

Towards Cross-Fertilization Between Data Mining and Constraints

Lakhdar Saïs

CRIL - CNRS UMR 8188
Université d'Artois, France

*Joint work with
Said Jabbour, Badran Raddaoui, Yakoub Salhi, Karim Tabia
Takeaki Uno*

<http://www.cril.univ-artois.fr/decMining/>

Journées GT CAVIAR, Paris, november 24th 2017

Data mining

Extracting **useful knowledge** from data

Mainly driven by **real-life applications** including biology (e.g. gene expression data), business intelligence (e.g. market basket), Web (e.g. XML data), ...

At the crossroad of many disciplines:

- ▶ Databases,
- ▶ Artificial Intelligence,
- ▶ Statistics
- ▶ and data analysis, combinatorics, algorithmic, ...

Data mining

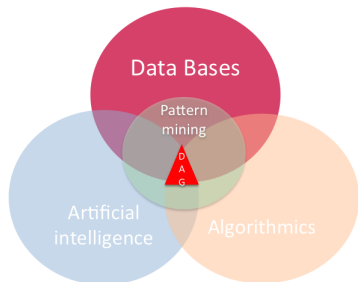
Extracting **useful knowledge** from data

Several research issues:

- ▶ **Pattern mining**
extracting regularities
- ▶ **Clustering & community detection**
extracting meaningful groups
- ▶ Machine learning
extracting predictive models
- ▶ Recommender systems
extracting preferences
- ▶ ...

ANR Défis - DAG (2009 - 2013)

Declarative Approaches for Enumerating Interesting Patterns



<http://liris.cnrs.fr/dag/>

- ▶ CRIL, University d'Artois, Lens
- ▶ LIMOS, University of Blaise Pascal, Clermont-Ferrand
- ▶ LIRIS, University Claude Bernard, Lyon 1, Lyon

Outline

Part I: Declarative approaches for pattern mining problems

Sequence, Itemset and association rules mining

Part II: Data mining \leftarrow AI

Symmetries in Itemset Mining

Part III: Data mining \rightarrow AI

Mining-based Compression Approach of Propositional
Formulae

Conclusion & Perspectives

Boolean Satisfiability Problem (SAT)

- ▶ A conjunction of clauses:

$$\overbrace{(x_1 \vee \cdots \vee x_l)}^{\text{clause}} \wedge (y_1 \vee \cdots \vee y_m) \wedge (z_1 \vee \cdots \vee z_n) \cdots$$

- ▶ Clause: a disjunction of literals (x , $\neg x$)

- ▶ Example :

$$\Phi = \overbrace{(p \vee \neg q \vee \neg r)}^1 \wedge \overbrace{(p \vee \neg q \vee s)}^1 \wedge \overbrace{p}^1 \wedge \overbrace{(r \vee \neg s)}^1$$

horn *unary* *binary*

$$\mathcal{M}(p) = 1 \text{ and } \mathcal{M}(r) = 1 \text{ (Model)}$$

Satisfiability: $\exists \mathcal{M}, \mathcal{M}(\Phi) = 1$ (NP-complete [Cook 71])

Boolean Satisfiability Problem (SAT)

- ▶ Spectacular progress → Modern SAT solvers
 - ▶ application instances with millions of variables and clauses
- ▶ Many applications
 - ▶ Formal Verification
 - ▶ Planning
 - ▶ Bioinformatics
 - ▶ Cryptography
 - ▶ ...
- ▶ CRIL Projects
 - ▶ Microsoft Research Cambridge (UK): 2008-2012
- ▶ CRIL Solvers : Glucose (Sequential), ManySAT (Parallel)
- ▶ Books:
 - ▶ Lakhdar Sais (eds), *Problème SAT : Progrès et Défis*, Hermes Publishing Ltd, pp.352, 2008
 - ▶ Youssef Hamadi and Lakhdar Sais (eds), **Handbook of Parallel Constraint Reasoning**, Springer, February 2018

Part I : SAT based approach for Sequences mining

- ▶ Alphabet: a set Σ
- ▶ Wildcard (or Joker): $\circ \notin \Sigma$
- ▶ Sequence S : word $S_1 S_2 \dots S_n$ in Σ^*
- ▶ Pattern P : word $P_1 P_2 \dots P_m$ in $(\Sigma \cup \{\circ\})^*$
 - ▶ $P_1 \neq \circ$ and $P_m \neq \circ$
 - ▶ Sequences are patterns

abbac, ab \circ c $\circ\circ$ d, ~~ab \circ c~~, ~~ab \circ e \circ~~

Frequent Patterns in a Sequence

Let $P = P_1 P_2 \dots P_m$ and $P' = P'_1 P'_2 \dots P'_n$

▶ $P \subseteq_p P'$ if $\forall i \in \{1, \dots, m\}$:

▶ either $P_i = P'_{p+i-1}$

▶ or $P_i = \circ$

▶ $P \subseteq P'$ if $\exists p$ st. $P \subseteq_p P'$

▶ $L_S(P) = \{p \mid P \subseteq_p S\}$

Example

$a \circ b \subseteq_2 \underline{aaabbaabab}$ $L_{aaabbaabab}(a \circ b) = \{2, 3, 6\}$

$a \circ \circ b \subseteq_1 \underline{a \circ ab \circ b}$ $a \circ \circ b \subseteq_3 a \circ \underline{ab \circ b}$

$a \circ bb \not\subseteq a \circ \circ b$

Definition (Finding frequent patterns in a sequence)

Input: a sequence S and an integer λ

Output: all patterns P st. $|L_S(P)| \geq \lambda$

Representing the Searched Pattern as Boolean Variables

Pattern $P = P_1P_2\dots P_m$

For each position i : $1 \leq i \leq m$:

- ▶ For each character $a \in \Sigma \cup \{o\}$
 - ▶ variable p_i^a is true iff $P_i = a$

$$\neg p_1^o \quad \wedge \quad \bigwedge_{i=1}^m \bigvee_{a \in \Sigma \cup \{o\}} p_i^a \quad \wedge \quad \bigwedge_{i=1}^m \bigwedge_{a,b \in \Sigma \cup \{o\}, a \neq b} (\neg p_i^a \vee \neg p_i^b) \quad (1)$$

Set trailing $p_{m'+1}^o \dots p_m^o$ to true to express pattern of size $m' < m$

Location and Support

The pattern P is found at position k in $S = S_1 \dots S_n$:

$$loc(k, P, S) = \bigwedge_{i=1}^m (p_i^\circ \vee p_i^{S_{i+k-1}})$$

Assuming $S_{i+k-1} = \circ$ when $i + k - 1 > n$

New variables $t_1 \dots t_n$, encoding $L_S(P)$

- ▶ t_k is true iff $P \subseteq_k S$

$$supp(P, S) = \bigwedge_{i=1}^n (t_k \Leftrightarrow loc(k, P, S)) \quad (2)$$

Example

Sequence:

a a a b b a a b a b

Pattern (max size 6):

a ○ *b* ○ ○ ○

Example

Sequence:

a a a b b a a b a b

Pattern (max size 6):

<i>a</i>	○	<i>b</i>	○	○	○
<i>p₁^a</i>	<i>p₂^a</i>	<i>p₃^a</i>	<i>p₄^a</i>	<i>p₅^a</i>	<i>p₆^a</i>
<i>p₁^b</i>	<i>p₂^b</i>	<i>p₃^b</i>	<i>p₄^b</i>	<i>p₅^b</i>	<i>p₆^b</i>
<i>p₁[○]</i>	<i>p₂[○]</i>	<i>p₃[○]</i>	<i>p₄[○]</i>	<i>p₅[○]</i>	<i>p₆[○]</i>

true, false

Example

Sequence:

<i>a</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>
<i>t</i> ₁	<i>t</i> ₂	<i>t</i> ₃	<i>t</i> ₄	<i>t</i> ₅	<i>t</i> ₆	<i>t</i> ₇	<i>t</i> ₈	<i>t</i> ₉	<i>t</i> ₁₀

$$L_S(P) = \{2, 3, 6\}$$

Pattern (max size 6):

<i>a</i>	○	<i>b</i>	○	○	○
<i>p</i> ₁ ^{<i>a</i>}	<i>p</i> ₂ ^{<i>a</i>}	<i>p</i> ₃ ^{<i>a</i>}	<i>p</i> ₄ ^{<i>a</i>}	<i>p</i> ₅ ^{<i>a</i>}	<i>p</i> ₆ ^{<i>a</i>}
<i>p</i> ₁ ^{<i>b</i>}	<i>p</i> ₂ ^{<i>b</i>}	<i>p</i> ₃ ^{<i>b</i>}	<i>p</i> ₄ ^{<i>b</i>}	<i>p</i> ₅ ^{<i>b</i>}	<i>p</i> ₆ ^{<i>b</i>}
<i>p</i> ₁ [○]	<i>p</i> ₂ [○]	<i>p</i> ₃ [○]	<i>p</i> ₄ [○]	<i>p</i> ₅ [○]	<i>p</i> ₆ [○]

true, false

Frequency constraint

$$\text{freq}(P, S, \lambda) = \sum_{i=1}^n t_k \geq \lambda \quad (3)$$

Several possible encodings of the boolean cardinality constraint:

- ▶ Transformation of 0/1 linear inequalities to CNF [Warners 1996]
- ▶ Cardinality networks [Asín et al. 2011]
- ▶ BDD encoding [Bailleux et al. 2003]

Polynomial Encoding of $\sum_{j=1}^n x_j \geq \lambda$ to CNF

$$\bigwedge_{k=1}^{\lambda} (\neg p_{ki} \vee x_i), \quad i = 1, \dots, n \quad (4)$$

$$\bigvee_{i=1}^n p_{ki}, \quad k = 1, \dots, \lambda \quad (5)$$

$$\bigwedge_{1 \leq k < k' \leq \lambda} (\neg p_{ki} \vee \neg p_{k'i}), \quad i = 1, \dots, n \quad (6)$$

(5) and (6) encode the pigeon hole problem PHP_n^λ

- ▶ p_{ki} expresses that pigeon k is in hole i
- ▶ x_i is true if the hole i contains one of the pigeons k for $k = 1, \dots, \lambda$

Complexity $O(\lambda \times n)$ vars and $O(n \times \lambda^2)$ clauses

With Symmetry breaking $\Rightarrow O(\lambda \times (n - \lambda))$ vars and clauses

Part I: SAT based approach for Itemsets Mining

- ▶ Transactions database \mathcal{D} over a set of items

$$\mathcal{I} = \{Camus, Djaout, Djebbar, Kateb, \dots, Mimouni\}$$

$T_{id}(\mathcal{D})$	itemset
000	<i>Djebbar, Djaout, Dib</i>
001	<i>Feraoun, Mimouni, Kateb</i>
002	<i>Djebbar, Dib</i>
003	<i>Camus, Mimouni, Kateb</i>
004	<i>Fanon, Haddad</i>
005	<i>Mimouni, Mammeri</i>

- ▶ Support: $\mathcal{S}(\{Mimouni, Kateb\}, \mathcal{D}) = |\{001, 003\}| = 2$

Frequent Itemset Mining problem

Compute $FIM(\mathcal{D}, \lambda) = \{I \subseteq \mathcal{I} \mid \mathcal{S}(I, \mathcal{D}) \geq \lambda\}$

Example

$$FIM(\mathcal{D}, 2) =$$

$\{\{Mimouni\}, \{Kateb\}, \{Mimouni, Kateb\}, \{Djebbar\},$
 $\{Dib\}, \{Djebbar, Dib\}\}$

Condensed Representations of Frequent Itemsets

Maximal frequent

$$\text{Max}(\mathcal{D}, \lambda) = \{I \in \text{FIM}(\mathcal{D}, \lambda) \mid \forall J \supset I, J \notin \text{FIM}(\mathcal{D}, \lambda)\}$$

Closed frequent itemsets

$$\text{CI}(\mathcal{D}, \lambda) = \{I \in \text{FIM}(\mathcal{D}, \lambda) \mid \forall J \supset I, \mathcal{S}(J, \mathcal{D}) < \mathcal{S}(I, \mathcal{D})\}$$

Example

$$\text{Max}(\mathcal{D}, 2) = \{\{Djebar, Dib\}, \{Mimouni, Kateb\}\}$$

$$\text{CI}(\mathcal{D}, 2) = \{\{Mimouni\}, \{Djebar, Dib\}, \{Mimouni, Kateb\}\}$$

Problem Statement

Mining Frequent Closed itemsets \mathcal{FCIM}_λ

- ▶ **Input:** $\mathcal{D} = \{(0, t_0), \dots, (n-1, t_{n-1})\}$ a transaction database over a set of items \mathcal{I} . λ a minimum support threshold.
- ▶ **Output:** all frequent closed itemsets

SAT-based Encoding for \mathcal{FCIM}_λ

- ▶ Associate to each item $a \in \mathcal{I}$ a boolean variable p_a .
 - ▶ Such boolean variables encode the candidate itemset $I \subseteq \mathcal{I}$, i.e., $p_a = \text{true}$ **iff** $a \in I$.
- ▶ $\forall i \in \{0, \dots, n-1\}$, associate to the i -th transaction a Boolean variable b_i .

SAT-based Encoding for \mathcal{FCIM}_λ

A constraint to capture all the transactions where the candidate itemset does not appear:

$$\bigwedge_{i=0}^{n-1} (b_i \leftrightarrow \bigvee_{a \in \mathcal{I} \setminus t_i} p_a) \quad (7)$$

A constraint to force the candidate itemset to be **closed**:

$$\bigwedge_{a \in \mathcal{I}} \left(\bigwedge_{i=0}^{n-1} \bar{b}_i \rightarrow a \in t_i \right) \rightarrow p_a \quad (8)$$

A constraint to consider only the frequent itemsets:

$$\sum_{i \in 0 \dots n-1} \bar{b}_i \geq \lambda \quad (9)$$

Note: for association rules and variants see our [IJCAI'2016, PAKDD'2017] papers

Outline

Part I: Declarative approaches for pattern mining problems

Sequence, Itemset and association rules mining

Part II: Data mining \leftarrow AI

Symmetries in Itemset Mining

Part III: Data mining \rightarrow AI

Mining-based Compression Approach of Propositional
Formulae

Conclusion & Perspectives

Itemset mining

Output of huge size

- ▶ Difficult to retrieve useful information.

- ▶ Reducing the size of the output is crucial for practical data mining
 - ▶ Search for condensed representations by exploiting the structure of the itemsets data
(e.g. closed, maximal, discriminative itemset patterns, etc.)

Symmetries

- ▶ Fundamental concept (structural knowledge) in Computer Science, Mathematics, Physics and many other domains.
 - ▶ Many human artifacts (e.g. classroom in a university, aircraft seats, circuit patterns) and entities in nature (e.g. plants, molecules, DNA sequences, atoms) exhibits symmetries.
 - ▶ \Rightarrow Useful for reasoning and understanding more complex entities and systems.

Symmetries in CP and SAT

- ▶ Symmetry resolution proof system [Krishnamurthy 1985]
- ▶ Dynamic symmetry detection and elimination in propositional calculus [Benhamou & Saïs 1992]
- ▶ Symmetry breaking predicates [Crawford 1992]
- ▶ Variable and value symmetries [Puget 1993]
- ▶ Many other contributions [Sakallah 2011, Walsh 2012...]

How to exploit symmetries in itemset mining?

1. by dynamic integration in Apriori-like algorithms for search space pruning.
2. by rewriting the transaction databases in a preprocessing step (items elimination).
 - ▶ → new transaction database + symmetry group.
 - ▶ → condensed representation of the output.

Symmetry in Frequent Itemset Mining

Definition (Transaction Renaming)

A renaming f over $\mathcal{T}_{id}(\mathcal{D})$ is a bijective mapping from $\mathcal{T}_{id}(\mathcal{D})$ to $\mathcal{T}_{id}(\mathcal{D})$.

We can extend a renaming f to \mathcal{D} as follows:

$$f(\mathcal{D}) = \{(f(t_i), I) | (t_i, I) \in \mathcal{D}\}.$$

Definition (Permutation)

A permutation σ over \mathcal{I} is a bijective mapping from \mathcal{I} to \mathcal{I} .

We extend a permutation σ to \mathcal{D} as follows:

$$\sigma(\mathcal{D}) = \{(t_i, \sigma(I)) | (t_i, I) \in \mathcal{D}\} \text{ where } \sigma(I) = \{\sigma(a) | a \in I\}.$$

Symmetry in Frequent Itemset Mining

Each permutation σ can be represented by a set of cycles $c_1 \dots c_n$ where each cycle $c_j = (a_1, \dots, a_k)$ is a list of elements of \mathcal{I} such that $\sigma(a_j) = a_{j+1}$ for $j = 1, \dots, k - 1$, and $\sigma(a_k) = a_1$.

Definition (Symmetry)

A symmetry of \mathcal{D} is a permutation $\sigma \in \mathcal{P}(\mathcal{I})$ such that there exists a transaction renaming f over $\mathcal{T}_{id}(\mathcal{D})$ where $\sigma(\mathcal{D}) = f(\mathcal{D})$ i.e. $f^{-1}(\sigma(\mathcal{D})) = \mathcal{D}$.

Proposition

Let σ a symmetry of \mathcal{D} , λ a minimal support threshold and I an itemset. $I \in \mathcal{FLM}(\mathcal{D}, \lambda)$ iff $\sigma(I) \in \mathcal{FLM}(\mathcal{D}, \lambda)$.

Symmetry in Frequent Itemset Mining

Example

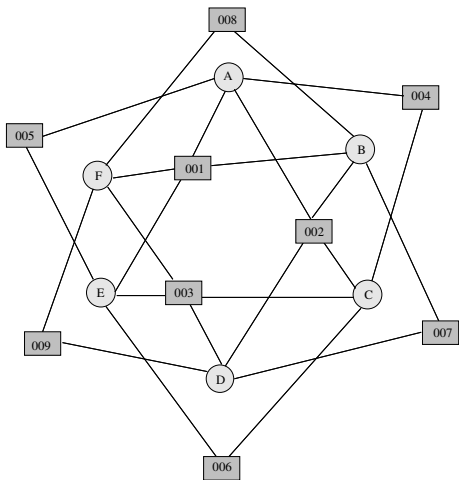
$\sigma = (C,E)(D,F)$ is a symmetry

t_i	itemset
001	A, B, E, F
002	A, B, C, D
003	C, D, E, F
004	A, C,
005	A, E,
006	C, E,
007	B, D,
008	B, F,
009	D, F,

$$f(t_i) = \begin{cases} 001 & \text{if } t_i=002 \\ 002 & \text{if } t_i=001 \\ 003 & \text{if } t_i=003 \\ 004 & \text{if } t_i=005 \\ 005 & \text{if } t_i=004 \\ 006 & \text{if } t_i=006 \\ 007 & \text{if } t_i=008 \\ 008 & \text{if } t_i=007 \\ 009 & \text{if } t_i=009 \end{cases}$$

Symmetry Detection in Transaction Databases

- ▶ Convert the original problem \mathcal{D} into a colored undirected graph \mathcal{G} , where vertices are labeled with colors.
- ▶ Look for the automorphism group of \mathcal{G} .
- ▶ Symmetries of \mathcal{D} are equivalent to the automorphisms of the colored undirected graph \mathcal{G} .
- ▶ Employ a general-purpose graph symmetry tool to uncover the symmetries [Mckay'81, Aloul'03].



t_i	itemset
001	A, B, E, F,
002	A, B, C, D
003	C, D, E, F
004	A, C
005	A, E
006	C, E
007	B, D
008	B, F
009	D, F

Symmetry Pruning

Integration in Apriori-like algorithm

→ proceeds by a level-wise search of the elements of $FIM(\mathcal{D}, \lambda)$.

1. Starts by computing the elements of $FIM(\mathcal{D}, \lambda)$ of size 1.
2. Assuming $FIM(\mathcal{D}, \lambda)$ of size n known, computes a set of candidates of size $n + 1$ so that l is a candidate if and only if all its subsets are in $FIM(\mathcal{D}, \lambda)$.
3. This procedure is iterated until no more candidate is found.

Symmetry-Based Pruning in Apriori-like algos

- ▶ Let \mathcal{D} be a transaction database such that $\mathcal{I}(\mathcal{D}) = \{A, B, C, D\}$ and σ is a symmetry such that $\sigma = (A, D)(B, C)$.
- ▶ Assume that the itemsets $\{A\}$, $\{B\}$, $\{C\}$ and $\{D\}$ are frequent. We also assume that in iteration 2, we find that the itemset $\{A, B\}$ is not frequent.

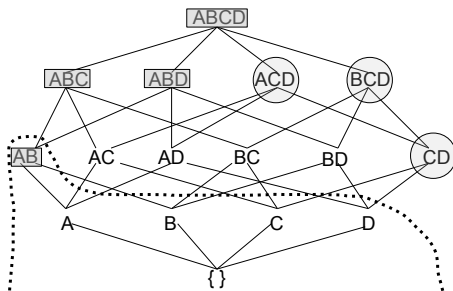


Figure: Symmetry Pruning

Symmetry Breaking

- ▶ Breaking symmetries in a preprocessing step.
 - ▶ Eliminate items from the original transaction database.
 - ▶ The frequent itemsets generated using the new transaction database together with the symmetry group can be used to retrieve the whole set of frequent itemsets of the original

Symmetry Breaking

Let \mathcal{D} a transaction database and $\sigma = (a, b)(c, d)$ a symmetry

$$FIM(\mathcal{D}, \lambda) = \{a, \dots\}, \{b, \dots\}, \{c, \dots\}, \{d, \dots\}, \{a, b, \dots\}, \\ \{a, c, \dots\}, \{a, d, \dots\}, \{b, c, \dots\}, \{b, d, \dots\}, \{c, d, \dots\}$$

$$\{a, \dots\} \rightarrow \{b, \dots\}$$

$$\{b, \dots\} \rightarrow \{a, \dots\}$$

$$\{a, d, \dots\} \rightarrow \{b, c, \dots\}$$

$$\{b, c, \dots\} \rightarrow \{a, d, \dots\}$$

$$\{a, c, \dots\} \rightarrow \{b, d, \dots\}$$

$$\{d, \dots\} \rightarrow \{c, \dots\}$$

$$\{a, b, \dots\} \rightarrow \{a, b, \dots\}$$

$$\{b, d, \dots\} \rightarrow \{a, c, \dots\}$$

Symmetry Breaking

Let \mathcal{D} a transaction database and $\sigma = (a, b)(c, d)$ a symmetry

$$FIM(\mathcal{D}, \lambda) = \{a, \dots\}, \{b, \dots\}, \{c, \dots\}, \{d, \dots\}, \{a, b, \dots\}, \\ \{a, c, \dots\}, \{a, d, \dots\}, \{b, c, \dots\}, \{b, d, \dots\}, \{c, d, \dots\}$$

$\{a, \dots\}$	\rightarrow	$\{b, \dots\}$	$\{b, \dots\}$	\rightarrow	$\{a, \dots\}$
$\{a, d, \dots\}$	\rightarrow	$\{b, c, \dots\}$	$\{b, c, \dots\}$	\rightarrow	$\{a, d, \dots\}$
$\{a, c, \dots\}$	\rightarrow	$\{b, d, \dots\}$	$\{d, \dots\}$	\rightarrow	$\{c, \dots\}$
$\{a, b, \dots\}$	\rightarrow	$\{a, b, \dots\}$	$\{b, d, \dots\}$	\rightarrow	$\{a, c, \dots\}$

Symmetry Breaking

Let \mathcal{D} a transaction database and $\sigma = (a, b)(c, d)$ a symmetry

$$FIM(\mathcal{D}, \lambda) = \{a, \dots\}, \{b, \dots\}, \{c, \dots\}, \{d, \dots\}, \{a, b, \dots\}, \\ \{a, c, \dots\}, \{a, d, \dots\}, \{b, c, \dots\}, \{b, d, \dots\}, \{c, d, \dots\} + \sigma$$

$\{a, \dots\} \rightarrow \{b, \dots\}$	$\{b, \dots\} \rightarrow \{a, \dots\}$
$\{a, d, \dots\} \rightarrow \{b, c, \dots\}$	$\{b, c, \dots\} \rightarrow \{a, d, \dots\}$
$\{a, c, \dots\} \rightarrow \{b, d, \dots\}$	$\{d, \dots\} \rightarrow \{c, \dots\}$
$\{a, b, \dots\} \rightarrow \{a, b, \dots\}$	$\{b, d, \dots\} \rightarrow \{a, c, \dots\}$

- ▶ \Rightarrow **b** can be removed from each $T \in \mathcal{D}$ if $\{a, b\} \not\subseteq T$
- ▶ \Rightarrow **d** can be removed from each $T \in \mathcal{D}$ if $\{a, d\} \not\subseteq T$ and $\{c, d\} \not\subseteq T$

Symmetry Breaking

Proposition

Let \mathcal{D} a transaction database and

$\sigma = (x_1, y_1)(x_2, y_2) \cdots (x_j, y_j) \cdots (x_n, y_n)$ a symmetry

$\Rightarrow y_j$ can be removed from each $T \in \mathcal{D}$ if $\{x_i, y_j\} \not\subseteq T, \forall i \leq j$

Remark

Symmetries can be broken independently

Symmetry Breaking: an example

t_j	itemset
001	A, B, E, F,
002	A, B, C, D
003	C, D, E, F
004	A, C
005	A, E
006	C, E
007	B, D
008	B, F
009	D, F

$$\sigma_1 = (A\ C)(B, D)$$

$$\sigma_2 = (A\ B)(C, D) (E\ F)$$

$$\sigma_3 = (C, E)(D, F)$$

t_j	itemset
001	A, B, E , F
002	A, B, C, D
003	C D E F
004	A, C,
005	A, E ,
006	C E
007	B D
008	B F
009	D F

Table: Itempair-based Symmetry Breaking approach

Symmetry Breaking: an example

t_i	itemset
001	A, B, E, F,
002	A, B, C, D
003	C, D, E, F
004	A, C
005	A, E
006	C, E
007	B, D
008	B, F
009	D, F

$$\sigma_1 = (A\ C)(B, D)$$

$$\sigma_2 = (A\ B)(C, D) (E\ F)$$

$$\sigma_3 = (C, E)(D, F)$$

t_i	itemset
001	A, B,
002	A, B, C, D
003	
004	A, C
005	A
006	
007	
008	
009	

Table: Itempair-based Symmetry Breaking approach

Outline

Part I: Declarative approaches for pattern mining problems

Sequence, Itemset and association rules mining

Part II: Data mining \leftarrow AI

Symmetries in Itemset Mining

Part III: Data mining \rightarrow AI

Mining-based Compression Approach of Propositional
Formulae

Conclusion & Perspectives

Motivation

Growing success obtained in solving real-world SAT problems highlights a real transition to industrial and commercial scale.

- ▶ increasing use of SAT technology to solve new real-world applications (bioinformatics, cryptography, etc.)
- ▶ a rapid growth in the size of the CNF instances encoding real-world problems.

→ **Challenge:** Design of new efficient models for representing and solving SAT instances of very large sizes ("Big" instances).

Modeling in SAT

- ▶ Knowledge representation using CNF formulae

		6	1		2	5		
	3	9				1	4	
				4				
9		2		3		4		1
	8						7	
1		3		6		8		9
				1				
	5	4				9	1	
		7	5		3	2		

8	4	6	1	7	2	5	9	3
1	3	9	6	5	8	1	4	2
5	2	1	3	4	9	7	6	8
9	6	2	8	3	7	4	5	1
4	8	5	9	2	1	3	7	6
1	7	3	4	6	5	8	2	9
2	9	8	7	1	4	6	3	5
3	5	4	2	8	6	9	1	7
6	1	7	5	9	3	2	8	4

- ▶ Example : $n \times n$ Sudoku

- ▶ Associate to each cell, n propositional variables
- ▶ Each cell contains at least one value:

$$\bigwedge_{l=1}^n \bigwedge_{c=1}^n (\bigvee_{v=1}^n p_{(l,c,v)}) \implies n^2 \text{ clauses of size } n$$

- ▶ Leads usually to formulae of huge size

Modeling in SAT: an example from formal verification

Name of the CNF instance : post-cbmc-zfcp-2.8-u2.cnf (BMC)

p cnf **11 483 525** (vars) **32 697 150** (clauses)

1 -3 0

2 -3 0 $x_3 = x_1 \wedge x_2$

1 -2 3 0

: 1million pages later

-11482897 -11483041 -11483523 0

11482897 11483041 -11483523 0

$x_3 \leftrightarrow x_4 \leftrightarrow x_5$

11482897 -11483041 11483523 0

-11482897 11483041 11483523 0

-11483518 -11483524 0

-11483519 -11483524 0

-11483520 -11483524 0

-11483521 -11483524 0

$x_6 = (x_7 \wedge x_8 \wedge x_9 \wedge x_{10} \wedge x_{11} \wedge x_{12})$

-11483522 -11483524 0

-11483523 -11483524 0

11483518 11483519 11483520 11483521 11483522 11483523 11483524 0

-8590303 -11483524 -11483525 0

8590303 11483524 -11483525 0

$x_{13} \leftrightarrow x_{14} \leftrightarrow x_{15}$

8590303 -11483524 11483525 0

-8590303 11483524 11483525 0

-11483525 0

Transformation - Extension principle [G. Tseitin 1965]

- ▶ Introduce new variables to represent truth value of sub-formulae

- ▶ Example : DNF \longrightarrow CNF

$$(x_1 \wedge y_1) \vee (x_2 \wedge y_2) \vee \cdots \vee (x_n \wedge y_n)$$

- ▶ Naïve approach: 2^n clauses and $n \times 2^n$ literals

$$(x_1 \vee \cdots \vee x_{n-1} \vee x_n) \wedge (x_1 \vee \cdots \vee x_{n-1} \vee y_n) \wedge \cdots \wedge (y_1 \vee \cdots \vee y_{n-1} \vee y_n)$$

- ▶ Tseitin approach: $2 \times n + 1$ clauses and $n + 2 \times 2 \times n$ literals

$$(z_1 \vee \cdots \vee z_n) \wedge (\neg z_1 \vee x_1) \wedge (\neg z_1 \vee y_1) \wedge \cdots \wedge (\neg z_n \vee x_n) \wedge (\neg z_n \vee y_n)$$

CNF formula as transactions database

- ▶ Goals : Reduce
 - ▶ *the size of the Formula* : reduce the number of literals using the frequent sets of literals and Tseitin extension principle
 - ▶ *the solving time* (bonus)
- ▶ Items: literals
- ▶ Transactions: clauses > 2

Example

$$(x_1 \vee \neg x_2 \vee \neg x_3) \wedge (x_1 \vee \neg x_2 \vee x_4) \wedge x_1 \wedge (x_3 \vee \neg x_4)$$

itemset
$x_1, \neg x_2, \neg x_3$
$x_1, \neg x_2, x_4$

Reduce the number of literals

- ▶ Introduce new Boolean variables:

$$(x_1 \vee \cdots \vee x_n \vee \alpha_1) \wedge \cdots \wedge (x_1 \vee \cdots \vee x_n \vee \alpha_k)$$

equivalent w.r.t. SAT

\Rightarrow

$$(y \vee \alpha_1) \wedge \cdots \wedge (y \vee \alpha_k) \wedge (\neg y \vee x_1 \vee \cdots \vee x_n)$$

- ▶ $n \times k$ literals substituted by $k + n + 1$ literals
- ▶ Size reduction: $n \times k - (k + n + 1) > 0 \rightarrow k > \frac{n+1}{n-1}$
- ▶ Minimum support threshold: $k \begin{cases} \geq 4 & \text{si } n = 2 \\ \geq 3 & \text{si } n = 3 \\ \geq 2 & \text{otherwise} \end{cases}$

Closed Vs. Maximal

- ▶ Maximal \subseteq Closed : more informations with closed

$$(x_1 \vee \dots \vee x_k \vee \dots \vee x_n \vee \alpha_1) \wedge \dots \wedge (x_1 \vee \dots \vee x_k \vee \dots \vee x_n \vee \alpha_m) \wedge \\ (x_1 \vee \dots \vee x_k \vee \beta_1) \wedge \dots \wedge (x_1 \vee \dots \vee x_k \vee \beta_{m'})$$

We suppose that the set of itemsets are frequent

$$\Rightarrow \text{Max} = \{\{x_1, \dots, x_n\}\}, \text{Clos} = \{\{x_1, \dots, x_k\}, \{x_1, \dots, x_n\}\}$$

- ▶ Use of $\{x_1, \dots, x_n\}$:

$$(y \vee \alpha_1) \wedge \dots \wedge (y \vee \alpha_m) \wedge \\ (x_1 \vee \dots \vee x_k \vee \beta_1) \wedge \dots \wedge (x_1 \vee \dots \vee x_k \vee \beta_{m'}) \wedge \\ (\neg y \vee x_1 \vee \dots \vee x_n)$$

- ▶ Use of $\{x_1, \dots, x_k\}$:

$$(y \vee \alpha_1) \wedge \dots \wedge (y \vee \alpha_m) \wedge \\ (z \vee \beta_1) \wedge \dots \wedge (z \vee \beta_{m'}) \wedge \\ (\neg y \vee z \vee x_{k+1} \vee \dots \vee x_n) \wedge (\neg z \vee x_1 \vee \dots \vee x_k)$$

Weighted Patterns

- ▶ The best:

- ▶ X if

- $|X| \times \mathcal{S}(X) - (\mathcal{S}(X) + |X| + 1) \geq |Y| \times \mathcal{S}(Y) - (\mathcal{S}(Y) + |Y| + 1)$

- ▶ Y otherwise

- ▶ Associates a weight to frequent itemsets:

$$|X| \times \mathcal{S}(X) - (\mathcal{S}(X) + |X| + 1)$$

Overlaps

- ▶ Problem with overlaps:
 - ▶ $\{x_1, x_2, x_3\}$ et $\{x_2, x_3, x_4\}$ two frequents itemsets s.t.
 $\mathcal{S}(\{x_1, x_2, x_3\}) = 3$, $\mathcal{S}(\{x_2, x_3, x_4\}) = 3$ and
 $\mathcal{S}(\{x_1, x_2, x_3, x_4\}) = 2$
 - ▶ Use of $\{x_1, x_2, x_3\} \rightarrow \mathcal{S}(\{x_2, x_3, x_4\}) = 1$

- ▶ **Overlap classes:**

- ▶ X overlaps with Y ($X \sim Y$): $X \cap Y \neq \emptyset$
- ▶ overlaps class (Overlap class): an equivalence class (transitive closure of \sim)

$$Y \in [X] \text{ iff } Y = Y_1 \sim Y_2 \sim \dots \sim Y_k = X$$

- ▶ **Overlap class = Connected Component** on $G = (V, E)$,
 - ▶ V the set of patterns \mathcal{P}
 - ▶ $E = \{\{P_i, P_j\} | P_i \cap P_j \neq \emptyset\}$.
- ▶ **Optimal solution** \rightarrow optimal solution in each overlaps class

Compression as an Optimisation Problem

The compression problem can be formulated as an optimisation problem

Problem : $Comp(\Phi, \mathcal{P})$

- ▶ **Input:** Φ a CNF formula, and \mathcal{P} a set of patterns
- ▶ **Output:** a compressed formula Φ of minimal size using \mathcal{P}

Compression as an Optimisation Problem

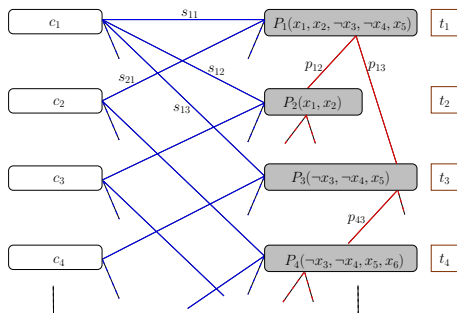
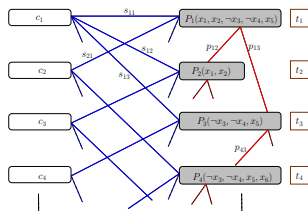


Figure: Compression using location problem

- ▶ $S = \{s_{ij} | P_j \subseteq C_i\}$. If $subst(c_i, P_j)$, then $s_{ij} = 1$; else $s_{ij} = 0$.
- ▶ $\mathcal{T} = \{t_j | 1 \leq j \leq m\}$. If $used(P_j)$, then $t_j = 1$, else $t_j = 0$
- ▶ $\mathcal{P} = \{p_{ij} | P_j \subseteq P_i\}$. If $subst(P_i, P_j)$, then $p_{ij} = 1$; else $p_{ij} = 0$.

Compression as an Optimisation Problem



A **formulation** as 0/1 linear program

$$\text{Max } \sum_{s_{ij} \in \mathcal{S}} (|P_j| - 1) \times s_{ij} + \sum_{p_{ij} \in \mathcal{P}} (|P_j| - 1) \times p_{ij} - (\sum_{j=1}^m (|P_j| + 1) \times t_j)$$

- $s_{ij} - t_j \leq 0 \quad s_{ij} \in \mathcal{S}$
- $p_{ij} - t_j \leq 0, p_{ij} - t_i \leq 0 \quad p_{ij} \in \mathcal{P}$
- $s_{ij} + s_{ik} \leq 1 \quad s_{ij} \in \mathcal{S}, s_{ik} \in \mathcal{S}, P_j \cap P_k \neq \emptyset$
- $s_{ij} \in \{0, 1\} \quad s_{ij} \in \mathcal{S}$
- $t_j \in \{0, 1\} \quad 1 \leq j \leq m$

Greedy Algorithm

- ▶ Search for frequent closed patterns (sub-clauses)
- ▶ Sort the patterns according to their weights (size reduction)
- ▶ Substitution of the patterns following the ordering

Algorithm

Require: A formula ϕ , an overlap class of closed frequent itemsets C

- 1: **while** $C \neq \emptyset$ **do**
- 2: $I \leftarrow C.MostInterestingElement();$
- 3: $\phi.replace(I, y);$
- 4: $\phi.Add(I, y);$
- 5: $C.remove(I);$
- 6: $C.replaceSubset(I, y);$
- 7: $C.removeUninterestingElements();$
- 8: $C.updateSupports();$
- 9: **end while**
- 10: **return** ϕ

Experiments: Industrial SAT instances

Instance	orig.	comp.	% red
1dlx_c.iq57_a	190 Mb	164 Mb	13.68 %
6pipe_6_ooo.*-as.sat03-413	11 Mb	7.7 Mb	30.00 %
9dlx_vliw_at_b_iq6.*-04-347	76 Mb	65 Mb	14.47 %
abb313GPIA-9-c.*.sat04-317	21 Mb	6.9 Mb	67.14 %
E05F18	3.7 Mb	2.2 Mb	40.54 %
eq.atree.braun.11.unsat	120 Kb	72 Kb	40.00 %
eq.atree.braun.12.unsat	144 Kb	88 Kb	38.88 %
k2mul.miter.*-as.sat03-355	1.5 Mb	1.3 Mb	13.33 %
korf-15	1.2 Mb	752 Kb	37.33 %
rbcl_xits_08_UNSAT	1.1 Mb	856 Kb	22.18 %
SAT_dat.k45	3.5 Mb	2.6 Mb	25.71 %
traffic_b.unsat	18 Mb	12 Mb	33.33 %
x1mul.miter.*-as.sat03-359	1.1 Mb	928 Kb	15.63 %
9dlx_vliw_at_b_iq3	19 Mb	15 Mb	21.05 %
9dlx_vliw_at_b_iq4	31 Mb	26 Mb	16.12 %
AProVE07-09	2.8 Mb	2.7 Mb	3.57 %
eq.atree.braun.10.unsat	96 Kb	56 Kb	41.66 %
goldb-heqc-frg1mul	348 Kb	328 Kb	5.74 %
minand128	7.7 Mb	2.6 Mb	66.23 %
ndhf_xits_09_UNSAT	2.6 Mb	2.1 Mb	19.23 %
velev-pipe-o-uns-1.1-6	5.5 Mb	4.4 Mb	20.00 %

Table: Results of Mining4SAT : a general approach

Application: A compact representation of 2-CNF

instance	#cls	#bin	(%) bin
velev-pipe-o-uns-1.1-6	304026	268354	88,26 %
9dlx_vliw_at_b.iq2	542253	500227	92,24 %
1dlx_c.iq57_a	8562505	7567948	88,38 %
7pipe_k	751116	722278	96,16 %
SAT_dat.k100.debugged	670701	523153	78,00 %
BM_FV_2004_rule_batch	445444	339588	76,23 %
sokoban-sequential-p145-*.040-*	1413816	1364160	96,48 %
openstacks-*.p30_1.085-*	1621926	1601145	98,71 %
aaai10-planning-ipc5-*.12-step16	1029036	991140	96,31 %
k2fix_gr_rcs_w8.shuffled	271393	270136	99,53 %
homer17.shuffled	1742	1716	98,50 %
gripper13u.shuffled-as.sat03-395	38965	35984	92,34 %
grid-strips-grid-y-3.045-*	2750755	2695230	97,98 %

Table: Ratio of binary clauses in some SAT instances

Application: A compact representation of 2-CNF

Example

Let us consider the following 2-CNF Φ :

$$\begin{aligned}\Phi = & (x_1 \vee x_2) \wedge (x_1 \vee x_3) \wedge (x_1 \vee x_4) \wedge (x_1 \vee x_5) \quad \wedge \\ & (x_1 \vee x_6) \wedge (x_1 \vee x_7) \wedge (x_2 \vee x_3) \wedge (x_2 \vee x_4) \quad \wedge \\ & (x_2 \vee x_5) \wedge (x_2 \vee x_6) \wedge (x_2 \vee x_7) \wedge (x_3 \vee x_4) \quad \wedge \\ & (x_3 \vee x_6) \wedge (x_3 \vee x_7) \wedge (x_3 \vee x_5) \wedge (x_4 \vee x_5) \quad \wedge \\ & (x_4 \vee x_6) \wedge (x_4 \vee x_7) \wedge (x_5 \vee x_6) \wedge (x_5 \vee x_7) \quad \wedge \\ & (x_6 \vee x_7)\end{aligned}$$

Definition (B-implication)

A *B-implication* is a Boolean formula of the following form :
 $x \vee \beta(x)$ where $\beta(x)$ is a conjunction of literals.

Application: A compact representation of 2-CNF

Using the complete order relation $x_1 \prec \dots \prec x_7$ over \mathcal{L}_Φ
rewrite Φ as set of B-implications $B_{[\vee(\wedge)]}^1(\Phi)$:

$$\begin{aligned} & \{[x_1 \vee (x_2 \wedge x_3 \wedge x_4 \wedge x_5 \wedge x_6 \wedge x_7)], \\ & [x_2 \vee (x_3 \wedge x_4 \wedge x_5 \wedge x_6 \wedge x_7)], \\ & [x_3 \vee (x_4 \wedge x_5 \wedge x_6 \wedge x_7)], \\ & [x_5 \vee (x_6 \wedge x_7)], \\ & [x_6 \vee (x_7)]\} \end{aligned}$$

tid	itemset					
tid_{x_1}	x_2	x_3	x_4	x_5	x_6	x_7
tid_{x_2}		x_3	x_4	x_5	x_6	x_7
tid_{x_3}			x_4	x_5	x_6	x_7
tid_{x_4}				x_5	x_6	x_7
tid_{x_5}					x_6	x_7
tid_{x_6}						x_7

Application: A compact representation of sets of 2-CNF

FIM process on the conjunctive part of $B_{\vee[\wedge]}^1(\Phi)$

Using $\{x_5, x_6, x_7\}$ a 4-frequent itemset, we can rewrite

$B_{\vee[\wedge]}^1(\Phi)$ as:

$$B_{\vee[\wedge]}^2(\Phi) = \{ [x_1 \vee (x_2 \wedge x_3 \wedge y)] , \\ [x_2 \vee (x_3 \wedge x_4 \wedge y)] , \\ [x_3 \vee (x_4 \wedge y)] , \\ [x_5 \vee (x_6 \wedge x_7)] , \\ [x_6 \vee (x_7)] , \\ [\neg y \vee (x_5 \wedge x_6 \wedge x_7)] \}$$

$\text{CNF}(B_{\vee[\wedge]}^2(\Phi)) =$

$$(x_1 \vee x_2) \wedge (x_1 \vee x_3) \wedge (x_1 \vee y) \quad \wedge$$

$$(x_2 \vee x_3) \wedge (x_2 \vee x_4) \wedge (x_2 \vee y) \quad \wedge$$

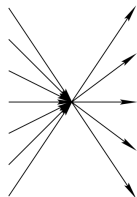
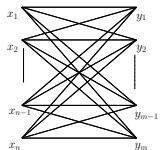
$$(x_3 \vee x_4) \wedge (x_3 \vee y) \quad \wedge$$

$$(x_5 \vee x_6) \wedge (x_5 \vee x_7) \quad \wedge$$

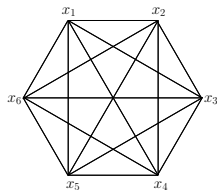
$$(x_6 \vee x_7) \quad \wedge$$

$$(\neg y \vee x_5) \wedge (\neg y \vee x_6) \wedge (\neg y \vee x_7)$$

Two particular cases: bi-cliques and cliques



$n \times m$ clauses $\Rightarrow n + m$ clauses and 1 new variable



$\mathcal{O}(n^2)$ clauses $\Rightarrow \mathcal{O}(n)$ clauses and $\mathcal{O}(n)$ new variables \Rightarrow
 $\sum_{i=1}^n x_i = 2$

More details on bi-cliques

Let $\Phi = [(x_1 \vee y_1) \wedge (x_1 \vee y_2) \wedge \cdots \wedge (x_1 \vee y_m)] \cdots [(x_n \vee y_1) \wedge (x_n \vee y_2) \wedge \cdots \wedge (x_n \vee y_m)]$

- ▶ Using a complete order relation defined by:

$$f(x_i) = i, f(y_j) = n + j.$$

- ▶ $B_{[\vee(\wedge)]}(\Phi)$ corresponds exactly to $\{(x_i \vee [y_1 \wedge y_2 \wedge \cdots \wedge y_m]) \mid 1 \leq i \leq n\}$

- ▶ Using a single closed frequent itemset $\{y_1, y_2, \dots, y_m\}$

$$\Phi' = [\wedge_{1 \leq i \leq n} (x_i \vee z)] \wedge [\wedge_{1 \leq j \leq m} (\neg z \vee y_j)].$$

Experiments: Industrial SAT instances

Instance	orig.	comp.	% red
velev-pipe-o-uns-1.1-6	5.5 Mb	3.2 Mb	41.81 %
9dlx_vliw_at_b.iq2	11 Mb	6 Mb	44.45 %
1dlx_c.iq57_a	190 Mb	124 Mb	34.73 %
7pipe_k	14 Mb	5.4 Mb	61.42 %
SAT_dat.k100.debugged	16 Mb	13 Mb	18.75 %
IBM_FV_2004_rule_batch _2_31_1_SAT_dat.k80.debugged	9.7 Mb	7.5 Mb	22.68 %
sokoban-sequential-p145-*.040-*	24 Mb	14 Mb	41.66 %
openstacks-*.p30_1.085-*	30 Mb	26 Mb	13.33 %
aaai10-planning-ipc5-*.12-step16	17 Mb	12 Mb	29.41 %
k2fix_gr_rcs_w8.shuffled	3.4 Mb	1.7 Mb	50.00 %
homer17.shuffled	20 Kb	16 Kb	20.00 %
gripper13u.shuffled-as.sat03-395	524 Kb	364 Kb	30.35 %
grid-strips-grid-y-3.045-*	52 Mb	42 Mb	19.23 %

Table: Results of Mining4Binary: a 2-CNF approach

Combining Binary and Non Binary Clauses

$$\begin{array}{ccc} x_0 \vee \neg x_4, & x_0 \vee \neg x_5, & x_0 \vee \neg x_6, \\ \neg x_3 \vee \neg x_4, & \neg x_3 \vee \neg x_5, & \neg x_3 \vee \neg x_6, \\ \hline \neg x_0 \vee x_1 & \vee & | x_4 \vee x_5 \vee x_6 |, \\ & & | x_4 \vee x_5 \vee x_6 |, \\ & & | x_4 \vee x_5 \vee x_6 |, \\ & & | x_4 \vee x_5 \vee x_6 | \\ \hline \end{array}$$

Suppose that $(x_4 \vee x_5 \vee x_6)$ is frequent

Combining Binary and Non Binary Clauses

$$(x_0 \vee [\neg x_4 \wedge \neg x_5 \wedge \neg x_6])$$
$$(\neg x_3 \vee [\neg x_4 \wedge \neg x_5 \wedge \neg x_6])$$

$$\begin{array}{rcl} \neg x_0 \vee x_1 & \vee & | x_4 \vee x_5 \vee x_6 |, \\ & & x_3 \vee | x_4 \vee x_5 \vee x_6 |, \\ \neg x_1 \vee x_2 & \vee & | x_4 \vee x_5 \vee x_6 |, \\ \neg x_2 \vee x_3 & \vee & | x_4 \vee x_5 \vee x_6 | \end{array}$$

Combining Binary and Non Binary Clauses

$$\begin{aligned} & (x_0 \vee \neg[x_4 \vee x_5 \vee x_6]) \\ & (\neg x_3 \vee \neg[x_4 \vee x_5 \vee x_6]) \end{aligned}$$

$$\begin{array}{rcc} \neg x_0 \vee x_1 & \vee & | x_4 \vee x_5 \vee x_6 |, \\ & & x_3 \vee | x_4 \vee x_5 \vee x_6 |, \\ \neg x_1 \vee x_2 & \vee & | x_4 \vee x_5 \vee x_6 |, \\ \neg x_2 \vee x_3 & \vee & | x_4 \vee x_5 \vee x_6 | \end{array}$$

Combining Binary and Non Binary Clauses

$$\begin{array}{l} x_0 \vee \neg \mathbf{y} \\ \neg x_3 \vee \neg \mathbf{y} \\ \neg x_0 \vee x_1 \quad \vee \quad \mathbf{y} \\ \quad \quad x_3 \quad \vee \quad \mathbf{y} \\ \neg x_1 \vee x_2 \quad \vee \quad \mathbf{y} \\ \neg x_2 \vee x_3 \quad \vee \quad \mathbf{y} \\ \quad \quad \neg \mathbf{y} \quad \vee \quad x_4 \vee x_5 \vee x_6 \\ \mathbf{y} \vee \neg x_4 \\ \mathbf{y} \vee \neg x_5 \\ \mathbf{y} \vee \neg x_6 \end{array}$$

Experiments

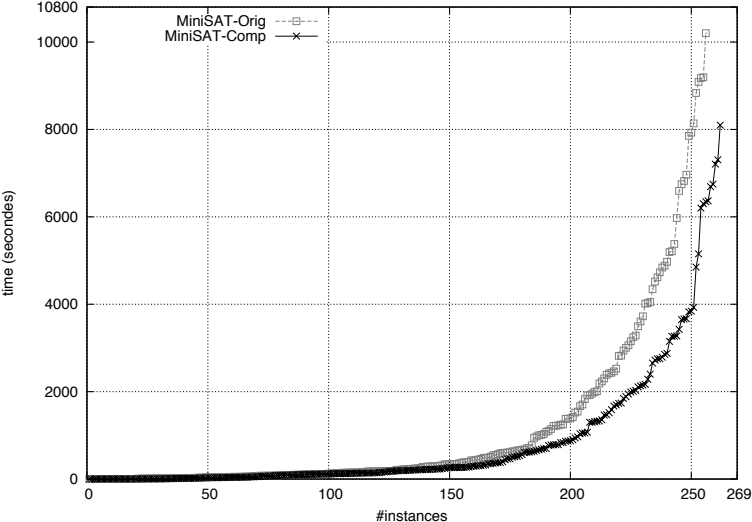


Figure: MiniSAT on SAT instances with and without compression

Application: A compact Graph Representation

For free, we can apply our approach for graphs.

- ▶ 2-CNF \leftrightarrow graphs
- ▶ Adjacency lists \leftrightarrow A set of B-implications
- ▶ $2 \rightarrow [4, 6, 8, 12] \leftrightarrow 2 \vee [4 \wedge 6 \wedge 8 \wedge 12]$

Outline

Part I: Declarative approaches for pattern mining problems

Sequence, Itemset and association rules mining

Part II: Data mining \leftarrow AI

Symmetries in Itemset Mining

Part III: Data mining \rightarrow AI

Mining-based Compression Approach of Propositional
Formulae

Conclusion & Perspectives

Perspectives

- ▶ Cross-fertilization between AI and Data mining
 - ▶ $DM \leftarrow AI$ (e.g. preferences, symmetries, knowledge compilation, etc.)
 - ▶ $DM \rightarrow AI$ (e.g. extracting structural knowledge, compression, etc.)
 - ▶ ...

Thank you for your attention