  thus have same domain.

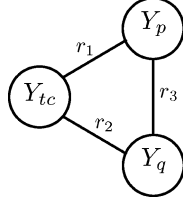


Figure 3. Dual representation of Example 8.

- For each pair of dual constrained variables  $(Y_{p.i}, Y_{q.j})$  such that literals  $p_i$  and  $q_j$  in  $\mathcal{C}$  share at least one variable in  $\mathcal{C}$ , a binary dual constraint  $r(Y_{p.i}, Y_{q.j})$  is defined. This dual constraint expresses the fact that the assignment of variables  $Y_{p.i}$  and  $Y_{q.j}$  must be consistent in the sense of Definition 1. (The particular case where a primal variable occurs more than once in a literal  $p_i$ , is also handled through a binary dual constraint  $r(Y_{p.i}, Y_{p.i})$ ).

Every dual constraint  $r(Y_{p.i}, Y_{q.j})$ , similar to the substitution graph in Scheffer, Herbrich, and Wysotzki (1997), is satisfied by all pairs  $(p', q')$  of literals in  $Ex$ , such that mapping  $\{p_i/p', q_j/q'\}$  is consistent according to Definition 1. In other words, any admissible assignment of (dual) variables  $Y_{p.i}$  and  $Y_{q.j}$  corresponds to a consistent assignment of the primal variables in  $\mathcal{C}$ .

*Example 9.* The  $\theta$ -subsumption problem presented in Example 8 is now binarized as follows (figure 3):

- The set of variables is  $\mathcal{X} = \{Y_{tc}, Y_p, Y_q\}$ ,
- The variable domains are:  $\text{dom}(Y_{tc}) = \{tc(a_0)\}$ ,  $\text{dom}(Y_p) = \{p(a_0, a_1), p(a_1, a_2)\}$ ,  $\text{dom}(Y_q) = \{q(a_0, a_2, a_3), q(a_0, a_1, a_3)\}$ ;
- The set of constraints is now:

$$\begin{aligned}
 r_1(Y_{tc}, Y_p) \text{ is associated to } \text{rel}(r_1) &= \{\langle tc(a_0), p(a_0, a_1) \rangle\}; \\
 r_2(Y_{tc}, Y_q) \text{ is associated to } \text{rel}(r_2) &= \{\langle tc(a_0), q(a_0, a_2, a_3) \rangle, \\
 &\quad \langle tc(a_0), q(a_0, a_1, a_3) \rangle\}; \\
 r_3(Y_p, Y_q) \text{ is associated to } \text{rel}(r_3) &= \{\langle p(a_0, a_1), q(a_0, a_2, a_3) \rangle, \\
 &\quad \langle p(a_0, a_1), q(a_0, a_1, a_3) \rangle\};
 \end{aligned}$$

### 4.3. Discussion

As mentioned above, dual constraints encode the same information as substitution graphs (Section 2.4, Scheffer, Herbrich, & Wysotzki (1997)). In both cases, the point is to check whether a consistent substitution can be derived from partial literal mappings.

A first difference is that dual constraints are more parsimonious, in the sense that they only consider literals  $p$  and  $q$  in  $\mathcal{C}$  that do share variables. In opposition, substitution graphs systematically consider all pairs of mappings  $\{p/p'\}$  and  $\{q/q'\}$ , where  $p$  and  $q$  (respectively  $p'$  and  $q'$ ) belong to  $\mathcal{C}$  (resp  $Ex$ ).

A second difference is that substitution graphs materially store the list of pairs ( $\{p/p'\}$ ,  $\{q/q'\}$ ) such that assignment  $\{p/p', q/q'\}$  is consistent. Although this list can be pruned by using graph contexts, it remains that the construction and exploitation of graph contexts is computationally expensive in some regions of the  $\theta$ -subsumption landscape (we shall return to this point in Section 6.3). In opposition, we do not store the domain of relation  $r(Y_p, Y_q)$ ; rather, this constraint is implemented as a function, checking for each pair  $p'$  and  $q'$  respectively in the domain of  $Y_p$  and  $Y_q$  whether assignment  $\{p/p', q/q'\}$  is consistent.

Dual constraints are exploited through arc consistency, i.e. with complexity  $\mathcal{O}(mL^2)$ , where  $m$  is the number of constraints and  $L$  the variable domain size. Arc consistency enables the pruning of the dual variables domains, thereby restricting the search space of consistent mappings.

To further restrict the search space, constraint propagation is activated when a dual variable domain is limited to a single value (literal in  $Ex$ ). Note that this is always the case for the head of  $\mathcal{C}$ , which is thus directly handled and propagated over all determinate literals (e.g. dual variable  $Y_p$  can only be assigned to  $p(a_0, a_1)$  in Example 9, due to dual constraint  $r_1(Y_{tc}, Y_p)$  and  $\text{dom}(Y_{tc}) = \{tc(a_0)\}$ ).

**Proposition.** *Let AC-Pol denote the set of  $\theta$ -subsumption problems which can be resolved with polynomial complexity using the above binarization combined with arc consistency, and let GC-Pol similarly denote the set of  $\theta$ -subsumption problems which can be resolved with polynomial complexity using graph contexts (Scheffer, Herbrich, & Wysotzki, 1997). Then, GC-Pol is a subset of AC-Pol.*

**Proof:** By induction, we show that  $i$ -GC test is equivalent to arc consistency with  $(i - 1)$  steps of propagation of the eliminated values.

The 1-GC test of a literal  $p$  in  $\mathcal{C}$  wrt a literal  $p'$  in  $Ex$ , built on the same predicate symbol as  $p$ , checks whether  $p'$  shares the same links as  $p$  (variables in  $\mathcal{C}$ , variables or constants in  $Ex$ ) with the other literals. Therefore, the 1-GC check succeeds if the arc-consistency check of  $p'$  in the value domain of  $p$  succeeds.

However, if the arc-consistency test fails, literal  $p'$  is removed from the domain of  $p$ —and its elimination is propagated on the adjacent literals. Since a GC of depth 2 is recursively defined as 1-GC on the neighbors of  $p$  and  $p'$ ,  $p'$  will be removed from the domain if a value that supports  $p'$  is not itself supported and thus, 2-GC test is equivalent to an arc-consistency check with a propagation limited to one step.

In the general case, a GC of depth  $i$  is thus equivalent to arc consistency with  $(i - 1)$  steps of propagation of the eliminated values.

This shows that *GC-Pol* is a subset of *AC-Pol*.

The strict inclusion is shown as follows. Let  $(\mathcal{C}, Ex)$  be a *GC-Pol* problem instance, and consider the addition of a new literal to  $Ex$ . In the GC-based approach, this literal can destroy the  $\theta$ -subsumption determinacy, entailing an exponential complexity of the GC-based resolution of  $\theta$ -subsumption.

In the dual binarization approach, this literal simply intervenes as an additional value in the (dual) constrained variable domain; this does not modify the complexity class of the problem, as the resolution complexity is quadratic in the variable domain size.  $\square$

This shows that the presented approach extends over (Scheffer, Herbrich, & Wysotzki, 1997) the set of problems which can be polynomially addressed.

#### 4.4. Signature

An additional data structure is defined to further restrict the search space.

To each literal in  $\mathcal{C}$  and  $Ex$  is associated a *signature*, similar to the 1-depth graph context (Section 2.3, Scheffer, Herbrich, & Wysotzki (1997)). The signature encodes the links of the literal with other literals through sharing (primal) variables in  $\mathcal{C}$  or constants in  $Ex$ . It is clear that a necessary condition for assigning a dual variable  $Y_p$  (literal in  $\mathcal{C}$ ) to a dual constant  $p'$  (literal in  $Ex$ ) is that the signature of  $Y_p$  be included into that of  $p'$ . Signatures thus enable the pruning of variable domain by arc consistency.

For instance in Example 4, the signature associated to  $p(X_0, X_1)$  states that the first variable appears in a literal built on symbol  $tc$ , position 1, and a literal built on  $p$ , position 1; and the second variable appears in a literal built on symbol  $p$ , positions 1 and 2.

Contrasting with the generality of graph contexts (Scheffer, Herbrich, & Wysotzki, 1997), signatures are restricted to 1-depth graph contexts, which allows for an optimized implementation.

Naturally, signatures suffer from the same limitations as 1-GC:

*Example 10.*

$$\begin{aligned} \mathcal{C}: & \quad p(X, Y), q(Y, X), \\ Ex: & \quad p(x_0, y_0), p(x_1, y_1), q(y_0, x_1), q(y_1, x_0) \end{aligned}$$

Though all literals in  $Ex$  have a signature including that of literals in  $\mathcal{C}$ , the  $\theta$ -subsumption test fails.

Additional signatures, termed 2-signatures, are thus defined. The 2-signature associated to a literal  $p$  describes the literals sharing at least *two* variables with  $p$ ; more precisely, it encodes the relative position of all variables in each literal.

For instance, the 2-signature associated to literal  $p(X, Y)$  in the above example is  $\{\langle q(1.2, 2.1) \rangle\}$ ; it expresses the fact that  $p(X, Y)$  shares two variables with a literal built on symbol  $q$ , such that the first variable in  $p(X, Y)$  is the second in the literal built on  $q$ , and the second variable in  $p(X, Y)$  appears as first variable in this same literal. Literals with several occurrences of a primal variable are also handled through signatures. For instance, the signature associated to the literal  $p(X, X)$  is  $\{\langle 1/p.2 \rangle, \langle 2/p.1 \rangle\}$ .

Like signatures, 2-signatures are used to prune the  $\theta$ -subsumption search space: a necessary condition for a literal  $p$  in  $\mathcal{C}$  to be mapped onto a literal  $p'$  in  $Ex$ , is that the 2-signature associated to  $p$  be included in that of  $p'$ .

## 5. Experimental goal and setting

This section describes the experimental setting chosen for the validation of the proposed  $\theta$ -subsumption approach. This setting is inspired from the now standard setting for validating constraint satisfaction algorithms, which has been imported to the ILP framework by Giordana and Saitta (2000).

The experimental goals are then discussed.

### 5.1. CSP validation paradigm

As mentioned in the introduction, the worst- and average-case complexity analysis of constraint satisfaction problems has been extended using a stochastic modeling.

Stochastic complexity modeling relies on the definition of order parameters on CSP problems (Hogg, 1996). Two main order parameters are considered: constraint density and constraint tightness. For the sake of clarity, only binary constraints will be considered in the rest of the section.

The constraint density  $p_1$  is defined as the fraction of variable pairs which are related by a constraint; letting  $n$  be the number of variables and  $m$  the number of constraints,

$$p_1 = \frac{m}{\frac{n(n-1)}{2}} = \frac{2m}{n(n-1)}$$

The constraint tightness  $p_2$  is the average fraction of value pairs ruled out by a constraint; letting  $L$  be the size of the variable domain, and  $N$  the size of the relation associated to a constraint,

$$p_2 = 1 - \frac{N}{L^2}$$

Stochastic computational complexity is viewed as a random variable depending on the constraint density and constraint tightness parameters. For a given density  $p_1 = v_1$  and tightness  $p_2 = v_2$ , this random variable is realized as the computational effort required to solve any CSP instance with density  $v_1$  and tightness  $v_2$ .

Experimentally, a sample of CSP instances associated to each value pair  $(v_1, v_2)$  is drawn; every CS algorithm is assessed from its average computational effort over the sample.

Stochastic analysis offers a more realistic perspective on CS complexity than worst- and average-case analysis. It accounts for the fact that, although the worst-case complexity is exponential ( $m \times L^n$ , where  $m$  is the number of constraints,  $n$  the number of variables and  $L$  the size of the domain variable), still the empirical cost of resolution is moderate for a vast majority of CSP instances (Cheeseman, Kanefsky, & Taylor, 1991).

More precisely, the CS complexity landscape defined by the order parameters  $p_1$  and  $p_2$  involves three regions:

- The so-called YES region, including under-constrained problems, is associated to a negligible stochastic complexity: a problem instance admits many solutions on average, and finding one of them is easy.



- The so-called NO region, including over-constrained problems, also is associated to a negligible stochastic complexity: a problem instance admits no solution on average, and showing that there is no solution is easy.
- The YES and NO regions are separated by a narrow region termed the *Phase Transition* region (PT). The probability of satisfiability abruptly drops from 1 (in the YES region) to 0 (in the NO region), and the stochastic complexity reaches a peak in the PT region; most hard-to-solve CS instances lie in this region. The height of the peak depends on the algorithm (see Section 6), while it is conjectured that the location of the peak depends on the distribution used to sample the problem instances.

It must be emphasized that stochastic complexity analysis only delivers general trends. Indeed, there exist satisfiable problems in the NO region and unsatisfiable problems in the YES region; and some problems in the PT can be solved with moderate computational effort, while some problems in the YES or NO regions are hard to solve (Gomes et al., 2000).

However, this stochastic modeling offers a principled framework to evaluate the performance and regions of competence of a CS algorithm. Further, this framework allows for a visual analysis of the algorithm competence, through a *competence map*, characterizing the average algorithmic performance with respect to the order parameters (see Section 5.3).

Actually, these competence maps can be used to decide which algorithm is best in a particular region. This facility is highly desirable as it is now widely recognized that no CS algorithm dominates over the others in all regions.

## 5.2. A stochastic complexity model for $\theta$ -subsumption

Stochastic modeling of CS problems has been ported to  $\theta$ -subsumption problems and Inductive Logic Programming by Giordana and Saitta (2000), as follows.

A  $\theta$ -subsumption problem instance is given as a pair (clause  $\mathcal{C}$ , example  $Ex$ ) defined as Horn clauses (after the learning by entailment ILP setting (De Raedt, 1997)). With no loss of generality, only the bodies of  $\mathcal{C}$  and  $Ex$  are considered thereafter (as noted earlier on, the head of  $\mathcal{C}$  can only be mapped onto the head of  $Ex$ ).

Three assumptions are made to facilitate the definition of order parameters; their implications are discussed further in Section 5.4:

- All literals in  $\mathcal{C}$  are built on distinct binary predicate symbols and clause  $\mathcal{C}$  is connected. The latter requirement prevents the  $\theta$ -subsumption problem from being decomposable into simpler problems; its effective complexity is therefore representative of the complexity attached to a given constraint density and tightness.
- Example  $Ex$  is a conjunction of ground literals. Each literal in  $Ex$  is built on a predicate symbol occurring in  $\mathcal{C}$  (other literals are irrelevant to the  $\theta$ -subsumption problem). The number of literals in  $Ex$  per predicate symbol is constant.
- All arguments of all predicates have same value domain.

After these assumptions, four order parameters for  $\theta$ -subsumption problems ( $\mathcal{C}$ ,  $Ex$ ) are defined:

- $n$ : The number of variables in  $\mathcal{C}$ ;
- $m$ : The number of literals in  $\mathcal{C}$  (=the number of distinct predicate symbols);
- $N$ : The number of literals built on each predicate symbol in  $Ex$ ;
- $L$ : The number of constants in  $Ex$ , i.e. the size of each variable domain.

Within the stochastic framework, the algorithmic complexity of  $\theta$ -subsumption ( $\mathcal{C}$ ,  $Ex$ ) is a random variable depending on the sizes of  $\mathcal{C}$  (the number  $m$  of predicates and the number  $n$  of variables) and  $Ex$  (the number  $L$  of constants and the number  $N$  of literals *per* predicate symbol; in total  $Ex$  includes  $N \times m$  literals).

### 5.3. The competence map format

The competence map describes the performance of the algorithm in a given region of interest, measured after a sampling mechanism.

The *region of interest* specifies the ranges for the order parameters  $n$ ,  $N$ ,  $m$  and  $L$ . To keep the total computational cost within reasonable limits, the number  $n$  of variables in  $\mathcal{C}$  is set to 10 and the number  $N$  of literals per predicate symbol is set to 100 (complementary experiments have been done with  $n = 12, 15$  and  $N = 120, 150$  to confirm the general trends reported in Section 6).

The number  $m$  of literals in  $\mathcal{C}$  varies in  $[10, 50]$  and the number  $L$  of constants in example  $Ex$  varies in  $[10, 50]$ . The region of interest thus is defined as  $n = 10$ ,  $N = 100$ ,  $m \in [10, 50]$  and  $L \in [10, 50]$ .

This region corresponds to larger  $\theta$ -subsumption problems than encountered in typical ILP benchmark problems. For instance in the mutagenesis problem (King, Srinivasan, & Sternberg, 1995), the clauses found involve four or five chemical atoms ( $n = 4$ ), with a few literals ( $m < 10$ ) besides the attributes of the molecule (e.g. determinate predicates such as `logp`); the examples involve two predicate symbols (*atm* and *bond*),  $L < 40$  chemical atoms, and  $N < 80$  literals per predicate.

However, from a multi-relational learning perspective (Dzeroski & Lavrac, 2001), this region corresponds to middle-sized problems ( $L \leq 50$  identifiers per example,  $m \leq 50$  tables with  $N = 100$  rows *per* table, and a target clause involving  $n = 10$  independent objects).

The competence map describes the *computational cost* of the algorithm under study. The regions in the competence map refer to the (algorithm-independent) probability of  $\theta$ -subsumption success, measured from the percentage of instances such that  $\mathcal{C}$   $\theta$ -subsumes  $Ex$  (figure 4; see also Giordana & Saitta (2000)). The YES region of under-constrained CS problems (left part of the  $\theta$ -subsumption landscape) corresponds to “short” clauses  $\mathcal{C}$  respectively to “long” examples  $Ex$ ; accordingly,  $\mathcal{C}$  almost surely subsumes  $Ex$ . The NO region of over-constrained problems (right part of figure 4) corresponds to “long” clauses respectively to “short” examples, where  $\mathcal{C}$  almost never subsumes  $Ex$ . In the PT region, the probability for  $\mathcal{C}$  to subsume  $Ex$  abruptly drops from 1 to 0.

In the rest of the paper, the competence map associated to an algorithm describes the average computational effort for solving the  $\theta$ -subsumption test.

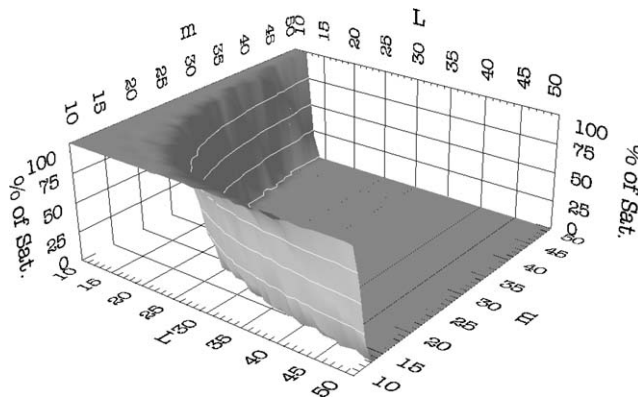


Figure 4. Percentage of  $\theta$ -subsumption tests satisfied on 100 pairs  $(\mathcal{C}, Ex)$ , where  $\mathcal{C}$  is uniformly generated with  $n = 10$  variables and  $m$  literals ( $m$  in  $[10, 50]$ ), and  $Ex$  is uniformly generated with  $N = 100$  literals built on each one of the  $m$  predicate symbols in  $\mathcal{C}$ , and  $L$  constants ( $L$  in  $[10, 50]$ ).

The competence map last depends on *the sampling mechanism*, used to select problem instances  $(\mathcal{C}, Ex)$  conditioned by the number of variables  $n$  and the number of literals  $m$  in  $\mathcal{C}$ , and the number of literals  $N \times m$  and the number of constants  $L$  in  $Ex$ . Following Giordana and Saitta (2000), we use a uniform sampling mechanism.

The  $2 \times m$  arguments in  $\mathcal{C}$  (since all  $m$  predicates are binary) are partitioned into  $n = 10$  subsets (the variables in  $\mathcal{C}$ ). Further, one ensures that  $\mathcal{C}$  is connected (see Section 5.4).

The  $N = 100$  literals  $p_i(a_j, a_k)$  built on each predicate symbol  $p_i$  in  $Ex$  are constructed by uniformly selecting without replacement the pairs  $(a_i, a_j)$ , where  $i, j$  are in  $[1, L]$ .

A sample of 100 problems instances  $(\mathcal{C}, Ex)$  is uniformly drawn for each value of  $(n = 10, N = 100, m \in [10, 50], L \in [10, 50])$ .

This relatively small sample size is justified as no heavy tailed distribution was observed for  $\theta$ -subsumption cost in the parameter range considered, although such distributions were reported for large size CSPs (Gomes et al., 2000).

#### 5.4. Discussion

Before presenting the experimental results in Section 6, we discuss the scope and limitations of our experimental setting. The reliability of the competence map depends on the representativity of the artificial problems generated as above. These artificial problems might differ from  $\theta$ -subsumption problems encountered in ILP benchmarks in several respects.

A first issue is the non-decomposability of clause  $\mathcal{C}$ . In the general CS case, the decomposition of the problem at hand into (scarcely related) subproblems is a major practical issue (Gyssens, Jeavons, & Cohen, 1994; Dechter & Pearl, 1989); similarly, the decomposition of  $\mathcal{C}$  into  $k$ -local components is a key issue in relational logic and ILP (Gottlob & Leitsch,

1985; Kietz & Lübbe, 1994). Such a decomposition is desirable as it entails exponential savings in the resolution cost.

However, on one hand decomposability is not often encountered in ILP problems. Scheffer, Herbrich, and Wysotzki (1997) shows empirically that determinate matching was rarely operational on the mesh problem (Dolsak & Muggleton, 1991; Costa et al., 2004) shows that the level of non-determinism is high, with a branching factor *circa* 30 in the mutagenesis (King, Srinivasan, & Sternberg, 1995) and carcinogenesis (Srinivasan, 1997) domains.

On the other hand, evaluating a  $\theta$ -subsumption algorithm on decomposable clauses introduces an optimistic bias, as the cost of simple problems is wrongly reported for complex ones (Giordana & Saitta, 2000). Such an optimistic bias reduces the reliability of the competence map, which is not desirable.

A second issue concerns the uniform distribution of the  $\theta$ -subsumption problems considered. Each predicate symbol occurs exactly once in the clause  $\mathcal{C}$ , and  $N = 100$  times in example  $Ex$ . Moreover, each predicate argument has same domain  $\{a_1, \dots, a_L\}$ .

In real-world examples, some predicate symbols occur more frequently than others; predicate arguments are typed, and their domains have diverse sizes. The use of types enables significant savings as it restricts the  $\theta$ -subsumption search space (Lloyd, 2003); furthermore, the diversity of domain sizes results in a better efficiency of CSP heuristics, e.g. constraint propagation or dynamic variable ordering. In this respect, the fact that  $Ex$  involves  $N$  literals for every predicate symbol, all instantiated from pairs of  $L$  constants, leads to a pessimistic estimate of the average  $\theta$ -subsumption cost. On the other hand, this restriction permits a simple and direct interpretation of the competence map.

The third issue concerns the arity of the predicate symbols, as the presented setting is limited to binary predicate symbols. However, the predicate arity does not intervene directly on computational complexity (note that the worst-case complexity  $L^n$  only depends on the variable domain size, exponentiated by the number of variables). More precisely, the key complexity factor depends on how many arguments in  $\mathcal{C}$  ( $m \times k$ , if  $k$  denotes the predicate arity) are grouped into how many subsets (the number  $n$  of variables). The  $\theta$ -subsumption equivalent for CS density is the probability for two literals in  $\mathcal{C}$  to share a variable (meaning that they cannot be independently mapped onto literals in  $Ex$ ).

This probability is  $1 - \frac{\binom{k(n-k)}{\binom{k}{n}}}{\binom{k}{n}}$  (two literals are independent iff the  $k$  arguments in the second literal, are mapped onto  $k$  among the  $n - k$  variables not involved in the first literal), of order  $\mathcal{O}(1 - (1 - \frac{k}{n})^k)$ . According to this analysis, increasing the predicate arity  $k$  by factor  $t$  can be likened to increasing the total number of variables  $n$  by approximately  $t^2$ . The competence map drawn for binary predicates can thus be used to estimate the  $\theta$ -subsumption algorithm behavior for non-binary predicates, by moving toward higher values of  $n$ .

Experimentally, increasing the number of variables  $n$  causes the phase transition to move toward longer hypotheses everything else being equal (toward the region on the right in figure 5), with exponential increase of the complexity peak (Giordana & Saitta, 2000).

We restrict therefore the experimental setting to binary predicates, to keep the total computational cost within reasonable limits.

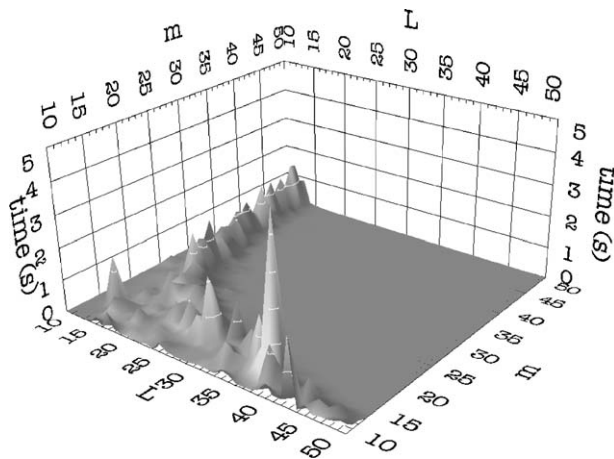


Figure 5. Competence map of *Django.V1*, reporting the average  $\theta$ -subsumption computation cost measured as in figure 4.

In conclusion, the assumptions made in Section 5.2 are meant to construct a reliable  $\theta$ -subsumption competence map, though it might be pessimistic for ILP problems with distributional diversity.

### 5.5. Experimental goals

As mentioned already, there is no universally efficient CSP algorithm. One primary goal of the paper thus is to determine which CSP algorithms and procedures are most efficient wrt  $\theta$ -subsumption problems in general. A secondary goal is to characterize where these procedures are particularly efficient.

To this aim, we studied 10 variants of the *Django* algorithm, involving various combinations of CSP procedures and summarized in Table 1. For each of these variants, a competence map is constructed as detailed in Section 5.3.

The baseline algorithm *Django.V1* combines arc-consistency checking (Mackworth, 1977) and forward checking (the propagation of the current assignment is restricted to the next variable domain).

The influence of variable ordering is investigated through *Django.V2* to *Django.V6*. These variants involve distinct heuristics, all based on the First Fail Principle. In *Django.V2*, variables with minimal domain are ranked first.<sup>2</sup> In *Django.V3*, variables subject to a maximal number of dual constraints are ranked first (favoring the literals in  $\mathcal{C}$  which are most connected to others literals). Both criteria are combined in *Django.V4* and *Django.V5*, with different priorities.

In *Django.V6*, the order criterion is heuristically defined as the domain size, divided by the number of neighbors in the graph; variables are sorted by increasing value of this criterion.

Table 1. *Django*, Variants **V1** to **V10**.

---

|            |   |
|------------|---|
|            | <i>Baseline</i>   |
| <b>V1</b>  | Arc consistency + simple Forward Checking (propagation of the current assignment wrt the next variable) |
|            | <i>Dynamic Variable Ordering</i>  |
| <b>V2</b>  | <b>V1</b> + DVO based on minimal domain (random choice in case of tie)                                  |
| <b>V3</b>  | <b>V1</b> + DVO based on maximal connectivity (random choice in case of tie)                            |
| <b>V4</b>  | <b>V1</b> + DVO based on (min. domain, max. connectivity) (minimal domain, then maximal connectivity)   |
| <b>V5</b>  | <b>V1</b> + DVO based on (max. connectivity, min. domain) (maximal connectivity then minimal domain)    |
| <b>V6</b>  | <b>V1</b> + DVO based on min. domain ÷ max. connectivity (random choice in case of tie)                 |
|            | <i>Forward Checking</i>   |
| <b>V7</b>  | <b>V6</b> + improved Forward Checking propagation of forced assignments (singleton candidate value)     |
|            | <i>Arc Consistency</i>  |
| <b>V8</b>  | <b>V7</b> + AC based on signatures  |
| <b>V9</b>  | <b>V7</b> + AC based on signatures and 2-signatures   |
|            | <i>Maintaining Arc Consistency</i>  |
| <b>V10</b> | <b>V6</b> + Full Arc Consistency on assignments + AC based on signatures and 2-signatures               |

---

The influence of Forward Checking is investigated in *Django.V7*; in addition to the 1-step propagation of the current assignment, forced assignments (singleton candidate value for any variable) are propagated.

The influence of arc consistency is investigated using signatures and 2-signatures. *Django.V8* differs from *Django.V7* as it considers the literal signatures; *Django.V9* considers both signatures and 2-signatures. Signatures and 2-signatures allow for pruning the domains of the literals in a preprocessing step, and redundant tests are avoided during the constraint propagation.

Finally, *Django.V10* implements a Maintaining Arc Consistency (MAC) procedure, where each assignment is fully propagated in each resolution step, using both signatures and 2-signatures.

## 6. Experimental validation

The CSP-inspired approach of  $\theta$ -subsumption implemented within the *Django* system is compared to several reference algorithms. These algorithms are introduced and the format of results is presented. The experimental results are then discussed and some interpretations are proposed.

### 6.1. Reference algorithms

*Django* (written in C++) is compared to three reference algorithms (written in C): SLD Prolog, determinate matching (Kietz & Lübbe, 1994) and graph contexts (Scheffer, Herbrich,

& Wysotzki, 1997). We used the reference algorithm implementations kindly given by T. Scheffer. Run times are measured on a PC-Pentium II, 300 MHz.

As might have been expected, Prolog SLD does not keep up when  $\mathcal{C}$  involves more than a few literals, and it had to be stopped for  $m > 5$  (note that example  $Ex$  involves  $100 \times m$  literals).

Determinate matching (Kietz & Lübbe, 1994) does significantly better than SLD for small size hypotheses; however, it runs out of resources for  $m > 10$ . Indeed, determinate matching is inappropriate to the uniform structure of the examples.

Graph contexts (Scheffer, Herbrich, & Wysotzki, 1997) also turned out to be not applicable, for efficiency reasons (discussed further in Section 6.3).

Finally only the maximal clique search (MCS) (Scheffer, Herbrich, & Wysotzki, 1997) could be used in the same problem range as *Django*.

## 6.2. Presentation of the results

The behavior of each algorithm is graphically displayed through a competence map (Section 5.3), reporting for each ( $n = 10, N = 100, m \in [10, 50], L \in [10, 50]$ ) the computational cost averaged over 100 pairs  $(\mathcal{C}, Ex)$  uniformly drawn, where  $\mathcal{C}$  involves  $n = 10$  variables and  $m$  literals, and  $Ex$  involves  $N = 100$  literals for each of the  $m$  predicate symbols, and  $L$  constants.

Figures 5 and 6 respectively display the competence maps associated to *Django.V1* and MCS.

With respect to the  $\theta$ -subsumption phase transition (figure 4), both competence maps show similar trends in the YES, NO and PT regions. The PT region coincides with the

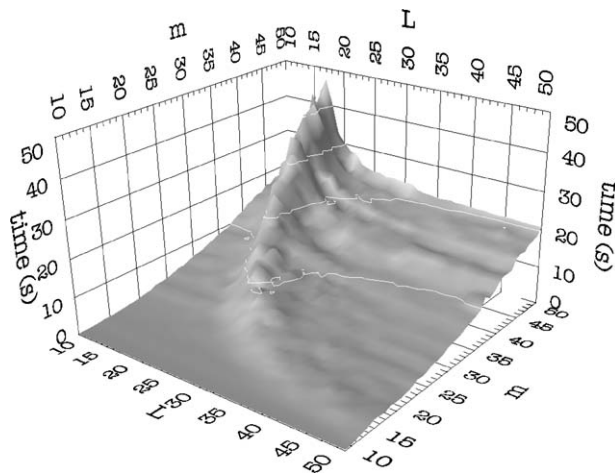


Figure 6.  $\theta$ -subsumption cost( $m, L$ ) for MCS, averaged on 100 pairs  $(\mathcal{C}, Ex)$ , with scaling factor  $\times 10$  wrt *Django.V1* (figure 5).

Table 2. *Django*: Average  $\theta$ -subsumption cost in milliseconds (Pentium II) in the YES, NO and PT regions, Artificial setting.

|                  | YES region | Phase transition | NO region |
|------------------|------------|------------------|-----------|
| MCS              | 4582.4     | 10509.5          | 12680.9   |
| <i>Django.V1</i> | 158.3      | 764.1            | 38.1      |
| <b>V2</b>        | 43.6       | 51.3             | 20.5      |
| <b>V3</b>        | 50.5       | 83.3             | 22.3      |
| <b>V4</b>        | 43.2       | 50.1             | 20.4      |
| <b>V5</b>        | 46.8       | 68.1             | 21.3      |
| <b>V6</b>        | 43.0       | 48.7             | 20.4      |
| <b>V7</b>        | 40.9       | 45.0             | 19.0      |
| <b>V8</b>        | 24.3       | 33.7             | 16.8      |
| <b>V9</b>        | 26.4       | 30.9             | 10.3      |
| <b>V10</b>       | 29.1       | 28.9             | 10.3      |

region of high complexity peaks, like a mountain chain of hyperbolic shape in plane  $(m, L)$ . The peak values increase with  $L$  as expected from the worst complexity analysis.

The NO region, on the left of the PT (high values of  $m$  and/or  $L$ ), is a low complexity region with moderate variance.

The YES region, on the right of the PT, is a moderate complexity region, with high variance.

For the sake of precise comparisons, the cost of each algorithm is also aggregated in the YES, NO and PT regions. Table 2 gives the average cost over all points  $(m, L)$  such that the  $\theta$ -subsumption probability of success is respectively greater than 90% (YES region), less than 10% (NO region), and in [10%, 90%] (PT region); the standard deviation is graphically displayed in figures 9 and 10.

### 6.3. Comparison with MCS

*Django.V1* significantly outperforms MCS; the gain factor is about 30 in the YES region, 15 in the PT, and 330 in the NO region (a scaling factor of 10 is used in figure 6 compared to figure 5).

The phase transition phenomenon is most marked for *Django.V1* (figure 5), the cost in the PT region being 5 times the cost in the YES region and 20 times the cost in the NO region. In comparison, the cost of MCS in the PT region is only twice the cost in the YES region, and the cost reaches its maximum in the NO region, close to the PT. This suggests that MCS does not detect early the inconsistencies during the search, and unnecessarily explores large portions of the search tree.

The lesser efficiency of MCS is blamed in part on the construction of the substitution graph, which must be completed before the maximal clique search takes place. This



construction is computationally heavy for small sized problem instances (in  $\mathcal{O}(m^2N^2)$ ). For large sized problem instances, the construction cost is negligible compared with the (exponential) maximal clique search and well worth the effort as the substitution graph significantly speeds up the maximal clique search. Unfortunately, the storage of the substitution graph then becomes unaffordable.

Also for complexity reasons, 2-graphs contexts could not be used for  $m > 10$  (as  $k$ -graphs are constructed and checked with exponential complexity in  $k$ ).

In contrast, *Django.V1* interleaves the search and the constraint propagation, avoiding the construction of the whole graph whenever a solution is found along the search.

The computational cost of MCS also shows a high variance compared to *Django.V1*. Practically, a time limit of one minute was set on every  $\theta$ -subsumption test to ensure the feasibility of the experiments with MCS; in preliminary runs, the resolution of some problems instances required more than 30 minutes on Pentium II. This variance is also blamed on the lack of adequate pruning mechanism within MCS, ensuring an early detection of inconsistencies.

#### 6.4. Comparison of CSP procedures

The impact of the diverse CSP procedures in *Django* is presented in Table 2. The first set of procedures is concerned with dynamic variable ordering (DVO, *Django.V2* to *Django.V6*). *Django.V2* significantly improves on *Django.V1*, with a dynamic variable ordering based on minimal domain, which can be viewed as a non-deterministic generalization of determinate matching (Kietz & Lübke, 1994). Sorting the variables (literals in  $\mathcal{C}$ ) based on their connectivity (*Django.V3* and *Django.V5*) is less efficient, though it also improves on *Django.V1*. The best sorting criterion is that of *Django.V6*, where literals are sorted by increasing value of their domain size divided by their number of 1-neighbors (literals sharing at least one variable).

Again, the efficiency of DVO depends on the region considered; the gain factor is less than 2 in the NO region, 3 in the YES region, and circa 15 in the PT region. The major effect of DVO is to smoothen the complexity peak in the PT region.

The addition of Forward Checking with propagation of singleton domains (*Django.V7*) very slightly improves on *Django.V6*, with a gain factor circa 5% in all three regions.

The effects of arc consistency (*Django.V8* uses signatures in a pre-processing step, and *Django.V9* likewise uses signatures and 2-signatures) are significantly beneficial compared to *Django.V7*. Note that *Django.V8* is more efficient in the YES region than *Django.V9* (with a gain factor 40% on *Django.V7*), and *Django.V9* is more efficient in the PT and NO regions (with a gain factor respectively 30% and 45% over *Django.V7*). This is interpreted as the construction and exploitation of 2-signatures is not worth the effort for small sized problem instances.

The efficiency of signature-based heuristics surprisingly contrasts with the relative inefficiency of graph contexts (Scheffer, Herbrich, & Wysotzki, 1997). Indeed, graph contexts and signatures alike proceed by reducing the variable domains. The difference is that all consistency checks are computed and stored in the substitution graph; in contrast, in *Django*

signatures are only meant to speed up consistency checks, which are performed as needed along constraint propagation; note that a domain not reduced through signatures does not undergo arc consistency. Furthermore, it appears that graph contexts suffer from a too general and thereby inefficient implementation, devised for handling  $k$ -GC though  $k \leq 2$  in practice.

Last, a Maintaining Arc Consistency (MAC) procedure is integrated in *Django.V10* (figure 8). This heuristic is mostly beneficial in the PT region: the additional constraint checks remove inconsistencies that would otherwise lead to backtracks (gain factor of 7% over *Django.V9*).

Globally, the extra cost of MAC is not worth the effort in the YES region, since few backtracks occur in this region anyway (loss factor of 10% over *Django.V9*). Further, MAC makes no difference in the NO region compared to *Django.V9*. This is explained as MAC is only called after arc consistency and forward checking, while arc consistency fail most of the times in the NO region.

Still, MAC brings a significant improvement in the PT region, increasing with the instance size. This improvement is explained as MAC avoids the dramatic extra cost due to the hardest problems handled by *Django.V9*. This way, the use of MAC smooths both the computational cost, and its variance, in the PT region. Figures 9 and 10 respectively show the standard deviation of the computational  $\theta$ -subsumption cost of *Django.V9* and *Django.V10*.

Finally, even though MAC only slightly improves on Forward Checking in terms of computational gain, *Django.V10* appears much more stable than *Django.V9* as it drastically decreases the cost of the hardest problems. From an operational perspective, this would recommend the use of *Django.V10* as more appropriate to real-time or manufacturing applications than *Django.V9*.

The *Django* system at its best (*Django.V9* and *Django.V10*), thus improves on MCS with a gain factor about 150 in the YES region, 360 in the PT region, and 1200 in the NO region (figures 7 and 8) on the artificial problems. Real-world problems will be considered in the next section.

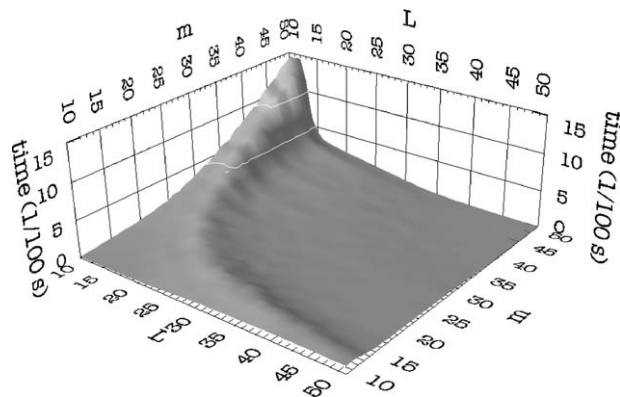


Figure 7.  $\theta$ -subsumption cost( $m, L$ ) for *Django.V9*, averaged on 100 pairs ( $C, Ex$ ). A factor  $\div 25$  is observed compared to *Django.V1* (figure 5).

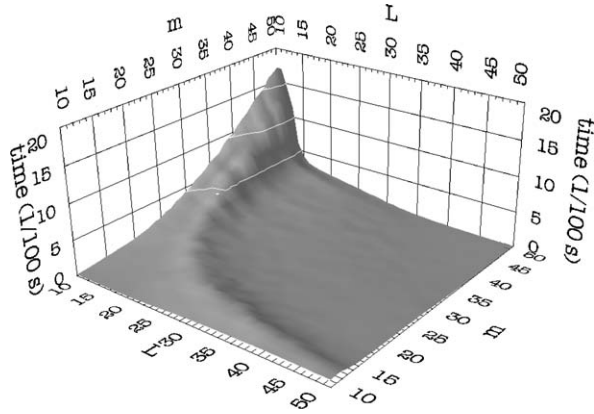


Figure 8.  $\theta$ -subsumption cost( $m, L$ ) for *Django.V10*, averaged on 100 pairs ( $C, Ex$ ). A factor  $\div 25$  is observed compared to *Django.V1* (figure 5).

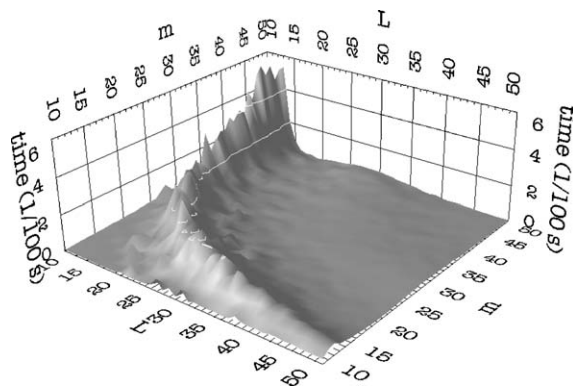


Figure 9. Standard deviation for *Django.V9* (average results reported in figure 7).

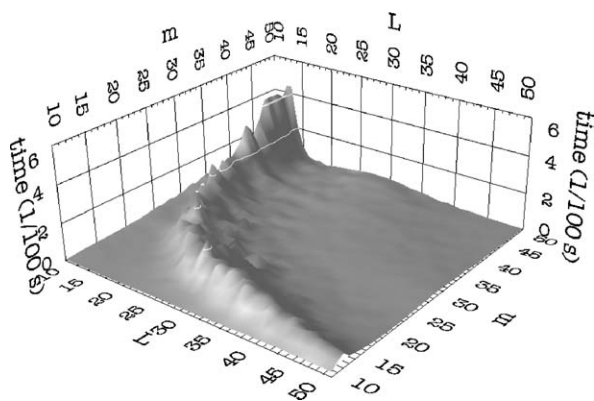


Figure 10. Standard deviation for *Django.V10* (average results reported in figure 8).

## 7. *Meta-Django*

This section is concerned with the choice of the best algorithm in *Django*. The goal is defined and discussed with respect to the Meta-learning goal (Pfahring, Bensusan, & Giraud-Carrier, 2000; Bensusan & Kalousis, 2001). Then the *Meta-Django* algorithm is presented, which exploits the competence maps built in the previous section. Experimental validation considers artificial problems and the real-world mutagenesis problem (King, Srinivasan, & Sternberg, 1995).

### 7.1. *Meta-learning*

The experimental results presented in Section 6 confirm that none of *Django.V9* and *Django.V10* dominates the other. This raises the question of how to choose the best algorithm depending on the problem instance.

This general question, recognized a key question in Machine Learning since the 90's, was formulated as a Meta-Learning problem in Brazdil, Gama, and Henery (1994), Pfahring, Bensusan, and Giraud-Carrier (2000) and Bensusan and Kalousis (2001): learn which algorithm is most suited to the problem instance at hand. Training examples for the meta-learning task are made of a pair (problem, algorithm), labeled with the success of the algorithm for the problem. Depending on whether the label is numerical (predictive accuracy, rank of the algorithm) or boolean (the algorithm reaches the best predictive accuracy), the meta-learning task is one of regression or classification. Another meta-learning setting considers examples made of a triplet (problem, algorithm1, algorithm2), labeled as *true* when algorithm1 dominates algorithm2 on this problem (Fürnkranz & Petrak, 2001).

The very elegant approach of meta-learning unfortunately depends on the quantity and quality of available training data, like all learning applications. A first difficulty regards the representation of the examples (problem, algorithm). Algorithms are represented by their name (categorical attribute), together with their parameters. Problems are described along many attributes, either static (size, number of missing values,...) or dynamic ("landmarking" the problem instance with the predictive accuracy of a  $k$  nearest neighbor, naive Bayes,...) (Pfahring, Bensusan, & Giraud-Carrier, 2000).

A second difficulty is the choice of the problem instances. In Bensusan and Kalousis (2001), the problem instances were constructed through principled perturbations of problems in the Irvine repository (Blake, Keogh, & Merz, 1998), (e.g. increasing the rate of missing values, incorporating irrelevant attributes,...).

However, the choice of the perturbations is a strong bias on Meta-learning; for instance, the importance of the "missing value rate" attribute in the meta-rules directly reflects its use in the perturbations (Kalousis, 2002).

Returning to the selection of the best  $\theta$ -subsumption algorithm, the stochastic complexity approach taken in this paper can be viewed as an alternative to meta-learning. The representation of the problem instances, set to the order parameters, is firstly defined; based on these order parameters, principled experiments yield the competence map of each algorithm; the competence maps can be used as look-up tables. These look-up tables encapsulate the same knowledge as sought for by meta-learning: they predict the behavior of the algorithms

conditioned by the general characteristics of the problem instance. This prediction supports the built-in selection of the best algorithm, achieved in *Meta-Django*.

According to experimental evidence, *Meta-Django* must select *Django.V9* (Forward checking) in the YES region and *Django.V10* (Maintaining Arc Consistency, MAC) in the NO and PT region. (From an implementation perspective, these procedures involve the same data structure; the selection is thus integrated seamlessly to the acquisition of the  $\theta$ -subsumption problem instance).

*Meta-Django* thus only requires to determine whether the current problem instance lies in the YES region.

## 7.2. Localizing the problem instance

To this aim, we borrow to the CSP literature the  $\kappa$  parameter defined by Gent et al. (1996):

$$\kappa = \frac{-\sum_{r \in \mathcal{R}} \log_2(1 - p_2)}{\log_2(|\text{dom}(\mathcal{X})|)}$$

where  $\mathcal{R}$  is the set of constraints,  $p_2$  is the constraint tightness (defined in Section 3 as the number of assignments ruled out by a constraint  $r$ ), and  $|\text{dom}(\mathcal{X})|$  the total size of the variable domains ( $L^n$  since all  $n$  variables have same domain of size  $L$ ). After Gent et al. (1996), problems in the PT region are associated with  $\kappa$  values between 0.75 and 1; the YES region is associated with  $\kappa < .75$  and the NO region is associated with  $\kappa > 1$ .

The  $\kappa$  parameter is first attuned to  $\theta$ -subsumption problems. Term  $(1 - p_2)$  corresponds to the probability of an assignment to be admissible wrt a given constraint. With respect to  $\theta$ -subsumption (before binarization), each constraint corresponds to some literal  $p$  in  $\mathcal{C}$ ; the number of admissible assignments is thus upper-bounded by the number of literals in  $Ex$  built on the same predicate symbol noted  $|\text{rel}(p)|$ ; in the meanwhile, the number of possible assignments is estimated as the total number of constants or variables in  $Ex$ , exponentiated by the arity of the current predicate symbol ( $|\text{arg}(Ex)|^{\text{arity}(p)}$ ).

Finally:

$$\begin{aligned} \kappa &= \frac{-\sum_{p \in \mathcal{C}} \log_2(|\text{rel}(p)|) - \text{arity}(p) \times \log_2(|\text{arg}(Ex)|)}{\log_2(|\text{dom}(\mathcal{X})|)} \\ &= \frac{-\sum_{p \in \mathcal{C}} \log_2(|\text{rel}(p)|) - \text{arity}(p) \times \log_2(|\text{arg}(Ex)|)}{n \log_2(|\text{arg}(Ex)|)} \\ &= \frac{m \times (2 \log_2 L - \log_2 N)}{n \log_2 L} \end{aligned}$$

where  $n$  denotes the number of variables in  $\mathcal{C}$ .

There is a good empirical agreement between parameter  $\kappa$  and the actual location of the  $\theta$ -subsumption problem instances. Practically, *Meta-Django* uses the FC procedure (Version **V9**) for problems such that  $\kappa < .7$  (Yes region) and the MAC procedure (Version **V10**) is used in all other cases.

Table 3. Meta-*Django*: Average  $\theta$ -subsumption cost in milliseconds (Pentium II) in the YES, NO and PT regions defined as in Table 2, Artificial setting.

|                     | YES region | Phase transition | NO region |
|---------------------|------------|------------------|-----------|
| <b>V9</b>           | 26.4       | 30.9             | 10.3      |
| <b>V10</b>          | 29.1       | 28.9             | 10.3      |
| Meta- <i>Django</i> | 26.4       | 28.9             | 10.3      |

### 7.3. Experimentation on artificial problems

Using the same experimental setting as in Section 5.3, by construction Meta-*Django* reaches the best cost among *Django*V9 and *Django*.V10 (Table 3), and the competence map of Meta-*Django* (figure 11) combines the best of respectively *Django*.V9 (in the YES region) and *Django*.V10 (in the NO and PT regions).

In the YES region, *Django*.V9 (FC) is selected; Meta-*Django* thus avoids checking and reducing the variable domains (as in MAC), which is in any case useless in the YES region.<sup>3</sup> However, the FC procedure is less efficient than MAC on exceptionally hard problem instances; this is demonstrated as the standard deviation of Meta-*Django* computational cost in the YES region is worse than for *Django*.V10 (figure 12).

*Django*.V10 (MAC) is selected in the PT region, with an additional propagation of assignments compared to FC. This propagation results in significant reductions of the variable domains, preventing further expensive backtracks.

*Django*.V10 is selected in the NO region for the sake of simplicity, as there is no actual difference in this region between MAC and FC. Actually, the resolution is mostly achieved during the preliminary arc consistency step, which is common to both FC and MAC procedures.

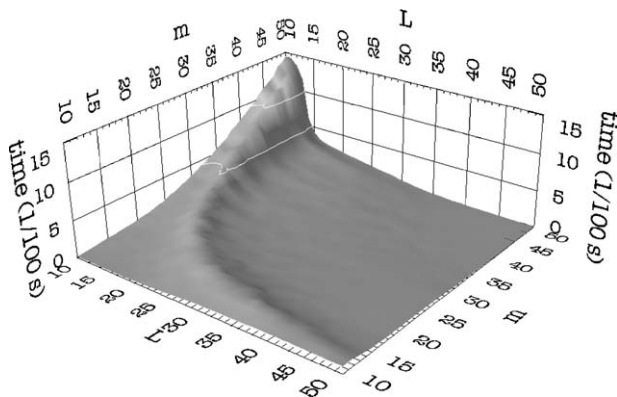


Figure 11. Meta-*Django*:  $\theta$ -subsumption cost( $m, L$ ) averaged on 100 pairs ( $C, Ex$ ). A factor  $\div 25$  is observed compared to *Django*.V1 (figure 5).

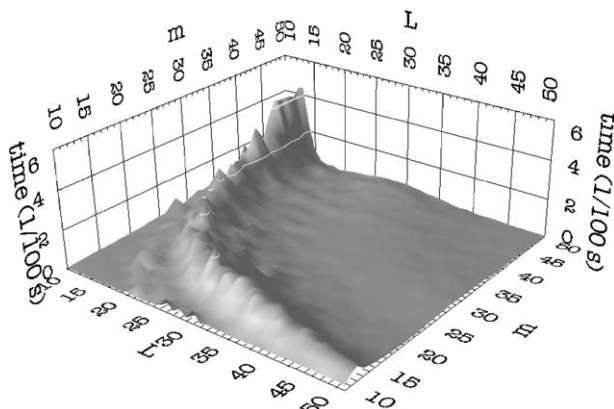


Figure 12. Meta-Django: Standard deviation on  $\text{cost}(m, L)$ .

#### 7.4. Experimentation on the mutagenesis domain

In order to validate *Django* and *Meta-Django* on a real-world like domain, they were integrated into an ILP learner inspired by Warmr (Dehaspe & De Raedt, 1997) and presented in Maloberti and Suzuki (2003), with the goal of finding the clauses with support  $\geq 10\%$  in the mutagenesis domain (King, Srinivasan, & Sternberg, 1995). However the focus is here on comparing the  $\theta$ -subsumption algorithms, rather than on the frequent clauses.

For the sake of computational feasibility, the following restrictions were made:

- The learner proceeds in a level-wise fashion (Agrawal & Srikant, 1994). At each level, only 10,000 clauses are retained out of the solution clauses. These retained clauses are specialized by addition of another literal.
- Clauses are limited to 10 literals. They are connected and they only involve variables (though *Django* and *Meta-Django* support constants); further, all variables in a literal are distinct.

Furthermore, the constants corresponding to the molecule identifier in the examples are omitted. In the learning from entailment setting (De Raedt, 1997) considered in this paper, the body of the example is made of all literals relevant to the description of this example; therefore there is no need for linking these literals explicitly through the example identifier.

Each candidate clause is tested against all examples in the dataset; the CPU time is measured on a Pentium IV 3 GHz. The  $\theta$ -subsumption cost reported in Table 4 gives the average cost of checking one candidate clause against all examples; the  $\theta$ -subsumption tests made to prune the logically equivalent solution clauses, or the candidate clauses that cover an infrequent clause, are not considered.

The graph context algorithm (Scheffer, Herbrich, & Wyszotzki, 1997) could not be used besides the third level and its results are omitted.

Experimental evidence in Table 4 shows some differences compared to the artificial setting (Table 2).

Table 4. MCS, *Django* and Meta-*Django*: Average  $\theta$ -subsumption cost in milliseconds (Pentium IV) in the YES, NO and PT regions defined as in Table 2, Mutagenesis domain.

|                     | YES region | Phase transition | NO region |
|---------------------|------------|------------------|-----------|
| MCS                 | 532.4      | 644.4            | 389.5     |
| <i>Django.V1</i>    | 35.1       | 370.3            | 408.0     |
| <b>V2</b>           | 33.5       | 246.9            | 59.7      |
| <b>V3</b>           | 34.3       | 413.5            | 61.6      |
| <b>V4</b>           | 33.6       | 239.6            | 45.3      |
| <b>V5</b>           | 34.1       | 351.1            | 57.5      |
| <b>V6</b>           | 33.6       | 237.9            | 45.2      |
| <b>V7</b>           | 35.7       | 303.9            | 53.3      |
| <b>V8</b>           | 17.5       | 274.3            | 8.2       |
| <b>V9</b>           | 18.3       | 269.7            | 8.0       |
| <b>V10</b>          | 21.7       | 264.8            | 7.0       |
| Meta- <i>Django</i> | 21.6       | 263.6            | 6.9       |

MCS is still outperformed by all *Django* variants (except for *Django.V1* in the NO region). However the gain factor is less than for the artificial problems. This is explained by the smaller size of the  $\theta$ -subsumption problems, reducing the cost of the substitution graphs and making MCS more efficient. However, MCS still does not detect inconsistencies as early as possible, and the best gain factor still occurs in the NO region.

Dynamic variable ordering based on minimal domains (*Django.V2*, *Django.V4* and *Django.V6*) likewise improves the results compared to *Django.V1*, and the gain factor is less than for the artificial problems. In fact, DVO appears to be beneficial mostly in the NO region.

Surprisingly, the forced propagation of singleton domains (*Django.V7*) appears to be at best useless and most often harmful in all regions. This might be explained as the connectivity of the literals is more regular in the mutagenesis domain than for artificial problems; the overhead of singleton propagation is not compensated for by the exploitation of highly constrained variables.

Quite the contrary, the use of signatures (*Django.V8* and *Django.V9*) appears to be very efficient, with a gain factor of 2 and 6 respectively in the YES and NO regions. As could have been expected, 2-signatures (*Django.V9*) does not improve significantly the performances, as literals share at most one variable in general (being reminded that the molecule identifier was omitted). But as the type of both *bond* and *atm* is coded by an integer, some literals happen to share two variables. Of course, in a real-world application the example description can be inspected to determine a priori whether the use of 2-signatures is relevant or should better be discarded.

Last, Maintaining Arc Consistency (*Django.V10*) slightly improves over *Django.V9* in the PT and NO regions, and does no good in the YES region. On these problems,



the computational overhead of full arc consistency penalizes *Django.V10* compared to *Django.V9*. As noted earlier on, the extra-cost of MAC is worth the effort on large-sized problems only. Complementary experiments show that MAC catches up and significantly outperforms FC for clauses with 12 literals or more. Regarding *Meta-Django*, since the algorithm selection is built in from the competence maps in Section 6, there is no wonder that it selects *Django.V10* in the *NO* and *PT* regions, though *Django.V8* was the best one in the *NO* region on the mutagenesis domain. What is more surprising, *Meta-Django* also tends to select *Django.V10* in the *YES* region. This is interpreted as the  $\kappa$  factor is not sufficient to localize the *PT* when dealing with a non-uniform example distribution. Other factors, such as the clause size, should also be considered.

As an alternative, *Meta-Django* can also be directly updated and refined from Table 4: clearly, *Django.V8* improves on *Django.V10* in the *YES* region for this problem.

Nevertheless, in its actual version *Meta-Django* improves on MCS by a gain factor ranging from 2.5 (in the *PT* region) to 24 (in the *YES* region) and 55 (in the *NO* region) on this real-world problem.

### 7.5. Conclusions regarding *Meta-Django*

Complementary experiments show that the comparative performance of *Django.V8* and *Django.V10* is sensitive to the size of the cache memory (128 Kb and 512 Kb). It appears that a small cache size hinders MAC (*Django.V10*) as frequent memory access might not be compensated for by effective pruning. However, for “sufficiently” large domain sizes *Django.V10* catches up and ultimately dominates *Django.V8*.

In summary, experiments on artificial problems and on the mutagenesis domain show three main trends. First of all, Dynamic Variable Ordering with minimal domain (*Django.V2*, *Django.V4* and *Django.V6*) significantly improves on *Django.V1*, particularly in the *PT* and *NO* region. Second, the use of signatures (*Django.V8*) also improves the performance though to a lesser extent, particularly in the *YES* region. Last, the use of MAC (*Django.V10*) is only beneficial for problems with large domain sizes.

## 8. Related works

The  $\theta$ -subsumption test was recently reconsidered in the ILP literature along several lines.

A novel generality relation, termed *subsumption under object identity* (OI) was defined by Esposito et al. (2000). OI requires substitution  $\theta$  to be injective, i.e. to map distinct variables onto distinct constants. OI is mostly based on semantic motivations (the roles of variables in a hypothesis should be distinct), though the additional requirement on substitutions also restricts the search space and thereby increases the efficiency of  $\theta$ -subsumption; note that *Django* would achieve OI  $\theta$ -subsumption by adding *diff constraints* on all pairs of variables. Along the same CSP lines, a novel  $\theta$ -subsumption algorithm has been proposed to find all OI-based successful substitutions in Mauro et al. (2003).

Another approach was proposed by Blockeel and De Raedt (1998) to limit the computational cost of hypothesis assessment. This approach proceeds by sharing the data structures

and cost of  $\theta$ -subsumption among several similar candidate clauses. Some modifications are needed (and planned for further research) in order to handle query packs in *Django*. Typically, the arc consistency and constraint propagation algorithms must effect a limited propagation of the domain reductions; the check of signatures and 2-signatures must handle in a different way the variables shared with the parent of the clauses.

Also, Costa et al. (2004) propose the combined use of several transformations to speed up  $\theta$ -subsumption, based on additional information on the variable types. This additional information allows for decomposing the clause (*cut-transformation*) into independent parts up to the instantiation of “grounded variables”. The *cut-transformation* is extended by the *once-transformation*, recursively decomposing the independent parts in independent subparts that share only variables instantiated in a previous step. In parallel, the *smartcall-transformation* prunes independent subsets of the clause for which a solution was found. However, it is noted that *once-transformation* is limited to acyclic clauses, where the instantiations of variables can be hierarchically ordered.

## 9. Conclusion and perspectives

The thesis advocated in this paper is that many results in the CSP literature are relevant and can lead to theoretical and algorithmic advances in Inductive Logic Programming, as first advocated in Giordana and Saitta (2000).

The main result of the paper is a novel  $\theta$ -subsumption system first presented in Maloberti and Sebag (2001), *Django*, which improves on earlier algorithms (Kietz & Lübke, 1994; Scheffer, Herbrich, & Wysotzki, 1997) by one or several orders of magnitude. Such a computational gain is good news as ILP systems routinely perform thousands of subsumption tests. This result is obtained by combining a specific re-description of  $\theta$ -subsumption problems into CSPs and data structures (signatures), and well-known CSP algorithms (dynamic variable ordering, forward checking and arc consistency).

Principled experiments are exploited to build a control layer on the top of *Django*, termed *Meta-Django*, which achieves the selection of the algorithm most suited to the problem at hand.

Further work is concerned with extending *Django*, using more sophisticated CSP algorithms (e.g. replacing the AC3 arc consistency with worst-case cubic complexity, by an AC6, AC7 or AC8, with worst-case quadratic complexity; adding a Conflict Based Jumping into the backtracking mechanism), or adapting other CSP algorithms (e.g. path consistency (Augier, 2000)). The extensions will be documented with their competence map and the *Meta-Django* algorithm will be extended accordingly.

Along the same lines, the use of Forward Checking on the primal and dual formulations of a  $\theta$ -subsumption problems will be investigated in more depth; as shown by Chen (2000), the FC efficiency might vary by an exponential factor between both formulations, depending on the CS problem at hand.

This work also opens several perspectives of research in ILP.

The CSP approach can be used to speed up other ILP operators than  $\theta$ -subsumption, such as for instance  $\theta$ -reduction (Hirata, 2003) or propositionalization (Alphonse & Rouveirol, 2000).  $\theta$ -reduction is used as post processing for the least general generalization (lgg)

operator defined by Plotkin (1970); it reduces the lgg output into a minimal unique (up to variable renaming) clause (Nienhuys-Cheng & de Wolf, 1997).  $\theta$ -reduction algorithms currently proceed by removing one literal at a time and checking whether the initial clause still subsumes the resulting one. By considering the CSP problem defined as the  $\theta$ -subsumption of the lgg by itself, one might exploit the indeterminacies of the  $\theta$ -subsumption solutions to detect groups of equivalent literals. Also, the decomposition of the lgg can be exploited to restrict the search space.

A second perspective concerns the choice of a representation for a particular domain application, and the assessment of hypothesis and example languages. As known by all ILPers, the predicate structure (arity, links) strongly influences learning performances. Better predicates can be introduced dynamically through refinement operators, as for instance in Lachiche and Flach (2003) and Peña-Castillo and Wrobel (2002). But the adequacy of a hypothesis language can also be experimentally assessed through systematic  $\theta$ -subsumption tests, prior to learning. By testing random clauses against the training examples, the phase transition region can be localized. This region likely includes the target clauses: the target clauses hardly belong to the YES and NO regions as they separate the training examples. Therefore the localization of the PT region indicates the complexity of the learning task, i.e. the number of variables and literals in the target clauses within the hypothesis language. This complexity estimate, which also reflects the characteristics of the training examples, can be used to compare several hypothesis languages.

### Acknowledgments

We gratefully acknowledge Lorenza Saitta, Attilio Giordana and Marco Botta, for many discussions about Phase transitions and its consequences on ILP. We also thank Tobias Scheffer, who kindly gave us his implementation of graph contexts and determinate matching.

Let Sébastien Augier and Fabien Torre be particularly thanked for their support and their many critics during the elaboration of *Django*.

We last thank the anonymous referees for their remarks and suggestions, that significantly helped us to improve the paper.

### Notes

1. The *Django* system is available at URL <http://www.lri.fr/~malobert/>.
2. Note that the determinate matching heuristic in the primal  $\theta$ -subsumption problem (Kietz and Lübke, 1994) corresponds to a particular case of the minimal domain heuristic applied on the dual CSP.
3. Incidentally, our last experiments show that removing the arc consistency step is beneficial in the YES region.

### References

- Alphonse, E., & Rouveirol, C. (2000). Lazy propositionalisation for relational learning. In W. Horn (Ed.), *Proceedings of the 14th European Conference on Artificial Intelligence, ECAI'00* (pp. 256–260). IOS Press.

- Agrawal, R., & Srikant, R. (1994). Fast algorithms for mining association rules. In J. B. Bocca, M. Jarke, & C. Zaniolo (Eds.), *Proceedings of the 20th Int. Conf. Very Large Data Bases, VLDB* (pp. 487–499). Morgan Kaufmann.
- Augier, S. (2000). Apprentissage Supervisé relationnel par algorithmes d'évolution. PhD thesis, Université de Paris 11 – Orsay.
- Blockeel, H., Dehaspe, L., Demoen, B., Janssens, G., Ramon, J., & Vandecasteele, H. (2000). Executing query packs in ILP. In *Proceedings of the 10th International Conference on Inductive Logic Programming, LNAI 1866* (pp. 60–77). Berlin: Springer-Verlag.
- Brazdil, P., Gama, J., & Henery, B. (1994). Characterizing the applicability of classification algorithms using meta-level learning. In *Proceedings of the 7th European Conference on Machine Learning* (pp. 83–102). Lecture Notes in Artificial Intelligence, Springer.
- Botta, M., Giordana, A., Saitta, L., & Sebag, M. (2003). Relational learning as search in a critical region. *Journal of Machine Learning Research*, 4, 431–463.
- Bensusan, H., & Kalousis, A. (2001). Estimating the predictive accuracy of a classifier. In P. Flach, & L. De Raedt (Eds.), *Proceedings of the 12th European conference on machine learning* (pp. 25–36). Springer Verlag.
- Blake, C., Keogh, E., & Merz, C.J. (1998). *UCI Repository of machine learning databases*. <http://www.ics.uci.edu/~mlern/MLRepository.html>.
- Bessière, C., & Régim, J.-C. (1996). Mac and combined heuristics: Two reasons to forsake FC (and CBJ ?) on hard problems. In *Proceedings of the 2nd Int. Conf. on Principles and Practice of Constraint Programming* (pp. 61–75).
- Blockeel, H., & De Raedt, L. (1998). Top-down induction of first order logical decision trees. *Artificial Intelligence*, 101:1/2, 285–297.
- Chen, X. (2000). A theoretical comparison of selected CSP solving and modeling techniques. PhD thesis, University of Alberta, Canada.
- Cheeseman, P., Kanefsky, B., & Taylor, W.M. (1991). Where the really hard problems are. In *Proceedings of Int. Joint Conf. on Artificial Intelligence, IJCAI'91* (pp. 331–337).
- Costa, V., Srinivasan, A., Camacho, R., Blockeel, H., Demoen, B., Janssens, G., Struyf, J., Vandecasteele, H., & Van Laer, W. (2004). Query transformations for improving the efficiency of ILP systems. *Journal of Machine Learning Research*, 4, 465–491.
- Dzeroski, S., & Lavrac, N. (Eds.). (2001). *Relational data mining*. Springer Verlag.
- Dolsak, D., & Muggleton, S. (1991). The application of ILP to finite element mesh design. In S. Muggleton (Ed.), *Proceedings of the First International Workshop on Inductive Logic Programming* (pp. 453–472).
- Dechter, R., & Pearl, J. (1989). Tree clustering for constraint networks. *Artificial Intelligence*, 38, 353–366.
- Dehaspe, L., & De Raedt, L. (1997). Mining association rules in multiple relations. In *Proceedings of the 7th International Workshop on Inductive Logic Programming, LNCS 1297* (pp. 125–132). Berlin: Springer-Verlag.
- Esposito, F., Fanizzi, N., Ferilli, S., & Semeraro, G. (2000). Ideal theory refinement under object identity. In P. Langley (Ed.), *Proceedings of the 17th International Conference on Machine Learning* (pp. 263–270). Morgan Kaufmann.
- Fürnkranz, J., & Petrak, J. (2001). An evaluation of landmarking variants. Technical report, Austrian Research Institute for Artificial Intelligence.
- Gomes, C.P., Selman, B., Crato, N., & Kautz, H.A. (2000). Heavy-tailed phenomena in satisfiability and constraint satisfaction problems. *Journal of Automated Reasoning*, 24:1/2, 67–100.
- Gyssens, M., Jeavons, P.G., & Cohen, D.A. (1994). Decomposing constraint satisfaction problems using database techniques. *Artificial Intelligence*, 66, 57–89.
- Gottlob, G., & Leitsch, A. (1985). On the efficiency of subsumption algorithms. *Journal of the Association for Computing Machinery*, 32:2, 280–295.
- Gottlob, G., Leone, N., & Scarcello, F. (2000). A comparison of structural CSP decomposition methods. *Artificial Intelligence*, 124:2, 243–282.
- Gent, I.P., MacIntyre, E., Prosser, P., & Walsh, T. (1996). The constrainedness of search. In *AAAI/IAAI*, vol. 1 (pp. 246–252).
- Giordana, A., & Saitta, L. (2000). Phase transitions in relational learning. *Machine Learning*, 2, 217–251.
- Hogg, T., Huberman, B.A., & Williams, C.P. (Eds.). (1996). *Artificial intelligence: Special issue on frontiers in problem solving: Phase transitions and complexity*, vol. 81(1/2), Elsevier.

- Hirata, K. (2003). On condensation of a clause. In T. Horvath, & A. Yamamoto (Eds.), *Proceedings of the 15th Int. Conf. on Inductive Logic Programming*, LNAI (pp. 164–179). Springer-Verlag.
- Hogg, T. (1996). Refining the phase transition in combinatorial search. *Artificial Intelligence*, 81, 127–154.
- Kalousis, A. (2002). Algorithm Selection via Meta-Learning. PhD thesis, Université de Genève, Centre Universitaire d'Informatique.
- Kietz, J.-U., & Lübke, M. (1994). An efficient subsumption algorithm for inductive logic programming. In W. Cohen, & H. Hirsh (Eds.), *Proceedings of the 11th International Conference on Machine Learning* (pp. 130–137). Morgan Kaufmann.
- Kramer, S., Lavrac, N., & Flach, P. (2001). Propositionalization approaches to relational data mining. In S. Dzeroski, & N. Lavrac (Eds.), *Relational data mining* (pp. 262–291). Springer-Verlag.
- King, R.D., Srinivasan, A., & Sternberg, M.J.E. (1995). Relating chemical activity to structure: An examination of ILP successes. *New Gen. Comput.*, 13.
- Lachiche, N. & Flach, P.A. (2003). A true first-order bayesian classifier. In S. Matwin (Ed.), *Proceedings of the 14th Int. Conf. on Inductive Logic Programming*, LNAI (pp. 133–148). Springer.
- Lloyd, J. (2003). *Logic for learning: Learning comprehensive theories from structured data*. Springer-Verlag.
- Mackworth, A.K. (1977). Consistency in networks of relations. *Artificial Intelligence*, 8, 99–118.
- Di Mauro, N., Basile, T.M.A., Ferilli, S., Esposito, F., & Fanizzi, N. (2003). An exhaustive matching procedure for the improvement of learning efficiency. In T. Horvath, & A. Yamamoto (Eds.), *Proceedings of the 15th Int. Conf. on Inductive Logic Programming*, LNAI (pp. 112–129). Springer-Verlag.
- Muggleton, S., & De Raedt, L. (1994). Inductive logic programming: Theory and methods. *Journal of Logic Programming*, 19, 629–679.
- Maloberti, J., & Sebag, M. (2001). Theta-subsumption in a constraint satisfaction perspective. In C. Rouveirol, & M. Sebag (Eds.), *Proceedings of the 13th Inductive Logic Programming*, LNAI 2157 (pp. 164–178). Springer-Verlag.
- Maloberti, J., & Suzuki, E. (2003). Improving efficiency of frequent query discovery by eliminating non-relevant candidates. In *Proceedings of the 6th International Conference on Discovery Science (DS 2003)*, LNAI 2843 (pp. 219–231). Berlin, Heidelberg, New York, Springer-Verlag.
- Muggleton, S. (1995). Inverse entailment and PROLOG. *New Gen. Comput.*, 13, 245–286.
- Nienhuys-Cheng, S., & de Wolf, R. (1997). *Foundations of inductive logic programming*. Springer-Verlag.
- Nédellec, C., Rouveirol, C., Adé, H., Bergadano, F., & Tausend, B. (1996). Declarative bias in ILP. In L. de Raedt (Ed.), *Advances in ILP* (pp. 82–103). IOS Press.
- Pfahring, B., Bensusan, H., & Giraud-Carrier, C. (2000). Meta-learning by landmarking various learning algorithms. In *Proceedings of the 17th International Conference on Machine Learning* (pp. 743–750). Morgan Kaufmann.
- Peña-Castillo, L., & Wrobel, S. (2002). Macro-operators in multirelational learning: A search-space reduction technique. In H. Mannila et al. (Eds.), *Proceedings of the 13th European Conference on Machine Learning*, LNCS 2430 (pp. 357–368). Berlin: Springer-Verlag.
- Plotkin, G. (1970). A note on inductive generalization. In B. Meltzer, & D. Michie (Eds.), *Machine Intelligence*, vol. 5 (pp. 153–163). Edinburgh University Press.
- Prosser, P. (1993). Hybrid algorithms for the constraint satisfaction problem. *Computational Intelligence*, 9, 268–299.
- Quinlan, J.R. (1990). Learning logical definitions from relations. *Machine Learning*, 5:3, 239–266.
- De Raedt, L. (1997). Logical settings for concept learning. *Artificial Intelligence*, 95:1, 187–201.
- De Raedt, L., & Bruynooghe, M. (1993). A theory of clausal discovery. In *Proceedings of Int. Joint Conf. on Artificial Intelligence, IJCAI'93* (pp. 1058–1063). Morgan Kaufmann.
- Read, R.C., & Corneil, D.G. (1977). The graph isomorph disease. *Journal of Graph Theory*, 1, 339–363.
- Robinson, J. (1965). A machine-oriented logic based on the resolution principle. *Journal of the ACM*, 12:1, 23–41.
- Scheffer, T., Herbrich, R., & Wyszotki, F. (1997). Efficient  $\theta$ -subsumption based on graph algorithms. In S. Muggleton (Ed.), *Proceedings Int. Workshop on Inductive Logic Programming* (pp. 212–228). Springer-Verlag.
- Sebag, M., & Rouveirol, C. (2000). Resource-bounded relational reasoning: Induction and deduction through stochastic matching. *Machine Learning*, 38, 41–62.

- Srinivasan, A. (1997). The predictive toxicology evaluation challenge. In *Proceedings of Int. Joint Conf. on Artificial Intelligence, IJCAI'97* (pp. 4–8). Morgan Kaufmann.
- Srinivasan, A. (2001). The Aleph manual, <http://web.comlab.ox.ac.uk/oucl/research/areas/machlearn/Aleph/>.
- Tsang, E. (1993). *Foundations of constraint satisfaction*. Academic Press.

Received June 18, 2002

Revised October 6, 2003

Accepted January 6, 2004

Final manuscript January 6, 2004