

Approche basée sur la Relaxation pour la Fouille de Motifs Fermés et Diversifiés

A. Hien^{1*} S. Loudni³ N. Aribi² Y. Lebbah² M. Laghzaoui² A. Ouali¹ A. Zimmermann¹

¹ Normandie Univ., UNICAEN, CNRS – UMR GREYC, France

² University of Oran1, Lab. LITIO, 31000 Oran, Algeria

³ TASC (LS2N-CNRS), IMT Atlantique, FR – 44307 Nantes, France

La fouille interactive de motifs a fait évoluer la fouille de données vers un modèle centré utilisateur [3, 6]. Il s'agit de prendre en compte les préférences de l'utilisateur afin de guider la recherche vers des motifs intéressants pour lui. Cela est rendu possible par l'introduction de mécanismes de *feedback* qui permettent à l'utilisateur de spécifier ses préférences sur les motifs qu'on lui présente [3]. Un aspect important de ce cadre est la nécessité de présenter à l'utilisateur des *résultats diversifiés* [1, 2].

Dans cet article, nous utilisons la programmation par contraintes pour extraire efficacement les motifs fréquents fermés diversifiés. La diversité est contrôlée par un seuil sur l'indice de Jaccard. Nous montrons que cette mesure n'a pas de propriété de monotonie, ce qui rend le processus d'extraction infaisable. Pour y remédier, nous proposons une nouvelle contrainte globale, CLO-SEDDIVERSITY, qui exploite une relaxation anti-monotone de l'indice de Jaccard pour élaguer les motifs non diversifiés. Une seconde relaxation, basée sur une borne supérieure, est exploitée via une nouvelle heuristique de branchement.¹

Notion de diversité entre motifs. Soit \mathcal{I} un ensemble de n items, un motif P est un sous-ensemble non vide de \mathcal{I} . Une base transactionnelle \mathcal{D} est un multi-ensemble de transactions sur \mathcal{I} , où chaque transaction t est un sous-ensemble de \mathcal{I} , i.e., $t \subseteq \mathcal{I}$. Un motif P apparaît dans une transaction t , ssi $P \subseteq t$. La couverture de P dans \mathcal{D} est l'ensemble des transactions qui le supportent : $t(P) = \{t \in \mathcal{D} \mid P \subseteq t\}$. Le support de P dans \mathcal{D} est le cardinal de sa couverture : $sup(P) = |t(P)|$. Un motif P est dit fréquent si son support dépasse un seuil de fréquence minimal θ , $sup(P) \geq \theta$. The clôture d'un motif P , notée $Clos(P)$, est l'ensemble des items communs à toutes les transactions dans $t(P)$: $Clos(P) = \{i \in \mathcal{I} \mid \forall t \in t(P), i \in t\}$. Un motif P est dit fermé ssi $Clos(P) = P$.

L'indice de Jaccard est une mesure de similarité classique sur les ensembles. Nous l'utilisons pour quantifier le chevauchement des couvertures entre deux motifs.

Définition 1 (Indice de Jaccard) Soient deux motifs P et Q , l'indice de Jaccard est le rapport entre le cardinal de l'intersec-

tion et le cardinal de l'union des couvertures des deux motifs : $Jac(P, Q) = \frac{|t(P) \cap t(Q)|}{|t(P) \cup t(Q)|}$.

Un indice de Jaccard plus petit est synonyme d'une faible similarité en termes de couverture entre motifs et peut donc être utilisé comme mesure de diversité entre des paires de motifs.

Définition 2 (Contrainte de Jaccard/Diversité) Soit P et Q deux motifs, Jac la mesure de Jaccard et J_{max} un seuil de diversité, la paire P et Q est dite diversifiée ssi $Jac(P, Q) \leq J_{max}$.

Notre objectif est d'exploiter la contrainte de Jaccard durant la recherche pour élaguer les motifs non-diversifiés. Toutefois, comme l'indique la proposition 1, la contrainte de Jaccard n'est ni monotone ni anti-monotone, ce qui implique un élagage limité lors de la recherche. Pour faire face à ce problème, nous proposons deux relaxations anti-monotones : (i) Une relaxation par la borne inférieure, permettant d'élaguer les motifs non-diversifiés lors de la recherche, (ii) une relaxation par la borne supérieure pour trouver les items menant vers des motifs diversifiés.

Proposition 1 Soit P , Q et P' trois motifs avec $P \subset P'$. $Jac(P, Q)$ peut être plus petit, égal ou supérieur à $Jac(P', Q)$.

Relaxations de la contrainte de Jaccard. L'idée clé de notre approche est d'exploiter une relaxation de la contrainte de Jaccard ayant des propriétés de monotonie afin de les exploiter pour la réduction de l'espace de recherche.

Définition 3 (Formulation du problème) Soit \mathcal{H} l'historique courant contenant les paires de motifs clos fréquents diversifiés, J_{max} un seuil de diversité, LB_J et UB_J les bornes inférieure et supérieure de l'indice de Jaccard, le problème relaxé consiste à extraire les motifs P tels que $\forall H \in \mathcal{H}, LB_J(P, H) \leq J_{max}$. Quand $UB_J(P, H) \leq J_{max}$, pour tout $H \in \mathcal{H}$, alors la contrainte Jaccard est satisfaite.

Nous présentons à présent une borne inférieure et supérieure de l'indice de Jaccard et montrons comment les exploiter. Pour cela, nous introduisons la notion de *couverture propre* d'un motif.

*Papier doctorant : A. Hien¹ est auteur principal.

1. La version originale de ce papier est publiée dans [4].

Définition 4 (Couverture propre) Soit P et Q deux motifs. La couverture propre de P par rapport à Q est définie par $t_Q^{pp}(P) = t(P) \setminus \{t(P) \cap t(Q)\}$.

Proposition 2 (Borne inférieure) Soit un motif $H \in \mathcal{H}$. Soit P un motif rencontré durant la recherche tel que $\text{sup}(P) \geq \theta$, et $t_H^{pp}(P)$ est la couverture propre de P par rapport à H . $LB_J(P, H) = \frac{\theta - |t_H^{pp}(P)|}{|t(P)| + |t(H)| + |t_H^{pp}(P)| - \theta}$ est un minorant de $Jac(P, H)$.

Proposition 3 (Monotonie du LB_J) Soit $H \in \mathcal{H}$ un motif. Pour toute paire de motifs P et Q avec $P \subseteq Q$, nous avons $LB_J(P, H) \leq LB_J(Q, H)$.

La propriété 3 définit une condition nécessaire d'élagage basée sur la monotonie de la borne inférieure : Si $LB_J(P, H) > J_{max}$, alors aucun motif $Q \supseteq P$ ne pourra satisfaire la contrainte de Jaccard (car LB_J est une borne inférieure), rendant la contrainte elle-même *anti-monotone*. Ainsi, le motif Q peut être élagué.

Comme notre relaxation est une sur-estimation des solutions de la contrainte Jaccard, on pourrait alors avoir des motifs P tels que $LB_J(P, H) < J_{max}$ mais $Jac(P, H) > J_{max}$. Pour traiter ce cas, nous définissons une borne supérieure sur l'indice Jaccard pour évaluer la satisfaction de la contrainte, c'est-à-dire les cas avec $UB_J(P, H) \leq J_{max}, \forall H \in \mathcal{H}$.

Proposition 4 (Borne supérieure) Soit un élément H de l'historique \mathcal{H} , et un motif P tel que $\text{sup}(P) \geq \theta$. $UB_J(P, H) = \frac{|t(H) \cap t(P)|}{|t_P^{pp}(H)| + \max\{\theta, |t(H)| \cap |t(P)|\}}$ est un majorant de $Jac(P, H)$.

Proposition 5 (Anti-monotonie de UB_J) Soit $H \in \mathcal{H}$ un motif. Pour toute paire de motifs P et Q avec $P \subseteq Q$, nous avons $UB_J(P, H) \geq UB_J(Q, H)$.

La contrainte globale CLOSEDDIVERSITY. Pour représenter un motif inconnu, nous utilisons un vecteur x de variables booléennes $(x_1, \dots, x_{|\mathcal{I}|})$, où x_i désigne la présence de l'item $i \in \mathcal{I}$ dans le motif. Nous utiliserons les notations suivantes : $x^+ = \{i \in \mathcal{I} \mid \text{dom}(x_i) = \{1\}\}$ les items présents, $x^- = \{i \in \mathcal{I} \mid \text{dom}(x_i) = \{0\}\}$ les items absents, et $x^* = \{i \in \mathcal{I} \mid i \notin x^+ \cup x^-\}$ l'ensemble des items non-assignés.

Définition 5 (CLOSEDDIVERSITY) Soit x une instanciation complète des variables d'items, \mathcal{H} un historique de motifs clos fréquents diversifiés (initialement vide), θ un seuil de fréquence, J_{max} un seuil de diversité et \mathcal{D} une base transactionnelle. La contrainte $CLOSEDDIVERSITY_{\mathcal{D}, \theta}(x, \mathcal{H}, J_{max})$ est satisfaite ssi : (1) $\text{Clos}(Px^+) = x^+$; (2) $\text{sup}(x^+) \geq \theta$; (3) $\forall H \in \mathcal{H}, LB_J(x^+, H) \leq J_{max}$.

Notre contrainte globale permet de mettre à jour \mathcal{H} de manière incrémentale durant la recherche. Le propagateur de $CLOSEDDIVERSITY$ exploite les règles de filtrage de la contrainte $CLOSEDPATTERNS$ [5]. Il utilise également notre relaxation LB pour supprimer les items qui ne peuvent mener à une solution contenant x^+ . Notons par x_{Freq}^- les items filtrés par la règle des items non fréquents et par x_{Div}^- les items filtrés par notre règle sur LB .

Proposition 6 (Filtrage de CLOSED DIVERSITY) Soit \mathcal{H} un historique de motifs clos fréquents diversifiés, x^+ une instanciation partielle sur x , et un item libre $i \in x^*, x^+ \cup \{i\}$ ne peut mener à un motif diversifié dans l'un des deux cas suivants : (1) si $\exists H \in \mathcal{H}$ s.t. $LB_J(x^+ \cup \{i\}, H) > J_{max}$, alors $1 \notin \text{dom}(x_i)$; (2) if $\exists k \in x_{Div}^-$ s.t. $t(x^+ \cup \{i\}) \subseteq t(x^+ \cup \{k\})$, alors $LB_J(x^+ \cup \{i\}, H) > LB_J(x^+ \cup \{k\}, H) > J_{max}$ et $1 \notin \text{dom}(x_i)$.

Exploiter la borne UB_J . On peut avoir $LB_J(x^+, H) \leq J_{max}$ mais $Jac(x^+, H) > J_{max}$. Nous proposons d'exploiter notre relaxation UB pour garantir la satisfaction de la contrainte Jaccard. Pendant la recherche, nous calculons incrémentalement $UB(x^+ \cup \{i\}, H)$ de toute extension de l'affectation partielle x^+ avec un item libre i . Si, $\forall H \in \mathcal{H}, UB_J(x^+ \cup \{j\}, H) \leq J_{max}$, alors la contrainte est satisfaite (i.e. propriété de *temoin*). De plus, grâce à l'anti-monotonie de UB_J (voir Proposition 5), tous les sur-ensembles de $x^+ \cup \{i\}$ satisferont également la contrainte de Jaccard. Pour accélérer la certification des solutions, nous proposons une nouvelle heuristique d'ordre sur les variables qui *sélectionne le premier item libre satisfaisant la propriété témoin*, sinon nous sélectionnons la variable ayant le plus petit support estimé (voir [4] pour plus de détails).

Résultats expérimentaux. Nous avons évalué les performances de notre contrainte globale $CLOSEDDIV$, ainsi que la qualité des résultats (en terme de diversité) sur les jeux de données UCI. Les résultats confirment l'efficacité de la règle de filtrage LB pour l'élagage des motifs non diversifiés. Ces résultats montrent aussi que la borne supérieure se rapproche beaucoup de la valeur de l'indice Jaccard, ce qui signifie que notre borne UB fournit une relaxation serrée. Ceci est révélateur de la qualité des motifs trouvés en termes de diversité et démontre l'intérêt et la force de notre règle de branchement pour obtenir des motifs diversifiés par rapport à la contrainte globale $CLOSEDPATTERNS$.

Dans ce papier, nous avons proposé une approche introduite sous la forme d'une contrainte globale appelée $CLOSEDDIV$. Cette contrainte exploite deux relaxations LB/UB (anti-)monotones, qui permettent d'avoir un filtrage et une règle de branchement efficaces, dynamisant l'ensemble du processus de recherche.

Références

- [1] G. BOSCH, J-F. BOULICAUT, Ch. RAÏSSI et M. KAYTOUE : Anytime discovery of a diverse set of patterns with monte carlo tree search. *Data mining and knowledge discovery*, 32(3):604–650, 2018.
- [2] B. BRINGMANN et A. ZIMMERMANN : The chosen few : On identifying valuable patterns. *In Proceedings of ICDM 2007*, pages 63–72, 2007.
- [3] V. DZYUBA et M. van LEEUWEN : Interactive discovery of interesting subgroup sets. *In International Symposium on Intelligent Data Analysis*, pages 150–161. Springer, 2013.
- [4] A. HIEN, S. LOUDNI, N. ARIBI, Y. LEBBAH, M. LAGHZAOU, A. OUALI et A. ZIMMERMANN : A relaxation-based approach for mining diverse closed patterns. *In Proceedings of ECML PKDD 2020*, volume 12457, pages 36–54. Springer, 2020.
- [5] N. LAZAAR, Y. LEBBAH, S. LOUDNI, M. MAAMAR, V. LEMIERE, C. BESSIERE et P. BOIZUMAUULT : A global constraint for closed frequent pattern mining. *In Proceedings of the 22nd CP*, pages 333–349, 2016.
- [6] M. van LEEUWEN : *Interactive Data Exploration Using Pattern Mining*, pages 169–182. Springer, 2014.

A Relaxation-based Approach for Mining Diverse Closed Patterns

Arnold Hien², Samir Loudni^{2,3}✉, Nouredine Aribi¹, Yahia Lebbah¹, Mohammed Laghzaoui¹, Abdelkader Ouali², and Albrecht Zimmermann²

¹ University of Oran1, Lab. LITIO, 31000 Oran, Algeria

² Normandie Univ., UNICAEN, CNRS – UMR GREYC, France

³ TASC (LS2N-CNRS), IMT Atlantique, FR – 44307 Nantes, France

Abstract. In recent years, pattern mining has moved from a slow-moving repeated three-step process to a much more agile iterative/user-centric mining model. A vital ingredient of this framework is the ability to *quickly* present a set of *diverse* patterns to the user. In this paper, we use constraint programming (well-suited to user-centric mining due to its rich constraint language) to efficiently mine a diverse set of closed patterns. Diversity is controlled through a threshold on the Jaccard similarity of pattern occurrences. We show that the Jaccard measure has no monotonicity property, which prevents usual pruning techniques and makes classical pattern mining unworkable. This is why we propose anti-monotonic lower and upper bound relaxations, which allow effective pruning, with an efficient branching rule, boosting the whole search process. We show experimentally that our approach significantly reduces the number of patterns and is very efficient in terms of running times, particularly on dense data sets.

1 Introduction

The original data analysis model based on pattern mining consists of three steps in a kind of *multi-waterfall cycle*: 1) a user chooses the values of one or several mining parameters, 2) an underlying engine extracts patterns (often taking not inconsiderable time to do so), and 3) the user sifts through a (potentially very large) set of result patterns and interprets them, using their insights to return to the first step and repeat the cycle.

Recently, this approach has been challenged by an increasing focus on *user-centered*, *interactive*, and *anytime* pattern mining [14]. This new paradigm stresses that users should be presented quickly with patterns likely to be interesting to them, and typically affect later iterations of the mining process by giving feedback. A powerful framework for taking a variety of user feedback into account is pattern mining via constraint programming (CP). Much of the current focus in this domain is on user-centered/interactive mining, particularly the ability to elicit and exploit user feedback [9, 14, 18]. An important aspect of requesting such feedback is that the user be quickly presented with *diverse* results. If patterns are too similar to each other, deciding which one to prefer can become challenging, and if they appear in several successive iterations, it eventually becomes a slog. Similarly, a method that produces diverse results but takes a long time to do so, risks that the user checks out of the process. Older work on diversity either post-process patterns derived from the process described above [5, 12, 21], use heuristics [20] or view

it purely from the point of view of speeding up the extraction process [8]. Recent work, on the other hand, pushes diversity constraints into the mining process itself [3, 4]. At the algorithmic level, additional user-specified constraints often require new implementations to filter out the patterns violating or satisfying the user’s constraints, which can be computationally infeasible for large databases.

In the last decade, data mining has been combined with constraint programming to model various data mining problems [2, 6, 13, 19]. The main advantage of CP for pattern mining is its declarativity and flexibility, which include the ability to incorporate new user-specified constraints without the need to modify the underlying system. Moreover, CP allows to define flexible search strategies.⁴ In this paper, we propose to add to the literature on *explicitly* taking the diversity of patterns (in terms of the data instances they describe) into account and to use an exhaustive process to find candidates for inclusion into a result set. To achieve this, we use the widely accepted Jaccard index to compare patterns and formulate a diversity constraint, which has no monotonicity property, implying limited pruning during search. To cope with this problem, we propose two anti-monotonic relaxations: (i) A lower bound relaxation, which allows to prune non-diverse items during search. This is integrated in our constraint programming based approach through a new global constraint taking into account diversity with its filtering algorithms (aka, propagators); (ii) An upper bound relaxation to find items ensuring diversity. This is exploited through a new branching rule, boosting the search process towards diverse patterns. We demonstrate the performance of our proposed method experimentally, comparing to the state-of-the-art in CP-based closed pattern mining.

2 Preliminaries

2.1 Itemset Mining

Let $\mathcal{I} = \{1, \dots, n\}$ be a set of n items, an *itemset* (or pattern) P is a non-empty subset of \mathcal{I} . The language of itemsets corresponds to $\mathcal{L}_{\mathcal{I}} = 2^{\mathcal{I}} \setminus \emptyset$. A transactional dataset \mathcal{D} is a bag (or multiset) of transactions over \mathcal{I} , where each *transaction* t is a subset of \mathcal{I} , i.e., $t \subseteq \mathcal{I}$; $\mathcal{T} = \{1, \dots, m\}$ a set of m transaction indices. An itemset P occurs in a transaction t , iff $P \subseteq t$. The *cover* of P in \mathcal{D} is the set of transactions in which it occurs: $\mathcal{V}_{\mathcal{D}}(P) = \{t \in \mathcal{D} \mid P \subseteq t\}$. The *support* of P in \mathcal{D} is the size of its cover: $sup_{\mathcal{D}}(P) = |\mathcal{V}_{\mathcal{D}}(P)|$. An itemset P is said to be *frequent* when its support exceeds a user-specified minimal threshold θ , $sup_{\mathcal{D}}(P) \geq \theta$. Given $S \subseteq \mathcal{D}$, $items(S)$ is the set of common items belonging to all transactions in S : $items(S) = \{i \in \mathcal{I} \mid \forall t \in S, i \in t\}$. The *closure* of an itemset P , denoted by $Clos(P)$, is the set of common items that belong to all transactions in $\mathcal{V}_{\mathcal{D}}(P)$: $Clos(P) = \{i \in \mathcal{I} \mid \forall t \in \mathcal{V}_{\mathcal{D}}(P), i \in t\}$. An itemset P is said to be *closed* iff $Clos(P) = P$. Constraint-based pattern mining aims at extracting all patterns P of $\mathcal{L}_{\mathcal{I}}$ satisfying a selection predicate c (called *constraint*) which is usually called *theory* [5]: $Th(c)$. A common example is the frequency measure leading to the minimal support constraint, which can be combined with the closure constraint to mine closed frequent itemsets.

⁴ Opposed to more rigid search in classical pattern mining algorithms, which often rely on exploiting the properties of a particular constraint.

Example 1. Figure 1 shows the itemset lattice derived from a toy dataset with five items and 100 transactions. As the figure shows, there exist 26 frequent closed itemsets with $\theta = 7$.

Most constraint-based mining algorithms take advantage of monotonicity which offers pruning conditions to safely discard non-promising patterns from the search space. Several frameworks exploit this principle to mine with a monotone or an anti-monotone constraint. Other classes of constraints have also been considered [15, 16]. However, for constraints that are not anti-monotone, pushing them into the discovery algorithm might lead to less effective pruning phases. Thus, we propose in this paper to exploit the *witness* concept introduced in [11] to handle such constraints. A witness is a single itemset on which we can test whether a constraint holds and derive information about properties of other itemsets.

Definition 1 (Witness). Let P, Q itemsets, and $C : \mathcal{I} \mapsto \{true, false\}$, then $W, P \subseteq W \subseteq P \cup Q$, is called a positive (negative) witness iff $\forall P', P \subseteq P' \subseteq P \cup Q : C(W) = true \Rightarrow C(P') = true$ ($C(W) = false \Rightarrow C(P') = false$).

2.2 Diversity of Itemsets

The Jaccard index is a classical similarity measure on sets. We use it to quantify the overlap of the covers of itemsets.

Definition 2 (Jaccard index). Given two itemsets P and Q , the Jaccard index is the relative size of the overlap of their covers : $Jac(P, Q) = \frac{|\mathcal{V}_D(P) \cap \mathcal{V}_D(Q)|}{|\mathcal{V}_D(P) \cup \mathcal{V}_D(Q)|}$.

A lower Jaccard indicates low similarity between itemset covers, and can thus be used as a measure of diversity between pairs of itemsets.

Definition 3 (Diversity/Jaccard constraint). Let P and Q be two itemsets. Given the Jac measure and a diversity threshold J_{max} , we say that P and Q are pairwise diverse iff $Jac(P, Q) \leq J_{max}$. We will denote this constraint by c_{Jac} .

Our aim is to push the Jaccard constraint during pattern discovery to prune non-diverse itemsets. To achieve this, we maintain a *history* \mathcal{H} of extracted pairwise diverse itemsets during search and constrain the next mined itemsets to respect a maximum Jaccard constraint with all itemsets already included in \mathcal{H} . This problem can be formalized as follows.

Definition 4 (k diverse frequent itemsets). Given a current history $\mathcal{H} = \{H_1, \dots, H_k\}$ of k pairwise diverse frequent closed itemsets, the Jac measure and a diversity threshold J_{max} , the task is to mine new itemsets P such that $\forall H \in \mathcal{H}, Jac(P, H) \leq J_{max}$.

Example 2. The lattice in Figure 1 depicts the set of diverse FCIs (marked with blue and green solid line circles) with $J_{max} = 0.19$ and $\mathcal{H} = \{BE\}$. ACE is a diverse FCI (i.e., $Jac(ACE, BE) = 0.147 < 0.19$).

Proposition 1. Let P, Q and P' be three itemsets s.t. $P \subset P'$. $Jac(P, Q)$ may be smaller, equal or greater than $Jac(P', Q)$.

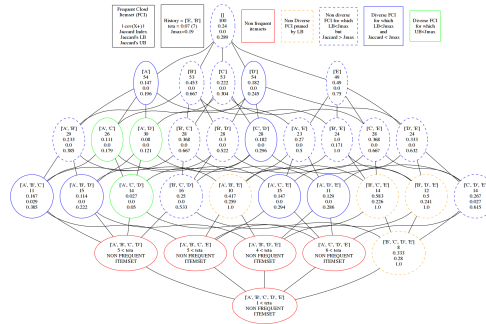


Fig. 1: The powerset lattice of frequent closed itemsets ($\theta = 7$) for the dataset \mathcal{D} of Example 1.

Based on the above proposition, the anti-monotonicity of the maximum Jaccard constraint does not hold, which disables pruning. Thus, instead of solving the problem of Definition 4 directly, we introduce bounds in Section 3 that allow us to prune the search space using a relaxation of the Jaccard constraint. The appeal of this approach is that we are able to infer monotone and anti-monotone properties from this relaxation.

2.3 Constraint Programming (CP)

Constraint programming [10] is a powerful paradigm which offers a generic and modular approach to model and solve combinatorial problems. A CP model consists of a set of variables $X = \{x_1, \dots, x_n\}$, a set of domains D mapping each variable $x_i \in X$ to a finite set of possible values $dom(x_i)$, and a set of constraints \mathcal{C} on X . A constraint $c \in \mathcal{C}$ is a relation that specifies the allowed combinations of values for its variables $X(c)$. An assignment on a set $Y \subseteq X$ of variables is a mapping from variables in Y to values in their domains. A solution is an assignment on X satisfying all constraints. Constraint solvers typically use backtracking search to explore the search space of partial assignments. Algorithm 1 provides a general overview of a CP solver. At each node of the search tree, procedure *Constraint-Search* selects an unassigned variable (line 8) according to user-defined heuristics and assigns it a value (line 9). It backtracks when a constraint cannot be satisfied, i.e. when at least one domain is empty (line 5). A solution is obtained (line 12) when each domain $dom(x_i)$ is reduced to a singleton and all constraints are satisfied. The main concept used to speed up the search is constraint propagation by *Filtering algorithms*. At each assignment, constraint filtering algorithms prune the search space by enforcing local consistency properties like *domain consistency*. A constraint c on $X(c)$ is domain consistent, if and only if, for every $x_i \in X(c)$ and every $v \in dom(x_i)$, there is an assignment satisfying c such that $(x_i = v)$. Global constraints are families of constraints defined by a relation on any number of variables [10].

2.4 A CP Model for Frequent Closed Itemset Mining

The first constraint programming model for frequent closed itemset mining (FCIM) was introduced in [6]. It is based on reified constraints to connect item variables to transac-

Algorithm 1: Constraint-Search(D)

```

1 In:  $X$  : a set of decision variables;  $C$  : a set of constraints;
2 InOut:  $D$  : a set of variable domains;
3 begin
4    $D \leftarrow \text{Filtering}(D, C)$ 
5   if there exists  $x_i \in X$  s.t.  $\text{dom}(x_i)$  is empty then
6     return failure
7   if there exists  $x_i \in X$  s.t.  $|\text{dom}(x_i)| > 1$  then
8     Select  $x_i \in X$  s.t.  $|\text{dom}(x_i)| > 1$ 
9     forall  $v \in \text{dom}(x_i)$  do
10      Constraint-Search( $\text{Dom} \cup \{x_i \rightarrow \{v\}\}$ )
11   else
12      output solution  $D$ 

```

tion variables. The first global constraint CLOSEDPATTERNS for mining frequent closed itemsets was proposed in [13]. The global constraint COVERSIZES for computing the exact size of the cover of an itemset was introduced in [19]. It offers more flexibility in modeling problems. We present the global constraint CLOSEDPATTERNS.

Global Constraint CLOSEDPATTERNS. Most declarative methods use a vector x of Boolean variables $(x_1, \dots, x_{|\mathcal{I}|})$ for representing itemsets, where x_i represents the presence of the item $i \in \mathcal{I}$ in the itemset. We will use the following notations: $x^+ = \{i \in \mathcal{I} \mid \text{dom}(x_i) = \{1\}\}$, $x^- = \{i \in \mathcal{I} \mid \text{dom}(x_i) = \{0\}\}$ and $x^* = \{i \in \mathcal{I} \mid i \notin x^+ \cup x^-\}$.

Definition 5 (CLOSEDPATTERNS). Let x be a vector of Boolean variables, θ a support threshold and \mathcal{D} a dataset. The global constraint $\text{CLOSEDPATTERNS}_{\mathcal{D}, \theta}(x)$ holds if and only if x^+ is a closed frequent itemset w.r.t. the threshold θ .

Definition 6 (Closure extension [22]). A non-empty itemset P is a closure extension of Q iff $\mathcal{V}_{\mathcal{D}}(P \cup Q) = \mathcal{V}_{\mathcal{D}}(Q)$.

Filtering of CLOSEDPATTERNS. [13] also introduced a complete filtering algorithm for CLOSEDPATTERNS based on three rules. The first rule filters 0 from $\text{dom}(x_i)$ if $\{i\}$ is a closure extension of x^+ (see Definition 6). The second rule filters 1 from $\text{dom}(x_i)$ if the itemset $x^+ \cup \{i\}$ is infrequent w.r.t. θ . Finally, the third rule filters 1 from $\text{dom}(x_i)$ if $\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\})$ is a subset of $\mathcal{V}_{\mathcal{D}}(x^+ \cup \{j\})$ where j is an absent item, i.e. $j \in x^-$.

To show the strength and the flexibility of the CP approach in taking into account user’s constraints, we formulate a CP model to extract more specific patterns using the following four global constraints: $\mathcal{C} = \{\text{CLOSEDP}_{\mathcal{D}, \theta}(X), \text{ATLEAST}(X, lb), \text{KNAPSACK}(X, z, w), \text{REGULAR}(X, DFA)\}$. The ATLEAST constraint enforces that at least lb variables in X are assigned to 1; the KNAPSACK constraint restricts a weighted linear sum to be no more than a given capacity z , i.e. $\sum_i w_i X_i \leq z$; the REGULAR constraint imposes that X is accepted by deterministic finite automaton (DFA), which recognizes a regular expression.

Example 3. Let us consider the example of Figure 1 using $lb = 2$, $w = \langle 8, 7, 5, 14, 16 \rangle$, $z = 20$, $\theta = 7$ and the regular expression $0^* 1^+ 0^*$ ensuring items’ contiguity. Solving this CP model provides the solution: $Th(\mathcal{C}) = \{\{AB\}, \{BC\}, \{CD\}, \{ABC\}\}$.

3 A CP Model for Mining Diverse Frequent Closed Itemsets

We present our approach for computing diverse FCIs. The key idea is to compute an approximation of the set of diverse FCIs by defining two bounds on the Jaccard index that allow us to reduce the search space. All the proofs are given in the Supp. material [1].

3.1 Problem Reformulation

Proposition 1 states that the Jaccard constraint is neither monotonic nor anti-monotonic. So, we propose to approximate the theory of the original constraint c_{Jac} by a larger collection corresponding to the solution space of its *relaxation* $c_{Jac}^r: Th(c_{Jac}) \subseteq Th(c_{Jac}^r)$. The key idea is to formulate a relaxed constraint having suitable monotonicity properties in order to exploit them for search space reduction. More precisely, we want to exploit *upper and lower bounding* operators to derive a monotone relaxation and an anti-monotone one of c_{Jac} .

Definition 7 (Problem reformulation). *Given a current history $\mathcal{H} = \{H_1, \dots, H_k\}$ of extracted k pairwise diverse frequent closed itemsets, a diversity threshold J_{max} , a lower bound LB_J and an upper bound UB_J on the Jaccard index, the relaxed problem consists of mining candidate itemsets P such that $\forall H \in \mathcal{H}, LB_J(P, H) \leq J_{max}$. When $UB_J(P, H) \leq J_{max}$, for all $H \in \mathcal{H}$, the Jaccard constraint is fully satisfied.*

3.2 Jaccard Lower Bound

Let us now formalize how to compute the lower bound and how to exploit it. To arrive at a lower bound for the Jaccard value between two itemsets, we need to consider the situation where the *overlap* between them has become as small as possible, while the coverage that is proper to each itemset remains as large as possible.

Definition 8 (Proper cover). *Let P and Q be two itemsets. The proper cover of P w.r.t. Q is defined as $\mathcal{V}_Q^{pr}(P) = \mathcal{V}_D(P) \setminus \{\mathcal{V}_D(P) \cap \mathcal{V}_D(Q)\}$.*

The lowest possible Jaccard would reduce the numerator to 0, which is however not possible under the minimum support threshold θ . The denominator, on the other hand, consists of $|\mathcal{V}_D(H)|$ (which cannot change) and the part of P 's coverage that does not overlap with H , i.e. $\mathcal{V}_H^{pr}(P)$.

Proposition 2 (Lower bound). *Consider a member pattern H of the history \mathcal{H} . Let P an itemset encountered during search such that $sup_D(P) \geq \theta$, and $\mathcal{V}_H^{pr}(P)$ be the proper cover of P w.r.t. H . $LB_J(P, H) = \frac{\theta - |\mathcal{V}_H^{pr}(P)|}{|\mathcal{V}_D(P)| + |\mathcal{V}_D(H)| + |\mathcal{V}_H^{pr}(P)| - \theta}$ is a lower bound of $Jac(P, H)$.*

The lower bound on the Jaccard index enables us to discard some non-diverse itemsets, i.e., those with an LB_J value greater than J_{max} are *negative witnesses*.

Example 4. The set of all non diverse FCIs with a lower bound value greater than J_{max} are marked in Figure 1 with orange line circles.

Proposition 3 (Monotonicity of LB_J). *Let $H \in \mathcal{H}$ be an itemset. For any two itemsets $P \subseteq Q$, the relationship $LB_J(P, H) \leq LB_J(Q, H)$ holds.*

Property 3 establishes an important result to define a pruning condition based on the monotonicity of the lower bound (cf. Section 3.4). If $LB_J(P, H) > J_{max}$, then no itemset $Q \supseteq P$ will satisfy the Jaccard constraint (because LB_J is a lower bound), rendering the constraint itself *anti-monotone*. So, we can safely prune Q .

3.3 Jaccard Upper Bound

As our relaxation approximates the theory of the Jaccard constraint, i.e. $Th(c_{Jac}) \subseteq Th(c_{Jac}^r)$, one could have itemsets P such that $LB_J(P, H) < J_{max}$ but $Jac(P, H) > J_{max}$ (see itemsets marked with blue dashed line circles in Figure 1). To tackle this case, we define an upper bound on the Jaccard index to evaluate the satisfaction of the Jaccard constraint, i.e., those with $UB_J(P, H) \leq J_{max}, \forall H \in \mathcal{H}$, are *positive witnesses*.

To derive the upper bound, we need to follow the opposite argument as for the lower bound: the highest possible Jaccard will be achieved if $\mathcal{V}_D(H) \cap \mathcal{V}_D(P)$ stays unchanged but the set $\mathcal{V}_H^{pr}(P)$ is reduced as much possible (under the minimum support constraint). If the intersection is greater than or equal θ , in the worst case scenario (leading to the highest Jaccard), a future P' covers only transactions in the intersection. If not, the denominator needs to contain a few elements of $\mathcal{V}_H^{pr}(P)$, $\theta - |\mathcal{V}_D(H) \cap \mathcal{V}_D(P)|$, to be exact.

Proposition 4 (Upper bound). *Given a member pattern H of the history \mathcal{H} , and an itemset P such that $sup_D(P) \geq \theta$. $UB_J(P, H) = \frac{|\mathcal{V}_D(H) \cap \mathcal{V}_D(P)|}{|\mathcal{V}_P^{pr}(H)| + \max\{\theta, |\mathcal{V}_D(H) \cap \mathcal{V}_D(P)\}}$ is an upper bound of $Jac(P, H)$.*

Example 5. The set of all diverse FCIs with UB_J values less than J_{max} are marked with green line circles in Figure 1.

Our upper bound can be exploited to evaluate the Jaccard constraint during mining. More precisely, in the enumeration procedure, if the upper bound of the current candidate itemset P is less than J_{max} , then c_{Jac} is fully satisfied. Moreover, if the upper bound is *monotonically decreasing* (or anti-monotonic), then all itemsets Q derived from P are also diverse (see Proposition 5).

Proposition 5 (Anti-monotonicity of UB_J). *Let H be a member pattern of the history \mathcal{H} . For any two itemsets $P \subseteq Q$, the relationship $UB_J(P, H) \geq UB_J(Q, H)$ holds.*

3.4 The Global Constraint CLOSEDDIVERSITY

This section presents our new global constraint CLOSEDDIVERSITY that exploits the LB relaxation to mine pairwise diverse frequent closed itemsets.

Algorithm 2: Filtering for CLOSED DIVERSITY

```

1 In:  $\theta, J_{max}$  : frequency and diversity thresholds;  $\mathcal{H}$  : history of solutions encountered during search;
2 InOut:  $x = \{x_1 \dots x_n\}$  : Boolean item variables;
3 begin
4   if ( $|\mathcal{V}_{\mathcal{D}}(x^+)| < \theta \vee !\mathcal{P}Growth_{LB}(x^+, \mathcal{H}, J_{max})$ ) then return false;
5   foreach  $i \in x^*$  do
6     if ( $|\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\})| < \theta$ ) then
7        $dom(x_i) \leftarrow dom(x_i) - \{1\}$ ;  $x_{Freq}^- \leftarrow x_{Freq}^- \cup \{i\}$ ;  $x^* \leftarrow x^* \setminus \{i\}$ ; continue;
8     if ( $|\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\})| = |\mathcal{V}_{\mathcal{D}}(x^+)|$ ) then
9        $dom(x_i) \leftarrow dom(x_i) - \{0\}$ ;  $x^+ \leftarrow x^+ \cup \{i\}$ ;  $x^* \leftarrow x^* \setminus \{i\}$ ;
10    if ( $!\mathcal{P}Growth_{LB}(x^+ \cup \{i\}, \mathcal{H}, J_{max})$ ) then
11       $dom(x_i) \leftarrow dom(x_i) - \{1\}$ ;  $x_{Div}^- \leftarrow x_{Div}^- \cup \{i\}$ ;  $x^* \leftarrow x^* \setminus \{i\}$ ; continue;
12    foreach  $k \in (x_{Freq}^- \cup x_{Div}^-)$  do
13      if ( $\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\}) \subseteq \mathcal{V}_{\mathcal{D}}(x^+ \cup \{k\})$ ) then
14         $dom(x_i) \leftarrow dom(x_i) - \{1\}$ 
15        if  $k \in x_{Freq}^-$  then  $x_{Freq}^- \leftarrow x_{Freq}^- \cup \{i\}$ ;
16        else  $x_{Div}^- \leftarrow x_{Div}^- \cup \{i\}$ ;
17         $x^* \leftarrow x^* \setminus \{i\}$ ; break;
18    return true;
19 Function  $\mathcal{P}Growth_{LB}(x, \mathcal{H}, J_{max})$ : Boolean
20   foreach  $H \in \mathcal{H}$  do
21     if ( $LB_J(x, H) > J_{max}$ ) then return false
22   return true

```

Definition 9 (CLOSED DIVERSITY). Let x be a vector of Boolean item variables, \mathcal{H} a history of pairwise diverse frequent closed itemsets (initially empty), θ a support threshold, J_{max} a diversity threshold and \mathcal{D} a dataset. The $CLOSED DIVERSITY_{\mathcal{D}, \theta}(x, \mathcal{H}, J_{max})$ global constraint holds if and only if: (1) x^+ is closed; (2) x^+ is frequent, $\sup_{\mathcal{D}}(x^+) \geq \theta$; (3) x^+ is diverse, $\forall H \in \mathcal{H}, LB_J(x^+, H) \leq J_{max}$.

Initially, the history \mathcal{H} is empty. Our global constraint allows to incrementally update \mathcal{H} with diverse FCIs encountered during search. Condition (3) expresses a necessary condition ensuring that x^+ is diverse. Indeed, one could have $LB_J(x^+, H) \leq J_{max}$ but $Jac(x^+, H) > J_{max}$. Thus, we propose in Section 4 to exploit our *UB* relaxation to guarantee the satisfaction of the Jaccard constraint.

The propagator for $CLOSED DIVERSITY$ exploits the filtering rules of $CLOSED PATTERNS$ (see Section 2.4). It also uses our *LB* relaxation to remove items i that cannot belong to a solution containing x^+ . We denote by x_{Freq}^- the set of items filtered by the rule of infrequent items and by x_{Div}^- the set of items filtered by our *LB* rule.

Proposition 6 (CLOSED DIVERSITY Filtering rule). Given a history \mathcal{H} of pairwise diverse frequent closed itemsets, a partial assignment on x , and a free item $i \in x^*$, $x^+ \cup \{i\}$ cannot lead to a diverse itemset if one of the two cases holds:

- 1) if $\exists H \in \mathcal{H}$ s.t. $LB_J(x^+ \cup \{i\}, H) > J_{max}$, then we remove 1 from $dom(x_i)$.
- 2) if $\exists k \in x_{Div}^-$ s.t. $\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\}) \subseteq \mathcal{V}_{\mathcal{D}}(x^+ \cup \{k\})$, then $LB_J(x^+ \cup \{i\}, H) > LB_J(x^+ \cup \{k\}, H) > J_{max}$, thus we remove 1 from $dom(x_i)$.

Algorithm. The propagator for CLOSED DIVERSITY is presented in Algorithm 2. It takes as input the variables x , the support threshold θ , the diversity threshold J_{max} and the current history \mathcal{H} of pairwise diverse frequent closed itemsets. It starts by computing the cover of the itemset x^+ and checks if x^+ is either infrequent or not diverse (see function $\mathcal{P}Growth_{LB}$), if so the constraint is violated and a fail is returned (line 4). Algorithm 2 extends the filtering rules of CLOSED PATTERNS (see Section 2.4) by examining the diversity condition of the itemset $x^+ \cup \{i\}$ (see Proposition 3). For each element $H \in \mathcal{H}$, the function $\mathcal{P}Growth_{LB}(x^+ \cup \{i\}, \mathcal{H}, J_{max})$ computes the value of $LB_J(x^+ \cup \{i\}, H)$ and tests if there exists an H s.t. $LB_J(x^+ \cup \{i\}, H) > J_{max}$ (lines 20-21). If so, we return *false* (line 21) because $x^+ \cup \{i\}$ cannot lead to a diverse itemset w.r.t. \mathcal{H} , remove 1 from $dom(x_i)$ (line 11), update x_{Div}^- and x^* and we continue with the next free item. Otherwise, we return true. Second, we remove 1 from each free item variable $i \in x^*$ such that its cover is a superset of the cover of an absent item $k \in (x_{Freq}^- \cup x_{Div}^-)$ (lines 12-17). The LB filtering rule associated to the case $k \in x_{Div}^-$ is a new rule taking its originality from the reasoning made on absent items.

Proposition 7 (Consistency and time complexity). *Algorithm 2 enforces Generalized Arc Consistency (GAC) (a.k.a. domain consistency [10]) in $\mathcal{O}(n^2 \times m)$.*

4 Using witnesses and the estimated frequency within the search

In this section, we show how to exploit the witness property and the estimated frequency so as to design a more informed search algorithm.

Positive Witness. During search, we compute incrementally the $UB(x^+ \cup \{i\}, H)$ of any extension of the partial assignment x^+ with a free item i . If, for each $H \in \mathcal{H}$, this upper bound is less or equal to J_{max} , then c_{Jac} is fully satisfied and $x^+ \cup \{i\}$ is a *positive witness*. Moreover, thanks to the anti-monotonicity of UB_J (see Proposition 5), all supersets of $x^+ \cup \{i\}$ will satisfy the Jaccard constraint.

Estimated frequency. The frequency of an itemset can be computed as the cardinality of the intersection of its items' cover: $sup_{\mathcal{D}}(x^+) = |\cap_{i \in x^+} \mathcal{V}_{\mathcal{D}}(i)|$, the intersection between 2 covers being performed by a bitwise-AND. To limit the number of intersections, we use an estimation of the frequency of each item $i \in \mathcal{I}$ w.r.t the set of present items x^+ , denoted $eSup_{\mathcal{D}}(i, x^+)$. This estimation constitutes a *lower bound* of $|\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\})|$. Interestingly, if $eSup_{\mathcal{D}}(i, x^+) \geq \theta$ then $|\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\})| \geq \theta$, meaning that the intersection between covers is performed only if $eSup_{\mathcal{D}}(i, x^+) < \theta$, thereby leading to performance enhancement. In addition, we argue that the estimated support is an interesting heuristic to reinforce the witness branching rule. Indeed, branching on the variable having the minimum estimated support (using the lower bound of the real support) will probably activate our filtering rules (see Algorithm 2), thus reducing the search space. It will be denoted as MINCOV variable ordering heuristic.

We propose Algorithm 3 as a branching procedure (returns the next variable to branch on). When the search begins, for each item $i \in x^*$, its estimated frequency is initialized to $eSup_{\mathcal{D}}(i, \emptyset) = |\mathcal{V}_{\mathcal{D}}(i)|$. Once an item j has been added to the partial solution, the estimated frequencies of unbound items must be updated (see lines 4-9). Thus, we first find the variable x^{es} having the minimal estimated support (line 4).

Algorithm 3: Branching for CLOSED DIVERSITY

```

1 In:  $J_{max}$  : diversity thresholds;  $\mathcal{H}$  : history of solutions ;
2 Out: First witness index or  $x^{es}$  as the item with the smallest estimated support
3 begin
4    $x^{es} \leftarrow \operatorname{argmin}_{i \in x^*} (eSup_{\mathcal{D}}(i, x^+))$ ;
5    $\text{diff} \leftarrow (|\mathcal{V}_{\mathcal{D}}(x^+)| - |\mathcal{V}_{\mathcal{D}}(x^+ \cup \{x^{es}\})|)$ ;
6   foreach  $i \in x^* \setminus \{x^{es}\}$  do
7      $eSup_{\mathcal{D}}(i, x^+ \cup \{x^{es}\}) \leftarrow eSup_{\mathcal{D}}(i, x^+) - \text{diff}$ ;
8     if  $(eSup_{\mathcal{D}}(i, x^+ \cup \{x^{es}\}) < \theta)$  then
9        $eSup_{\mathcal{D}}(i, x^+ \cup \{x^{es}\}) \leftarrow |\mathcal{V}_{\mathcal{D}}(x^+ \cup \{x^{es}\}) \cap \mathcal{V}_{\mathcal{D}}(i)|$ ;
10
11   foreach  $i \in x^*$  do
12     if  $(\mathcal{P}Growth_{UB}(x^+ \cup \{i\}, \mathcal{H}, J_{max}))$  then
13       return  $\langle i, \text{true} \rangle$ ;
14   return  $\langle x^{es}, \text{false} \rangle$ 
15
16 Function  $\mathcal{P}Growth_{UB}(x^+ \cup \{j\}, \mathcal{H}, J_{max})$  : Boolean
17   foreach  $H \in \mathcal{H}$  do
18     if  $(UB_J(x^+ \cup \{j\}, H) > J_{max})$  then
19       return false
20   return true

```

Next, each item $i \in x^* \setminus \{x^{es}\}$ may lose some support, but no more than $|\mathcal{V}_{\mathcal{D}}(x^+)| - |\mathcal{V}_{\mathcal{D}}(x^+ \cup \{x^{es}\})|$, since some removed transactions may not contain i (line 5). Using this upper bound (denoted by *diff*), the estimated frequency of i is updated and set to $eSup_{\mathcal{D}}(i, x^+) - \text{diff}$ (lines 6-9). As indicated above, if $eSup_{\mathcal{D}}(i, x^+) \geq \theta$ then $|\mathcal{V}_{\mathcal{D}}(x^+ \cup \{i\})| \geq \theta$. Otherwise, we have to compute the right support by performing the intersection between covers (line 9). It is important to stress that the branching variable x^{es} will be returned (line 13) only if no positive witness is found (lines 10-12). Finally, the function $\mathcal{P}Growth_{UB}(x^+ \cup \{i\}, \mathcal{H}, J_{max})$ allows to test whether the current instantiation x^+ can be extended to a witness itemset using the free item $\{i\}$. It returns true if the upper bound of the current itemset x^+ when adding one item $\{i\}$ is less than J_{max} for all $h \in \mathcal{H}$ (lines 15-17). Here, the Jaccard constraint is fully satisfied and thus, we return the item $\{i\}$ with the witness flag set to *true*. This information will be supplied to the search engine (line 12) to accelerate solutions certification. We will denote by FIRSTWITCOV, our variable ordering heuristic that branches on the first free item satisfying the witness property.

Exploring the witness subtree. Let N be the node associated to the current itemset x^+ extended to a free item $\{i\}$. When the node N is detected as a positive witness during the branching, all supersets derived from N will also satisfy the Jaccard constraint. As these patterns are more likely to have similar covers, so a rather high Jaccard between them, we propose a simple strategy which avoids a complete exploration of the witness sub-tree rooted at N . Thus, we generate the first closed diverse itemset from N , add it to the current history and continue the exploration of the remaining search space using the new history. With a such strategy we have no guarantee that the closed itemset added to the history have the best Jaccard. But this strategy is fast.

5 Related work

The question of mining sets of diverse patterns has been addressed in the recent literature, both to offer more interesting results and to speed up the mining process. Van Leeuwen *et al.* propose populating the beam for subgroup discovery not purely with the *best* partial patterns to be extended but to take coverage overlap into account [20]. Beam search is heuristic, as opposed to our exhaustive approach and since they mine all patterns at the same time, diverse partial patterns can still lead to a less diverse final result. Dzyuba *et al.* propose using *XOR* constraints to partition the transaction set into disjoint subsets that are small enough to be efficiently mined using either a CP approach or a dedicated itemset miner [8]. Their focus is on efficiency, which they demonstrate by approximating the result set of an exhaustive operation. While they discuss pattern sets, they limit themselves to a strict non-overlap constraint on coverages. In [4], the authors propose using Monte Carlo Tree Search and upper confidence bounds to direct the search towards interesting regions in the lattice given the already explored space. While MCTS is necessarily randomized, it allows for anytime mining. The authors of [3] consider sets of subgroup descriptions as *disjunctions* of such patterns. Using a greedy algorithm exploiting upper bounds, the authors propose to iteratively extract up to k subgroup descriptions (similarly to our work). Notably, this approach requires a target attribute *and* a target value to focus on while our approach allows for unsupervised mining.

Earlier work has treated reducing redundancy as a post-processing step, e.g. [12] where a number of redundancy measures such as entropy are exploited in exhaustive search and the number of patterns in the set limited, [7] where the constraint-based itemset mining constraint is adapted to the pattern set settings, [5], which exploit bounds on predicting the presence of patterns from the patterns already included in \mathcal{H} in a heuristic algorithm, or [21], which exploits the MDL principle to minimize redundancy among itemsets (and, in later work, sequential patterns). All of those methods require a potentially rather costly first mining step, and none exploits the Jaccard measure. As discussed in Section 2.1, there exist a number of constraint properties that allow for pruning, and Kifer *et al.*'s witness concept unifies them and discusses how to deal with constraints that do not have monotonicity properties [11]. The way to proceed in such a case is establishing *positive* and *negative* witnesses for the constraint, something we have done for the maximum pairwise Jaccard constraint. A rarely discussed aspect is that witnesses are closely related to CP since every witness enforces/forbids the inclusion of certain domain values.

6 Experiments and Results

The experimental evaluation is designed to address the following questions: (1) How (in terms of CPU-times and # of patterns) does our global constraint (denoted CLOSED-DIV) compare to the CLOSEDPATTERNS global constraint (denoted CLOSEDP) and the approach of Dzyuba *et al.* [8] (denoted FLEXICS)? (2) How do the resulting diverse FCIs compare qualitatively with those resulting from CLOSEDP and FLEXICS? (3) How far is the distance between the Jaccard index and the upper/lower bounds.

Dataset	$ \mathcal{I} \times \mathcal{T} $ $\rho(\%)$	#Patterns		Time (s)		#Nodes	
		$\theta(\%)$	(1)	(2)	(1)	(2)	(2)
CHESS 75 × 3196 49.33%	20	22,808,625	96	2838.30	5.87	45,617,249	436
	15	50,723,131	393	5666.03	75.40	101,446,261	1,855
	10	OOM	4,204	OOM	3825.29	OOM	18,270
HEPATITIS 68 × 137 50.00%	30	83,048	12	9.64	0.09	166,095	29
	20	410,318	87	42.00	0.57	820,635	162
	10	1,827,264	2,270	169.59	76.91	3,654,527	5,256
KR-VS-KP 73 × 3196 49.32%	30	5,219,727	17	682.94	0.74	10,439,453	82
	20	21,676,719	96	2100.79	5.64	43,353,437	448
	10	OOM	4,120	OOM	3035.49	OOM	17,861
CONNECT 129 × 67557 33.33%	30	460,357	18	1666.14	14.81	920,713	77
	18	2,005,476	197	5975.44	573.66	4,010,951	900
	15	3,254,780	509	9534.07	1989.35	6,509,559	2,188
HEART-CLEVELAND 95 × 296 47.37%	10	12,774,456	3,496	1308.63	257.39	25,548,911	7,977
	8	23,278,687	12,842	2278.97	2527.38	46,357,373	28,221
	6	43,588,346	58,240	4126.84	46163.06	87,176,691	124,705
SPLICE1 287 × 3190 20.91%	10	1,606	422	6.55	25.25	3,211	843
	5	31,441	8,781	117.15	5616.47	62,881	17,594
	2	589,588	-	1179.55	-	1,179,175	-
MUSHROOM 112 × 8124 18.75%	5	8,977	727	10.02	60.70	17,953	1,704
	1	40,368	12,139	34.76	12532.95	80,735	25,154
	0.5	62,334	27,768	50.05	64829.06	124,667	56,873
T4010D100K 942 × 100000 4.20%	8	138	127	75.91	447.20	275	253
	5	317	288	331.47	1561.34	633	575
	1	65,237	7,402	5574.31	58613.88	130,473	14,887
PUMSB 2113 × 49046 3.50%	40	-	4	-	57.33	-	16
	30	-	15	-	267.72	-	64
	20	-	52	-	852.39	-	250
T1014D100K 870 × 100000 1.16%	5	11	11	1.73	6.31	21	21
	1	386	361	434.25	3125.06	771	722
	0.5	1,074	617	881.31	7078.90	2,147	1,257
BMS1 497 × 59602 0.51%	0.15	1,426	609	11362.71	68312.38	2,851	1,220
	0.14	1,683	668	11464.93	68049.00	3,365	1,339
	0.12	2,374	823	13255.79	79704.88	4,747	1,651
RETAIL 16470 × 88162 0.06%	5	17	13	10.74	33.44	33	25
	1	160	111	297.21	1625.73	319	227
	0.4	832	528	6073.53	31353.23	1,663	1,093

Table 1: CLOSED DIV ($J_{max} = 0.05$) vs CLOSED P. For columns #Patterns and #Nodes, the values in bold indicate a reduction more than 20% of the total number of patterns and nodes. “-” is shown when time limit is exceeded. OOM : Out Of Memory. (1): CLOSED P (2): CLOSED DIV

Experimental protocol. Experiments were carried out on classic UCI data sets, available at the FIMI repository (fimi.ua.ac.be/data). We selected several real-world data sets, their characteristics (name, number of items $|\mathcal{I}|$, number of transactions $|\mathcal{T}|$, density ρ) are shown in the first column of Table 1. We selected data sets of various size and density. Some data sets, such as Hepatitis and Chess, are very dense (resp. 50% and 49%). Others, such as T10 and Retail, are very sparse (resp. 1% and 0.06%). The implementation of the different global constraints and their constraint propagators were carried out in the Choco solver [17] version 4.0.3, a Java library for constraint programming. The source code is publicly available.⁵ Experiments were conducted on AMD Opteron 6174, 2.2 GHz with a RAM of 256 GB and a time limit of 24 hours. The default maximum heap size allowed by the JVM is 30 GB. We have selected for every data set frequency thresholds to have different numbers of frequent closed itemsets ($|Th(c)| \leq 15000$, $30000 \leq |Th(c)| \leq 10^6$, and $|Th(c)| > 10^6$). The only exception are the very large and sparse data sets Retail and Pumsb, where we do not find a large number of solutions. We used the CLOSED P CP model as a baseline to determine suitable thresholds used with the CLOSED DIV CP model. To evaluate the quality of a set of patterns in terms of diversity, we measured the average ratio of exclusive pattern coverages: $ECR(P_1, \dots, P_k) = avg_{1 \leq i \leq k} \left(\frac{sup_D(P_i) - |V_D(P_i) \cap \bigcup_{j \neq i} V_D(P_j)|}{sup_D(P_i)} \right)$.

(a) Comparing CLOSED DIV with CLOSED P and FLEXICS. Table 1 compares the performance of the two CP models for various values of θ on different data sets. Here, we report the CPU time (in seconds), the number of extracted patterns, and the number of nodes explored during search. This enables to evaluate the amount of inconsistent values pruned by each approach (filtering algorithm). We use MINCOV as variable ordering heuristic. The maximum diversity threshold J_{max} is set to 0.05. First, the results highlight the great discrepancy between the two models with a distinctly lower number of patterns generated by CLOSED DIV (in the thousands) in comparison to CLOSED P (in the millions). On dense and moderately dense data sets (from CHESS to MUSHROOM),

⁵ <https://github.com/lobnury/ClosedDiversity>

the discrepancy is greatly amplified, especially for small values of θ . For instance, on CHES, the number of patterns for CLOSEDIV is reduced by 99% (from $\sim 50 \cdot 10^6$ solutions to 393) for θ equal to 15%. The density of the data sets provides an appropriate explanation for the good performance of CLOSEDIV. As the number of closed patterns increases with the density, redundancy among these patterns increases as well. On very sparse data sets, CLOSEDIV still outputs fewer solutions than CLOSEDP but the difference is less pronounced. This is explainable by the fact that on these data sets, where we have few solutions, almost all patterns are diverse.

Second, regarding runtime, CLOSEDIV exhibits different behaviours. On dense data sets ($\rho \geq 30\%$), CLOSEDIV is more efficient than CLOSEDP and up to an order of magnitude faster. On CHES (resp. CONNECT), the speed-up is 1455 (resp. 112) for $\theta = 30\%$. For instances resulting in between 500 and 5000 diverse FCIs, the speed-up is up to 5. This good performance of CLOSEDIV is mainly due to the strength of the *LB* filtering rule that provides the CP solving process with more propagation to remove more inconsistent values in the search space. In addition, the number of nodes explored by CLOSEDIV is always small comparing to CLOSEDP. These results support our previous observations. The only exception is HEART-CLEVELAND for which CLOSEDIV is slower (especially for values of $\theta \leq 8\%$). This is mainly due to the relative large number of diverse patterns (≥ 12000), which induces higher lower bound computational overhead. We observe the same behaviour on the two moderately dense data sets SPLICE1 and MUSHROOM. On sparse data sets, CLOSEDIV can take significantly more time to extract all diverse FCIs. This can be explained by the fact that on these instances almost all FCIs are diverse w.r.t. lower bound (on average about 70% for RETAIL and 39% for BMS1, see Table 1). Thus, non-solutions are rarely filtered, while the lower bound overhead greatly penalizes the CP solving process. On the very large PUMSB data set, finally, our approach is very efficient while CLOSEDP fails to complete the extraction.

Finally, Figure 2a compares CLOSEDIV with FLEXICS (two variants) for various values of θ on different data sets: GFLEXICS, which uses CP4IM [6] as an oracle to enumerate the solutions, and EFLEXICS, a specialized variant, based on ECLAT [23]. We run WEIGHTGEN with values of $\kappa \in \{0.1, 0.5, 0.9\}$ [9]. For each instance, we fixed the number of samples to the number of solutions returned by CLOSEDIV. We report results corresponding to the best setting of parameter κ . First, CLOSEDIV largely dominates GFLEXICS, being more than an order of magnitude faster. Second, while EFLEXICS is faster than GFLEXICS, our approach is almost always ranked first, illustrating its usefulness for mining diverse patterns in an anytime manner.

(b) Impact of varying J_{max} . We varied J_{max} from 0.1 to 0.7. The minimum support θ is fixed for each data set (indicated after '-'). Figure 2 shows detailed results. As expected, the greater J_{max} , the longer the CPU time. In fact, the size of the history \mathcal{H} grows rapidly with the increase of J_{max} . This induces significant additional costs in the lower and upper bound computations. Moreover, when J_{max} becomes sufficiently large, the *LB* filtering of CLOSEDIV occurs rarely since the lower bound is almost always below the J_{max} value (see Figure 3). Despite the hardness of some instances (i.e., $J_{max} \geq 0.35$), our CP approach is able to complete the extraction for almost all values of J_{max} . The only exception are the large and dense data sets, where CLOSEDIV fails

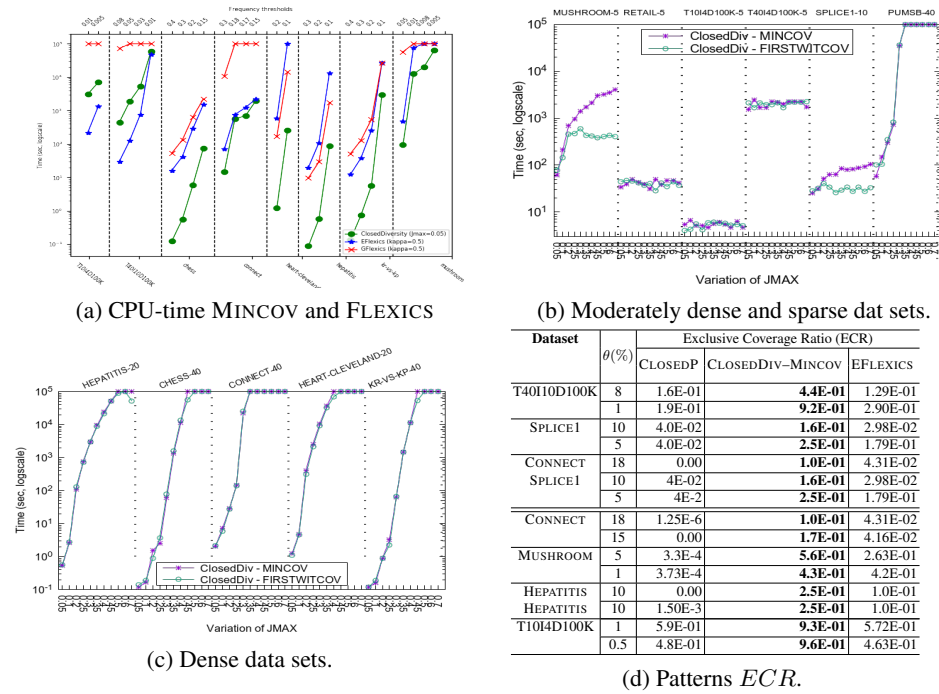
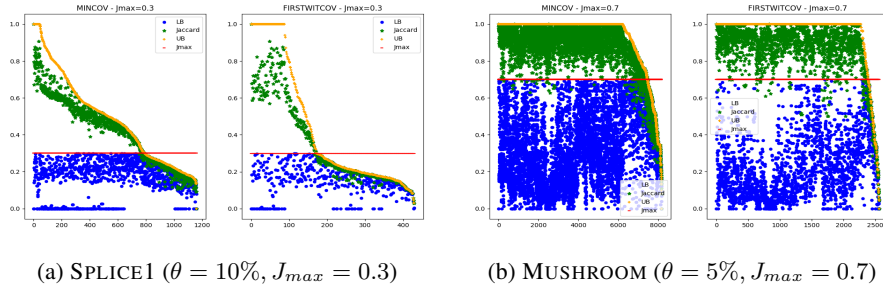


Fig. 2: CPU-time analysis (MINCOV vs FIRSTWITCOV and CLOSEDDIV vs FLEXICS) and patterns discrepancy analysis.

to complete the extraction within the time limit for $J_{max} \geq 0.45$. However, in practice, the user will only be interested in small values of J_{max} because the diversity of patterns is maximal and the number of patterns returned becomes manageable.

Figures 2b and 2c also compare the resolution time of our CP model using the two variable ordering heuristics MINCOV and FIRSTWITCOV. First, on dense data sets, both heuristics perform similarly, with a slight advantage for FIRSTWITCOV. On these data sets, the number of witness patterns mined remains very low (≤ 100), thus the benefits of FIRSTWITCOV is limited (see Supp. material). On moderately dense data sets (MUSHROOM and SPLICE1), FIRSTWITCOV is very effective; on MUSHROOM it is up 10 times faster than MINCOV for J_{max} equal to 0.7. On these data sets, the number of witness patterns extracted is relatively high compared to dense ones. In this case, FIRSTWITCOV enables to guide the search to find diverse patterns more quickly. On sparse data sets, no heuristic clearly dominates the other. When regarding the number of diverse patterns generated (see Supp. material), we observe that FIRSTWITCOV returns less patterns on moderately dense and sparse data sets, while on dense data sets the number of diverse patterns extracted remains comparable.

(c) Qualitative analysis of the proposed relaxation. In this section, we shed light on the quality of the relaxation of the Jaccard constraint. Figure 3a shows, for a particular instance SPLICE1 with $J_{max} = 0.3$, the evolution of the LB_J and UB_J of the solutions

Fig. 3: Qualitative analysis of the LB and UB relaxations.

found during search. Here, the solutions are sorted according to their UB_J . Concerning the lower bound, one can observe that the LB_J values are always below the J_{max} value. This shows how frequently the LB filtering rule of CLOSEDIV occurs. This also supports the suitability of the LB filtering rule for pruning non-diverse FCIs. With regard to the upper bound, it is interesting to see that it gets very close to the Jaccard value, meaning that our Jaccard upper bounding provides a tight relaxation. Moreover, a large number of solutions have UB_J values either below or very close to J_{max} . This is indicative of the quality of the patterns found in terms of diversity. We recall that when $UB_J < J_{max}$, all partial assignments can immediately be extended to diverse itemsets, thanks to the anti-monotonicity property of our UB (see Proposition 5). We observe the same behaviour on MUSHROOM with $J_{max} = 0.7$ (see Figure 3b). Finally, we can see that FIRSTWITCOV allows to quickly discover solutions of better quality in terms of UB_J and Jaccard values compared to MINCOV. This demonstrates the interest and the strength of our UB_J branching rule to get diverse patterns.

(d) Qualitative analysis of patterns. Figure 2d compares CLOSEDIV with CLOSEDP and EFLEXICS in terms of the ECR measure, which should be as high as possible. Due to the huge number of patterns generated by CLOSEDP, a random sample of $k = 10$ solutions of all patterns is considered. Reported values are the average over 100 trials. ECR penalises overlap, and thus having two similar patterns is undesirable. According to ECR , leveraging Jaccard in CLOSEDIV clearly leads to pattern sets with more diversity among the patterns. This is indicative of patterns whose coverage are (approximately) mutually exclusive. This should be desirable for an end-user tasked with exploring and interpreting the set of returned patterns.

7 Conclusions

In this paper, we showed that mining diverse patterns using a maximum Jaccard constraint cannot be modeled using an anti-monotonic constraint. Thus, we have proposed (anti-)monotonic lower and upper bound relaxations, which allow to make pruning effective, with an efficient branching rule, boosting the whole search process. The proposed approach is introduced as a global constraint called CLOSEDIV where diversity is controlled through a threshold on the Jaccard similarity of pattern occurrences. Experimental results on UCI datasets demonstrate that our approach significantly reduces

the number of patterns, the set of patterns is diverse and the computation time is lower compared to CLOSEDP global constraint, particularly on dense data sets.

References

1. Supplementary Material (June 2020), <https://github.com/lobnury/ClosedDiversity>.
2. Belaid, M., Bessiere, C., Lazaar, N.: Constraint programming for mining borders of frequent itemsets. In: Proceedings of IJCAI 2019, Macao, China. pp. 1064–1070 (2019)
3. Belfodil, A., Belfodil, A., Bendimerad, A., Lamarre, P., Robardet, C., Kaytoue, M., Plantevit, M.: Fssd-a fast and efficient algorithm for subgroup set discovery. In: Proceedings of DSAA. pp. 91–99 (2019)
4. Bosc, G., Boulicaut, J.F., Raïssi, C., Kaytoue, M.: Anytime discovery of a diverse set of patterns with monte carlo tree search. *Data mining and knowledge discovery* **32**(3), 604–650 (2018)
5. Bringmann, B., Zimmermann, A.: The chosen few: On identifying valuable patterns. In: Proceedings of ICDM 2007. pp. 63–72
6. De Raedt, L., Guns, T., Nijssen, S.: Constraint programming for itemset mining. In: 14th ACM SIGKDD. pp. 204–212 (2008)
7. De Raedt, L., Zimmermann, A.: Constraint-based pattern set mining. In: 7th SIAM SDM. pp. 237–248. SIAM (2007)
8. Dzyuba, V., van Leeuwen, M., De Raedt, L.: Flexible constrained sampling with guarantees for pattern mining. *Data Mining and Knowledge Discovery* **31**(5), 1266–1293 (2017)
9. Dzyuba, V., van Leeuwen, M.: Interactive discovery of interesting subgroup sets. In: International Symposium on Intelligent Data Analysis. pp. 150–161. Springer (2013)
10. Hoeve, W., Katriel, I.: Global constraints. In: Handbook of Constraint Programming, pp. 169–208. Elsevier Science Inc. (2006)
11. Kifer, D., Gehrke, J., Bucila, C., White, W.: How to quickly find a witness. In: Constraint-Based Mining and Inductive Databases. pp. 216–242. Springer Berlin Heidelberg (2006)
12. Knobbe, A.J., Ho, E.K.: Pattern teams. In: Proceedings of ECML-PKDD. pp. 577–584. Springer (2006)
13. Lazaar, N., Lebbah, Y., Loudni, S., Maamar, M., Lemièrre, V., Bessiere, C., Boizumault, P.: A global constraint for closed frequent pattern mining. In: Proceedings of the 22nd CP. pp. 333–349 (2016)
14. van Leeuwen, M.: Interactive Data Exploration Using Pattern Mining, pp. 169–182. Springer Berlin Heidelberg, Berlin, Heidelberg (2014)
15. Ng, R.T., Lakshmanan, L.V.S., Han, J., Pang, A.: Exploratory mining and pruning optimizations of constrained association rules. In: Proceedings of ACM SIGMOD. pp. 13–24 (1998)
16. Pei, J., Han, J., Lakshmanan, L.V.S.: Mining frequent item sets with convertible constraints. In: Proceedings of ICDE. pp. 433–442 (2001)
17. Prud’homme, C., Fages, J.G., Lorca, X.: Choco Solver Documentation (2016)
18. Puolamäki, K., Kang, B., Lijffijt, J., De Bie, T.: Interactive visual data exploration with subjective feedback. In: Proceedings of ECML PKDD. pp. 214–229. Springer (2016)
19. Schaus, P., Aoga, J.O.R., Guns, T.: Coversize: A global constraint for frequency-based itemset mining. In: Proceedings of the 23rd CP 2017. pp. 529–546 (2017)
20. Van Leeuwen, M., Knobbe, A.: Diverse subgroup set discovery. *Data Mining and Knowledge Discovery* **25**(2), 208–242 (2012)
21. Vreeken, J., Van Leeuwen, M., Siebes, A.: Krimp: mining itemsets that compress. *Data Mining and Knowledge Discovery* **23**(1), 169–214 (2011)

22. Wang, J., Han, J., Pei, J.: CLOSET+: searching for the best strategies for mining frequent closed itemsets. In: Proceedings of the Ninth KDD. pp. 236–245. ACM (2003)
23. Zaki, M., Parthasarathy, S., Ogihara, M., Li, W.: New algorithms for fast discovery of association rules. In: Proceedings of KDD 1997, Newport Beach, California, USA, August 14-17, 1997. pp. 283–286. AAAI Press (1997)